



UNIVERSIDADE DE BRASÍLIA
INSTITUTO DE GEOCIÊNCIAS – IG
PROGRAMA DE PÓS-GRADUAÇÃO EM GEOLOGIA

MACHINE LEARNING APLICADO NA CARACTERIZAÇÃO DA ASSINATURA
PETROFÍSICA, ESPECTRAL E GEOQUÍMICA DOS DEPÓSITOS AURÍFEROS DA
SERRA DE JACOBINA, CRÁTON SÃO FRANCISCO

Guilherme Ferreira da Silva

Tese de Doutorado Nº 187

Brasília, DF
Abril de 2022

Universidade de Brasília
Instituto de Geociências – IG
Programa de Pós-graduação em Geologia – PPGG
Área de Concentração: Prospecção e Geologia Econômica

Título: *MACHINE LEARNING* APLICADO NA CARACTERIZAÇÃO DA ASSINATURA PETROFÍSICA, ESPECTRAL E GEOQUÍMICA DOS DEPÓSITOS AURÍFEROS DA SERRA DE JACOBINA, CRÁTON SÃO FRANCISCO

Tese submetida ao Programa de Pós-Graduação em Geologia da Universidade de Brasília, como cumprimento parcial dos requerimentos para a outorga do grau de Doutor em Geologia.

Guilherme Ferreira da Silva

Tese de Doutorado N° 187

Orientadora: Profa. Dra. Adalene Moreira Silva

Banca examinadora

Prof. Dr. Álvaro Penteado Crosta (IG-Unicamp) – Titular
Prof. Dr. José Carlos Sicoli Seoane (DEGEO-UFRJ) – Titular
Profa. Dra. Susanne Taina Ramalho Maciel (FUP-UnB) – Titular
Prof. Dr. Augusto César Bitencourt Pires (IG-UnB) – Suplente
Profa. Dra. Roberta Mary Vidotti (IG-UnB) – Suplente

Brasília, DF

Abril de 2022



CIP - Catalogação na Publicação

Ferreira da Silva, Guilherme
FS586m MACHINE LEARNING APLICADO NA CARACTERIZAÇÃO DA
ASSINATURA PETROFÍSICA, ESPECTRAL E GEOQUÍMICA DOS
DEPÓSITOS AURÍFEROS DA SERRA DE JACOBINA, CRÁTON
SÃO FRANCISCO / Guilherme Ferreira da Silva. --
Brasília, 2022.
208 p. Ilustrado.

Orientadora: Adalene Moreira Silva.
Tese (doutorado) - Universidade de Brasília,
Instituto de Geociências, Programa de Pós-
Graduação em Geologia, 2022.

1. Integração de dados Multifonte. 2. Petrofísica.
3. Geoquímica. 4. Espectrorradiometria. 5.
Depósito de ouro em paleoplacer modificado. I.
Moreira Silva, Adalene, orient. II. Título.



Para Jacqueline e Ulisses.

Com todo meu carinho.



“All models are approximations. Essentially, all models are wrong, but some are useful. However, the approximate nature of the model must always be borne in mind”.

“Todos os modelos são aproximações. Essencialmente, todos os modelos estão errados, mas alguns são úteis. Entretanto, a natureza aproximada dos modelos deve sempre ser levada em conta”.

- George P. E. Box (1919–2013)

Estatístico britânico responsável por diversos avanços nas áreas de transformação de dados, controle de qualidade e inferência bayesiana.



AGRADECIMENTOS

Este trabalho é fruto do apoio e colaboração de família, amigos e colegas de trabalho, sem o qual não seria finalizado. Trago uma tentativa de síntese da minha gratidão, na esperança de não me esquecer de nenhum nome.

Eu inicio agradecendo o apoio dos meus pais, Antonio Marcos Ferreira da Silva e Eliete Barbosa de Brito: por me ensinarem a sempre buscar aprendizado. Sem seu direcionamento e incentivo a minha jornada certamente seria outra. Não me furto de agradecer ao apoio do meu padrasto, Gilmar Elias Rodrigues, que também contribuiu neste mesmo sentido, e da minha irmã, Jéssica Ferreira, que me hospedou em sua casa no primeiro ano do doutorado.

Agradeço ao Serviço Geológico do Brasil pelo apoio logístico, financiamento dos trabalhos de campo e pela liberação para capacitação através de processo seletivo interno. Direciono este agradecimento ao ex-chefe da Divisão de Geologia Econômica, Felipe Mattos Tavares e ao Gerente de Geologia e Recursos Minerais da Superintendência de Salvador, Valter Rodrigues Sobrinho, que tantas vezes apoiaram o meu trabalho.

Agradeço a Universidade de Brasília pela sua importância em mais um estágio da minha formação profissional, pela luta contínua sob a bandeira da educação pública, gratuita e de qualidade, mantendo-se como um dos bastiões da pesquisa acadêmica em tão difíceis tempos. O presente trabalho foi realizado com apoio da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Código de Financiamento 001

Agradeço a equipe da Jacobina Mineração e Comércio (Yamana Gold), mencionando os nomes de Artur Areal Braga, Juliano José de Souza e Wagner Pocay, além de toda a equipe de técnicos e geólogos. Agradeço por serem receptivos e generosos, viabilizando a coleta de amostras necessárias para este trabalho.

Agradeço a minha orientadora, Adalene Moreira Silva, por aceitar me acompanhar nesta caminhada e por me propor um desafio que contribuíra enormemente para meu crescimento técnico. Estendo o agradecimento aos meus coorientadores Catarina Labouré Bemfica Toledo, Evandro Luiz Klein e Farid Chemale Junior, pela convivência amigável e contribuições das mais diversas formas, dentro das possibilidades de cada um.

Agradeço a João Henrique Larizzatti, meu tutor, e Joseneusa Brilhante Rodrigues, ambos membros da Câmara Técnico Científica do Serviço Geológico do Brasil, responsáveis pela viabilização e acompanhamento da pesquisa. Agradeço a ambos pela paciência, amizade e pelos conselhos durante estes quatro anos.



Agradeço aos geólogos, amigos e colegas: Anderson Dourado, Carina Lopes, Kotaro Uchigasaki, Matheus Ferreira (meu irmão), Guilherme Teles, Fernando Almeida, Pedro Costa e Anderson Matias pelo apoio prestado durante os trabalhos de campo e levantamento de dados em laboratório.

Agradeço aos colegas de sala, Marcos Vinícius Ferreira e Iago Lima Costa, pelo incentivo contínuo no exercício das análises quantitativas em geociências, pela amizade e parceria.

Agradeço ao geólogo José Leonardo Andriotti Silva que involuntariamente redirecionou minha carreira, motivo pelo qual sou imensamente grato.

Finalizo essa seção agradecendo a minha esposa, Jaqueline, e filho, Ulisses, a quem dedico esse trabalho. Pelo apoio durante a elaboração desta tese e pela compreensão nos momentos de ausência. Sem fugir do clichê: somos o resultado dos livros que lemos, das viagens que fazemos e das pessoas que amamos.



RESUMO

Neste trabalho foram adquiridas variáveis categóricas e numéricas relacionadas a propriedades físicas das rochas, tais como densidade, susceptibilidade magnética, condutividade elétrica, concentração de radioelementos, reflectância entre outras propriedades químicas. As análises químicas de rocha foram obtidas através de medidas *in situ* de fluorescência de raios-X portátil (pXRF). Adicionalmente, foram analisadas descrições petrográficas e análises de química mineral quantitativas e semiquantitativas em amostras chave para a compreensão do sistema mineral. Ao todo, foram processadas 1950 análises de pXRF, 2484 medidas de espectrorradiometria, 7490 medidas de susceptibilidade magnética, 5720 medidas de condutividade elétrica, 598 medidas de densidade, 541 análises de química mineral (ablação de laser de espectrômetro de massa, LA-ICP-MS) e 304 medidas de radioelementos, além de 20 análises petrográficas por microscópio óptico e 5 análises por microscópio eletrônico. Utilizamos abordagens supervisionadas para fazer previsões e fornecer informações sobre as mineralizações auríferas em rochas do Grupo Jacobina, Cráton do São Francisco, usando os parâmetros petrofísicos e litogeoquímicos em escala de amostra. Um modelo de aprendizado de máquina baseado no algoritmo *Random Forests* foi aplicado para prever a mineralização em amostras de testemunho de sondagem. As acurácias médias foram de 0,87 para treinamento de validação cruzada, 0,91 para os dados de teste e 0,86 para previsão de todas as amostras. O resultado permitiu estimar a importância das variáveis de entrada para a predição e essas estimativas foram validadas por uma interpretação petrográfica de microscopia óptica e eletrônica de varredura, que foram realizadas para esclarecer a relação entre minerais de diferentes estágios com a mineralização do ouro. Paralelamente, utilizamos abordagens não-supervisionadas para extrair informações sobre a estruturação das amostras nos dados de LA-ICP-MS e de reflectância espectral. Usamos métodos de Agrupamento Aglomerativo (*Hierarchical Clustering*) para avaliar os padrões de elementos traços de acordo com o tipo de pirita (detritica ou epigenéticas) e níveis estratigráficos. Em seguida, implementamos a técnica *Uniform Manifold Approximation and Projection* (UMAP) para reduzir a dimensionalidade avaliada para uma projeção bidimensional buscando inspecionar a estrutura interna dos dados. Elementos como Cu, Zn, Ag, Sb, Te, Au, Pb e Bi são mobilizados durante a alteração mineral e foram cristalizados em minerais recém-formados, como calcopirita, pirrotita e esfalerita, que estão espacialmente associados à pirita epigenética e ouro. O padrão das piritas do Grupo Jacobina parece não variar ao longo da estratigrafia, o que sugere uma manutenção da fonte de sedimento ao longo da história de sedimentação ou um posterior reequilíbrio químico. Relativo



às análises de reflectância espectral, aplicamos o algoritmo de *Self-Organizing Maps* (SOM) para segmentar dados em vários agrupamentos baseados na matriz de distância das unidades e, em seguida, usamos a projeção UMAP para compactar a estrutura de dados para um gráfico bidimensional, mantendo os principais padrões de dados e comparando com os espectros de minerais conhecidos descritos nos metaconglomerados. Assim, estimamos a composição mineral com base na distância de cada medição dos minerais conhecidos e validamos essa inferência usando dados geoquímicos. Os resultados da inferência mineral corresponderam ao esperado pela análise geoquímica, validando a estimativa da composição mineral das amostras. Baseado nestes resultados, separamos as assinaturas das propriedades físicas e químicas nas zonas mineralizada, proximal e estéril e indicamos critérios que podem ser utilizados para a prospecção de ouro em zonas de paleoplacer modificado, como presença de calcopirita, esfalerita e outros sulfetos na matriz, além de pirita, teores de cromo, potássio e enxofre, susceptibilidade magnética, densidade e a presença de argilominerais.

PALAVRAS-CHAVE: Machine-learning aplicado a geociências; Prospecção mineral, Ouro em paleoplacer modificado; Integração de dados multifonte.



ABSTRACT

This thesis aims to characterize the signature of gold mineralization of the Serra do Córrego Formation, the basal unit of the Jacobina Group, using multisource data (petrophysics, spectroradiometrics, geochemistry, and mineral chemistry) through data integration and pattern verification using machine learning. Categorical and numerical variables related to the physical properties of rocks were acquired, such as density, magnetic susceptibility, electrical conductivity, the concentration of radioelements, and reflectance, among other chemical properties. Rock chemical analyzes were obtained by in situ portable X-ray fluorescence (pXRF) measurements. Petrographic descriptions and quantitative and semi-quantitative mineral chemistry analyses were also considered in samples for understanding the mineral system. Altogether, 1950 pXRF analyses, 2484 spectroradiometric measurements, 7490 magnetic susceptibility measurements, 5720 electrical conductivity measurements, 598 density measurements, 541 mineral chemistry analyses (mass spectrometer laser ablation, LA-ICP-MS), and 304 measurements of radio elements, in addition to 20 petrographic analyzes by optical microscope and 5 analyzes by electron microscope. We use supervised approaches to make predictions and provide information on gold mineralizations in rocks of the Jacobina Group, São Francisco Craton, using sample-scale petrophysical and litho-geochemical parameters. A machine learning model based on the Random Forests algorithm was applied to predict mineralization in drill core samples. Average accuracies were 0.87 for cross-validation training, 0.91 for testing, and 0.86 for all-sample prediction. The result allowed us to estimate the importance of the input variables for the prediction. These estimates were validated by a petrographic interpretation of optical and scanning electron microscopy, which were performed to understand better the relationship between minerals of different stages of gold mineralization. In parallel, we used unsupervised approaches to extract information about sample structuring from LA-ICP-MS and spectral reflectance data. We used Hierarchical Clustering methods to evaluate trace element patterns according to pyrite type (detrital or epigenetic) and stratigraphic levels. Then, we implemented the Uniform Manifold Approximation and Projection (UMAP) technique to reduce the evaluated dimensionality to a two-dimensional projection, seeking to inspect the internal structure of the data. Elements such as Cu, Zn, Ag, Sb, Te, Au, Pb, and Bi are mobilized during mineral alteration and crystallized into newly formed minerals such as chalcopyrite, pyrrhotite, and sphalerite, which are spatially associated with epigenetic pyrite and gold. The multivariate pattern of the pyrites of the Jacobina Group does not seem to vary along with the stratigraphy, which suggests maintenance of the sediment source throughout the



sedimentation history or a subsequent chemical rebalancing. Concerning spectral reflectance analyses, we apply the Self-Organizing Maps (SOM) algorithm to segment data into various groupings based on the best unit machine distance matrix. We then use the UMAP algorithm to compress the data structure into a two-dimensional graph, maintaining the main data patterns and comparing them with the spectra of known minerals described in the metaconglomerates. Thus, we estimate the mineral composition based on the distance of each measurement from known minerals and validate this inference using geochemical data. The results of the lithochemistry validate the estimate of the mineral composition of the samples. Based on all presented results, we separated the signatures of the physical and chemical properties in the mineralized, proximal and sterile zones. We indicated criteria that can be used for prospecting for gold in modified paleoplacer zones, such as chalcopyrite, sphalerite, and other sulfides in the matrix and pyrite, besides Cr, K, and S contents, magnetic susceptibility, density, and the presence of clay minerals.

KEYWORDS: Machine-learning applied to geosciences; Mineral prospecting; Gold in modified paleoplacer; Multisource data integration.



Sumário

1	INTRODUÇÃO	1
1.1	Apresentação da tese e justificativas.....	1
1.2	Materiais e métodos	6
1.2.1	Amostragem	7
1.2.2	Propriedades físicas de rocha	8
1.2.3	Química de rocha	11
1.2.4	Espectrorradiometria	12
1.2.5	Petrografia	13
1.2.6	Pré-processamento e análise exploratória de dados	14
2	PETROFÍSICA, GEOQUÍMICA E PREDIÇÃO DE MINERALIZAÇÃO.....	19
2.1	Introduction.....	21
2.2	Geology and gold mineralization in the Serra de Jacobina	23
2.2.1	Geological setting	23
2.2.2	Deformation, metamorphism, and hydrothermal alterations	26
2.2.3	Gold mineralization.....	28
2.3	Materials and methods	29
2.3.1	Drill core samples	29
2.3.2	Petrophysics	30
2.3.3	X-Ray Fluorescence	32
2.3.4	Machine learning analysis (MLA): Random Forests	33
2.4	Results and data analysis	38
2.4.1	Petrophysics and lithochemistry	38
2.4.2	Mineralization prediction	40
2.4.3	Probabilistic prediction approach.....	42
2.5	Discussion.....	43
2.5.1	Mineral Targeting	43
2.5.2	Drill core prediction	47



2.6	Conclusions.....	48
3	QUÍMICA MINERAL DE PIRITAS E ASSOCIAÇÃO GEOQUÍMICA DO OURO....	60
3.1	Introduction.....	63
3.2	Geological context	64
3.2.1	The Jacobina Group	64
3.2.2	Gold mineralization.....	67
3.3	Materials and methods	68
3.3.1	Sampling	68
3.3.2	Pyrite grains	69
3.4	Pyrite LA-ICP-MS data	70
3.4.1	Data processing	71
3.5	Results.....	75
3.5.1	LA-ICP-MS results	75
3.5.2	Agglomerative clustering.....	76
3.5.3	Dimensionality Reduction and data visualization.....	78
3.6	Discussions	79
3.6.1	Pyrite chemistry patterns controlled by texture and stratigraphy	79
3.6.2	Gold associations in detrital and epigenetic pyrite	80
3.7	Conclusions.....	83
4	MODELAGEM DE DADOS ESPECTRORRADIOMÉTRICOS	95
4.1	Introduction.....	98
4.2	Materials and methods	98
4.2.1	Drill core sampling and descriptions	98
4.2.2	Reflectance data acquisition.....	99
4.2.1	Spectral library	100
4.2.2	Lithogeochemistry	102
4.2.3	Dimensionality reduction.....	102
4.3	Results and Discussions.....	103
4.4	Concluding remarks	111



4.5	Computer Code Availability	112
5	DISCUSSÕES E CONSIDERAÇÕES FINAIS	128
	REFERÊNCIAS BIBLIOGRÁFICAS	131
	APÊNDICE A – Lista de Publicações	135
	APÊNDICE B – Sistemas Minerais e sua aplicação na exploração mineral	136
	Sistemas Minerais	136
	Modelagem de Potencial Mineral	139
	Aprendizagem de máquina para integração de dados	144
	APÊNDICE C – Código utilizado no Artigo 01	146
	APÊNDICE D – Código utilizado no Artigo 02	165
	APÊNDICE F –Análises de EDS	184

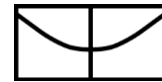


Lista de Figuras

Figura 1-1 – Profundidade de cobertura para as principais descobertas de minerais na Austrália de 1850-2010.	1
Figura 1-2 – Localização da área estudada.....	5
Figura 1-3 – Ilustração das ferramentas utilizadas para obtenção das variáveis.....	7
Figura 1-4 – Diagrama de caixa (boxplot) para as propriedades físicas adquiridas.....	10
Figura 1-5 – Cartas de controle de Shewhart.	12
Figura 1-6 – Diagrama Quantil-Quantil (QQ-plot) para os elementos selecionados (valores de concentração bruta).....	17
Figura 1-7 - Diagrama Quantil-Quantil (QQ-plot) para os elementos selecionados (valores de concentração transformada para razão logaritmica centralizada).....	18
Figura 2-1- Localization of São Francisco Craton in South America.	24
Figura 2-2 – Mineralized drill core samples and photomicrographs.....	27
Figura 2-3 – Conceptual framework describing the behavior of various physical properties..	31
Figura 2-4 – Hyperparameters evaluation and definition of optimum values.....	36
Figura 2-5 – ROC and AUC diagrams for SMOTE-balanced and imbalanced.	37
Figura 2-6 – Selected variables transformed to natural log and centered-log ratio distributions and schematic graphics at the lower portion.....	39
Figura 2-7 – Alluvial validation diagram for model prediction	41
Figura 2-8 – Ore probability analysis for all samples based on the mineralization status of the test dataset.....	42
Figura 2-9 – Color-coded strip log according to the lithologies for drill cores.....	44
Figura 2-10 – Mineral paragenesis flowchart.....	45
Figura 2-11– Variable importance rank.....	47
Figura 3-1 – Simplified geological map and localization in the São Francisco Craton of the Jacobina Group and surrounding domains.....	67
Figura 3-2 – Photomicrography of pyrites and related minerals from the Jacobina Group. ..	70
Figura 3-3 – Data processing fluxogram	71
Figura 3-4 – List of elements ordered by the percentage of non-missing data.	72
Figura 3-5 – Dendrograms and distance matrices for pyrite grains according to grain texture	77
Figura 3-6 – UMAP configuration for data with samples classified according to pyrite texture and stratigraphic level.....	79



Figura 3-7 – Linked dendrograms for Detrital and Epigenetic pyrites.....	82
Figura 3-8 – Epigenetic pyrite grains associated to sulphides in inclusions and in grain borders	83
Figura 3-9 – Quantile-plot for the concentration of V51 on a logarithmic scale according to different distributions.....	86
Figura 3-10 – Comparison of correlated log-transformed elements.....	87
Figura 4-1 – Metaconglomerates samples main aspects	99
Figura 4-2 – Stacked reflectance spectra of selected refence minerals	101
Figura 4-3 – SOM training and results.	105
Figura 4-4 – Reflectance values grouped by each assigned cluster.(Fe/S)	106
Figura 4-5 – Reflectance values grouped by each assigned cluster. (Cr/Fe).....	107
Figura 4-6 – Projection for the ultra-dimensional spectral data into a two-dimensional plot	108
Figura 4-7 – Mineral concentration estimation according to the relative position of the metaconglomerate sample in each the drill core.....	109
Figura 4-8 – Scatterplots of infered mineral contents and geochemical data.....	110
Figura 5-1 - Quadro resumo dos resultados obtidos na tese a respeito das propriedades físicas e químicas das rochas e minerais analisados..	130



1 INTRODUÇÃO

1.1 Apresentação da tese e justificativas

O avanço do conhecimento sobre o conceito de sistemas minerais (Wyborn et al., 1994) vem ao encontro a um desafio da indústria mineral no mundo, ou seja, o de mapear depósitos abaixo da cobertura e em maiores profundidades. A taxa de descobertas de novos depósitos tem diminuído na última década e novas tecnologias têm sido empregadas para se mapear terrenos potenciais em diferentes escalas. Estas tecnologias associam os fatores geológicos que controlam a geração de depósitos minerais e sua preservação ao longo da história geológica (Figura 1-1).

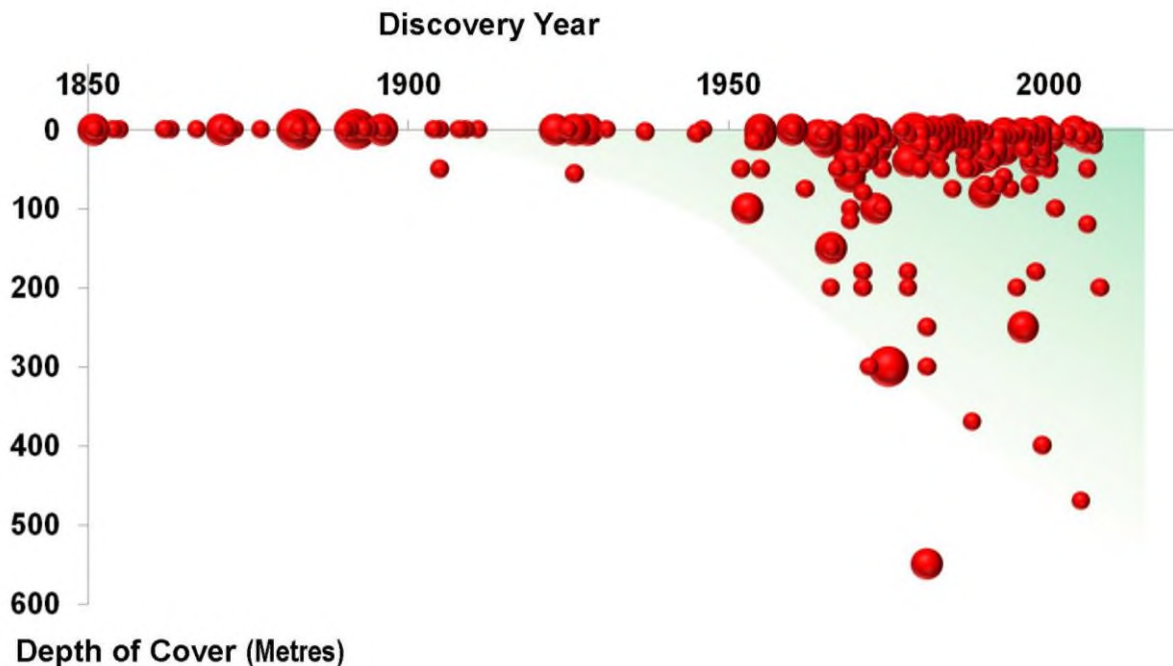
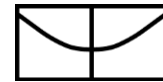


Figura 1-1 – Profundidade de cobertura para as principais descobertas de minerais na Austrália de 1850-2010 (Fonte: Richard Schodde, MinEx Consulting). Os resultados mostram a necessidade de se ampliar a utilização geofísica na prospecção sob a cobertura (“*undercover*”) e que seria a base para a compreensão das respostas geofísicas na ausência de controle geológico.

O conceito de sistema mineral sugere um novo conjunto de objetivos de exploração. Ou seja, a mudança do paradigma convencional que se baseia no mapeamento e na detecção do ambiente mineralizado em escala de depósito para uma compreensão mais abrangente do

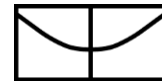


sistema levando em conta o ambiente tectônico, as fontes de metais, os caminhos percorridos pelos fluidos e, por exemplo, os paleo-reservatórios (Wyborn et al., 1994).

O fato é que as mineralizações estão associadas a interações fluido-rocha, em maior ou menor escala, e para avançarmos com o entendimento de dados geofísicos que mapeiam abaixo da cobertura é essencial a aplicação de petrofísica para uma melhor compreensão das respostas de processos geológicos e, em especial, da alteração hidrotermal (Dentith et al., 2020; Dentith & Mudge, 2014).

Vários estudos têm sido desenvolvidos para diferentes sistemas minerais, mas as mineralizações auríferas-uraníferas hospedadas em metaconglomerados representam ainda um desafio, devido à alta complexidade dos depósitos, além do recorrente debate sobre a formação das mineralizações. Depósitos de Au (U) em metaconglomerados são conhecidos em crátons pré-cambrianos em todo o mundo. Dentre várias mineralizações, destacam-se os depósitos associados à bacia arqueana de Witwatersrand no Cráton Kaapvaal, África do Sul. Estas minas são responsáveis por quase um terço da produção global de ouro (Frimmel, 2019, 2014; Frimmel et al., 2019). Outros exemplos incluem Tarkwa no Cráton da África Ocidental (Pigois et al., 2003), o Fortescue Group, no Cráton Pilbara, Austrália (Hennigh, 2016), e o Supergrupo Huronian na Província Superior do Canadá (Whymark & Frimmel, 2018), juntamente com vários outros exemplos com maior ou menor relevância econômica.

No Brasil, os exemplares de mineralizações mais conhecidas hospedadas em metaconglomerados auríferos e uraníferos encontram-se no Supergrupo Minas, na porção meridional do Cráton do São Francisco (Guimarães et al., 2019; Minter et al., 1990), na porção basal do Grupo Jacobina, Bloco Gavião, porção setentrional do Cráton do São Francisco (Ledru et al., 1997; Milési et al., 2002; Teixeira et al., 2001) e nos metassedimentos da Formação Castelo dos Sonhos, Província Tapajós, do Cráton Amazônico (Klein et al., 2017). A área objeto

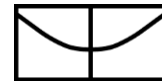


de estudo deste trabalho compreende os depósitos auríferos da Formação Serra do Córrego, unidade basal do Grupo Jacobina (Figura 1-2).

O nível de conhecimento e a quantidade de informações disponíveis em províncias geológicas historicamente conhecidas, permitem que novas abordagens sejam desenvolvidas visando a melhor compreensão da gênese dos depósitos e dos detalhes prospectivos que auxiliarão na procura por novos alvos ou em prolongar a vida útil de distritos maduros. Para tanto, precisamos de informações geoquímicas e mineralógicas para entender estas respostas e avançar nos processos de modelagem e integração de dados. É chave associar o estudo com dados de propriedades químicas, dados mineralógicos e geoquímicos "quantitativos" que irão auxiliar no entendimento do contexto geológico.

Nas últimas décadas, o uso de ferramentas de estatística computacional, *Machine Learning* e inteligência artificial (*latu sensu*) para integração de dados em geociências tem crescido, devido em parte à revolução digital e à consequente maior disponibilidade de equipamentos com boa capacidade de processamento e à tendência crescente de geração e armazenamento de dados das diversas naturezas (Davies, 2002; Flemming et al. 2021). Entretanto, a integração e interpretação conjunta deste grande volume de dados com perspectivas de gerar aplicações práticas na pesquisa mineral ainda é necessária.

Desta forma, esta tese se propõe a avaliar o uso de métodos de *Machine Learning* para integração de dados multifonte em geologia para fins de extração de informações para a pesquisa mineral, baseado em predições explícitas e no conhecimento gerado a partir da verificação da estrutura dos dados geoquímicos, espectrorradiométricos, petrofísicos e mineralógicos em várias escalas de trabalho. Assim, apresentamos os resultados na forma dos artigos técnico-científicos intitulados “*Predicting mineralization and targeting exploration criteria based on a machine learning approach in the Paleoproterozoic Jacobina quartz-pebble metaconglomerate Au-(U) deposit, northern of São Francisco Craton*”, no Capítulo 2,



“*Machine learning applied to the analysis of mineral chemistry in pyrite grains from the Jacobina gold deposits, São Francisco Craton, Brazil: geochemical patterns and implications to mineral exploration*”, no Capítulo 3, e os resultados parciais de um artigo em desenvolvimento denominado “*Unmixing spectral signal and estimating the mineral composition of metaconglomerates using dimensionality reduction and relative distance concepts*” no Capítulo 4. As considerações finais da tese são apresentadas no Capítulo 5.

1.2 Objetivos

O objetivo central desta pesquisa é caracterizar a assinatura dos da mineralização aurífera da Serra de Jacobina através do mapeamento do *footprint* petrofísico, mineralógico e geoquímico, bem como a geração de modelos de vetorização mineral através da integração de dados utilizando aprendizado de máquina.

Objetivos específicos incluem:

1. Gerar um banco de dados geofísicos, geoquímicos e de composição mineral das rochas hospedeiras e mineralizadas de amostras da Formação Serra do Córrego, base do Grupo Jacobina;
2. Adquirir dados petrofísicos qualitativos e quantitativos em furos de sondagem chaves ao longo da Formação Serra do Córrego;
3. Caracterizar a assinatura espectral das alterações hidrotermais associadas à mineralização de ouro na Formação Serra do Córrego;
4. Gerar de informações de prospecção mineral através da integração de dados multifonte que envolvam o reconhecimento de padrões e inteligência artificial (*machine learning*).

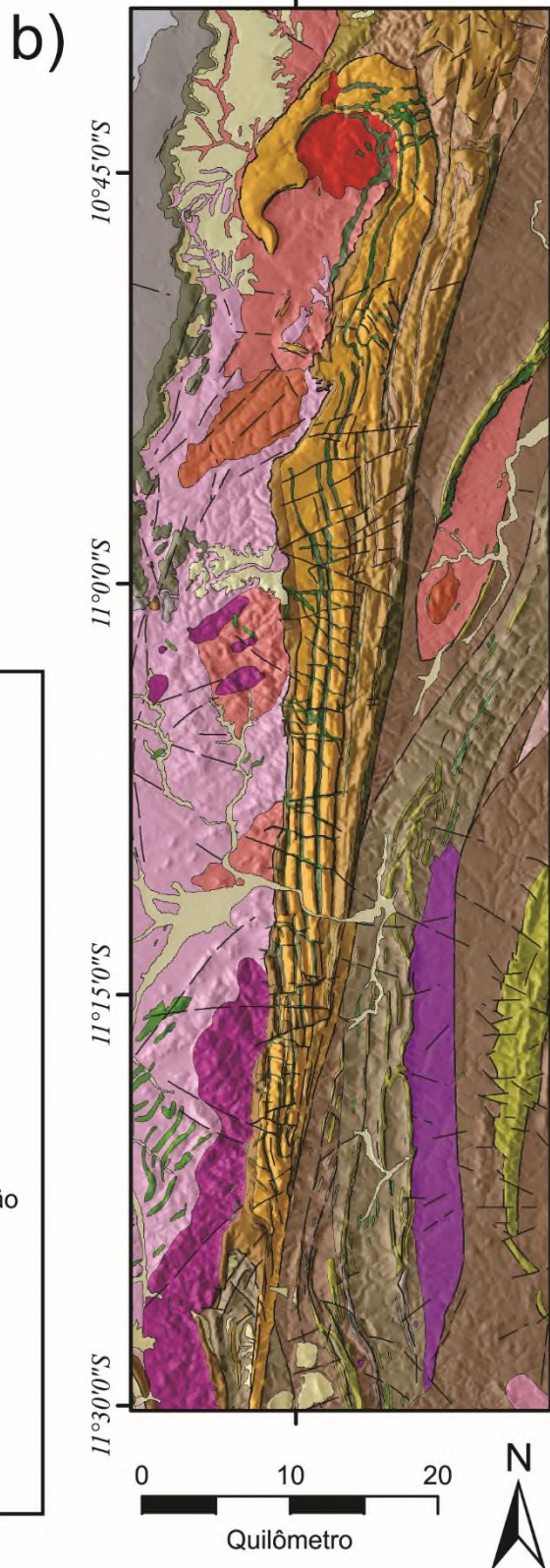
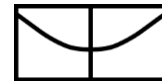


Figura 1-2 – Localização da área estudada, a) posição do cráton São Francisco dentro da plataforma Sul-Americana. O retângulo vermelho indica a localização da área de estudo, na porção setentrional do Cráton São Francisco; b) mapa geológico simplificado da Serra de Jacobina e arredores (modificado de Teles et al. 2015, Reis et al. 2020).



1.3 Materiais e métodos

Os materiais utilizados nesta pesquisa compreendem dados multifonte, como fluorescência de raios-X portátil, susceptímetro magnético e condutímetro elétrico, gamaespectrômetro, balança de precisão para cálculo de densidade, com estação acoplada para medida emersa em fase líquida e espectrorradiômetro de precisão (Figura 1-3)

Com a finalidade de alcançar os objetivos propostos neste projeto, foram desenvolvidas as seguintes etapas e métodos de trabalho, discutidas em maior detalhe nas seções seguintes:

- coleta e descrição de amostras de rocha em testemunhos de sondagem representativos dos níveis mineralizados da Formação Serra do Córrego;
- aquisição de dados de propriedades físicas de rocha;
- aquisição de dados de química de rocha seguindo os protocolos de controle de qualidade indicados;
- aquisição de dados de espectrorradiometria;
- descrição petrográficas em amostras chave para melhor compreensão da dinâmica de formação e inter-relação dos minerais associados à mineralização;
- pré-processamento e análise exploratória de dados;
- integração de dados e escolha do modelo supervisionado adequado

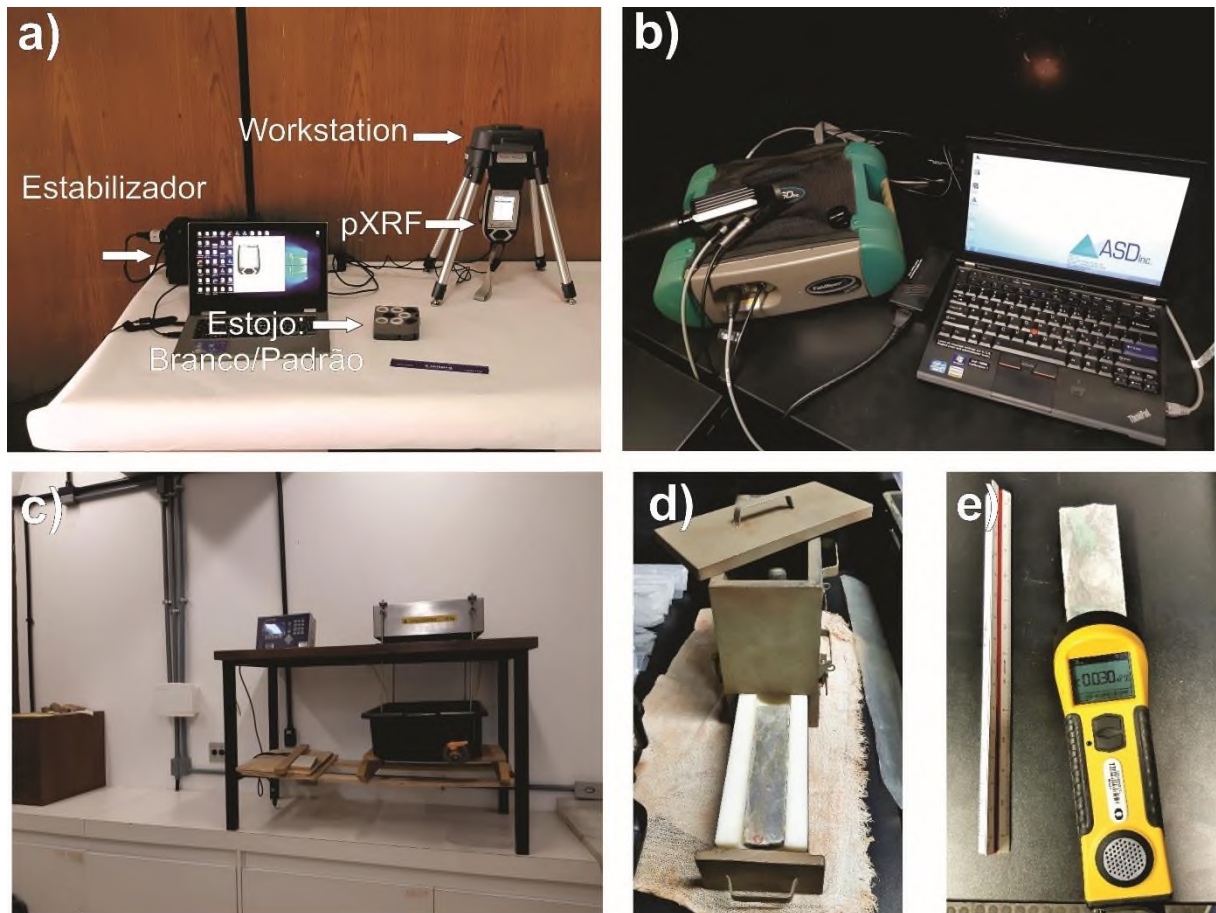
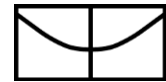


Figura 1-3 – Ilustração das ferramentas utilizadas para obtenção das variáveis numéricas: a) Fluorescência de Raios-X Portátil (pXRF) acoplada a estação de trabalho, b) espectrorradiômetro, c) balança de precisão com plataforma para medida submersa, d) caixa de chumbo para medição de radioelementos, com adaptação para amostras em testemunho de sondagem; Susceptibilímetro KT-10 plus com bobina circular.

1.3.1 Amostragem

Neste trabalho, enfocamos os depósitos de paleoplacer modificados que ocorrem na unidade inferior do Grupo Jacobina, denominada Formação Serra do Córrego. Coletamos 557 amostras de rocha de quatro diferentes testemunhos de sondagem, denominados CANIF-27, JBA-722, CAN-120 e CAN-144. Os furos interceptam toda Formação Serra do Córrego.

As amostras foram obtidas em testemunhos serrados ao meio, coletados em intervalos aproximados de 4 metros, onde foram selecionadas as amostras mais representativas. Em média, as amostras têm 15 centímetros de comprimento, variando de 8 a 50 centímetros, dependendo da representatividade e do tamanho do grão. Predominam amostras de quartzitos e



metaconglomerados, mas também estão descritas amostras de xistos, brechas e rochas meta-ultramáficas.

1.3.2 Propriedades físicas de rocha

As propriedades físicas consideradas nesta pesquisa são densidade, susceptibilidade magnética e condutividade elétrica (Figura 1-4). Os métodos de aquisição dos dados são descritos abaixo.

1.3.2.1 Densidade

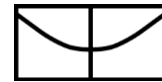
As medidas de densidade foram realizadas no Laboratório de Geofísica Aplicada (LGA), da Universidade de Brasília, com base no método de pesagem hidrostática. Os valores de densidade foram calculados após a obtenção do peso da amostra em balança padrão e após a obtenção do peso imerso em água à temperatura ambiente. Em seguida, a densidade foi calculada com base no Princípio de Arquimedes.

A água no recipiente de medição foi trocada a cada 50 amostras ou após a troca da matriz da amostra. Ao final de cada rodada de medições, algumas amostras foram escolhidas aleatoriamente para repetir o teste de densidade. Em caso de repetição, foi considerada a média dos valores. Se a diferença for significativa (ou seja, superior a $0,2 \text{ g/cm}^3$), repetíamos o procedimento até a estabilização.

Foram coletadas 598 medidas de densidade, considerando as replicatas.

1.3.2.2 Susceptibilidade magnética e condutividade elétrica

As propriedades de suscetibilidade magnética (unidade 10^{-6} SI) e condutividade elétrica (unidade S/m) foram avaliadas com base em um medidor de suscetibilidade magnética Terraplus KT-10 S/C e condutividade com uma bobina circular.



Essas propriedades foram medidas pelo menos dez vezes por amostra, em diferentes pontos ao longo da amostra, e o valor mediano de cada propriedade foi considerado a medida da amostra mais representativa e foi posteriormente usado para a modelagem.

Foram coletadas 7490 medidas de susceptibilidade magnética e 5720 medidas de condutividade elétrica. Para cada amostra, a mediana foi tomada como valor mais representativo, sendo o valor considerado para as transformações subsequentes e análises multivariadas.

A distribuição dos valores das propriedades susceptibilidade magnética e condutividade elétrica tendem a se aproximar de uma distribuição logarítmica (ou seja, valores baixos são mais frequentes do que valores altos). Neste caso, ambas as propriedades foram analisadas em escala logarítmica, com o objetivo de facilitar as interpretações e reduzir a assimetria das distribuições.

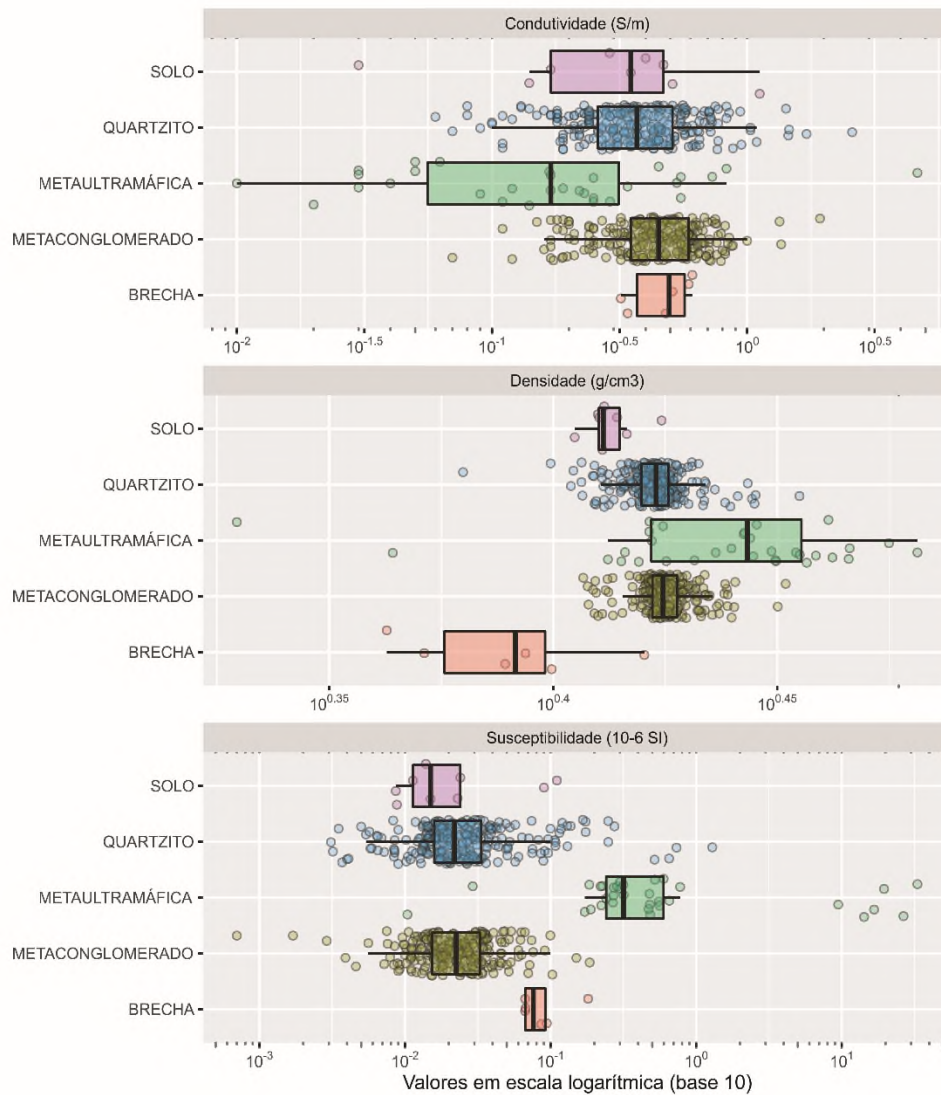
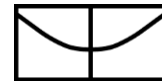


Figura 1-4 – Diagrama de caixa (*boxplot*) para as propriedades físicas adquiridas com transformação para escala logarítmica (base 10). As diferentes cores representam os litotipos considerados na descrição das amostras. Os pontos indicam a posição transformada de cada medida

1.3.2.3 Gamaespectrometria

O conteúdo de radioelementos foi medido em ensaios com uso de um gamaespectrômetro portátil RS-125, tomando os valores da amostra em medições obtidas em uma câmara de chumbo isolada por 300 segundos. Foram coletadas medidas de radioelementos em 304 amostras distintas, compreendendo parcialmente os furos CANIF-27 e JBA-722.

Contudo, os valores obtidos estão geralmente próximos do limite de detecção inferior e estas análises não foram consideradas. Uma interpretação possível é que essa resposta se deve ao baixo volume das amostras e, conseqüentemente, ao baixo nível de radiação, insuficiente



para excitar o cristal do instrumento de maneira adequada. Portanto, apesar de haver alguns minerais radioativos na assembleia, os valores de radioelementos não puderam ser considerados nas modelagens multiparamétricas utilizadas neste trabalho.

1.3.3 *Química de rocha*

Para avaliar as concentrações dos elementos nas amostras, usamos um analisador Thermo Scientific Niton XL3t Gold + XRF, com tubo anódico Au de 2W, 50kV Au e um detector de deriva de grande área geometricamente otimizado. O instrumento foi acoplado em bancada estacionária durante as medições, onde foram colocadas as amostras. Cada medição durou 120s, com 60s de duração para cada feixe. Ao todo, foram realizadas 1950 medidas com o pXRF, considerando as leituras realizadas em padrões certificados, utilizados para o controle de qualidade.

As medições foram realizadas em amostras de testemunhos de sondagem cortados ao meio, usando o modo de ensaio “*point and shoot*”. Os procedimentos adotados de controle de qualidade seguiram as sugestões apresentadas por Fisher et al. (2014) e Piercey (2014). No caso de análises cuja dissimilaridade ultrapassaram 10% do valor de referência para aquele elemento, foi realizada a retificação do valor através da aplicação de um coeficiente calculado a partir da razão da medida inacurada sobre o valor de referência. Este coeficiente então é aplicado em todas as leituras subsequentes (Figura 1.5).

Para verificar a representatividade das informações, coletamos uma segunda medição de cada amostra em um local diferente. Apesar de alguns valores atípicos, a distribuição principal é mantida tanto no primeiro quanto no segundo pontos analisados. Calculamos a média da primeira e da segunda medidas e, em seguida, pegamos o valor médio para a análise dos dados.

Os elementos foram selecionados para a análise multivariada baseado na proporção de valores não censurados (i.e., acima do limite de detecção). As variáveis selecionadas



apresentaram ao menos 75% de valores não censurados. Complementarmente, realizou-se a imputação dos valores censurados baseado na substituição por 50% dos respectivos limites de detecção inferiores.

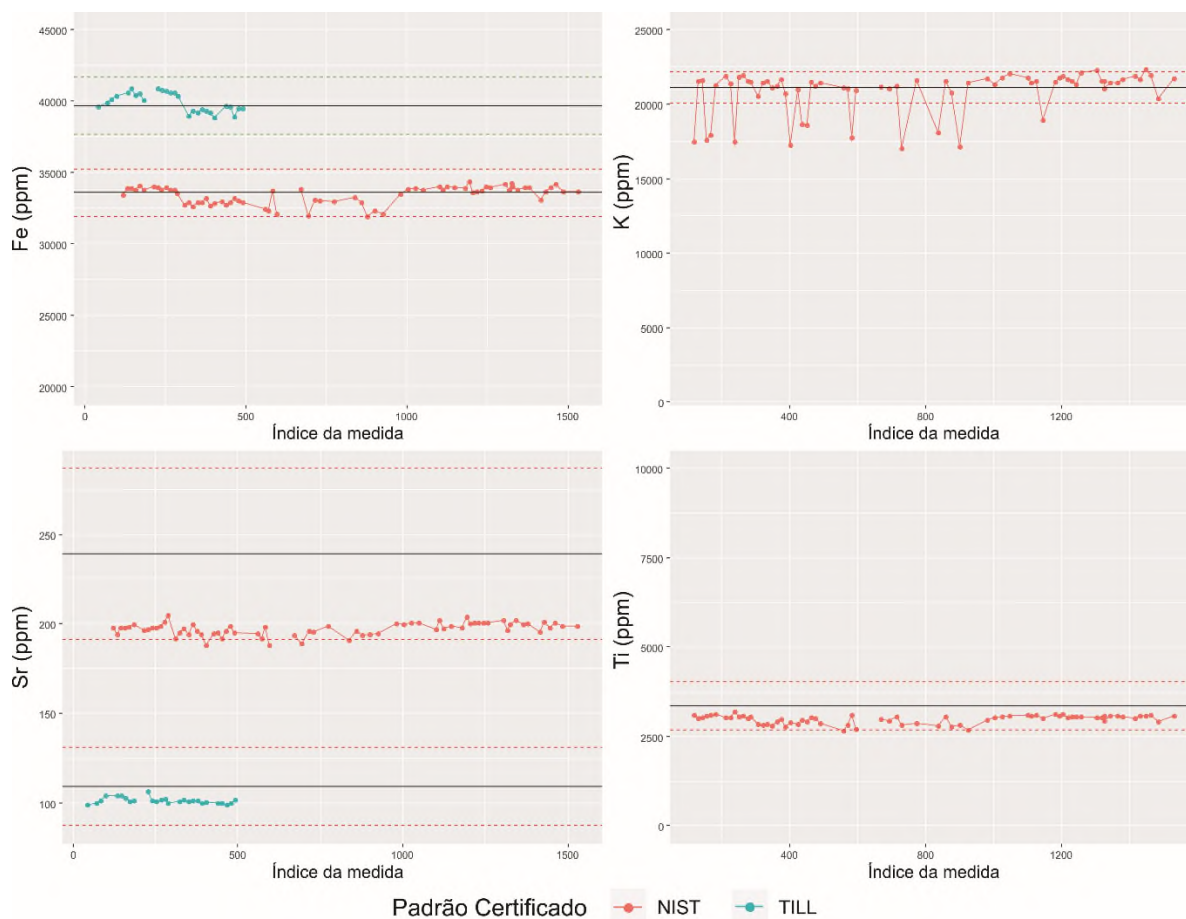


Figura 1-5 – Cartas de controle de Shewhart para as medidas do padrão certificado NIST e TILL, realizadas em determinados durante a análise. A linha pontilhada indica o limite de $\pm 10\%$ do valor referenciado.

1.3.4 Espectrorradiometria

As análises de assinatura espectral foram realizadas com o equipamento ASD FieldSpec 2, instalado em uma sala com controle de luz externa. As análises abrangem faixas do espectro eletromagnético desde o campo da luz visível (350-750 nm), infravermelho-próximo (751-1000 nm) e infravermelho de ondas curtas (1001-2500 nm).

Através da análise das assinaturas espectrais das amostras dos diversos litotipos, pretende-se investigar variações nos padrões relacionados às fases de alteração hidrotermal identificadas nas amostras.



Os dados de reflectância foram processados através de scripts próprios escritos na linguagem R (<http://r-project.org/>). Inicialmente, realizamos a segmentação do dado baseado no algoritmo de *Self-Organizing Maps* (SOM; Kohonen, 1998) que permitiu a redução da dimensionalidade para redução do custo de processamento, que é proporcional ao número de dimensões analisadas. O SOM consiste em uma grade regular bidimensional de nós (também chamados de “unidades” ou “neurônios”). A grade é automaticamente organizada baseado na estruturação dos dados de modo a agrupar as unidades semelhantes (Kohonen, 1998). Para esta finalidade, utilizamos uma grade de 400 unidades dispostas em uma malha quadrada (20 x 20). Testamos o treinamento através de diversas épocas, porém o treinamento atingiu a estabilidade em torno de 2500 iterações

Para agrupar o sinal espectral da reflectância para curvas de reflectância similares, aplicamos um agrupamento hierárquico baseado na distância euclidiana das unidades do SOM. Esta etapa foi necessária devido aos padrões de reflectância altamente complexos, que dificultam a interpretação direta.

1.3.5 Petrografia

Selecionamos 20 amostras mineralizadas para análise microscópica para descrever e verificar as relações texturais e minerais com a mineralização. Estas descrições permitiram o entendimento da relação textural entre as diversas fases minerais encontradas no depósito, inclusive a distinção entre pirritas detríticas (da fase de sedimentação) e epigenéticas (formadas durante o metamorfismo ou evento hidrotermal subsequente).

Adicionalmente, cinco amostras foram analisadas pelo método de espectroscopia de elétron/energia dispersiva retro espalhada (BSE-EDS) para melhor avaliação das texturas petrográficas e composições dos minerais (ver Apêndice F).



1.3.6 *Pré-processamento e análise exploratória de dados*

Considera-se pertencente à etapa de pré-processamento e análise exploratória de dados, todos os passos de organização, estruturação e transformação de dados, testes estatísticos e geração de figuras de exploratórias sobre a distribuição e comportamento dos dados.

Todos estes estágios foram desenvolvidos através da linguagem de programação R (versão 4.1.2). As etapas de manipulação dos dados e geração de gráficos finais foram realizadas através dos pacotes disponíveis na coleção Tidyverse (Wickham, 2014).

As transformações dos dados descritas neste tópico e nos anteriores foram realizadas em linha de comando em ambiente R, através do uso da biblioteca “geoquimica”, desenvolvida durante a execução desta pesquisa pelo candidato e disponibilizada através de repositório virtual (<https://github.com/gferrsilva/geoquimica>).

1.3.6.1 *Verificação da distribuição dos dados*

Uma etapa essencial na análise exploratória de dados é a determinação do tipo de distribuição. Se os dados forem paramétricos, a média e os desvios padrão são estimativas razoáveis para o centro e a dispersão dos dados e várias ferramentas podem ser usadas para analisar e inferir parâmetros populacionais. Caso contrário, os dados devem ser tratados de forma distinta, com base em métodos que não dependem de parâmetros geométricos.

Existem quase 40 testes disponíveis para verificação de normalidade da distribuição, mas vários autores (Razali & Wah, 2011; Saculinggan & Balase, 2013; Yap & Sim, 2011) mostram que o teste de Shapiro-Wilk é o teste preferido para a maioria dos tipos de distribuições e tamanho amostral. O teste de Shapiro-Wilk (Shapiro & Wilk, 1965) foi inicialmente definido para pequenas amostras ($n < 50$), e então foi aprimorado por Royston (1982) que expandiu o teste para uma faixa maior de valores ($3 \leq n \leq 5000$)

Neste trabalho, este teste foi realizado para cada elemento selecionado no banco de dados, considerando os valores da média de cada amostra. A hipótese nula consiste na

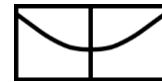


equivalência entre a distribuição dos dados e a distribuição gaussiana, enquanto a hipótese alternativa sustenta que a distribuição dos dados e a distribuição gaussiana não são equivalentes.

Para um nível de significância de 5% e analisando o parâmetro p.value, todos os elementos selecionados apresentaram p.value $\ll 0,01$ (Tabela 1.1), o que significa a rejeição da hipótese nula, e conseqüentemente a respectiva distribuição para cada elemento não é equivalente à distribuição gaussiana, o que se confirmada no diagrama QQ Plot (Figura 1.6).

Tabela 1-1 – Parâmetro estatístico e p.value para o teste de Shapiro-Wilk, com formulação de hipótese nula para a congruência com a distribuição gaussiana

Elemento	Dado bruto (ppm)		Dado Transformado (clr)	
	Estatística	p.value	Estatística	p.value
Zr	0.432434674	1.33E-38	0.914233291	2.77E-17
Sr	0.531424409	4.95E-36	0.971842006	6.60E-09
Cu	0.604893159	8.03E-34	0.868293075	2.45E-21
Ni	0.605874891	8.64E-34	0.818577858	1.16E-24
Fe	0.431073972	1.23E-38	0.962760516	1.02E-10
Cr	0.336519148	9.10E-41	0.982159738	2.32E-06
Ti	0.364852955	3.72E-40	0.974306805	2.35E-08
Ca	0.184832368	1.03E-43	0.951186247	1.17E-12
K	0.38565616	1.08E-39	0.913317044	2.23E-17
S	0.427819613	1.03E-38	0.959805637	3.01E-11
Ba	0.584526537	1.82E-34	0.816863879	9.17E-25
Cs	0.714989228	8.59E-30	0.807393671	2.59E-25
Te	0.734048894	5.67E-29	0.837710298	1.78E-23
Sb	0.713592928	7.51E-30	0.824527425	2.64E-24
Sn	0.750783251	3.25E-28	0.858613897	4.74E-22
Cd	0.758331722	7.36E-28	0.895298162	4.12E-19
Pd	0.73804041	8.53E-29	0.911212683	1.36E-17
Nd	0.870360028	3.53E-21	0.943938652	1.01E-13
Pr	0.891936733	2.07E-19	0.944743494	1.31E-13
Ce	0.758785709	7.74E-28	0.921160061	1.53E-16
La	0.862151281	8.56E-22	0.935599702	7.81E-15
P	0.623349799	3.25E-33	0.846132596	6.41E-23
Si	0.798403208	8.17E-26	0.667709196	1.19E-31
Cl	0.628795993	4.96E-33	0.781690128	1.06E-26



1.3.6.2 Transformação de dados composicionais

Como os dados geoquímicos são considerados variáveis composicionais (Aitchison, 1986), todas as análises geoquímicas selecionadas foram interpretadas após uma transformação da razão logarítmica centrada (*centered log-ratio* abreviada como clr). Essa transformação é uma etapa essencial para analisar dados composicionais com estatísticas multivariadas (Grunsky, 2001). Ainda assim, os dados das variáveis selecionadas não possuem conformidade com a distribuição gaussiana, e uma abordagem não paramétrica se faz necessária para a integração (Figura 1.7).

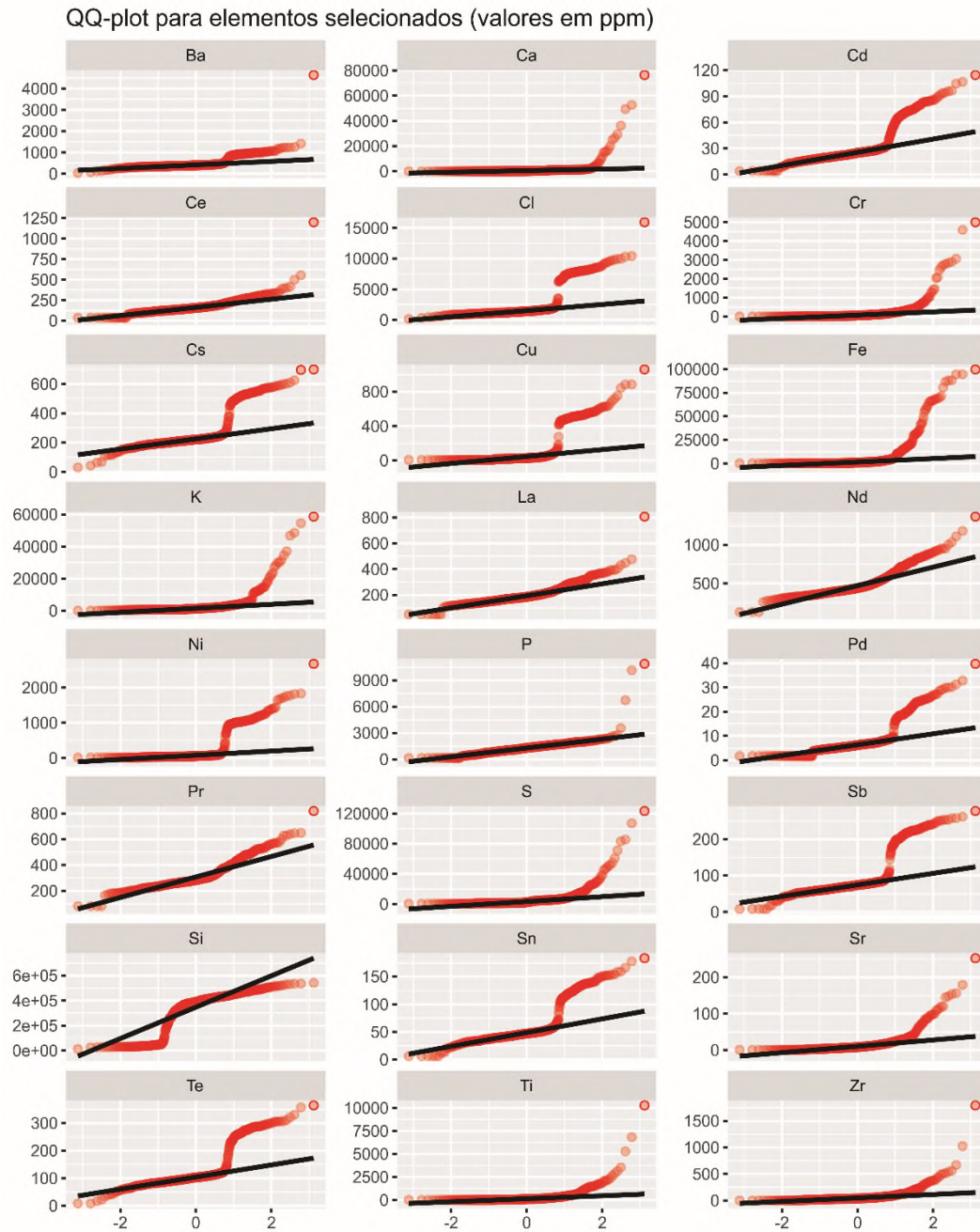
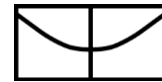
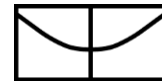


Figura 1-6 – Diagrama Quantil-Quantil (QQ-plot) para os elementos selecionados. A linha preta indica o comportamento esperado para uma distribuição gaussiana. A dispersão dos pontos nos valores extremos da distribuição sugere a presença de valores anômalos, principalmente na porção superior da distribuição. Estes valores anômalos não foram removidos da distribuição justamente por representarem um enriquecimento químico e potencialmente compor parte do objetivo desta pesquisa. A quebra de regularidade na direção de agrupamento dos pontos sugere mudanças na distribuição, o que em alguns elementos caracteriza uma distribuição polimodal. Isto é interpretado como as diferentes distribuições obtidas em litotipos distintos.



QQ-plot para elementos selecionados (transformação por razão logarítmica centralizada)

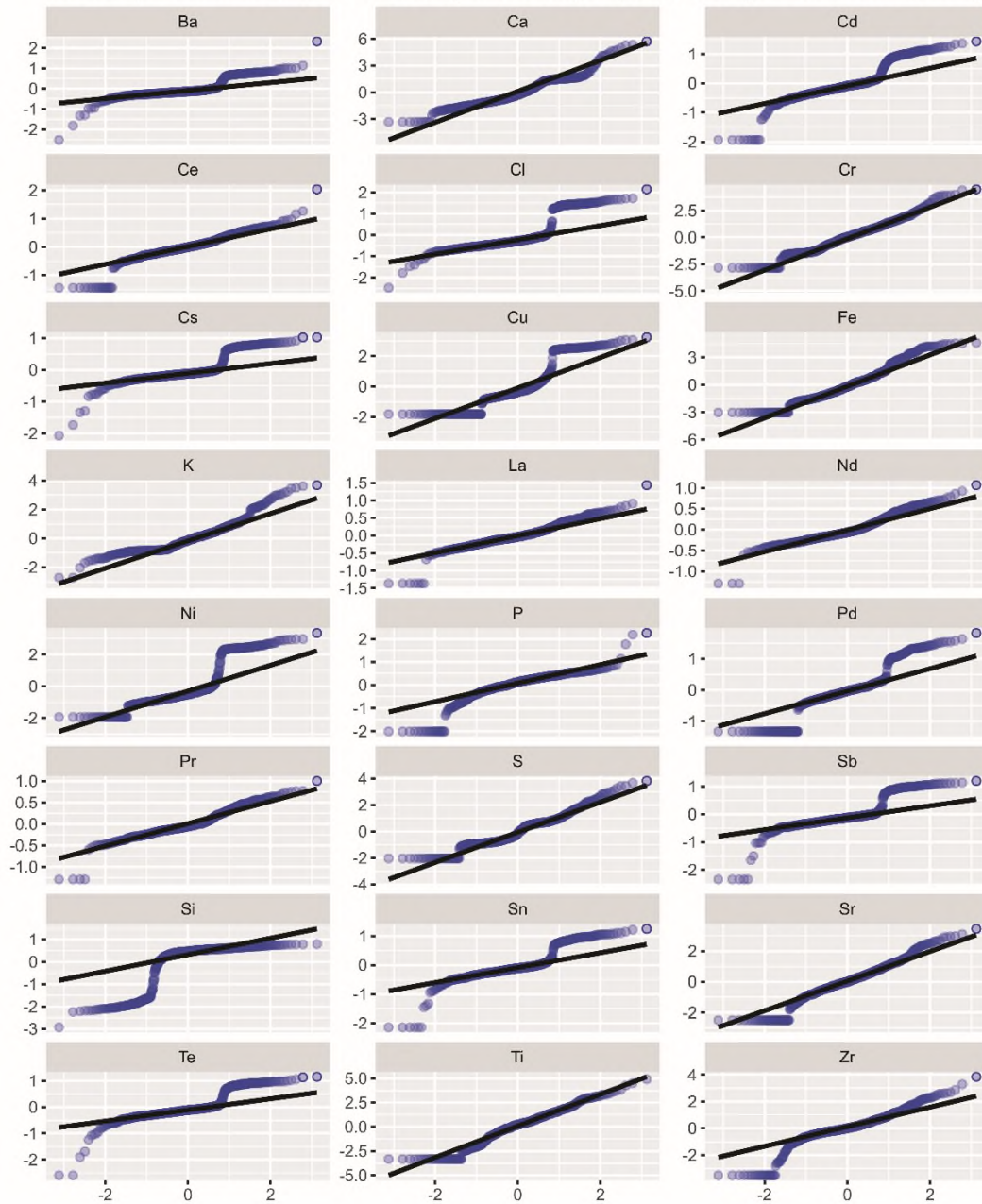
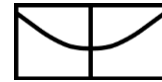


Figura 1-7 - Diagrama Quantil-Quantil (QQ-plot) para os elementos selecionados com valores transformados pelo método da razão logarítmica centralizada. Mesmo após a transformação dos dados, as distribuições não tendem a seguir a distribuição gaussiana (referenciada no gráfico pela linha preta diagonal). Da mesma forma, os valores limítrofes ainda estão em destaque em ambos os lados da distribuição, na maioria dos elementos selecionados e as quebras de tendência de direção de agrupamento também estão presentes.



2 PETROFÍSICA, GEOQUÍMICA E PREDIÇÃO DE MINERALIZAÇÃO

O presente capítulo está apresentado na forma de um artigo, intitulado “*Predicting mineralization and targeting exploration criteria based on machine-learning in the Serra de Jacobina quartz-pebble-metaconglomerate Au-(U) deposits, São Francisco Craton, Brazil*”. O artigo foi aceito para publicação na *Journal of South American Earth Sciences*, aceito para revisão e devolvido com solicitação de correções menores sob o número SAMES-D-21-00534R1.

Este trabalho discorre sobre um modelo de predição de amostras de rocha mineralizadas do Grupo Jacobina, Cráton São Francisco, através da interpretação conjunta de dados petrofísicos e geoquímicos, utilizando modelos de *Random Forests* para a integração dos dados, e geração de um modelo capaz de predizer o status de mineralização das amostras avaliadas.

Apresentamos ainda abordagens inéditas dentro do escopo de predição na prospecção mineral em escala de amostras, como a análise probabilística dos resultados de predição das *Random Forests*. Os resultados foram organizados pela posição das amostras nos respectivos furos. Adicionalmente, discriminamos ainda assembleias minerais existentes em ao menos duas fases hidrotermais presentes na história do Grupo Jacobina e sua relação com a mineralização. O código utilizado na elaboração deste trabalho está disponível no Apêndice C

1 **Predicting mineralization and targeting exploration criteria based on machine-**
2 **learning in the Serra de Jacobina quartz-pebble-metaconglomerate Au-(U) deposits,**
3 **São Francisco Craton, Brazil**

4 Guilherme Ferreira da Silva (Corresponding author) *

5 ¹ Programa de Pós-graduação em Geologia, Instituto de Geociências,
6 Universidade de Brasília, Brasília, DF, Brazil.

7 ² Geological Survey of Brazil (SGB/CPRM), SBN, Quadra 2, Bloco H, Edifício
8 Central Brasília, 2º andar, Brasília, DF, Brazil.

9 E-mail: guilherme.ferreira@cprm.gov.br,

10 OCID: <https://orcid.org/0000-0002-3675-7289>

11
12 Adalene Moreira Silva

13 ¹ Programa de Pós-graduação em Geologia, Instituto de Geociências,
14 Universidade de Brasília, Brasília, DF, Brazil.

15 E-mail: adalene@unb.br

16 ORCID: <https://orcid.org/0000-0001-6290-2374>

17
18 Catarina Labouré Bemfica de Toledo

19 ¹ Programa de Pós-Graduação em Geologia, Instituto de Geociências,
20 Universidade de Brasília, Brasília, DF, Brazil.

21 E-mail: catarinatoledo@unb.br

22
23 Farid Chemale Junior

24 ³ Programa de Pós-Graduação em Geologia, Universidade do Vale do Rio dos
25 Sinos, São Leopoldo, RS, Brazil.

26 E-mail: faridcj@unisinós.br

27 ORCID: <https://orcid.org/0000-0001-5003-5824>

28
29 Evandro Luiz Klein

30 ² Geological Survey of Brazil (SGB/CPRM), SBN, Quadra 2, Bloco H, Edifício
31 Central Brasília, 1º andar, Brasília, DF, Brazil.

32 ⁴ Grupo de Pesquisa em Geologia Econômica, Programa de Pós-graduação em
33 Geologia e Geoquímica, Universidade Federal do Pará, Belém, PA, Brazil

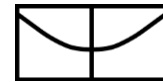
34 E-mail: evandro.klein@cprm.gov.br

35 ORCID: <https://orcid.org/0000-0003-4598-9249>

36
37 * **Corresponding author.**

38 Permanent address: Serviço Geológico do Brasil (SGB/CPRM), SBN, Quadra 2, Bloco
39 H, Edifício Central Brasília, 2º andar, Brasília, DF, Brazil. CEP: 70040-904. E-mail:
40 guilherme.ferreira@cprm.gov.br.

41



42 **Abstract**

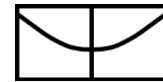
43

44 Defining mineral exploration criteria is a laborious, time-consuming, and generally an
45 empirical task often biased and limited to expert knowledge. To address this problem with a
46 different approach, we used data-driven analysis to make predictions and provide insights about
47 gold mineralization in rocks of the Jacobina Group, São Francisco Craton. The input variables
48 were petrophysical parameters (density, magnetic susceptibility, and electric conductivity) and
49 lithogeochemistry data obtained by X-Ray Fluorescence assays. A machine learning model
50 based on the Random Forests algorithm was applied to predict mineralization in drill core
51 samples. The database used for algorithm training was balanced using the Borderline-SMOTE
52 technique to provide approximately the same numbers of samples of the two classes in the
53 mineral status parameter (i.e., ore and barren samples). The quality of the predictions was
54 assessed with different datasets (i.e., training, testing, each drill core separately, and all
55 samples) and by parameters. The average accuracies were 0.87 for cross-validation training,
56 0.91 for testing, and 0.86 for all samples. Also, the model allowed us to estimate and rank the
57 importance of the input variables to the prediction. These estimates were validated by an
58 interpretation of optical and scanning electron microscopy petrographic analysis, which were
59 carried out to more clearly understand the relationship between minerals of different stages and
60 gold mineralization. As this approach can be easily replicated in mineral exploration, it is
61 feasible to put models like this in production based on numerical and categorical variables
62 obtained routinely.

63 **Keywords:** Gold-bearing quartz-pebble conglomerate; Modified Paleoplacer; Supervised
64 Machine Learning; Hard rock petrophysics; Random Forests.

65 **2.1 Introduction**

66 As more modern exploration techniques are developed, managing, processing, and
67 interpreting all information generated during the exploration workflow presents a significant

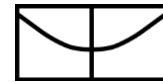


68 challenge. This does not change after the confirmation of positive results for mineralization.
69 With the massive amount of data generated in all stages of the mineral industry, defining, and
70 validating mineral exploration criteria turns into a laborious, time-consuming, and generally
71 empirical task (also often biased and confined to expert knowledge).

72 However, in the past few years, various machine learning algorithms (MLA) have been
73 used to deal with geological datasets and recurring tasks in several branches within geosciences,
74 including lithology prediction and segmentation (Bérubé et al., 2018; da Silva et al., 2022; Hall,
75 2016; Saporetti et al., 2018), semi-automated geological mapping (Carneiro et al., 2012; Costa
76 et al., 2019; Harris and Grunsky, 2015), mineral formula calculation (da Silva et al., 2021; Li
77 et al., 2020) and mineral potential modeling (Carranza and Laborte, 2016, 2015a; Prado et al.,
78 2020; Rodriguez-Galiano et al., 2015). Nevertheless, the majority of works using MLA in the
79 mineral exploration concerns the selection of areas with potential for mineralization (i.e.,
80 finding new targets), instead of focusing on integrating information to aid in the explanation of
81 a previously known mineralization and the ranking of mineral exploration criteria (Carranza
82 and Laborte, 2015b, 2015a; Chen, 2015; Ford, 2019; Niiranen et al., 2019; Saljoughi and
83 Hezarkhani, 2018; Yousefi and Nykänen, 2016; Zuo, 2017).

84 Therefore, this article aims to assess the quality of a predictive model based on MLA
85 built with quantitative variables collected in split drill-core samples to predict the mineralization
86 status (i.e., the sample prediction as Ore or Barren). We also discuss the footprint signature of
87 the Au-(U) mineralized samples of the quartz-pebble metaconglomerates of the Jacobina Group
88 through the analysis of petrophysical geochemical data by the variable's importance rank
89 obtained by the data-driven model.

90 Furthermore, we evaluate the mineralization through a probabilistic approach,
91 presenting the odds of a mineralized sample. This method can provide insights into probability



92 distribution across the drill core samples ordered by its depth. This probabilistic approach
93 allows investigating the problem in a binary way.

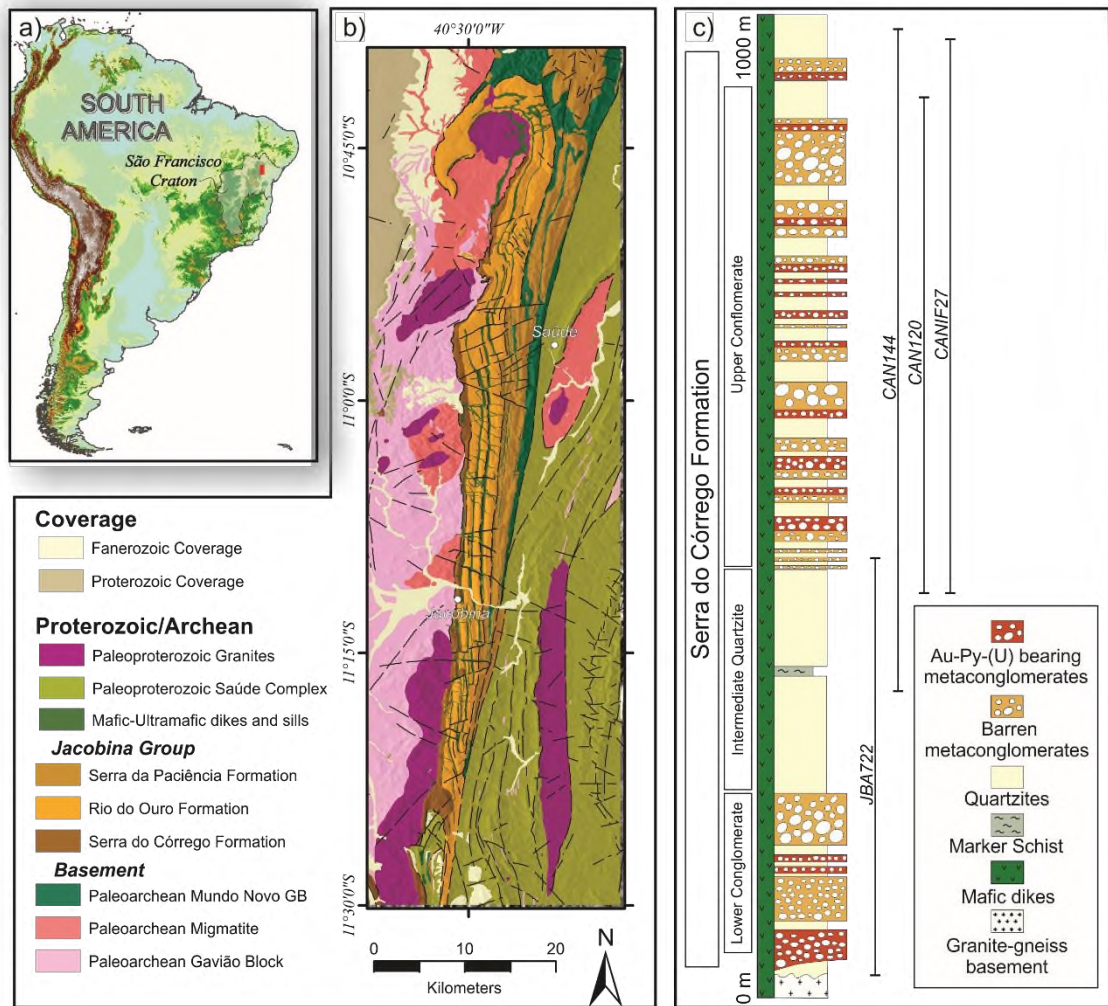
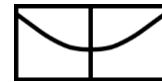
94 Our work brings an example of the prediction of mineralization status based on
95 quantitative petrophysical and lithochemical variables on gold-bearing rocks from the Jacobina
96 Group. Additionally, the MLA helps estimating exploration criteria based on statistical
97 parameters evaluated upon the evaluation metrics (e.g., model accuracy). Furthermore, those
98 parameters were validated after field descriptions and laboratory analysis of mineralized
99 samples.

100 **2.2 Geology and gold mineralization in the Serra de Jacobina**

101 *2.2.1 Geological setting*

102 The Serra de Jacobina (or Jacobina Range) contains Au deposits and is the
103 geomorphological expression of the Jacobina Group, with more than 170 km long and up to 12
104 km wide, set in the northern portion of the São Francisco Craton in eastern Brazil (Figura 2-1).
105 The Jacobina Group “lies on the eastern edge of the Paleoproterozoic Gavião Block, close to the
106 suture zone derived from the Paleoproterozoic collision with the surrounding terrains” (Barbosa
107 and Sabaté, 2004; Heilbron et al., 2017; Santos et al., 2019; Teixeira et al., 2017).

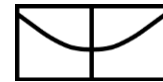
108 According to many authors (Alkmin et al., 1993; Barbosa and Barbosa, 2017; Santos et
109 al., 2019; Leite et al., 2007; Leite, 2002; Teixeira et al., 2001; Teles et al., 2020) the buildup of
110 the Serra de Jacobina occurred during the Paleoproterozoic orogeny developed by the
111 amalgamation of the Gavião, Serrinha, and Jequié paleoplates, between 2.1 Ga and 1.91 Ga.



112

113 Figura 2-1 a) Localization of São Francisco Craton in South America. The red rectangle indicates the position of
114 the Jacobina Group. b) Simplified geological map of the Jacobina Range and its surroundings (modified after
115 Santos et al., 2019; Reis et al., 2021; Teles et al., 2015 and the references therein); c) Serra do Córrego Formation
116 stratigraphy and an indication of the drill core approximate position (after Teles et al., 2015)

117 The Jacobina Group was deposited in a pre-GOE stage (> 2.3 Ga.), and the source of
118 the sediments are TTG (Tonalite-Trondhjemite-Granodiorite) with a contribution of mafic-
119 ultramafic rocks as indicated by high concentrations of Cr in lithochemistry analysis (Teles
120 et al., 2015). Teles et al. (2015, 2020) also describe rounded sedimentary pyrite grains
121 associated with other detrital minerals, consistent with a pre-GOE deposition. In addition, the
122 source rocks have Paleoproterozoic age (between 3.3 and 3.4 Ga), indicated by the U-Pb analysis
123 of several rounded zircons crystals (Teles et al., 2015). In the past, some authors considered the
124 possibility of foreland basin, as the tectonic setting for the Jacobina Group (Ledru et al., 1997;

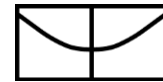


125 Leite and Marinho, 2012), despite that many works attributed it to inverted rift setting (Santos
126 et al., 2019; Pearson et al., 2005; Teixeira et al., 2001; Teles et al., 2015, 2020, and others).

127 There is some controversy in the literature about the extension of the Jacobina Group,
128 as some authors consider the Bananeira and Cruz das Almas formations as part of the Jacobina
129 Group, and not consider the Serra da Paciência (Leite et al., 2007; Leite, 2002; Leo, 1964;
130 Mascarenhas et al., 1998; Miranda et al., 2021a; Reis et al., 2021). The interpretation of
131 Jacobina Group adopted in this work comprises quartz-pebble metaconglomerates, quartzites,
132 and schists, from the base to the top of the stratigraphy, and includes the Serra do Córrego, Rio
133 do Ouro and Serra da Paciência formations, following Santos et al. (2018), Teles et al. (2015)
134 and Teles et al. (2020). Numerous mafic-ultramafic dikes and sills intersect those rocks, mainly
135 metamorphosed and partially serpentinized.

136 The Serra do Córrego Formation, base unit of the Jacobina Group, consists of an
137 alluvial-fluvial formation and comprises an association of quartz-pebble metaconglomerates
138 and quartzites (Santos et al., 2019). The metaconglomerates are mainly oligomictic, composed
139 primarily of quartz pebbles, but some polymictic varieties are found in the Upper Conglomerate
140 Unit with clasts consisting of lithic fragments (granite, quartzite, and metachert of different
141 colors). The metaconglomerates are commonly green-colored due to the presence of fuchsite
142 but can also be yellow-greyish and red-colored, depending on the amount of fuchsite or the
143 degree of oxidation (Mascarenhas et al., 1998).

144 The Rio do Ouro Formation is exposed in the central part of Serra da Jacobina and is
145 locally 2000 m thick (Teles et al., 2015). The unit consists primarily of high-purity fine-to-
146 medium quartzite, but also thin layers of metaconglomerates are presented in the base, making
147 gradational contact with the lower Serra do Córrego Formation (Pearson et al., 2005; Teles et
148 al., 2015).



149 The Serra da Paciência Formation is exposed along the eastern margin of the Jacobina
150 Range. It consists of thick packages of fine-to-coarse-grained quartzite, conglomeratic
151 quartzites, subordinate metaconglomerates with blue quartz grains of possible volcanic origin,
152 and local andalusite-quartz-graphite schist layers (Pearson et al., 2005; Teles et al., 2015).

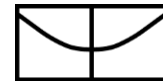
153 *2.2.2 Deformation, metamorphism, and hydrothermal alteration*

154 Two main deformation phases were identified in the Jacobina Group. The first
155 deformation phase (D₁) is compressional, representing the tectonic transport from east to west,
156 which progressively evolved to a sinistral transpressive phase (D₂; Santos et al., 2019). The
157 second deformational phase (D₂) occurred due to the rotation of the compressive vector to a
158 SE-NW orientation (Santos et al., 2019). In addition, the Jacobina Group is separated from the
159 Saúde Complex by the Pindobaçu Fault System, that has transcrustal dimensions and conforms
160 to the Pindobaçu Suture (Santos et al., 2019).

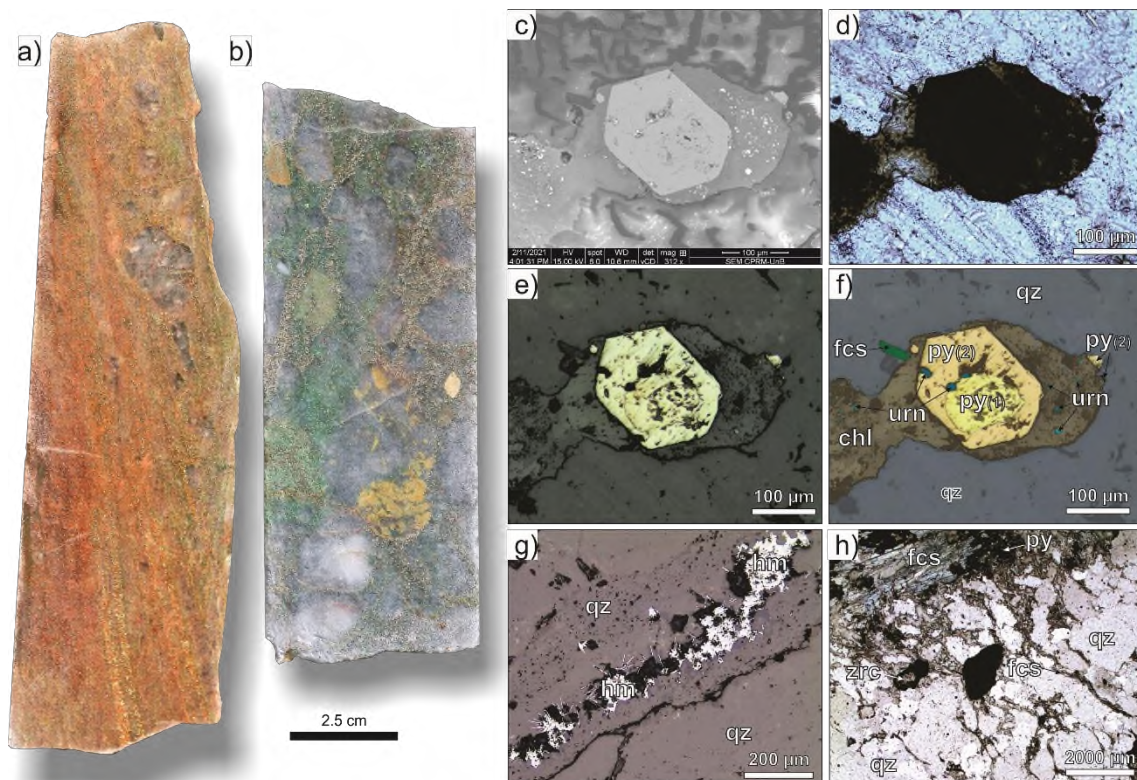
161 Several metamorphic minerals are described in the rocks from the Jacobina Group and
162 in the cross-cutting mafic-ultramafic dikes (Milési et al., 2002; Miranda et al., 2021; Santos et
163 al., 2019; Teixeira et al., 2001; Teles et al., 2015). Fuchsite, chlorite, epidote, uraninite (the last
164 three are more common at the base of the Jacobina Group), serpentine (in the mafic-ultramafic
165 dikes), andalusite, and graphite (at the top of the Serra da Paciência Formation; Pearson et al.,
166 2005).

167 Additionally, sigmoidal-shaped clasts and pressure-shadows (mainly involving
168 inclusion-rich quartz, fuchsite, pyrite, or chlorite crystals) are observed in hand samples and
169 thin sections, suggesting that the metamorphic minerals were produced during a ductile
170 deformation phase (either D₁ or D₂), as shown in Figura 2-2a to f.

171 Despite da Costa et al. (2020) suggestion for the sedimentary related Witwatersrand
172 Group, Teles et al. (2020) found epigenetic pyrite crystals associated with the metamorphism
173 in the Jacobina Group. In addition, Teles et al. (2020) suggest that epigenetic pyrite grains,

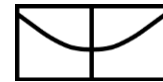


174 locally associated with remobilized gold, are the most common variety of crystals at the Serra
175 do Córrego and Rio do Ouro formations. Teles et al. (2020) associate this phenomenon with the
176 Paleoproterozoic orogeny (2.1 Ga. To 1.9 Ga) that affected the Jacobina Group. The other
177 hypothesis presented by Teles et al. (2020) is that the epigenetic pyrite could be formed by the
178 influence of the emplacement of granitic bodies in the same era. Both events may have provided
179 metamorphic and/or hydrothermal fluids and the heat required for the recrystallization of pre-
180 existing pyrite. Epigenetic pyrite (Py₂) with inclusions of uraninite that overgrow detrital pyrite
181 (Py₁) is here presented. The Py₂ are spatially associated with chlorite corona texture in a
182 shadow-pressure shape. In this texture, inclusion-bearing crystals of uraninite and tiny crystals
183 epigenetic pyrite occur (Figura 2-2c to f).



184

185 Figura 2-2 – Mineralized drill core samples and photomicrographs: a) pyrite-bearing metaconglomerate, with
186 deformed clasts of quartz and a pervasive iron-oxide alteration (orange and reddish colors), including among the
187 pyrite layers (light yellow); b) polymictic (clasts of quartzite, fuchsite quartzite, metachert, and lithic fragments)
188 pyrite-bearing metaconglomerate with fuchsite (green mica); c) to f) sigmoidal-shaped chlorite in a corona texture
189 around authigenic pyrite (Py₂) overgrown on detrital pyrite (Py₁), with uraninite inclusions in the chlorite and
190 authigenic pyrite (tiny white crystals in the Back Scattered Electron image - BSE). BSE image, photomicrography
191 (transmitted polarized light), photomicrograph (reflected polarized light) and drawn interpretations (respectively);
192 g) euhedral hematite in a quartz vein, locally replaced by goethite; h) crystals of fuchsite (green mica) among the



193 recrystallized quartz grains with a euhedral aggregate of pyrite (opaque minerals) spatially associated to a large
194 crystal of fuchsite. Abbreviations: qz – quartz, fcs – fuchsite, py – pyrite, urn – uraninite, chl – chlorite and hm –
195 hematite.

196 The last alteration described is oxidation. Despite being associated with rock
197 weathering, as suggested by the replacement of pyrite crystals by goethite and other iron oxides,
198 this alteration also forms euhedral hematite crystals, disseminated along with some layers or
199 crystallized within quartz veins (Figura 2-2g). Pearson et al. (2005) firstly described hematite
200 alteration as a late alteration stage in the Serra do Córrego Formation.

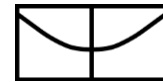
201 This alteration may vary the intensity depending on the proximity of some brittle
202 structures, and in some cases, breccias and other brittle structures seem to control the alteration.
203 In addition, euhedral crystals of hematite and goethite are more likely to be found closer to fault
204 zones. Despite this, the relation of this alteration with gold mineralization is uncertain in the
205 absence of more evidence.

206 2.2.3 *Gold mineralization*

207 Gold in the Serra do Córrego Formation quartz-pebble metaconglomerates occurs as
208 fine-grained native gold with pyrite or hematite, predominantly in the matrix of coarser
209 metaconglomerates (Pearson et al., 2005). There are two mineralized reefs (following the
210 terminology used in mining) within this formation, named lower and upper conglomerate
211 (Figure 1c; Teixeira et al., 2001).

212 Exploration in the Jacobina Range has occurred since the early 18th century, with
213 numerous artisanal miners (“*garimpeiros*”; Pearson et al., 2005). In modern days, three mines
214 of the Serra de Jacobina produced approximately 700,000 ounces of gold during the years of
215 1983-1998, and approximately 2 million of ounces during the years of 2003-2019 (Yamana,
216 2020).

217 In addition, there are small, clearly hydrothermal gold deposits associated with quartz
218 veins hosted in rocks of the Rio do Ouro Formation (Teixeira et al., 2001; Miranda et al., 2021).
219 These deposits may occur in different contexts but typically contain free gold spatially



220 associated with pyrite, chlorite, epidote, tourmaline, chalcopyrite, and other minor phases,
221 hosted in meta-ultramafic rocks, quartzites, and metaconglomerates.

222 Many works have discussed the origin of the gold in the Jacobina Group in the past few
223 decades (Ledru et al., 1997; Mascarenhas et al., 1998; Milési et al., 2002; Miranda et al., 2021;
224 Pearson et al., 2005; Reis et al., 2021; Teixeira et al., 2001; Teles et al., 2020). Although some
225 authors advocate for a pure hydrothermal origin (Ledru et al., 1997; Milési et al., 2002; Pearson
226 et al., 2005), recent evidence observed in rocks of the Serra do Córrego Formation, such as age
227 of the zircons, presence of detrital pyrite, and isotopic chemistry (Teles et al., 2015, 2020)
228 favored the “modified placer model” (Frimmel, 2019; Teixeira et al., 2001; Ledru et al., 1997;
229 Frimmel et al., 1993; Robb and Meyer, 1991), in which placer mineralization was followed by
230 mobilization of ore components by post-depositional fluids.

231 **2.3 Materials and methods**

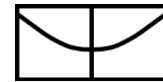
232 *2.3.1 Drill core samples*

233 We collected 557 rock samples from four drill cores intersecting the Serra do Córrego
234 Formation (please refer to Fig. 1c). Quartzites and metaconglomerates predominate, but there
235 are also schists, breccias, and meta-ultramafic rocks samples (Tabela 2-1).

236 Tabela 2-1 – Lithology distribution of samples through the drill cores.

Hole	Breccia (n: 6)	Meta- Conglomerate (n: 251)	Quartzite (n: 266)	Schist (n: 02)	Meta- Ultramafic (n: 32)
CAN120	0	62	51	0	3
CAN144	0	60	78	1	3
CANIF27	0	91	27	0	11
JBA722	6	38	110	1	15
Total	6	251	266	2	32

237



238 The samples were obtained on split drill cores and systematically collected with 4 meters
239 intervals, where the most representative samples were selected. On average, samples have 15
240 centimeters in length, and however, some samples vary from 8 to 50 centimeters in length,
241 depending on the representativeness and grain size.

242 The mineralization status of samples was classified based on the discretization of the
243 gold content in the original drill core's assays. All samples were labeled as "Ore" or "Barren"
244 based on the given threshold of 1 ppm (please refer to Tables A1 and A2 in the supplementary
245 data online for more details).

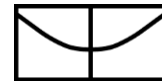
246 We described textural and mineral assemblages related to mineralization optical
247 microscopy. Furthermore, the selected samples were analyzed by the Back-scattered
248 Electron/Energy Dispersive Spectroscopy (BSE-EDS, FEI QUANTA 450) method to
249 complement the petrographic texture and mineral composition study.

250 *2.3.2 Petrophysics*

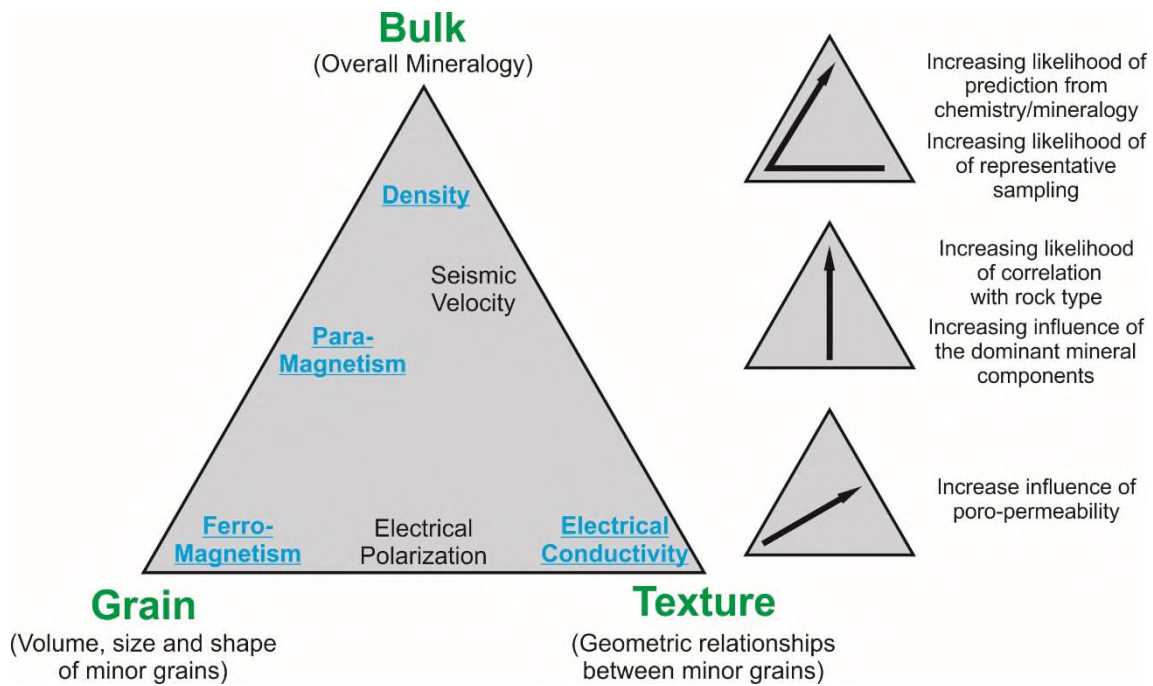
251 The mineral exploration industry often uses petrophysical data to evaluate and
252 characterize geological variations, including mineral changes caused by mineralization (Dentith
253 et al., 2020).

254 According to Dentith et al. (2020), to understand geological controls on physical
255 properties in hard rocks environment, it is necessary to analyze petrophysical data in terms of
256 the properties of different rock types. Still, it is also required to interpret data based on rock and
257 mineral characteristics, as alteration, metamorphism, and strain.

258 As descriptions used in geological exploration are usually categorical data, the measure
259 of rock properties can provide resources to numerically quantify the variations across the host
260 rocks and target mineralization. Still, the necessity of relying on quantitative variables related
261 to the target phenomena turns the MLA into a practical approach to build the model and assess
262 the accuracy of the predictions.



263 We acquired properties related to the three types of petrophysical behavior that can
264 respond to changes in overall mineralogy, texture, and grain-size variations (Figura 2-3). Thus,
265 we measured the density, magnetic susceptibility, and electrical conductivity, all of them
266 described in the following section.

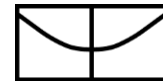


268 Figura 2-3 – Conceptual framework describing the behavior of various physical properties. The physical properties
269 highlighted in the graph were considered for this work (based on Dentith et al., 2020).

270 2.3.2.1 Density

271 The density measurements were taken in the Laboratory of Applied Geophysics (LGA),
272 at the University of Brasília, based on the hydrostatic weighing method. First, the density values
273 were calculated after obtaining the sample weight on a standard scale and after obtaining the
274 weight immersed in ambient temperature water. Then, the density was calculated based on the
275 Principle of Archimedes.

276 The water in the measuring vessel was changed every 50 samples or after changing the
277 sample matrix. At the end of each round of measurements, some samples were chosen randomly
278 to repeat the density test. In the case of repetition, the average of the values was considered.



279 We repeated the procedure until stabilization if the difference was not significant (i.e., higher
280 than 0.2 g/cm³).

281 2.3.2.2 *Magnetic susceptibility and electrical conductivity*

282 The magnetic susceptibility (unit 10⁻⁶ SI) and electric conductivity (unit S/m) properties
283 were assessed using a Terraplus KT-10 S/C magnetic susceptibility and conductivity meter with
284 a circular coil.

285 These properties were measured at least ten times per sample, at different points
286 throughout the sample, and the median value of each property was considered the most
287 representative sample's measure and was used as input for the modeling.

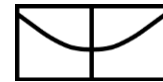
288 Due to the nature of the distribution of the values (i.e., low values are more frequent
289 than high values), both properties were transformed to a logarithm scale, aiming to ease the
290 interpretations and reduce the asymmetry of data distribution.

291 2.3.3 *X-Ray Fluorescence*

292 To assess the concentrations of elements in the samples, we used a Thermo Scientific
293 Niton XL3t Gold+ XRF analyzer, with 2W, 50kV Au anode tube, and a geometrically
294 optimized large area drift detector. The instrument was coupled on a stationary test stand during
295 the measurements, where the samples were placed. Each measurement took 120s, with 60s of
296 duration for each beam.

297 We collected a second sample measurement in a different spot in each sample to check
298 the representativeness of information. Despite some outlier values, the main distribution is
299 maintained in the first and second analyzed spots. We calculated the average of the first and
300 second measurements, and then we took the average value for the data analysis.

301 The QA/QC adopted procedures that followed the suggestions presented by Fisher et al.
302 (2014) and Piercey (2014). The measurements were performed on the sawn surface of split drill
303 core samples, using the “point and shoot” assay mode. The reference material RM 180-646 was



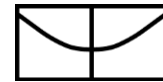
304 read between every ten samples (or 20 spot measurements). At the beginning of every analysis
305 batch, the instrument was set to read continually for at least 40 minutes to correct the
306 instrumental drift caused by variations on the cathode temperature (Ida, 2004; Thermo-
307 Scientific, 2013). We collected a total of 1114 measurements, excluding the reference material
308 analysis.

309 Geochemistry data are considered compositional variables (Aitchison, 1986), and all
310 selected geochemical variables were interpreted after a centered-log ratio transformation. This
311 transformation is essential for analyzing compositional data with multivariate statistics
312 (Grunsky, 2001).

313 2.3.4 *Machine learning analysis (MLA): Random Forests*

314 One of the most employed MLA in geoscience prediction problems is the Random
315 Forests (RF - Breiman, 2001). RF combines several independent estimators (decision trees) to
316 build classification or regression models through bootstrap aggregation. We processed data and
317 built the RF models in the R programming language, using the Tidyverse collection of packages
318 for data wrangling (Wickham, 2014) and the randomForest package ([https://cran.r-](https://cran.r-project.org/web/packages/randomForest)
319 [project.org/web/packages/randomForest](https://cran.r-project.org/web/packages/randomForest)) for modeling.

320 RF relies on the bootstrap principle (random sampling methods with reposition) and the
321 law of large numbers statistical principle, which tells that as the sample size gets large enough,
322 its mean gets closer to the actual average of the whole population (Breiman, 2001). That
323 statistical principle indicates that the RF algorithm does not overfit by considering more
324 estimators. Also, another of RF most significant advantages is that this algorithm has a high
325 performance combined with the low required numbers of hyperparameters, easy to tune. Several
326 geoscientific researchers have shown that the RF outperformed the other MLA, such as support
327 vector machines, artificial neural networks, and logistic regression (e.g., Kuhn et al., 2018;
328 McKay and Harris, 2016; Rodriguez-Galiano et al., 2015). These characteristics make RF



329 widely and effectively used (e.g., Carranza and Laborte, 2015b; Costa et al., 2019; Ford, 2019;
330 Hariharan et al., 2017; Harris et al., 2015).

331 In this work, we employ the RF algorithm to predict the mineralized samples, presented
332 in section 2.4.2, to evaluate the ore probability across the drill core samples, presented in section
333 2.4.3, and to build insights about mineralization targeting and signature by the analysis of
334 variable importance rank, discussed in section 2.5.2.

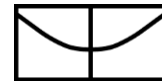
335 2.3.4.1 *Data balancing, train, and test splits*

336 Imbalanced data is a common issue in the implementation of predictive models. This
337 issue considerably reduces the capacity of the model to perform predictions, especially for the
338 minority classes, where the recognition rate decreases considerably (Japkowicz and Stephen,
339 2002). Therefore, resampling data configures a mandatory pre-processing step for a successful,
340 high-performance machine learning model (Koziarski et al., 2020).

341 This problem was previously discussed by (Prado et al., 2020) under the geoscience
342 scope, and the resolution of this issue, or the implementation of a data balancing stage,
343 significantly increased the accuracy of the discussed mineral potential model.

344 A reasonable solution to this problem was first introduced by Chawla et al. (2002),
345 presenting the SMOTE algorithm (Synthetic Minority Over-sampling Technique). This
346 technique over-samples the minority class by creating synthetic examples rather than over-
347 sampling by replacement. The new generated synthetic data is obtained by the linear
348 combination of previous data, considering the k -nearest neighbors of actual samples, until the
349 number of minority and dominant classes is achieved.

350 In this work, we used a modified version of the SMOTE algorithm, called Borderline-
351 SMOTE (Chawla et al., 2002; Han et al., 2005; Koziarski et al., 2020; Prado et al., 2020), to
352 balance samples classified as Ore ($n = 49$) and Barren samples ($n = 523$). The Borderline-
353 SMOTE method priorly classifies the minority data sample into two different subsets, named

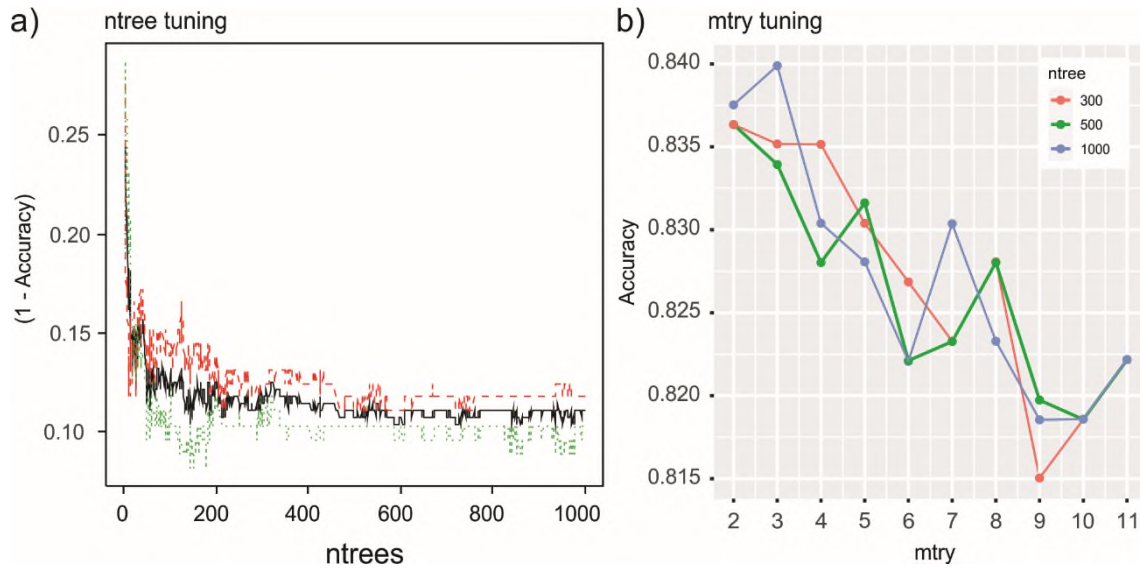
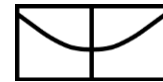


354 DANGER and SAFE samples. This approach assumes that the samples near the limit between
355 minority and dominant classes (DANGER zone) are more easily misclassified than those far
356 from the borderline (SAFE zone, Chawla et al., 2002; Han et al., 2005; Koziarski et al., 2020;
357 Prado et al., 2020). Thus, the algorithm first identifies the borderline between minority
358 examples, and then synthetic examples are generated and added to the original dataset to
359 strengthen the border.

360 To avoid overfitting, the model by resampling bias, we first randomly sample the data
361 evaluated by the Borderline-SMOTE approach, with a 300% oversampling rate that virtually
362 equalized the number of minority dominant samples. Then, we split the generated dataset by
363 the proportion of 0.7 of original data to separate the data into a train (281 samples) and test
364 dataset (120 samples), considering the original distribution of minority and dominant samples.
365 These data were used to train and assess the model's performance according to sections 2.3.4.2
366 and 2.3.4.3, respectively.

367 2.3.4.2 Model tuning

368 The model tuning is a prior stage to select the optimum hyperparameters of the chosen
369 machine learning model. For RF, the mandatory hyperparameters to tune are the number of
370 estimators (*ntree*, or the number of considered trees) and the number of features randomly taken
371 in each tree (*mtry*), as shown in Figura 2-4.



372

373 Figura 2-4 – Hyperparameters evaluation and definition of optimum values: a) ntree (number of estimators) search
 374 according to the error rate. The error starts to stabilize for ntree greater than 500. The red curve represents the error
 375 rate for the barren samples, while the green curve represents the ore samples, and the black curve is the average
 376 error; b) feature values and mean accuracy on a grid-search optimization in the cross-validation training step. Each
 377 curve represents a given number of trees. The highest accuracy value was taken with the parameters mtry (number
 378 of features taken randomly) set to 3 and ntree set to 1000.

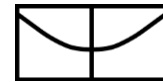
379 One of the possible approaches to defining the optimum values is when an exhaustive
 380 search is performed (e.g., “grid search,” Bérubé et al., 2018; Prado et al., 2020; Rodriguez-
 381 Galiano et al., 2014) and the values which the higher obtained accuracy was selected. For the
 382 model in consideration, we took the number of features randomly taken (mtry) and the number
 383 of predictors (ntree) values as 3 and 1000, respectively.

384 2.3.4.3 Performance evaluation

385 The model’s performance is assessed based on quantitative parameters calculated across
 386 some subsets of data. For this work, we tested the performance for the train data split, test data
 387 split, and for each one of the drill core’s samples separated (Tabela 2-2). We also compared the
 388 accuracy obtained for the unbalanced and SMOTE-balanced models, supporting the idea that a
 389 well-balanced model can make better predictions (Figura 2-5).

390 Tabela 2-2 – Evaluation parameters for the random forests model implemented in this work for each data subset.

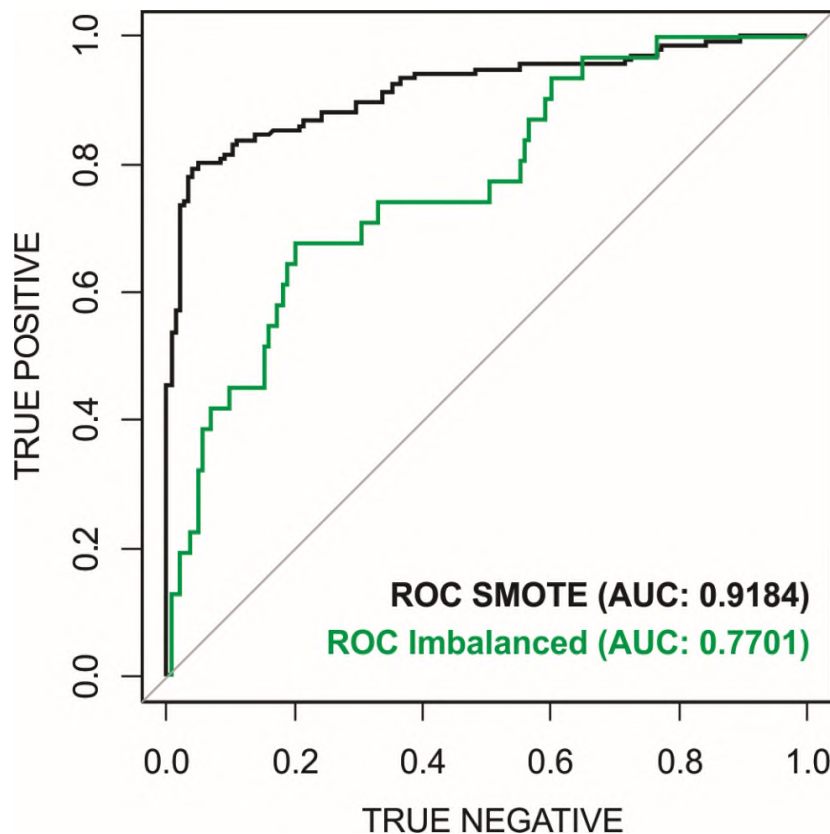
Parameters	Train Split	Test Split	CAN144 Drill Core	CAN120 Drill Core	CANIF27 Drill Core	JBA722 Drill Core	All Samples
Number of Samples	281	120	142	116	129	170	557



NIR	0.527	0.500	0.880	0.836	0.884	0.847	0.860
Raw Accuracy	0.854	0.917	0.958	0.922	0.977	0.871	0.926
Balanced Accuracy	0.854	0.917	0.849	0.805	0.989	0.577	0.769
AUC	0.918	0.917	0.942	0.894	0.917	0.934	0.905
Final Accuracy	0.875	0.917	0.916	0.874	0.960	0.794	0.867

391 Abbreviations; NIR: No Information Rate, AUC: Area Under the Curve.

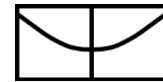
392 Raw accuracy is the proportion of the correct predictions, usually based on the analysis
 393 of previously labeled samples. We applied cross-validation accuracy for the train dataset in this
 394 work, with five repetitions with a 3-times folded analysis. Then, the accuracy obtained in the
 395 training dataset is the average of the accuracy obtained in all five repetitions. For the other data
 396 subsets, the raw accuracy is calculated as the ratio of correct predictions over the number of
 397 samples.



398

399 Figura 2-5 – ROC and AUC diagrams for SMOTE-balanced (AUC: 0.91) and imbalanced (AUC: 0.77) training
 400 datasets.

401 The “No Information Rate” (NIR) parameter is calculated based on the proportion of
 402 the dominant class in the dataset. The model is considered helpful for each data split if the



403 calculated raw accuracy is higher than the NIR value, and this value is closer to 0.5 if the dataset
404 is balanced. The NIR value was calculated based on the proportion of the barren samples in
405 each subset of data. In all considered subsets, the raw accuracy values were higher than the
406 NIR, attesting by this criterium that the models are valid.

407 The balanced accuracy is calculated as the average of the specificity and sensitivity
408 parameters (Zhu et al., 2010). In this paper, specificity is considered the proportion of correctly
409 identified mineralized samples, and sensitivity is the proportion of correctly identified barren
410 samples. The balanced accuracy approach is helpful to evaluate unbalanced datasets (i.e., with
411 a different number of instances for the considered classes).

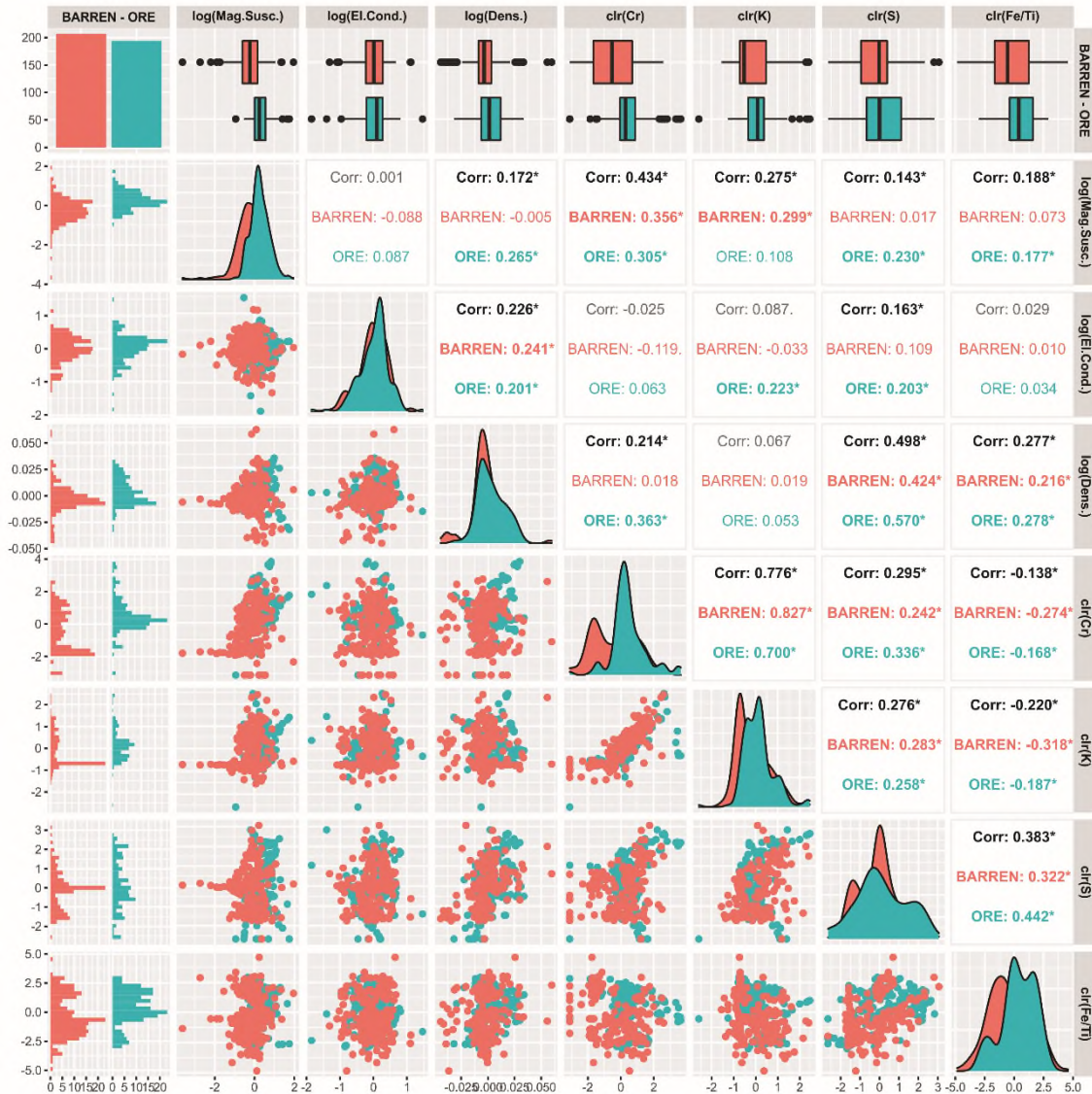
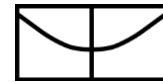
412 The Area Under the Curve (AUC) parameter is calculated based on the Receiver
413 Operating Characteristic (ROC), which estimates the trade-off of True positives and False
414 Positive rates (Torppa et al., 2019). The AUC parameter is calculated based on the ROC curve
415 and ranges from 0.5 (i.e., a completely random model) and 1.0 (perfectly accurate model).

416 To level and consider all described parameters, we calculated the final accuracy as the
417 average of the raw accuracy, balanced accuracy, and AUC values. The final accuracy values
418 range from 0.7937 (for the JBA 722 drill core subset) to 0.9600 (for the CANIF27 drill core
419 subset). The train and test final accuracy were 0.8755 and 0.9168. Thus, the model does not
420 overfit, as the test and train accuracies are close to each other.

421 **2.4 Results and data analysis**

422 *2.4.1 Petrophysics and lithochemistry*

423 Figura 2-6 presents the variation among the measured physical properties and some
424 selected elements using a combined graphs strategy (histograms, density plots, and boxplots)
425 and the bivariate analysis. The graphs below are color-coded according to the mineralization
426 status (i.e., Ore or Barren).

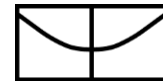


427

428 Figura 2-6 – Selected variables transformed to natural log and centered-log ratio distributions and schematic
 429 graphics at the lower portion. Histograms, scatterplots, and density plots for each variable and combinations of
 430 two variables coded according to mineralization status. At the upper portion, boxplot diagram, colored according
 431 to mineralization label and Spearman ranked correlations calculated based on barren samples (red), ore samples
 432 (blue), and both (black). The correlations values marked with asterisks were validated for a significance test at the
 433 level of 5%. Abbreviations: Mag.Suscep. – Magnetic susceptibility; EI.Cond. – Electric conductivity; Dens. –
 434 Density; clr – Centered log-ratio.

435 We calculated Spearman’s correlation for each selected variable pair to numerically assess their
 436 relationship. Additionally, we validated the correlation values by a statistical significance
 437 symmetrical t-test, with a tolerance level of 5%.

438 The Spearman’s correlation is an index calculated according to the rank of samples and
 439 used when data follows a multimodal or non-parametrical distribution. For the dataset, the
 440 absolute value of significant correlations ranges from 0.143 (among the Magnetic Susceptibility



441 to with the clr(S) content, considering all samples) to 0.827 (among the values of clr(K) and
442 clr(Cr), considering barren samples).

443 More important than the correlation strength between the variables is the correlation
444 between ore and barren populations to predict mineralization status correctly. In that
445 perspective, we point out the correlations between the following pairs of variables: log(Density)
446 and clr(Cr), log(Magnetic Susceptibility) and clr(S), log(Density) and log(Magnetic
447 Susceptibility), clr(Fe/Ti) and log(Magnetic Susceptibility), and clr(K) and log(Magnetic
448 Susceptibility).

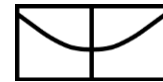
449 Also, by the analysis of the boxplot diagram, it is possible to notice that the position of
450 quantiles and median values according to the mineralization status has significant contrast for
451 some of the selected variables, mentioning the log(Magnetic Susceptibility), log(Density),
452 clr(Cr), clr(K), and clr(Fe/Ti) variables.

453 2.4.2 Mineralization prediction

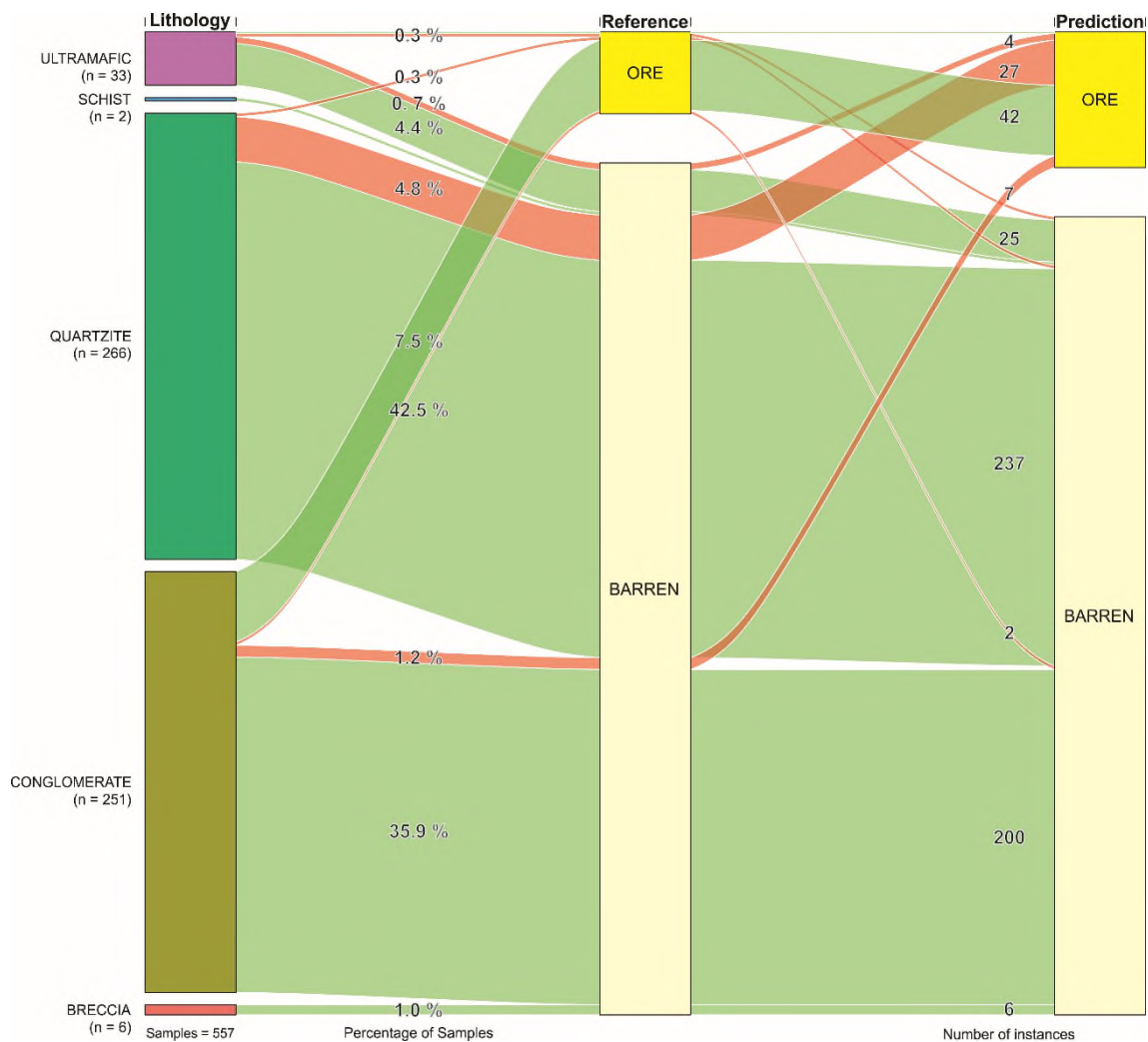
454 We predicted the mineralization status of the samples based on the Random Forests
455 classification model described in section 2.3.4.

456 We prepared an alluvial plot to visualize better the relations between predictions and
457 references across the lithology variations (Figura 2-7). In this type of plot, the reference labels
458 are disposed of in the columns, with the inner subdivisions expressed by their proportion inside
459 the columns. Then, a link (or alluvial) is related to each bar's portions and indicates if a set of
460 samples were labeled as ore or barren by the reference and prediction columns. Additionally,
461 each link is color-coded to aid the visualization of a good prediction or a disagreement with the
462 reference.

463 It is possible to observe that most misclassifications (i.e., poor mineralization
464 predictions) are associated with quartzite samples predicted as Ore but labeled as Barren
465 samples (False Positive Type, or FP, Figura 2-7). Twenty-seven misclassified samples were

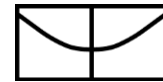


466 recognized, corresponding to a 5% error rate in the raw accuracy parameter considering all
 467 samples. Some other samples were wrongly classified as Barren, while the reference indicated
 468 mineralized (False Negative Type, or FN). Furthermore, the FN proportion is almost negligible
 469 and is restricted to five conglomerates, quartzite, or ultramafic mineralized samples.
 470 Considering the FP, the model predicted 50% mineralized samples than the reference, but most
 471 of these samples (n: 27) are from quartzites from the JBA-722 drill core.



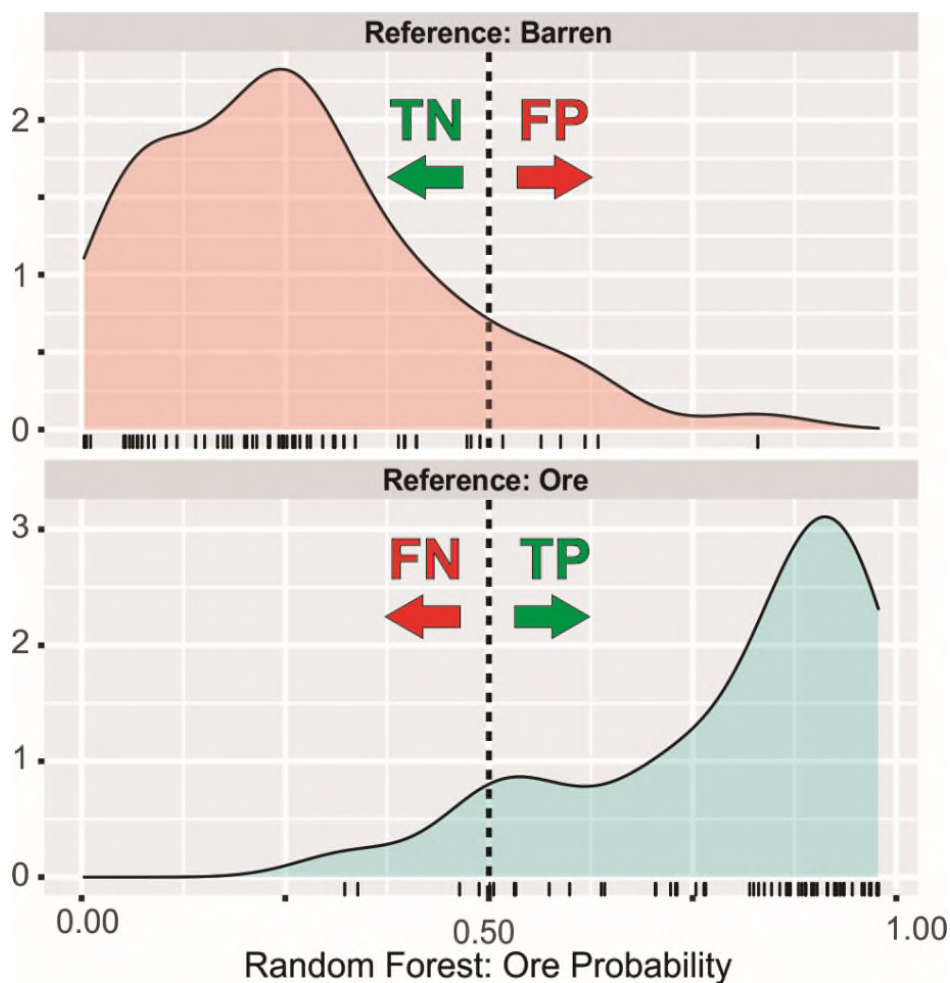
472

473 Figura 2-7 – Alluvial validation diagram for model prediction. This diagram uses column bars to identify the
 474 proportion of lithology, mineralization status for reference, and model prediction. Each link represents a relation
 475 between the three columns, and the links colored green and red represent correct and incorrect predictions,
 476 respectively. The proportion of each link through the sample space and the respective number of samples are
 477 indicated inside the chart.



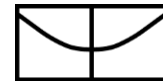
478 2.4.3 Probabilistic prediction approach

479 The mineralization status was predicted considering the average of the votes of all
480 estimators (decision trees in the forest). If the proportion of trees that classified the sample as
481 mineralized (i.e., Ore) is higher than 0.5, the sample is labeled in this way. Then, by analyzing
482 both the reference classification values and the proportion of samples that voted to the
483 considered class, we can build insights about the probability of mineralization of a given
484 sample, even then it is classified as non-mineralized (i.e., Barren, Figura 2-8).



485

486 Figura 2-8 – Ore probability analysis for all samples based on the mineralization status of the test dataset (barren
487 samples are represented in the red curve and ore samples in the blue curve). Most barren samples took a low Ore
488 probability, and the mineralized samples got the highest probabilities. The fields of True Negative (TN, i.e., barren
489 samples predicted as non-mineralized), False Positive (FP, barren samples predicted as mineralized), False
490 Negative (FN, ore samples predicted as non-mineralized), and True Positive (TP, ore samples predicted as
491 mineralized) are indicated in the plot. The ticks at the bottom of each plot indicate the calculated probability for
492 each test dataset sample.



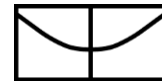
493 Most samples are grouped in the TP and TN fields for all the drill core data. Also, the
494 field with the highest density probability of samples referenced as ore is placed around the
495 probability of 0.9, which suggests that these samples were estimated adequately with a high
496 probability by the model. On the other hand, for the samples classified as barren in the
497 mineralization status, the highest density probability is concentrated around 0.1 to 0.4,
498 indicating that some samples may have any of the considered features indicating a chance of
499 being mineralized.

500 A schematic strip log color-coded mapping the lithology, mineral status, validation, and
501 ore probabilities parameters provide a better understanding of the mineralization (Figura 2-9).

502 **2.5 Discussion**

503 **2.5.1 Mineral Targeting**

504 Pearce et al. (2005) stated that metaconglomerates with blue-gray quartz pebbles and
505 fine-grained disseminated pyrite or hematite usually host gold at the Serra de Jacobina. They
506 noted that the same type of rock with the same pebble size, packing styles, percentage of the
507 matrix, white quartz pebbles, and fuchsite in the matrix tends to have less gold. So, the
508 description mentioned above can be used as a mineral footprint assemblage. A mineral
509 paragenesis schema representing the assemblage, mineral abundance, and the stage of the
510 hydrothermal alteration events is built up using the petrographic and microtexture analysis
511 (Figura 2-10).



Key Legend:

- BRECCIA
- CONGLOMERATE
- QUARTZITE
- SCHIST
- ULTRAMAFIC



512

513 Figura 2-9 – Color-coded strip log according to the lithologies for drill cores studied in this work and respective
 514 validation column. The calculated Ore probability is indicated as a bar beside the sample position for each sample
 515 and drill core, and the threshold of probability is indicated as the red dashed line. The circle color shows the
 516 reference values of the mineralization status at the end of the probability bar. The validation column is color-coded
 517 according to the verification of the predicted and reference mineralization status.



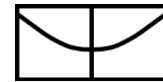
Mineral Phase	Sedimentation	Hydrothermal Alteration (E1)	Hydrothermal Alteration (E2)
<i>Pyrite</i>	████████████████████	████████████████████	
<i>Gold</i>	████████████████████	-----?--?--?	
<i>Uraninite</i>	-----	████████████████████	
<i>Fuchsite</i>		-----	
<i>Ilmenite</i>		-----	
<i>Chlorite</i>			-----
<i>Hematite</i>			-----
<i>Quartz</i>	████████████████████	████████████████████	████████████████████

518

519 Figura 2-10 – Mineral paragenesis flowchart. Pyrite, gold, and uraninite are present at the sedimentation and first
 520 hydrothermal alteration (E1). Fuchsite, ilmenite, chlorite, and hematite were only encountered at the first
 521 hydrothermal alteration event (E1). Hematite, as euhedral to subhedral crystals, is only encountered on the second
 522 hydrothermal alteration event (E2), in quartz veins, in the middle of the recrystallized matrix, or substituting
 523 previous crystals of pyrite.

524 The first hydrothermal alteration event (E₁) relates to the Paleoproterozoic deformation,
 525 either D1 or D2, due to the association with pressure-shadow texture, deformation pattern, and
 526 associated metamorphic minerals (e.g., fuchsite, chlorite, remobilized uraninite, or epigenetic
 527 pyrite). In addition, some authors described remobilized gold associated with this assemblage,
 528 occurring both as free gold and spatially associated with sulfides (Pearson et al., 2005; Teles et
 529 al., 2020). This alteration is interpreted here as a product of the first modification of the
 530 paleoplacer deposit. It may be due to the syntectonic metamorphism event during the
 531 Paleoproterozoic, with or without the participation of fluids derived from the Paleoproterozoic
 532 granitic intrusions (Teles et al., 2020). Despite that, the $\Delta^{33}\text{S}$ values do not show a disturbance
 533 in the isotopic system, which does not favor the entry of external fluids on the basin during
 534 metamorphism (Teles et al., 2020).

535 The second hydrothermal alteration event, also described as hematitization (E₂),
 536 affected the Serra do Córrego Formation rocks and is more prominent close to large brittle



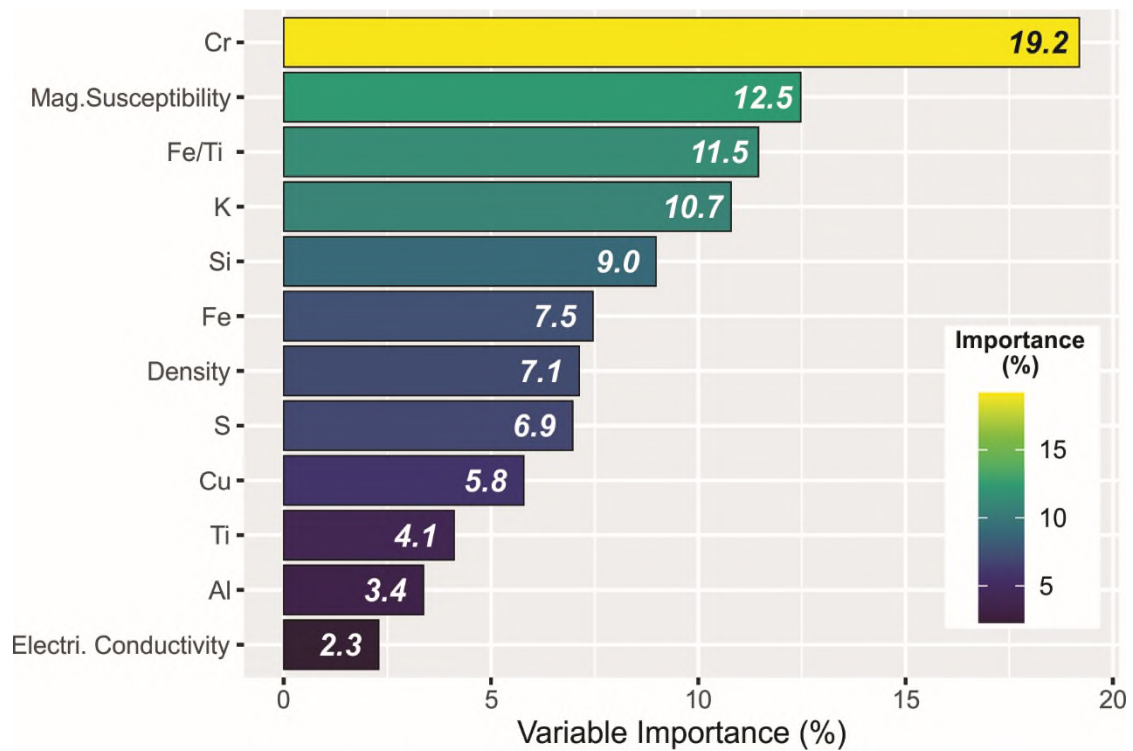
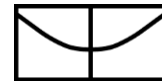
537 structures. Despite this fact, evidence on the relation between hematitization and gold
538 mineralization is still lacking.

539 RF models estimate the importance of the variables to predict a target by analyzing the
540 mean accuracy decrease for each variable. The raw accuracy of the prediction is compared with
541 the accuracy obtained for the estimators that did not evaluate the variable in question to
542 calculate the Importance parameter (Breiman, 2001). Thus, the average accuracy decrease is
543 calculated, and the variables are ranked according to the degree of importance. The variables
544 were normalized for better comparison. The Variable Importance rank is based on the mean
545 accuracy decrease parameter, and the ranking is driven from the data signature (Breiman, 2001).
546 In addition, as the variable in an estimator is taken by chance, the model bias is not significant.

547 For the model in question, the rank of the variables shows that properties such as Cr,
548 Magnetic Susceptibility, Fe/Ti ratio, K, Si, Fe, and S, density are the variables of the most
549 significant importance for the assertiveness of the predictions (Figura 2-11).

550 According to the minerals presented in Figure 10, the variables discussed here can be
551 traced back to the signature of minerals related to the mineralization event (i.e., mineralization
552 footprint). Cr and K are elements present in the mineral structure of the fuchsite, and Fe and S
553 are present in the mineral structure of pyrite and pyrrhotite.

554 Thus, based on Figures 2-10 and 2-11, it is possible to infer that the mineral contents
555 such as fuchsite $K(Al,Cr)_2(AlSi_3O_{10})(OH)_2$, pyrite FeS_2 , intermediate magnetic susceptibility
556 values (possibly associated with the presence of pyrrhotite, Fe_7S_8 , or ilmenite, $FeTiO_3$, observed
557 here in thin sections and the literature) and density (associated with the presence of the sulfides
558 and hematite) have an essential role in the quantitative prediction of mineralization through the
559 Random Forests model. All minerals listed above are associated with the hydrothermal
560 alteration event E_1 .



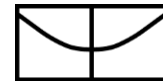
561

562 Figura 2-11– Variable importance rank based on the Mean Decrease Accuracy parameter normalized to a
563 percentage distribution. The bars are filled with the coded colors according to their respective importance.
564 Abbreviations: Mag.Susceptibility – Magnetic Susceptibility; Electri. Conductivity – Electric Conductivity.

565 Variables such as electrical conductivity, and Cu, Ti, and Al content do not significantly
566 influence the model’s accuracy and are not determinant to the predictions. This behavior may
567 occur either because these variables do not differentiate barren and ore samples or because they
568 were not involved in the mineralization event. For example, even though the sample has a
569 relevant sulfide content, included in the rock matrix on the metaconglomerates, the electrical
570 conductivity is not favored if the conductive minerals are dispersed and do not show the
571 continuous distribution in the rocks. This interpretation could explain why electrical
572 conductivity values did not play an essential role in the predictions.

573 2.5.2 Drill core prediction

574 We also evaluated the model’s performance across the samples grouped by the drill core
575 (see Figure 9). The final accuracy of the models significantly changes from the cores CAN120,
576 CAN144, and CANIF27 (ranging from 0.87 to 0.96) to the core JBA722 (0.79).



577 The drill cores CAN120, CAN144, and CANIF27, intersects the Upper Conglomerate
578 Unit and the Intermediate Quartzite, and the core JBA722 intersects from the Intermediate
579 Quartzite to the Lower Conglomerate Unit (see Figura 2-1c).

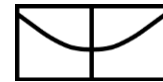
580 This interpretation may imply that Upper Conglomerate and Lower Conglomerate units
581 may differ in the mineralization styles. Nevertheless, this statement should be more carefully
582 investigated in the future, and previous assumptions based on the investigation of drill core
583 samples do not show a relevant distinction between them. In addition, even though the final
584 accuracy is lower in the core samples from JBA722 than in the others, it stills performs
585 satisfactorily.

586 Upon analyzing the probabilistic approach on all drill core samples (see Figura 2-9), we
587 observe that the barren samples surrounding the ore samples continuously increase their Ore
588 Probability in some cases. For example, this can be observed in core CANIF27 samples, close
589 to positions 75 and 120, and core JBA722, close to positions 15 and 140. So, we conclude that
590 the mineralization affects the surrounding samples, and it suggests that the observed variables
591 are mapping the footprint signature of the mineralization.

592 **2.6 Conclusions**

593 We presented in this work a supervised machine learning approach used to predict the
594 gold mineralization in the quartz-pebble metaconglomerate samples from the Serra do Córrego
595 Formation of the Jacobina Group. The implemented predictive model was based on the Random
596 Forests Algorithm.

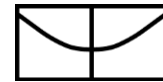
597 The model predicted the mineralization satisfactorily, showing an accuracy of 0.87 and
598 0.91, respectively, for the train and test datasets. We also ran the model to predict samples
599 grouped by the drill core label, and it resulted in minor but significant accuracy differences
600 between samples in cores from different positions in the Serra do Córrego Formation
601 stratigraphy.



602 Despite the satisfactory achieved accuracy, there are some ways to further improve the
603 machine learning model performance results. A possible follow-up suggestion for this work,
604 regarding the machine learning implementation, might be to separate samples into three groups
605 (ore, barren and proximal/altered). It is expected that this might solve the false positive
606 classification issue, as those samples might be barely altered. Also, a new modeling attempt
607 would benefit from a larger (and originally balanced) dataset, eliminating the need to generate
608 synthetic samples.

609 Regarding the interpretation of the model, the Variables Importance analysis showed
610 that properties such as Cr, Magnetic Susceptibility, Fe/Ti ratio, K, Fe, S content, and density
611 are the most significant. Petrographic evidence combined with probabilistic analysis (supported
612 by the Random Forests algorithm) made it possible to explain the relevance of the variables for
613 predicting mineralization status. Therefore, we can infer some mineral targeting criteria to
614 understand the ore formation phenomenon, as the role of the mineral assemblage on the
615 hydrothermal phases described in this work.

616 Additionally, petrographic information supported by back-scattered electron images and
617 energy dispersion spectroscopy semi-quantitative analysis allowed the interpretation that the
618 rocks from the Serra do Córrego Formation present at least two assemblages of hydrothermal
619 alteration, i.e., a first alteration stage composed by epigenetic pyrite, fuchsite, and uraninite,
620 with minor presence of other sulfides, and a second alteration stage with hematite and quartz,
621 associated with ductile and brittle deformation. There is enough evidence of the gold and
622 uranium remobilization during the second hydrothermal alteration stage. However, the relation
623 between the mineralization and the second described alteration stage still needs more
624 investigation. Therefore, the role of secondary hematite in the gold mineralization must be
625 further investigated, as this behavior may predominate at certain mineralized levels and not be
626 the rule.



627 This data-driven method is an alternative way to approach mineralization targeting and
628 provides valuable insights in different mineral exploration stages. In addition, petrophysical
629 measurements and geochemistry data can be obtained in the mineral exploration industry with
630 relatively low costs, and their evaluation may be beneficial, whenever the quality control
631 procedures are followed (avoiding undesirable bias). Thus, the use of machine learning
632 algorithms for aiding in the understanding of a complex mineralization is feasible and can bring
633 important practical insights for mineral exploration in many scenarios.

634 Further, to increase statistical representativeness and put the model in production under
635 a more diverse scenario, the machine learning approach presents the advantage that more
636 samples could be added to the training dataset. Thus, the updated machine learning model could
637 “learn” some new information and theoretically evaluate and understand even minor
638 particularities of a mineral deposit if enough data and contrasting variables are provided.

639 **Author contribution**

640 **GFS** Conceptualization, Sampling, Investigation, Data Curation, Methodology, Visualization,
641 Writing - Original Draft. **AMS** Supervision, Investigation, Validation, Writing - Review &
642 Editing. **CBT** Supervision, Investigation, Validation, Writing - Review & Editing. **FCJ**
643 Resources, Validation, Writing - Review & Editing. **ELK**: Sampling, Investigation, Validation,
644 Writing - Review & Editing

645 **Acknowledgments**

646 We want to thank Yamana Gold Inc. for its assistance during the collection of drill-core
647 samples. Also, we acknowledge the Geological Survey of Brazil (SGB/CPRM) for financially
648 supporting the fieldwork campaigns. Thanks are extended to the geologists Anderson Dourado,
649 Carina Lopes, Kotaro Ushigasaki, João Larizzatti, Joseneusa Rodrigues, Daniel Miranda, and
650 Valter Sobrinho for aid with field logistics, sampling, or during laboratory analysis. Lastly, the
651 authors thank the improvements made on the original manuscript by the two reviewers, R.S.D.
652 and P.M.P.G. The scripts used in this work may be provided by contacting the corresponding



653 author. This work was financed by the Coordenação de Aperfeiçoamento de Pessoal de Nível
654 Superior - Brasil (CAPES) - Finance Code 001. A.M. Silva, C.L.B. Toledo, F.C.J., and ELK
655 acknowledge the Brazilian National Council for Scientific and Technological Development
656 (CNPq) for their respective research grants.

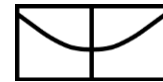
657 **Supplementary data**

658 Please refer to the online version to access the supplementary files.

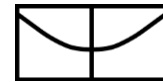
- 659 • Table A1 – Gold content (Fire assay) by drill-core
- 660 • Table A2 – Petrophysics and XRF data

661 **References**

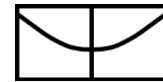
- 662 Aitchison, J., 1986. The Statistical Analysis of Compositional Data. *Stat. Anal. Compos. Data*
663 44, 139–177. <https://doi.org/10.1007/978-94-009-4109-0>
- 664 Alkmin, F.F., Brito Neves, B.B., Alves, J.A.C., 1993. Arcabouço tectônico do Cráton do São
665 Francisco: Uma revisão, in: Dominguez, J.M.L., Misi, A. (Eds.), *O Cráton Do São*
666 *Francisco*. SBG - Sociedade Brasileira de Geociências, Salvador, BA, pp. 45–62.
- 667 Barbosa, J.S.F., Sabaté, P., 2004. Archean and Paleoproterozoic crust of the São Francisco
668 Craton, Bahia, Brazil: geodynamic features. *Precambrian Res.* 133, 1–27.
669 <https://doi.org/10.1016/j.precamres.2004.03.001>
- 670 Bérubé, C.L., Olivo, G.R., Chouteau, M., Perrouy, S., Shamsipour, P., Enkin, R.J., Morris,
671 W.A., Feltrin, L., Thiémonge, R., 2018. Predicting rock type and detecting hydrothermal
672 alteration using machine learning and petrophysical properties of the Canadian Malartic
673 ore and host rocks, Pontiac Subprovince, Québec, Canada. *Ore Geol. Rev.* 96, 130–145.
674 <https://doi.org/10.1016/j.oregeorev.2018.04.011>
- 675 Breiman, L., 2001. Random forests. *Mach. Learn.* 56, 5–32.



- 676 Carneiro, C.D.C., Fraser, S.J., Crósta, A.P., Silva, A.M., Barros, C.E. de M., 2012. Semi-
677 automated geologic mapping using self-organizing maps and airborne geophysics in the
678 Brazilian Amazon. *GEOPHYSICS* 77, K17–K24. <https://doi.org/10.1190/geo2011-0302.1>
- 679 Carranza, E.J.M., Laborte, A.G., 2016. Data-Driven Predictive Modeling of Mineral
680 Prospectivity Using Random Forests: A Case Study in Catanduanes Island (Philippines).
681 *Nat. Resour. Res.* 25, 35–50. <https://doi.org/10.1007/s11053-015-9268-x>
- 682 Carranza, E.J.M., Laborte, A.G., 2015. Data-driven predictive mapping of gold prospectivity,
683 Baguio district, Philippines: Application of Random Forests algorithm. *Ore Geol. Rev.* 71,
684 777–787. <https://doi.org/10.1016/j.oregeorev.2014.08.010>
- 685 Chawla, N. V, Bowyer, K.W., Hall, L.O., Kegelmeyer, W.P., 2002. SMOTE: Synthetic
686 Minority Over-sampling Technique. *J. Artif. Intell. Res.* 16, 321–357.
687 <https://doi.org/10.1613/jair.953>
- 688 Chen, Y., 2015. Mineral potential mapping with a restricted Boltzmann machine. *Ore Geol.*
689 *Rev.* 71, 749–760. <https://doi.org/10.1016/j.oregeorev.2014.08.012>
- 690 Costa, I., Tavares, F., Oliveira, J., 2019. Predictive lithological mapping through machine
691 learning methods: a case study in the Cinzento Lineament, Carajás Province, Brazil. *J.*
692 *Geol. Surv. Brazil* 2, 26–36. <https://doi.org/10.29396/jgsb.2019.v2.n1.3>
- 693 da Costa, G., Hofmann, A., Agangi, A., 2020. A revised classification scheme of pyrite in the
694 Witwatersrand Basin and application to placer gold deposits. *Earth-Science Rev.* 201,
695 103064. <https://doi.org/10.1016/j.earscirev.2019.103064>
- 696 da Silva, G.F., Ferreira, M.V., Costa, I.S.L., Bernardes, R.B., Mota, C.E.M., Cuadros Jiménez,
697 F.A., 2021. Qmin – A machine learning-based application for processing and analysis of
698 mineral chemistry data. *Comput. Geosci.* 157, 104949.
699 <https://doi.org/10.1016/j.cageo.2021.104949>



- 700 da Silva, G.F., Larizzatti, J.H., da Silva, A.D.R., Lopes, C.G., Klein, E.L., Uchigasaki, K., 2022.
701 Unsupervised drill core pseudo-log generation in raw and filtered data, a case study in the
702 Rio Salitre greenstone belt, São Francisco Craton, Brazil. *J. Geochemical Explor.* 232,
703 106885. <https://doi.org/10.1016/j.gexplo.2021.106885>
- 704 Dentith, M., Enkin, R.J., Morris, W., Adams, C., Bourne, B., 2020. Petrophysics and mineral
705 exploration: a workflow for data analysis and a new interpretation framework. *Geophys.*
706 *Prospect.* 68, 178–199. <https://doi.org/10.1111/1365-2478.12882>
- 707 Fisher, L., Gazley, M.F., Baensch, A., Barnes, S.J., Cleverley, J., Duclaux, G., 2014. Resolution
708 of geochemical and lithostratigraphic complexity: A workflow for application of portable
709 X-ray fluorescence to mineral exploration. *Geochemistry Explor. Environ. Anal.* 14, 149–
710 159. <https://doi.org/10.1144/geochem2012-158>
- 711 Ford, A., 2019. Practical Implementation of Random Forest-Based Mineral Potential Mapping
712 for Porphyry Cu – Au Mineralization in the Eastern Lachlan Orogen , NSW , Australia.
713 *Nat. Resour. Res.* <https://doi.org/10.1007/s11053-019-09598-y>
- 714 Frimmel, H.E., 2019. The Witwatersrand Basin and Its Gold Deposits, in: Kröner, A., Hofmann,
715 A. (Eds.), *The Archaean Geology of the Kaapvaal Craton, Southern Africa*. Springer-
716 Verlag, pp. 255–275. https://doi.org/10.1007/978-3-319-78652-0_10
- 717 Frimmel, H.E., Le Roex, a P., Knight, J., Minter, W.E.L., 1993. A Case Study of the Post-
718 depositional Alteration. *Econ. Geol.* 88, 249–265.
- 719 Grunsky, E., 2001. Aspects of multivariate statistical analysis in geology, *Computers &*
720 *Geosciences.* [https://doi.org/10.1016/s0098-3004\(00\)00094-7](https://doi.org/10.1016/s0098-3004(00)00094-7)
- 721 Hall, B., 2016. Facies classification using machine learning. *Lead. Edge* 35, 906–909.
722 <https://doi.org/10.1190/tle35100906.1>



- 723 Han, H., Wang, W.Y., Mao, B.H., 2005. Borderline-SMOTE: A new over-sampling method in
724 imbalanced data sets learning. *Lect. Notes Comput. Sci.* 3644, 878–887.
725 https://doi.org/10.1007/11538059_91
- 726 Hariharan, S., Tirodkar, S., Porwal, A., Bhattacharya, A., Joly, A., 2017. Random Forest-Based
727 Prospectivity Modelling of Greenfield Terrains Using Sparse Deposit Data: An Example
728 from the Tanami Region, Western Australia. *Nat. Resour. Res.* 26, 489–507.
729 <https://doi.org/10.1007/s11053-017-9335-6>
- 730 Harris, J.R., Grunsky, E., Behnia, P., Corrigan, D., 2015. Data- and knowledge-driven mineral
731 prospectivity maps for Canada's North. *Ore Geol. Rev.* 71, 788–803.
732 <https://doi.org/10.1016/j.oregeorev.2015.01.004>
- 733 Heilbron, M., Cordani, U.G., Alkmim, F.F., 2017. São Francisco Craton, Eastern Brazil:
734 Tectonic Genealogy of a Miniature Continent. *Springer* 326. [https://doi.org/10.1007/978-](https://doi.org/10.1007/978-3-319-01715-0)
735 [3-319-01715-0](https://doi.org/10.1007/978-3-319-01715-0)
- 736 Ida, H., 2004. X-ray fluorescence analysis with portable instruments. Kyoto University.
- 737 Japkowicz, N., Stephen, S., 2002. The class imbalance problem: A systematic study1. *Intell.*
738 *Data Anal.* 6, 429–449. <https://doi.org/10.3233/IDA-2002-6504>
- 739 Koziarski, M., Woźniak, M., Krawczyk, B., 2020. Combined Cleaning and Resampling
740 algorithm for multi-class imbalanced data with label noise. *Knowledge-Based Syst.* 204,
741 106223. <https://doi.org/10.1016/j.knosys.2020.106223>
- 742 Kuhn, S., Cracknell, M.J., Reading, A.M., 2018. Lithologic mapping using Random Forests
743 applied to geophysical and remote-sensing data: A demonstration study from the Eastern
744 Goldfields of Australia. *GEOPHYSICS* 83, B183–B193. [https://doi.org/10.1190/geo2017-](https://doi.org/10.1190/geo2017-0590.1)
745 0590.1



- 746 Ledru, P., Milési, J.P., Johan, V., Sabaté, P., Maluski, H., 1997. Foreland basins and gold-
747 bearing conglomerates: a new model for the Jacobina Basin (São Francisco province,
748 Brazil). *Precambrian Res.* 86, 155–176.
- 749 Leite, C. de M.M., Barbosa, J.S.F., Nicollet, C., Sabaté, P., 2007. Evolução
750 metamórfica/metassomática paleoproterozóica do Complexo Saúde, da Bacia Jacobina e
751 de leucogranitos peraluminosos na parte norte do Cráton do São Francisco. *Rev. Bras.*
752 *Geociências* 37, 777–797. <https://doi.org/10.25249/0375-7536.2007374777797>
- 753 Leite, C.M.M., 2002. A evolução geodinâmica da orogênese paleoproterozóica nas regiões de
754 Capim Grosso - Jacobina e Pintadas - Mundo Novo (Bahia, Brasil): metamorfismo,
755 anatexia crustal e tectônica. Universidade Federal da Bahia.
- 756 Leite, C.M.M., Marinho, M.M., 2012. Serra de Jacobina e Contendas-Mirante, in: Barbosa,
757 J.S.F. (Ed.), *Geologia Da Bahia: Pesquisa e Atualização*. CBPM - Companhia Baiana de
758 Pesquisa Mineral, pp. 397–441.
- 759 Li, X., Zhang, C., Behrens, H., Holtz, F., 2020. Lithos Calculating amphibole formula from
760 electron microprobe analysis data using a machine learning method based on principal
761 components regression. *LITHOS* 362–363, 105469.
762 <https://doi.org/10.1016/j.lithos.2020.105469>
- 763 Mascarenhas, J.F., Ledru, P., Souza, S.L., Conceição-Filho, V.M., Melo, L.F.A., Lorenzo, C.L.,
764 Milési, J.P., 1998. *Geologia e recursos minerais do Grupo Jacobina e da parte sul do*
765 *Greenstone Belt de Mundo Novo. Série Arquivos Abertos*, vol. 13. CBPM - Companhia
766 Baiana de Pesquisa Mineral, Salvador, Brazil.
- 767 McKay, G., Harris, J.R., 2016. Comparison of the Data-Driven Random Forests Model and a
768 Knowledge-Driven Method for Mineral Prospectivity Mapping: A Case Study for Gold
769 Deposits Around the Huritz Group and Nueltin Suite, Nunavut, Canada. *Nat. Resour. Res.*
770 25, 125–143. <https://doi.org/10.1007/s11053-015-9274-z>



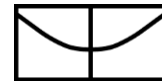
- 771 Milési, J., Ledru, P., Marcoux, E., Mougeot, R., Johan, V., Lerouge, C., Sabaté, P., Bailly, L.,
772 Respaut, J., Skipwith, P., 2002. The Jacobina Paleoproterozoic gold-bearing
773 conglomerates, Bahia, Brazil: a “hydrothermal shear-reservoir” model. *Ore Geol. Rev.* 19,
774 95–136. [https://doi.org/10.1016/S0169-1368\(01\)00038-5](https://doi.org/10.1016/S0169-1368(01)00038-5)
- 775 Miranda, D.A., Misi, A., Klein, E.L., Castro, M.P., Queiroga, G., 2021. A mineral system
776 approach on the Paleoproterozoic Au-bearing quartz veins of the Jacobina Range,
777 northeastern of the São Francisco Craton, Brazil. *J. South Am. Earth Sci.* 106.
778 <https://doi.org/10.1016/j.jsames.2020.103080>
- 779 Niiranen, T., Nykänen, V., Lahti, I., 2019. Scalability of the mineral prospectivity modelling –
780 An orogenic gold case study from northern Finland. *Ore Geol. Rev.* 109, 11–25.
781 <https://doi.org/10.1016/j.oregeorev.2019.04.002>
- 782 Pearson, W., Macêdo, P.M.M., Rúbio, A., Lorenzo, C.L., Karpeta, P., 2005. Geology and gold
783 mineralization of the Jacobina Mine and Bahia Gold Belt, Bahia, Brazil and comparison to
784 Tarkwa and Witwatersrand., in: *Proceedings, Geological Society of Nevada Symposium*,
785 Vol. 1. Geological Society of Nevada Symposium, Reno, Nevada, USA., pp. 757–786.
- 786 Piercey, S.J., 2014. Modern analytical facilities 2. A review of quality assurance and quality
787 control (qa/qc) procedures for lithogeochemical data. *Geosci. Canada* 41, 75–88.
788 <https://doi.org/10.12789/geocanj.2014.41.035>
- 789 Prado, E.M.G., de Souza Filho, C.R., Carranza, E.J.M., Motta, J.G., 2020. Modeling of Cu-Au
790 prospectivity in the Carajás mineral province (Brazil) through machine learning: Dealing
791 with imbalanced training data. *Ore Geol. Rev.* 124, 103611.
792 <https://doi.org/10.1016/j.oregeorev.2020.103611>
- 793 Reis, C., Menezes, R.C.L., Miranda, D.A., Santos, F.P. dos, Santos, R.S.V. dos, Menezes, A.R.,
794 2021. Áreas de Relevante Interesse Mineral (Arim) Integração Geológica E Avaliação do



- 795 Potencial Metalogenético da Serra de Jacobina e do Greenstone Belt Mundo Novo. Serviço
796 Geológico do Brasil - CPRM, Salvador, Brazil.
- 797 Robb, L.J., Meyer, F.M., 1991. A contribution to recent debate concerning epigenetic versus
798 syngenetic mineralization processes in the Witwatersrands Basin. *Econ. Geol.* 86, 396–
799 401.
- 800 Rodriguez-Galiano, V., Sanchez-Castillo, M., Chica-Olmo, M., Chica-Rivas, M., 2015.
801 Machine learning predictive models for mineral prospectivity: An evaluation of neural
802 networks, random forest, regression trees and support vector machines. *Ore Geol. Rev.* 71,
803 804–818. <https://doi.org/10.1016/j.oregeorev.2015.01.001>
- 804 Saljoughi, B.S., Hezarkhani, A., 2018. A comparative analysis of artificial neural network
805 (ANN), wavelet neural network (WNN), and support vector machine (SVM) data-driven
806 models to mineral potential mapping for copper mineralizations in the Shahr-e-Babak
807 region, Kerman, Iran. *Appl. Geomatics* 10, 229–256. [https://doi.org/10.1007/s12518-018-](https://doi.org/10.1007/s12518-018-0229-z)
808 0229-z
- 809 Santos, F.P. dos, Chemale Junior, F., Meneses, A.R.A.S., 2019. The nature of the
810 Paleoproterozoic orogen in the Jacobina Range and adjacent areas, northern São Francisco
811 Craton, Brazil, based on structural geology and gravimetric modeling. *Precambrian Res.*
812 332, 105391. <https://doi.org/10.1016/j.precamres.2019.105391>
- 813 Saporetti, C.M., da Fonseca, L.G., Pereira, E., de Oliveira, L.C., 2018. Machine learning
814 approaches for petrographic classification of carbonate-siliciclastic rocks using well logs
815 and textural information. *J. Appl. Geophys.* 155, 217–225.
816 <https://doi.org/10.1016/j.jappgeo.2018.06.012>
- 817 Teixeira, J.B.G., De Souza, J.A.B., Da Silva, M. da G., Leite, C.M.M., Barbosa, J.S.F.S.F.,
818 Coelho, C.E.S., Abram, M.B., Filho, V.M.C., Iyer, S.S.S., 2001. Gold mineralization in the



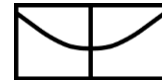
- 819 Serra de Jacobina region, Bahia Brazil: tectonic framework and metallogenesis. *Miner.*
820 *Depos.* 36, 332–344. <https://doi.org/10.1007/s001260100174>
- 821 Teixeira, W., Oliveira, E.P., Marques, L.S., 2017. Nature and evolution of the Archean crust of
822 the São Francisco Craton, in: Heilbron, M., Alkmin, F.F., Cordani, U.G. (Eds.), *The São*
823 *Francisco Craton and Its Margins, Eastern Brazil*. Springer-Verlag, pp. 29–56.
- 824 Teles, G., Jr, F.C., Oliveira, C.G. De, Chemale, F., de Oliveira, C.G., 2015. Paleoafrican record
825 of the detrital pyrite-bearing, Jacobina Au-U deposits, Bahia, Brazil. *Precambrian Res.* 256,
826 289–313. <https://doi.org/10.1016/j.precamres.2014.11.004>
- 827 Teles, G.S., Chemale, F., Ávila, J.N., Ireland, T.R., Dias, A.N.C., Cruz, D.C.F., Constantino,
828 C.J.L., 2020. Textural and geochemical investigation of pyrite in Jacobina Basin, São
829 Francisco Craton, Brazil: Implications for paleoenvironmental conditions and formation of
830 pre-GOE metaconglomerate-hosted Au-(U) deposits. *Geochim. Cosmochim. Acta* 273,
831 331–353. <https://doi.org/10.1016/j.gca.2020.01.035>
- 832 Thermo-Scientific, 2013. *Mining and exploration: Solutions from early-stage discovery*
833 *through mineral processing*. San Jose, California.
- 834 Torppa, J., Nykänen, V., Molnár, F., 2019. Unsupervised clustering and empirical fuzzy
835 memberships for mineral prospectivity modelling. *Ore Geol. Rev.* 107, 58–71.
836 <https://doi.org/10.1016/j.oregeorev.2019.02.007>
- 837 Wickham, H., 2014. Tidy Data. *J. Stat. Softw.* 59. <https://doi.org/10.18637/jss.v059.i10>
- 838 Yamana Gold, 2020. *Annual Report 2020 - NI43-101*, 176p.
- 839 Yousefi, M., Nykänen, V., 2016. Data-driven logistic-based weighting of geochemical and
840 geological evidence layers in mineral prospectivity mapping. *J. Geochemical Explor.* 164,
841 94–106. <https://doi.org/10.1016/j.gexplo.2015.10.008>



842 Zhu, W., Zeng, N., Wang, N., 2010. Sensitivity, specificity, accuracy, associated confidence
843 interval and ROC analysis with practical SAS® implementations. Northeast SAS Users Gr.
844 2010 Heal. Care Life Sci. 1–9.

845 Zuo, R., 2017. Machine Learning of Mineralization-Related Geochemical Anomalies: A
846 Review of Potential Methods. Nat. Resour. Res. 26, 457–464.
847 <https://doi.org/10.1007/s11053-017-9345-4>

848



5 DISCUSSÕES E CONSIDERAÇÕES FINAIS

Os resultados apresentados no capítulo 2 mostram que propriedades como suscetibilidade magnética, densidade e teores de Cr, K, Fe, e S são as variáveis mais significativas para a predição de amostras mineralizadas (não nessa ordem de importância). Evidências petrográficas combinadas com análises probabilísticas (derivadas das inferências estatísticas utilizando o algoritmo *Random Forests*) permitiram explicar a relevância das variáveis para predição da mineralização. Portanto, podemos inferir alguns critérios de direcionamento mineral para entender o fenômeno de formação do minério, como o papel da assembleia mineral nas fases hidrotermais descritas neste trabalho.

Assim, sugerimos essa abordagem na prospecção local, com a ressalva que os modelos foram construídos para serem representativos, porém a amostragem realizada pode não abranger toda a variância da população mineralizada, assim mais amostras podem ser adicionadas ao conjunto de dados de treinamento e o modelo deve ser atualizado. Desse modo, o modelo de aprendizado de máquina poderia “aprender” novas informações e avaliar particularidades de um depósito mineral se dados suficientes e variáveis contrastantes forem fornecidos.

Os resultados apresentados no capítulo 3 mostram que a investigação da assinatura de elementos traço em pirita ao longo das unidades do Grupo Jacobina sugere uma falta de contraste que pode implicar na manutenção da fonte de sedimentos durante a formação da bacia Jacobina, ou a posterior equilíbrio químico durante os estágios de alteração metamórficos/hidrotermais. No entanto, para elucidar essas questões, recomendamos a avaliação de mais amostras e análises em trabalhos futuros, pois isso pode ajudar a reduzir o viés amostral.

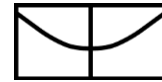
Para fins de exploração mineral, nossos resultados apontam para um importante papel de alguns minerais acessórios anteriormente negligenciados nas partes alteradas pelos eventos epigenéticos dos depósitos da Serra de Jacobina, como esfalerita, calcopirita, pirrotita e outros.



Esses minerais são considerados farejadores da alteração epigenética nos depósitos, e podem estar relacionados ao ouro livre, conforme indicado pela análise de dendrogramas e confirmado em lâminas delgadas. No entanto, concluímos que a alteração epigenética nos depósitos pode ter resultados positivos ou negativos sobre o conteúdo metálico (*endowment*) do depósito, pois uma forte modificação e consequente mobilização de ouro não canalizada em um mecanismo de concentração eficaz poderia espalhar o ouro em várias pequenas frações.

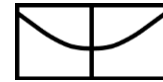
Os resultados parciais de reflectância espectrorradiométrica e geoquímica, apresentados no capítulo 4, indicam que a maioria das amostras mineralizadas possuem mais de 50% de pirita na composição das frações analisadas, sendo a fuchcita e a illita minerais secundários. Goethita, hematita e muscovita ocorrem como fases acessório na assembleia mineral e, exceto nas amostras mineralizadas do testemunho JBA-722, os minerais de óxido de ferro não são dominantes nos metaconglomerados auríferos. Além disso, os dados de litoquímicos validaram os valores inferidos de abundância dos minerais, destacando os teores de Cr, K, Al (no caso de fuchcita) e Fe e S (no caso de minerais de óxido de ferro e pirita).

Assim sendo, utilizando os resultados obtidos pelas ferramentas descritas nessa tese, e observando a dimensão das amostragens e as escalas de trabalho, resumimos os critérios de prospecção encontrados na Figura 5-1. Separamos as informações de acordo com a assinatura das propriedades físicas e químicas nas zonas mineralizada, proximal e estéril. Finalizando, sugerimos para trabalhos futuros que as abordagens aqui apresentadas sejam traduzidas em critérios de exploração para modelamento de potencial mineral e de prospecção local para que a abordagem seja validada (ver Apêndice B).



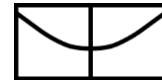
Formação Serra do Córrego		
Paleoplacer modificado		
Metaconglomerados		
Zona Mineralizada	Zona Proximal	Zona Estéril
+ <i>Fuchcita</i> na matriz + <i>Suscep.</i> <i>Magnética</i>	± <i>Fuchcita</i> na matriz <i>Poucas</i> <i>alterações</i>	+ <i>Goethita</i> e <i>Hematita</i>
+ <i>Cr em rocha</i> - <i>Pb, Zn e Cu</i> na <i>pirita</i>	+ <i>Esfalerita,</i> <i>Galena e pirrotita</i> na matriz	- <i>Cr em rocha</i>
+ <i>Densidade</i>	+ <i>Clorita</i>	+ <i>Condutividade</i> <i>Elétrica</i>
+ <i>Pirita</i> <i>epigenética</i>	± <i>Ilita?</i>	<i>Densidade</i> < 2.65 <i>SM</i> < 0.2×10^{-3} <i>SI</i>

Figura 5-1: Quadro resumo dos resultados obtidos na tese a respeito das propriedades físicas e químicas das rochas e minerais analisados, no escopo das amostras de metaconglomerados com mineralização do tipo paleoplacer modificado da Formação Serra do Córrego, Grupo Jacobina. Os controles do componente sedimentar da mineralização não foram abordados neste trabalho.

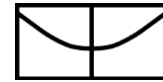


REFERÊNCIAS BIBLIOGRÁFICAS

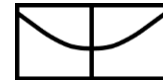
- Abedi, M., Norouzi, G.H., Bahroudi, A., 2012. Support vector machine for multi-classification of mineral prospectivity areas. *Comput. Geosci.* 46, 272–283. <https://doi.org/10.1016/j.cageo.2011.12.014>
- Agterberg, F.P., Bonham-Carter, G.F., 2005. Measuring the Performance of Mineral-Potential Maps. *Nat. Resour. Res.* 14, 1–17. <https://doi.org/10.1007/s11053-005-4674-0>
- Bérubé, C.L., Olivo, G.R., Chouteau, M., Perrouy, S., Shamsipour, P., Enkin, R.J., Morris, W.A., Feltrin, L., Thiémonge, R., 2018. Predicting rock type and detecting hydrothermal alteration using machine learning and petrophysical properties of the Canadian Malartic ore and host rocks, Pontiac Subprovince, Québec, Canada. *Ore Geol. Rev.* 96, 130–145. <https://doi.org/10.1016/j.oregeorev.2018.04.011>
- Bonham-Carter, G.F., 1994. *Geographic Information Systems for Geoscientists*, 1st Ed. ed. Pergamon.
- Breiman, L., 2001. Random forests. *Mach. Learn.* 56, 5–32.
- Breiman, L., Cutler, A., Forests, R., Ho, T.K., Labs, B., Kleinberg, E., Breiman, L., 1995. Random forest 1–5.
- Carranza, E.J.M. (Emmanuel J.M., 2009. *Geochemical anomaly and mineral prospectivity mapping in GIS*. Elsevier.
- Carranza, E.J.M., Laborte, A.G., 2016. Data-Driven Predictive Modeling of Mineral Prospectivity Using Random Forests: A Case Study in Catanduanes Island (Philippines). *Nat. Resour. Res.* 25, 35–50. <https://doi.org/10.1007/s11053-015-9268-x>
- Carranza, E.J.M., Laborte, A.G., 2015. Random forest predictive modeling of mineral prospectivity with small number of prospects and data with missing values in Abra (Philippines). *Comput. Geosci.* 74, 60–70. <https://doi.org/10.1016/j.cageo.2014.10.004>
- Costa, I., Tavares, F., Oliveira, J., 2019. Predictive lithological mapping through machine learning methods: a case study in the Cinzento Lineament, Carajás Province, Brazil. *J. Geol. Surv. Brazil* 2, 26–36. <https://doi.org/10.29396/jgsb.2019.v2.n1.3>
- Cover, T.M., Hart, P.E., 1967. Nearest Neighbor Pattern Classification. *IEEE Trans. Inf. Theory* 13, 21–27. <https://doi.org/10.1109/TIT.1967.1053964>
- Davies, J. C. 2002. *Statistics and Data Analysis in Geology*, 3rd Edition. John Wiley & Sons, New York –USA. 656 pages.



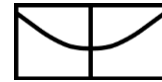
- Dentith, M., Enkin, R.J., Morris, W., Adams, C., Bourne, B., 2020. Petrophysics and mineral exploration: a workflow for data analysis and a new interpretation framework. *Geophys. Prospect.* 68, 178–199. <https://doi.org/10.1111/1365-2478.12882>
- Dentith, M.C., Mudge, S.T., 2014. *Geophysics for the mineral exploration geoscientist*, 1st ed, Cambridge University Press. Cambridge University Press, Cambridge, UK.
- Fleming, S.W., Watson, J.R., Ellenson, A., Canon, A.J., Vesselinov, V.C. 2021. Machine learning in Earth and environmental science requires education and research policy reforms. *Nature Geoscience.* 14, 878-880(2021). <https://doi.org/10.1038/s41561-021-00881-3>
- Frimmel, H.E., 2019. The Witwatersrand Basin and Its Gold Deposits, in: Kröner, A., Hofmann, A. (Eds.), *The Archaean Geology of the Kaapvaal Craton, Southern Africa*. Springer-Verlag, pp. 255–275. https://doi.org/10.1007/978-3-319-78652-0_10
- Frimmel, H.E., 2014. Chapter 10 A Giant Mesoarchean Crustal Gold-Enrichment Episode : Possible Causes and Consequences for Exploration 209–234.
- Frimmel, H.E., Groves, D.I., Kirk, J., Ruiz, J., Chesley, J., Minter, W.E.L., 2019. The Formation and Preservation of the Witwatersrand Goldfields, the World’s Largest Gold Province. *One Hundredth Anniv. Vol.* 769–797. <https://doi.org/10.5382/av100.23>
- Guimarães, F.S., de Freitas, M.E., Rios, F.J., Pedrosa, T.A., 2019. Mineralogical characterization and origin of uranium mineralization in Witwatersrand-like metaconglomerate of the Moeda Formation, Quadrilátero Ferrífero, Brazil. *Ore Geol. Rev.* 106, 423–445. <https://doi.org/10.1016/j.oregeorev.2019.01.016>
- Hagemann, S.G., Lisitsin, V.A., Huston, D.L., 2016. Mineral system analysis: Quo vadis. *Ore Geol. Rev.* 76, 504–522. <https://doi.org/10.1016/j.oregeorev.2015.12.012>
- Hennigh, Q., 2016. Conglomerate-Hosted Gold Mineralization in the Pilbara, Western Australia, in: *Association for Mineral Exploration British Columbia, Roundup 2016, Abstracts: Core Shack*. pp. 49–50.
- Hronsky, J.M.A., Kreuzer, O.P., 2019. Applying spatial prospectivity mapping to exploration targeting: Fundamental practical issues and suggested solutions for the future. *Ore Geol. Rev.* 107, 647–653. <https://doi.org/10.1016/j.oregeorev.2019.03.016>
- Joly, A., Porwal, A., McCuaig, T.C., 2012. Exploration targeting for orogenic gold deposits in the Granites-Tanami Orogen: Mineral system analysis, targeting model and prospectivity analysis. *Ore Geol. Rev.* 48, 349–383. <https://doi.org/10.1016/J.OREGEOREV.2012.05.004>



- Klein, E.L., Rodrigues, J.B., Queiroz, J.D.S., Oliveira, R.G., Guimarães, S.B., Chaves, C.L., 2017. Deposition and tectonic setting of the Palaeoproterozoic Castelo dos Sonhos metasedimentary formation, Tapajós Gold Province, Amazonian Craton, Brazil: age and isotopic constraints. *Int. Geol. Rev.* 59, 864–883. <https://doi.org/10.1080/00206814.2016.1237311>
- Ledru, P., Milési, J.P., Johan, V., Sabaté, P., Maluski, H., 1997. Foreland basins and gold-bearing conglomerates: a new model for the Jacobina Basin (São Francisco province, Brazil). *Precambrian Res.* 86, 155–176.
- Mccuaig, T.C., Beresford, S., Hronsky, J., 2010. Translating the mineral systems approach into an effective exploration targeting system. *Ore Geol. Rev.* 38, 128–138. <https://doi.org/10.1016/j.oregeorev.2010.05.008>
- McCuaig, T.C., Hronsky, J., 2014. The mineral systems concept: the key to exploration targeting. *Soc. Econ. Geol. - Spec. Publ.* 18, 153–175. <https://doi.org/10.1080/03717453.2017.1306274>
- Milesi, J., Ledru, P., Marcoux, E., Mougeot, R., Johan, V., Lerouge, C., Sabaté, P., Bailly, L., Respaut, J., Skipwith, P., 2002. The Jacobina Paleoproterozoic gold-bearing conglomerates, Bahia, Brazil: a “hydrothermal shear-reservoir” model. *Ore Geol. Rev.* 19, 95–136. [https://doi.org/10.1016/S0169-1368\(01\)00038-5](https://doi.org/10.1016/S0169-1368(01)00038-5)
- Minter, W.E.L., Renger, F.E., Siegers, A., 1990. Early Proterozoic gold placers of the Moeda Formation within the Gandarela Syncline, Minas Gerais, Brazil. *Econ. Geol.* 85, 943–951. <https://doi.org/10.2113/gsecongeo.85.5.943>
- Pigois, J.-P., Groves, D.I., Fletcher, I.R., McNaughton, N.J., Snee, L.W., 2003. Age constraints on Tarkwaian palaeoplacer and lode-gold formation in the Tarkwa-Damang district, SW Ghana. *Miner. Depos.* 38, 695–714. <https://doi.org/10.1007/s00126-003-0360-5>
- Razali, N.M., Wah, Y.B., 2011. Power comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors and Anderson-Darling tests. *J. Stat. Model. Anal.* 2, 21–33.
- Rodriguez-Galiano, V., Sanchez-Castillo, M., Chica-Olmo, M., Chica-Rivas, M., 2015. Machine learning predictive models for mineral prospectivity: An evaluation of neural networks, random forest, regression trees and support vector machines. *Ore Geol. Rev.* 71, 804–818. <https://doi.org/10.1016/j.oregeorev.2015.01.001>
- Royston, J.P., 1982. An Extension of Shapiro and Wilk’s W Test for Normality to Large Samples. *Appl. Stat.* 31, 115–124. <https://doi.org/10.2307/2347973>



- Saculinggan, M., Balase, E.A., 2013. Empirical power comparison of goodness of fit tests for normality in the presence of outliers. *J. Phys. Conf. Ser.* 435. <https://doi.org/10.1088/1742-6596/435/1/012041>
- Shapiro, S.S., Wilk, M.B., 1965. An analysis of variance test for normality (complete samples). *Biometrika* 52, 591–611. <https://doi.org/10.1093/biomet/52.3-4.591>
- Teixeira, J.B.G., De Souza, J.A.B., Da Silva, M. da G., Leite, C.M.M., Barbosa, J.S.F.S.F., Coelho, C.E.S., Abram, M.B., Filho, V.M.C., Iyer, S.S.S., 2001. Gold mineralization in the Serra de Jacobina region, Bahia Brazil: tectonic framework and metallogenesis. *Miner. Depos.* 36, 332–344. <https://doi.org/10.1007/s001260100174>
- Whymark, W.E., Frimmel, H.E., 2018. Regional gold-enrichment of conglomerates in Paleoproterozoic supergroups formed during the 2.45 Ga rifting of Kenorland. *Ore Geol. Rev.* 101, 985–996. <https://doi.org/10.1016/j.oregeorev.2017.04.003>
- Wickham, H., 2014. Tidy Data. *J. Stat. Softw.* 59. <https://doi.org/10.18637/jss.v059.i10>
- Wyborn, L.A.I., Heinrich, C.A., Jaques, A.L., 1994. Australian Proterozoic mineral systems: essential ingredients and mappable criteria. *Aust. Inst. Min. Metall. Publ. Ser.* 109–115.
- Yap, B.W., Sim, C.H., 2011. Comparisons of various types of normality tests. *J. Stat. Comput. Simul.* 81, 2141–2155. <https://doi.org/10.1080/00949655.2010.520163>
- Zuo, R., 2017. Machine Learning of Mineralization-Related Geochemical Anomalies: A Review of Potential Methods. *Nat. Resour. Res.* 26, 457–464. <https://doi.org/10.1007/s11053-017-9345-4>
- Zuo, R., Zhang, Z., Zhang, D., Carranza, E.J.M., Wang, H., 2015. Evaluation of uncertainty in mineral prospectivity mapping due to missing evidence: A case study with skarn-type Fe deposits in Southwestern Fujian Province, China. *Ore Geol. Rev.* 71, 502–515. <https://doi.org/10.1016/j.oregeorev.2014.09.024>



APÊNDICE A – Lista de Publicações

Trabalho Publicado em Eventos Científicos

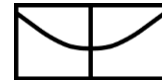
Silva, G.F., Silva, A.D.R., Lopes, C.G., Klein, E.L., Silva, A.M., Toledo, C.L.B. 2019. *k-Nearest Neighbors algorithm applied to lithology prediction based on handheld geochemistry and petrophysics analysis – a case study of metasedimentary and intrusive rocks from the Jacobina Range*. Anais do IV Simpósio Brasileiro de Metalogenia, p125-126. Gramado – RS.

Artigo científico aceitos e em preparação

Silva, G. F., Silva, A.M., Toledo, C.L.B., Chemale-Junior, F., Klein, E.L. 2022. *Predicting mineralization and targeting exploration criteria based on machine-learning in the Serra de Jacobina quartz-pebble-metaconglomerate Au-(U) deposits, São Francisco Craton, Brazil*". Manuscrito submetido ao periódico *Journal of South American Earth Sciences*, e aceito para revisão sob o código SAMES-D-21-00534R2.

Silva, G. F., Silva, A.M., Toledo, C.L.B., Teles, G.S., Chemale-Junior, F., Klein, E.L., Braga, A.A. 2022. *Machine learning applied to the analysis of mineral chemistry in pyrite grains from the Jacobina gold deposits, São Francisco Craton, Brazil: geochemical patterns and implications to mineral exploration*. Manuscrito submetido ao periódico *Journal of Geochemical Exploration*, e aceito para revisão sob o código GEXPLO-D-22-00084.

Silva, G. F., Silva, A.M., Toledo, C.L.B., Chemale-Junior, F. 2022. *Unmixing spectral signal and estimating the mineral composition of metaconglomerates using dimensionality reduction and relative distance concepts*. Manuscrito em preparação.



APÊNDICE B – Sistemas Minerais e sua aplicação na exploração mineral

Sistemas Minerais

O conceito de sistemas minerais, inicialmente desenvolvido por Wyborn et al. (1994), deriva do conceito análogo de sistemas petrolíferos, muito difundido na indústria do petróleo a partir de meados da década de 1980. Os sistemas petrolíferos categorizam todos os processos e elementos geológicos necessários para a formação e armazenamento de óleo e gás. Da mesma forma, os sistemas minerais foram definidos como “todos os fatores geológicos que controlam a geração e preservação de depósitos minerais, cujos processos são relacionados à mobilização do minério desde a fonte até a região de concentração, transporte e acumulação, e sua posterior preservação ao longo da história geológica” (Wyborn et al., 1994).

Ainda segundo Wyborn et al. (1994), a maior parte dos corpos de minério possuem menos do que 1km² de expressão, não consistindo portanto de um alvo significativo para a prospecção mineral. Felizmente, apesar de os depósitos não terem uma área muito expressiva e serem resultados de uma coincidência excepcional de determinados processos geológicos, estes processos podem ser mapeáveis em escalas regionais, sendo chave importante para o processo de prospecção. Em outras palavras, apesar de o depósito consistir em áreas de centenas de metros, o sistema total de interação entre fluido-rocha encaixante-mineralização pode se estender por áreas de até poucas dezenas de quilômetros, sendo detectáveis por determinadas ferramentas (Wyborn et al., 1994).

Os fatores críticos para a caracterização de qualquer sistema mineral incluem:

1. fonte do fluido mineralizadores e dos compostos ligantes;
2. fonte dos metais e outros componentes da mineralização;
3. caminhos de migração (*pathways*) do fluido mineralizante;
4. gradiente térmico;
5. Fonte de energia (por vezes relacionado ao item 4);
6. Estruturas ou mecanismos de concentração;
7. Condições químicas e físicas para a deposição da mineralização.

McCuaig et al., (2010) avaliam que o conceito de sistema mineral evoluiu à medida que foi aceito gradualmente pela indústria e academia nos 15 anos anteriores, apesar de que fora pouco utilizado na rotina das empresas de mineração. É enfatizado que a linha estruturante do conceito de sistema mineral está associada à compreensão dos vários processos geológicos que operam em todas as escalas, ao invés de focar na compreensão das características particulares de depósitos específicos em sua escala local. Desta forma, ressaltando que há certa dificuldade



em traduzir essa mudança de paradigma para a realidade da prospecção. McCuaig et al., (2010) propõem uma sistemática de quatro passos para guiar o *pipeline* de exploração através da ótica dos Sistemas Mineraiis.

Ainda nesta linha, McCuaig & Hronsky (2014) também contribuem para a evolução do conceito de Sistema Mineral. Estes autores postulam que a existência de um depósito mineral está condicionada a sobreposição de pelo menos quatro fatores críticos ao longo da história geológica, sendo eles fertilidade do terreno, favorabilidade geodinâmica, arquitetura litosférica e a posterior preservação das zonas de concentração de minério em um sistema mineral.

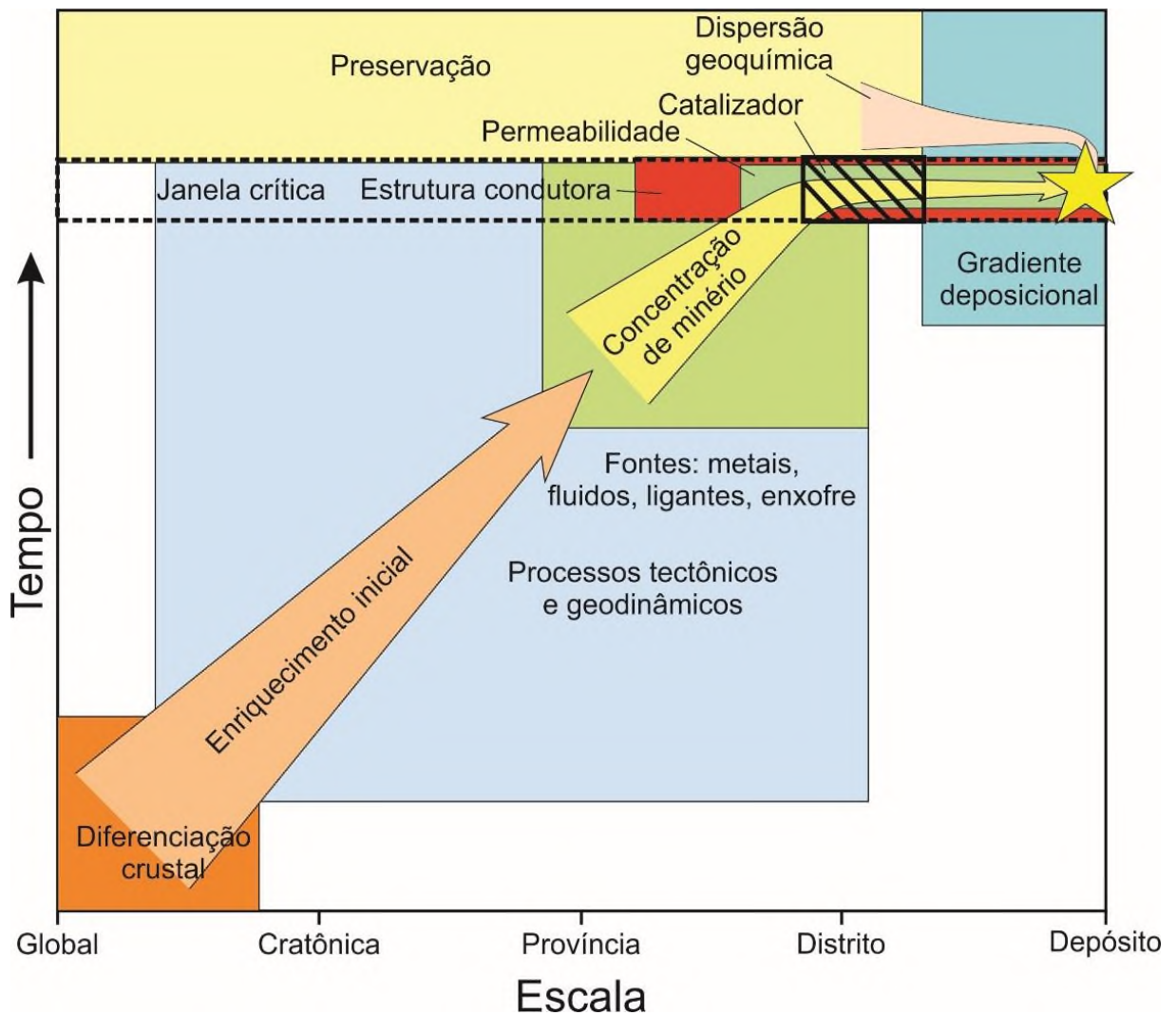


Figura B- 1 – Síntese dos processos envolvidos nos Sistemas Mineraiis nas suas variadas escalas, desde a diferenciação crustal até a dispersão do fluido em superfície (Fonte: www.ga.gov.au.)

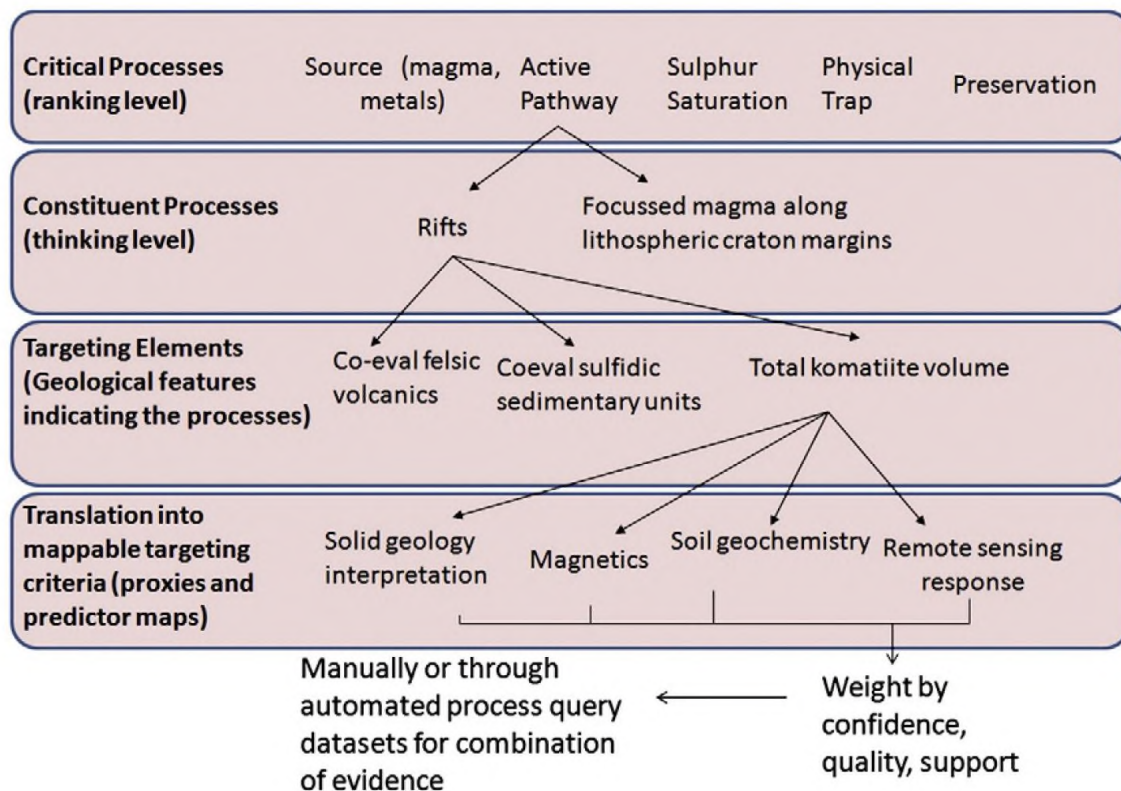
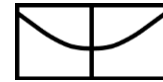


Figura B- 2 – Principais vertente de um sistema mineral do tipo Ni-Cu em komatiitos proposto por McCuaig et al., (2010). Observe o encadeamento dos processos geológicos desde a escala crustal até o mapeamento dos alvos potenciais em escala de depósito.

A fertilidade do terreno é definida como a tendência de uma região ou de uma época geológica de ser mais favorável para a formação de depósitos minerais do que outras (McCuaig & Hronsky, 2014). A fertilidade varia de acordo com a evolução da crosta em diversos momentos tais como processos de rifteamento e colisão subsequente, formação de supercontinentes, dentre outros.

A arquitetura litosférica tem relação com os padrões estruturais associados com as mineralizações, como tendências estruturais ou *orshoots*, assim como com as estruturas de dimensão crustais, que cortam do embasamento às sequências de topo, e por muitas vezes são utilizadas como condutos pelos fluidos mineralizadores em seu processo de migração. Este critério é muito importante em processos hidrotermais, porém é também relevante na formação de depósitos minerais associados a magmatismo, como em pórfiros, *greisens* ou outros depósitos associados a intrusões (McCuaig & Hronsky, 2014).

Com a evolução das técnicas de datação e a difusão destes métodos, foi possível perceber que os grandes depósitos minerais ocorrem em momentos muito restritos da história da terra (McCuaig & Hronsky, 2014). Desta forma, a condição geodinâmica pode ser notada em depósitos distantes entre si por centenas de quilômetros, mas formados em uma mesma



época geológica e com condições similares. Os autores destacam três ambientes geodinâmicos favoráveis para a formação de grandes depósitos:

1. estágios iniciais de eventos extensionais, com ascensão de magma mantélico e/ou instalação de uma pluma sob a litosfera;
2. compressão intermitente, importante na formação de depósitos tipo pórfiro;
3. variações na direção da deformação que resultem em um campo de tensões neutro, com ausência de estruturas geradoras de permeabilidade secundária.

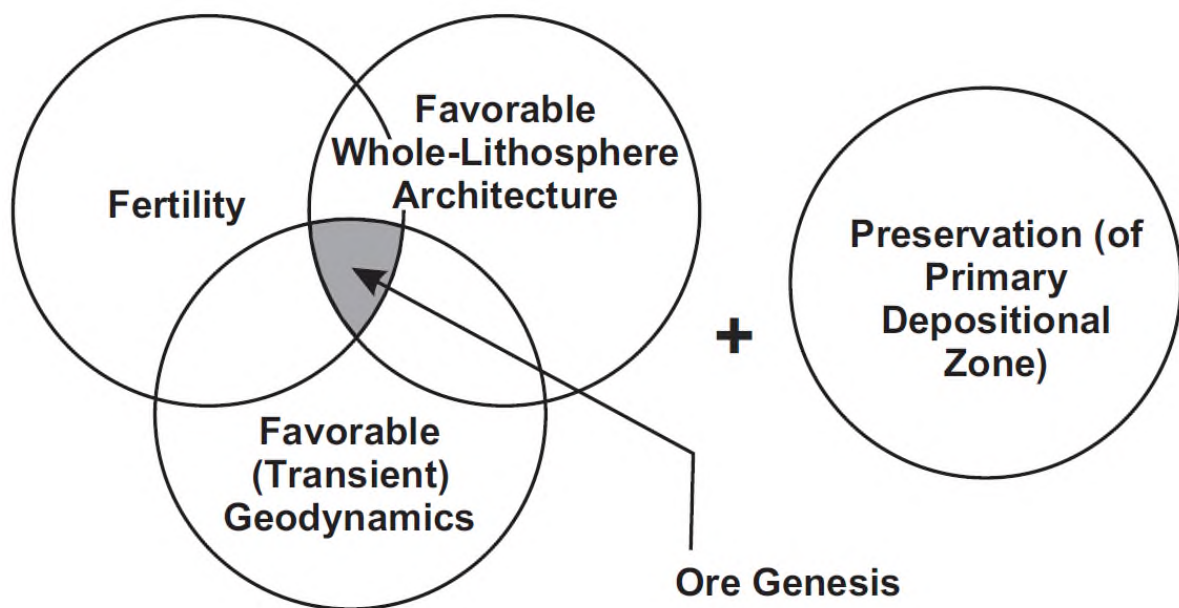
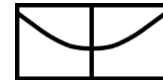


Figura B- 3 – Fatores Críticos para a formação de um depósito mineral (McCuaig & Hronsky, 2014).

Modelagem de Potencial Mineral

O desenvolvimento de uma abordagem reprodutível para identificar locais com alto potencial para exploração de uma *commodity* mineral é o objetivo central dos estudos da prospectividade mineral (Joly et al., 2012). Para tanto, se faz necessário um grande conjunto de informações consistentes em uma perspectiva multiparamétrica, aplicável à escala de interesse. Assim, os modelos de prospecção são tentativas de emular os processos formadores e dispersores de mineralizações, a fim de detectar novos alvos baseados em critérios do sistema mineral estabelecido.

O modelo de prospectividade pode ser construído com base em informações pré-definidas através de modelos orientados pelo conhecimento (ou *knowledge-driven model*) ou com base na assinatura de um depósito mineral conhecido por modelos orientados pelos dados (ou *data-driven model*). Desta forma, pode-se desenvolver uma prospecção mineral através de uma concepção regional de modelo geológico para a mineralização ou através da similaridade



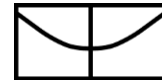
de condições com um depósito de minério conhecido previamente (Carranza, 2009; Joly et al., 2012).

Os modelos orientados pelo conhecimento prévio são muitas vezes adequados para áreas onde o conhecimento geológico é limitado (*greenfield*) ou as ferramentas disponíveis são escassas. Os intérpretes se baseiam em modelos clássicos de mineralizações, e tentam estimar o tipo de resposta esperada em cada ferramenta utilizada. Se a conjunção de respostas for como esperada, gera-se um alvo a ser verificado, caso contrário, troca-se o modelo de mineralização ou a área é descartada. Estes modelos têm a vantagem de serem mais versáteis, porém são altamente tendenciosos e dependem exclusivamente da experiência da equipe de prospecção e da similaridade entre os depósitos (Agterberg & Bonham-Carter, 2005; Bonham-Carter, 1994; Carranza, 2009).

Já os modelos orientados por dados, ou também modelos empíricos, são mais complexos, sendo construídos numericamente pelo conjunto dados fornecidos e são usados para prever, não para explicar as mineralizações. Um modelo empírico simula uma função matemática que captura a tendência dos dados. Essa tendência, ou padrão, pode ser inferida pela análise e interpretação direta, ou com a ajuda de lógica computacional pré-compilada (como Algoritmos de Aprendizado de Máquina ou *Machine Learning Algorithms* – MLA), que podem dar uma resposta não tendenciosa do problema envolvido (Agterberg & Bonham-Carter, 2005; Bonham-Carter, 1994; Carranza, 2009).

Para maior eficácia, os fatores críticos de formação da mineralização deverão corresponder a pelo menos um guia prospectivo, representados como vetores em um espaço N-dimensional, onde N é o número total de variáveis envolvidas no modelo desenvolvido. Este método seria mais adequado em províncias e distritos minerais razoavelmente maduros, onde o volume de informação gerada e sistematizada poderia satisfazer estes pré-requisitos. Porém, devido ao grande volume de dados, é esperado algum problema para a integração e processamento de toda a informação gerada, uma vez que os dados podem ser coletados com metodologias, momentos e escalas distintos. Assim, integração eficaz de dados passaria a ser um novo desafio.

Para Hagemann et al. (2016), a vantagem da abordagem de sistema mineral em relação a descrição taxonômica de depósitos. Ela enfocaria os processos geológicos críticos necessários para formar grandes mineralizações, além de incluir todos os elementos descritivos de um estilo específico de mineralização. Sistemas minerais podem explicar a coexistência espacial e temporal de depósitos minerais dentro de uma província mineral específica. Também pode



explicar famílias de sistemas minerais coevos que potencialmente se formaram nos mesmos terrenos ou províncias adjacentes. Apesar disto, estes autores relacionam uma série de desafios a serem superados para que esta linha de pesquisa.

Apesar de uma aparente atração de bases científicas sólidas e consistência metodológica interna, traduzir o entendimento teórico de sistemas minerais em modelos efetivos de prospectividade mineral e alvos específicos de exploração ainda é um grande desafio em vários níveis conceituais (Hagemann et al., 2016). A dificuldade mais séria (além da compreensão inevitavelmente incompleta e evolutiva dos sistemas minerais) permanece definindo critérios mapeáveis que representam adequadamente vários elementos do sistema mineral críticos para formar, expor e preservar províncias minerais férteis e zonas dotadas e campos minerais dentro deles. Isso é particularmente problemático porque as indicações diretas desses elementos críticos e correspondentes processos geológicos de grande escala, como uma fonte de metal fértil e uma arquitetura de subsolo pré-existente favorável, geralmente não são observáveis ou têm apenas expressões muito sutis em um nível crustal de formação de depósitos minerais.

Vários desenvolvimentos recentes nas geociências podem ajudar a melhorar a eficácia da segmentação de exploração preditiva em escalas regionais (província a distrito), embora algumas delas ainda devem ser totalmente apreciados e adequadamente abordados na prática comum de direcionamento de exploração. Esses desenvolvimentos incluem o reconhecimento de: (i) dependência de escala de elementos e processos do sistema mineral; (ii) composição da arquitetura tectônica como os principais fatores que controlam a fertilidade e a operação do sistema mineral na província para a escala de distrito metalogênico; (iii) limites de domínios crustais profundos e outras estruturas pré-existentes no subsolo com alguma expressão em superfície como controle fundamental do sistema mineral na província e nas imediações dos depósitos; (iv) efeitos da incerteza na tomada de decisão sobre o objetivo da exploração (Hagemann et al., 2016).

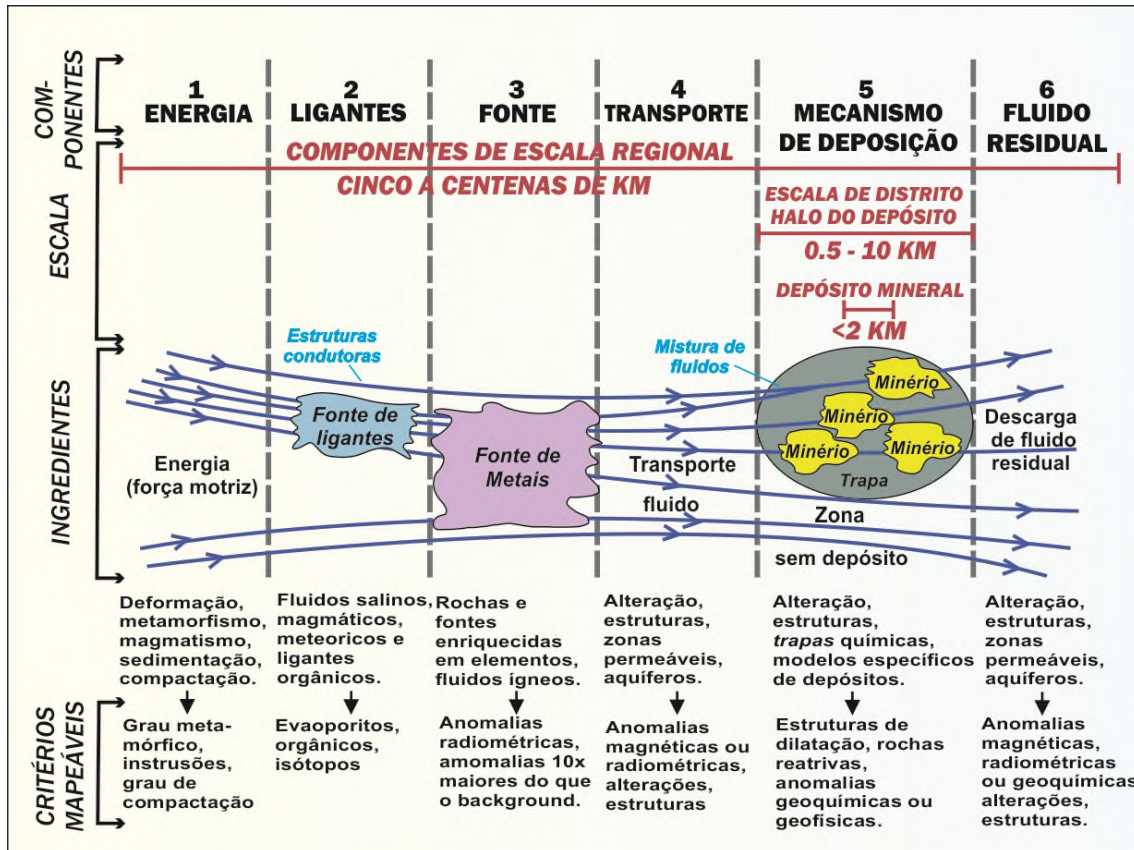
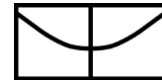
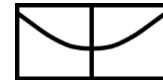


Figura B- 4 – Modelo genérico de Sistemas Minerais, com indicação de processos, escala de observação e possíveis critérios mapeáveis em levantamentos geológicos, geofísicos ou geoquímicos de exploração (basedo em Hagemann et al., 2016).

Em seu trabalho mais recente Hronsky & Kreuzer (2019) alegam que apesar de muitas décadas de desenvolvimento, a modelagem de prospectividade ainda não é amplamente utilizada ou aceita mundialmente em toda a indústria de exploração mineral, à exceção de equipes de empresas como DeBeers, Newmont e Kenex (Hronsky & Kreuzer, 2019). Uma crítica comum ao método é que ele não é praticamente útil porque tem sido utilizado para amadurecer regiões já conhecidas e muitas vezes gera áreas excessivamente grandes de alta prospectividade. Estes autores sugerem que a razão para isso não esteja primariamente relacionada a limitações nos algoritmos de mapeamento de prospectividade, mas a questões relativas ao uso de conjuntos de dados de entrada.

De acordo com Hronsky & Kreuzer (2019), críticas comuns das tentativas de modelagem prospectiva no setor de exploração relacionam-se a duas questões principais. Em primeiro lugar, a técnica é muito mais eficaz em encontrar os depósitos já conhecidos do que gerar novos alvos válidos. Isso obviamente se relaciona, pelo menos em parte, ao fato de que, na maioria das metodologias de modelagem de prospectividade comumente adotadas, os locais de depósito conhecidos são uma entrada de modelo chave. As exceções são sobreposição de



lógica difusa, sobreposição ponderada e modelos de “Sistemas de Inferência Fuzzy”, que não exigem locais de depósito como entrada (Bonham-Carter, 1994; Carranza, 2009). A segunda grande crítica é que, além de áreas adjacentes a depósitos conhecidos (e, portanto, óbvias), o próximo nível de domínios de alta prospecção gerados é tipicamente grande, relativo à área de interesse. Esta é uma questão muito importante porque, para ser praticamente útil na exploração mineral, qualquer técnica de direcionamento deve produzir pelo menos uma redução de ordem de grandeza na área de foco (Hronsky & Kreuzer, 2019) .

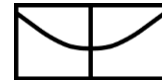
Finalmente, os autores apresentam algumas sugestões a serem tomadas a fim de avançar esta linha de pesquisa e torna-la mais adequada a realidade da indústria (Hronsky & Kreuzer, 2019):

- Desenvolvimento de novos fluxos de trabalho e metodologias de modelagem de prospectividade efetivos que estejam fortemente alinhados com a prática no mundo real da exploração mineral. Isso não deve ocorrer se esse desenvolvimento for deixado somente para os geoestatísticos. Ao invés disso, deve-se encorajar a formação de equipes colaborativas multidisciplinares com experiência em modelagem prospectiva e exploração mineral.

- Aplicação de uma abordagem “híbrida”, focada na amplificação da inteligência e não na inteligência artificial (IA), que aproveita poder combinado da mente humana para reconhecer, mapear e extrapolar padrões com o rigor de algoritmos baseados em máquinas. A chave para essa abordagem é a compilação dos principais mapas geológicos interpretativos que se situam intermediários entre os dados de entrada primários e os algoritmos de mapeamento de prospectividade, como os MLA. Em seguida, estes dados seriam usados como as principais entradas no processo de modelagem de prospectividade, adicionalmente a quaisquer conjuntos de dados disponíveis sistematicamente amostrados.

- Os Serviços Geológicos devem fornecer juntamente com as demais camadas de dados pré-competitivas as suas melhores interpretações de estruturas em larga escala que não são facilmente observáveis na geologia de superfície. O uso de dados gravimétricos regionais de alta resolução podem ser uma ferramenta muito importante para auxiliar nestas interpretações.

Ainda segundo Hronsky & Kreuzer (2019), embora os problemas discutidos limitem fortemente a aplicação da técnica de modelagem de prospectiva como atualmente praticada pela maioria, eles não são barreiras para a implementação bem-sucedida dessa tecnologia no futuro. Sugere-se que o método mais eficaz possa ser um híbrido de interpretação geológica humana subjetiva e análise objetiva com suporte em algoritmos, que capture os melhores aspectos dessas abordagens alternativas. Isso exigiria uma boa integração dos dados geológicos, geofísicos e



geoquímicos básicos disponíveis em camadas interpretativas que forneçam corretamente as entradas primárias para a análise de modelagem de prospectividade.

Aprendizagem de máquina para integração de dados

Os modelos orientados por dados são inferências matemáticas baseado em medidas de parâmetros quantitativos e na relação espacial destes resultados com os depósitos conhecidos (Agterberg & Bonham-Carter, 2005; Bonham-Carter, 1994; Carranza, 2009; Rodriguez-Galiano et al., 2015). Historicamente, os modelos “*data-driven*” desenvolvidos são baseados em regressões probabilísticas ou em lógica Bayesiana, o que resulta em inferências estatísticas complexas, por vezes operadas em ferramentas complexas em ambiente SIG, o que inibe o seu uso. Com o advento da popularização de técnicas de inteligência artificial para gerenciamento de dados, como os algoritmos de aprendizagem de máquina, o processamento de dados para modelagem *data driven* tem se tornado mais simples e rápido.

Algoritmos de aprendizagem de máquina (MLA na sigla em inglês, *Machine Learning Algorithms*) como redes neurais artificiais (*Artificial Neural Networks – ANN*), árvores de regressão (*Regression Tree – RT*), florestas aleatórias (*Random Forests – RF*) e máquinas de vetores de suporte (*Support Vector Machine – SVM*) são métodos *data-driven* eficientes que podem ser usados para identificar padrões em um conjunto de dados, visando modelagem de prospectividade mineral (Abedi et al., 2012; Bérubé et al., 2018; Breiman, 2001; Breiman et al., 1995; Carranza & Laborte, 2016, 2015a; Costa et al., 2019; Rodriguez-Galiano et al., 2015; Zuo, 2017; Zuo et al., 2015). Via de regra, os métodos de MLA para classificação automática necessitam reconhecer parte dos dados como forma de “calibrar” (ou treinar) as predições.

Rodriguez-Galiano et al. (2015) comparou a performance de modelos gerados por ANN, RT, RG e SVM no Distrito Mineiro de Rodalquilar, no sul da Espanha, uma área com ocorrências de ouro epitermal. Os dados de entrada para elaboração dos modelos são baseados em dados geológicos, geoquímicos, geofísicos e de sensoriamento remoto hiperespectral publicados em estudos anteriores (Rodriguez-Galiano et al., 2015).

A análise comparativa dos métodos MLA para modelagem mineral prospectividade foi realizada a partir de diferentes perspectivas: facilidade de aplicação e eficácia, sensibilidade à configuração dos parâmetros do modelo e redução de dados, precisão do mapeamento das classificações, e transparência e coerência dos modelos.

Rodriguez-Galiano et al. (2015) avaliam que os modelos apresentam uma dificuldade diferente em sua formulação (ou treinamento). Algoritmos baseados em árvore de decisão (RT e RF) envolvem são mais facilmente treinados do que os demais. Os autores afirmam ainda que



o desempenho dos métodos é extremamente dependente do conjunto de dados utilizados para a calibração inicial, sendo que o método que obteve melhor desempenho com um conjunto de dados de calibração mínimo foi o RF, sendo, portanto, o método mais indicado para uso em áreas onde há pouca informação geológica prévia.



1 APÊNDICE C – Código utilizado no Artigo 01

2 NOTEBOOK

3 This is the source code used in the manuscript:

4 “Predicting mineralization and targeting exploration criteria based on machine-learning in the
5 Serra de Jacobina quartz-pebble-metaconglomerate Au-(U) deposits, São Francisco Craton,
6 Brazil”

7 authored by: “Guilherme Ferreira” date: “02/03/2022”

8 The following code was written in R (3.5.6)

9 ABSTRACT

10 Defining mineral exploration criteria is a laborious, time-consuming, and generally an
11 empirical task often biased and limited to expert knowledge. To address this problem with a
12 different approach, we used data-driven analysis to make predictions and provide insights
13 about gold mineralization in rocks of the Jacobina Group, São Francisco Craton. The input
14 variables were petrophysical parameters (density, magnetic susceptibility, and electric
15 conductivity) and lithochemistry data obtained by X-Ray Fluorescence assays. A machine
16 learning model based on the Random Forests algorithm was applied to predict mineralization
17 in drill core samples. The database used for algorithm training was balanced using the
18 Borderline-SMOTE technique to provide approximately the same numbers of samples of the
19 two classes in the mineral status parameter (i.e., ore and barren samples). The quality of the
20 predictions was assessed with different datasets (i.e., training, testing, each drill core
21 separately, and all samples) and by parameters. The average accuracies were 0.87 for cross-
22 validation training, 0.91 for testing, and 0.86 for all samples. Also, the model allowed us to
23 estimate and rank the importance of the input variables to the prediction. These estimates
24 were validated by an interpretation of optical and scanning electron microscopy petrographic
25 analysis, which was carried out to understand the relationship between minerals of different
26 stages and gold mineralization. Thus, the techniques used in this work could help to decrease
27 the time spent in data integration and interpretation, as mineral exploration teams can easily
28 replicate this approach.



29 DATA WRANGLING

30 Dependencies

```
31 library(tidyverse) # ggplot2, tidyr, dplyr
32 library(readxl) # open XLSX data
33 library(geoquimica) # Data wrangling
34 library(caret) # Machine Learning
35 library(doParallel) # Parallel Processing
36 library(randomForest) # RF
37 library(randomForestExplainer) # RF
38 library(pROC) # ROC and AUC
39 library(smotefamily) # Smote
```

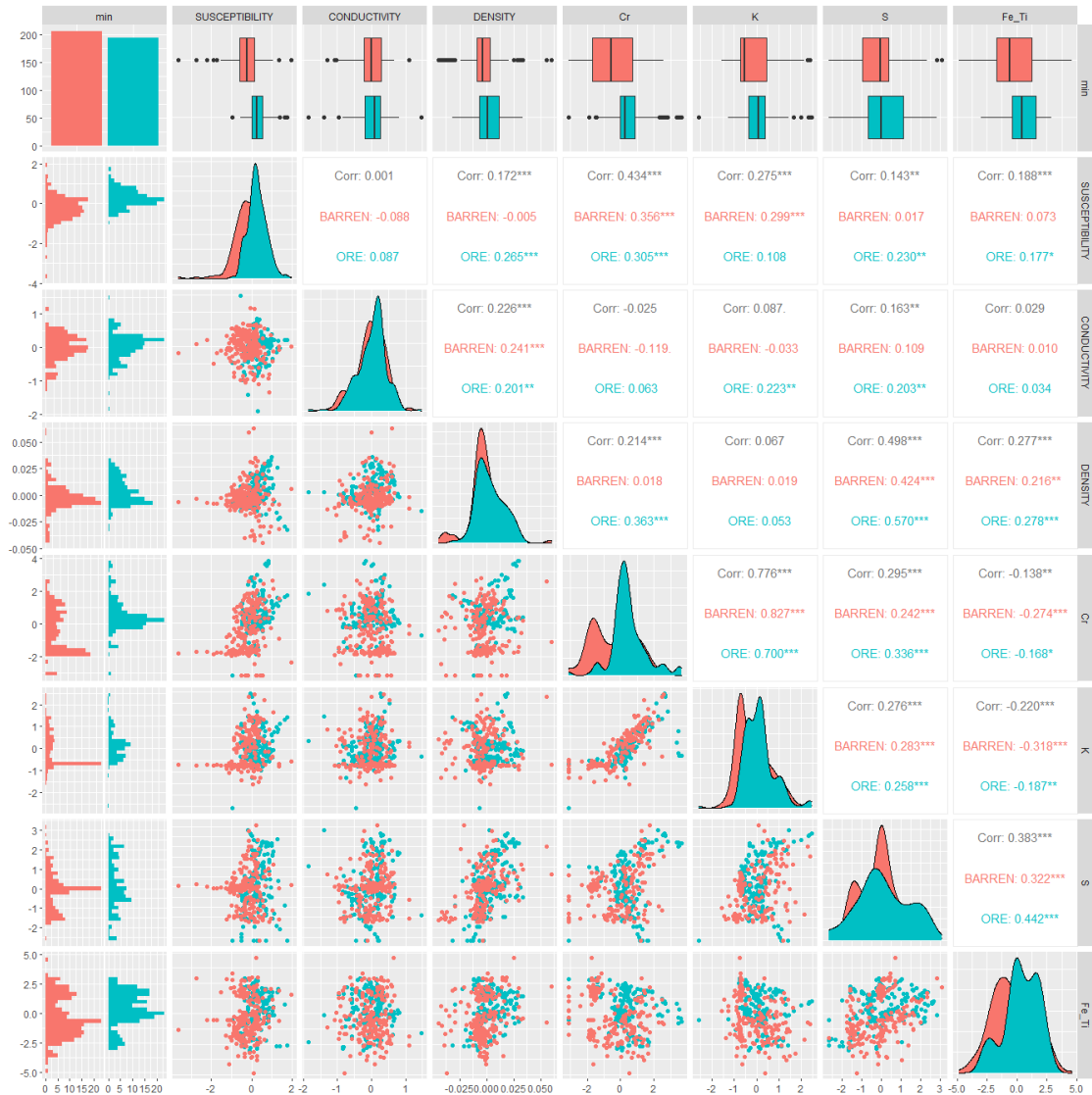
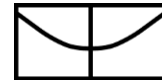
40 Data preparation

```
41 set.seed(0)
42
43 setwd('~/.GitHub/jacobina/data/pphy')
44
45
46 xrf <- read_xlsx(path = '~/.GitHub/jacobina/data/xrf/pXRF_Jacobina_SAMPLES.xlsx',
47                sheet = 1)
48
49 # Petrophysics data ----
50 files <- list.files(pattern = '.xlsx$', path = '~/.GitHub/jacobina/data/pphy')
51
52 phy <- lapply(files, read_xlsx, sheet = 1) %>%
53   bind_rows() %>%
54   mutate(HOLE = as.factor(HOLE),
55          ID = as.factor(ID),
56          FROM = as.numeric(FROM),
57          TO = as.numeric(TO),
58          LITHO = as.factor(LITHO),
59          MINERALIZATION = as.factor(MINERALIZATION),
60          SUSCEPTIBILITY = as.numeric(SUSCEPTIBILITY),
61          CONDUCTIVITY = as.numeric(CONDUCTIVITY),
62          DENSITY = as.numeric(DENSITY),
63          COMMENTS = as.character(COMMENTS)) %>%
64   arrange(ID)
65
66 phy <- phy %>%
67   mutate(SAMPLE = paste(HOLE, formatC(x = phy$FROM, flag = '0',
68                                     width = 6,
69                                     digits = 2,
70                                     format = 'f'), sep = '-'),
71          ROCK = case_when(phy$LITHO %in% c('GRIT', 'LMPC', 'LVLPC', 'MLPC', 'MPC', 'MSPC',
72                                           'SMPC', 'SPC', 'VSPC') ~ 'CONGLOMERATE',
73                            phy$LITHO %in% c('QTO', 'QTO_SX', 'QZ_VEIN') ~ 'QUARTZITE',
74                            phy$LITHO %in% c('ITV', 'UMF') ~ 'ULTRAMAFIC',
75                            phy$LITHO %in% c('BRX') ~ 'BRECCIA',
76                            phy$LITHO %in% c('XISTO') ~ 'SCHIST',
77                            phy$LITHO %in% c('SOLO') ~ 'SOIL',
78                            TRUE ~ as.character(phy$LITHO)))
79 # Merging dataset ----
80
81 df <- phy %>%
82   left_join(xrf, by = 'SAMPLE') %>%
83   mutate(min = factor(ifelse(test = phy$MINERALIZATION == 1 | phy$MINERALIZATION == 1000,
84                             yes = 'ORE', no = 'BARREN')),
85          Fe_Ti = Fe/Ti)
```



86 SMOTE

```
87 set.seed(0)
88
89 conglomerate <- as.data.frame(df) %>%
90   filter(ROCK == 'CONGLOMERATE') %>%
91   dplyr::select(min,7:9, Cu, Fe, Cr, Ti, K, Al, Si, S, Fe_Ti) %>%
92   na.omit()
93
94 # Split data to smote
95 index <- caret::createDataPartition(conglomerate$min,
96                                     p = 1,
97                                     list = FALSE,
98                                     times = 1)
99
100 toSmote <- conglomerate[index,]
101
102 fromSmote <- BLSMOTE(C = 5, dupSize = 0,
103                    X = as.data.frame(toSmote[, -1]),
104                    K = 5,
105                    target = as.data.frame(toSmote[, 'min']))
106
107 ## [1] "Borderline-SMOTE done"
108
109 fromSmote$data %>%
110   rename(min = class) %>%
111   select(min,1:3,Cr, K, S, Fe_Ti) %>%
112   elem_norm(method = 'clr') %>%
113   GGally::ggpairs(mapping=ggplot2::aes(colour = min), progress = FALSE)
114
115 df_smote <- fromSmote$data %>%
116   rename(min = class) %>%
117   mutate(min = as.factor(min))
```

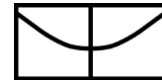
115

116 Fig. C. 1: Exploratory data analysis for selected variables after the BLSMOTE balancing. Data are color coded
 117 according to the Mineralization Status (i.e., Ore or Barren). The asterisk indicates the level of significance of the
 118 correlations. * for alpha = 0.15, ** for alpha = 0.05, and *** for alpha = 0.01.

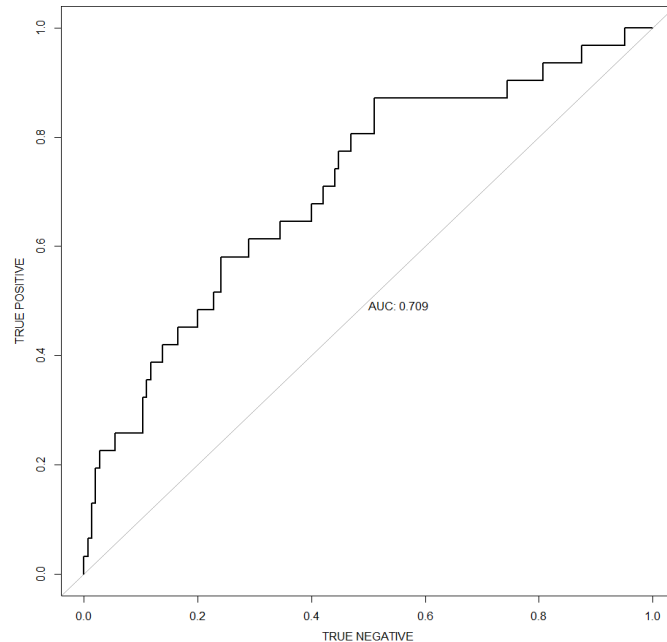
119 RANDOM FORESTS

```

120 # Imbalanced Model
121
122 splitIndex <- caret::createDataPartition(conglomerate$min,
123                                         p = .7,
124                                         list = FALSE,
125                                         times = 1)
126
127 trainSplit <- conglomerate[splitIndex,]
128 testSplit <- conglomerate[-splitIndex,]
129
130
131 imbal_minModel <- randomForest(min ~ .,
132                                trainSplit,
133                                proximity = TRUE,
134                                ntree = 1000,
135                                localImp = TRUE)
136
  
```



```
137 par(pty = 's')
138 roc1 <- plot.roc(trainSplit$min, imbal_minModel$votes[,1], legacy.axes=TRUE,
139                percent = FALSE, print.auc = TRUE, xlab="TRUE NEGATIVE",
140                ylab = 'TRUE POSITIVE')
141 par(pty = 'm')
```



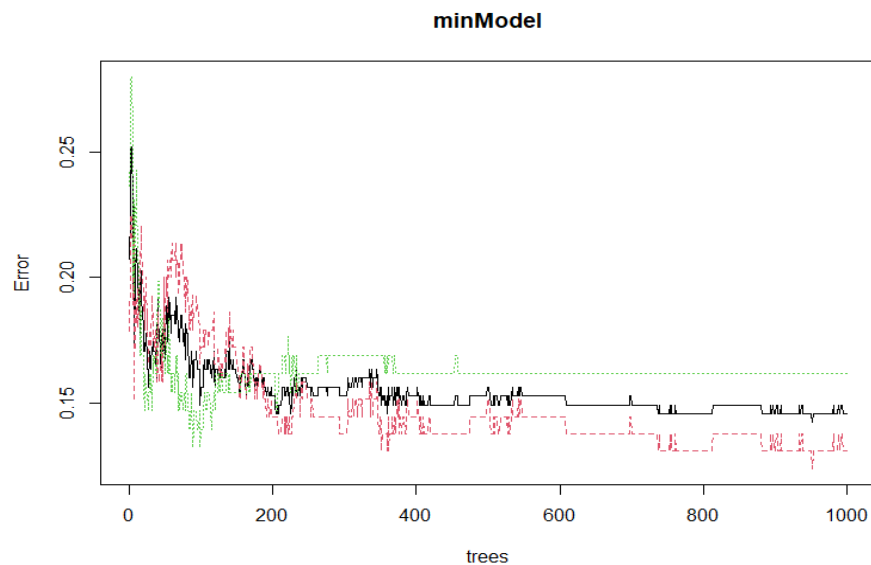
142

143 Fig. C. 2: AUC for imbalanced dataset

```
144 # Balanced Model
145
146 splitIndex <- caret::createDataPartition(df_smote$min, p = .7,
147                                         list = FALSE,
148                                         times = 1)
149
150 trainSplit <- df_smote[splitIndex,]
151 testSplit <- df_smote[-splitIndex,]
152
153 round(prop.table(table(trainSplit$min)), digits = 3)
154
155 ##
156 ## BARREN ORE
157 ## 0.516 0.484
158
159 round(prop.table(table(testSplit$min)), digits = 3)
160
161 ##
162 ## BARREN ORE
163 ## 0.517 0.483
164
165 doParallel::registerDoParallel(cores = 3)
166 ctrl <- caret::trainControl(method = 'repeatedcv',
167                             number = 5,
168                             repeats = 3, search = 'grid',
169                             allowParallel = TRUE,
170                             verboseIter = TRUE)
171
172 minMod_fram <- min_depth_distribution(minModel)
173
174 print(minModel)
```



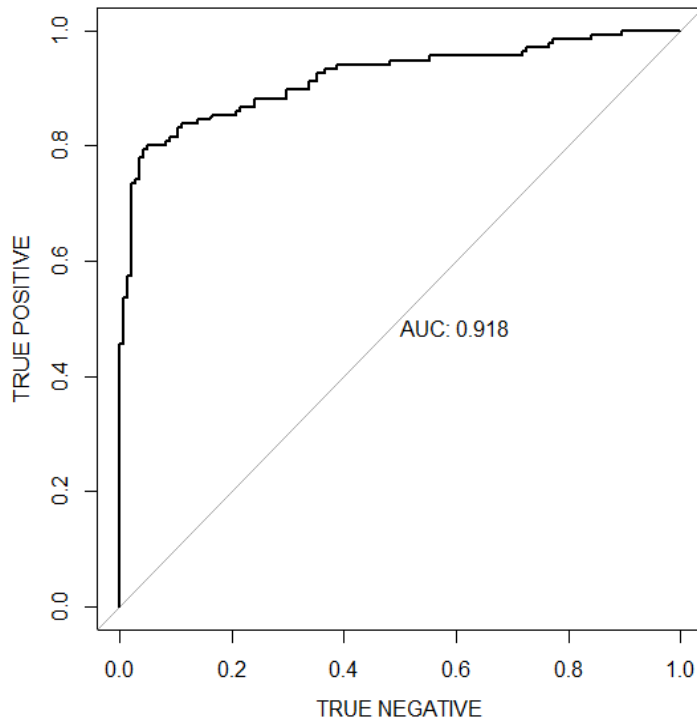
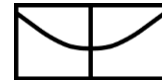
```
171 ##
172 ## Call:
173 ## randomForest(formula = min ~ ., data = trainSplit, proximity = TRUE,      ntree = 1000,
174 localImp = TRUE, mtry = tuning$bestTune[[1, 1]])
175 ##           Type of random forest: classification
176 ##           Number of trees: 1000
177 ## No. of variables tried at each split: 3
178 ##
179 ##           OOB estimate of error rate: 14.59%
180 ## Confusion matrix:
181 ##           BARREN ORE class.error
182 ## BARREN      126  19  0.1310345
183 ## ORE          22 114  0.1617647
184 plot(minModel)
```



185

186 Fig. C. 3: Error rate (1 – Accuracy) and number of estimators (trees) for the Random Forests models

```
187 par(pty = 's')
188 plot.roc(trainSplit$min, minModel$votes[,1], legacy.axes=TRUE,
189          percent = FALSE, print.auc = TRUE, xlab="TRUE NEGATIVE",
190          ylab = 'TRUE POSITIVE')
191 par(pty = 'm')
```

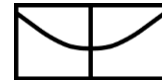


192

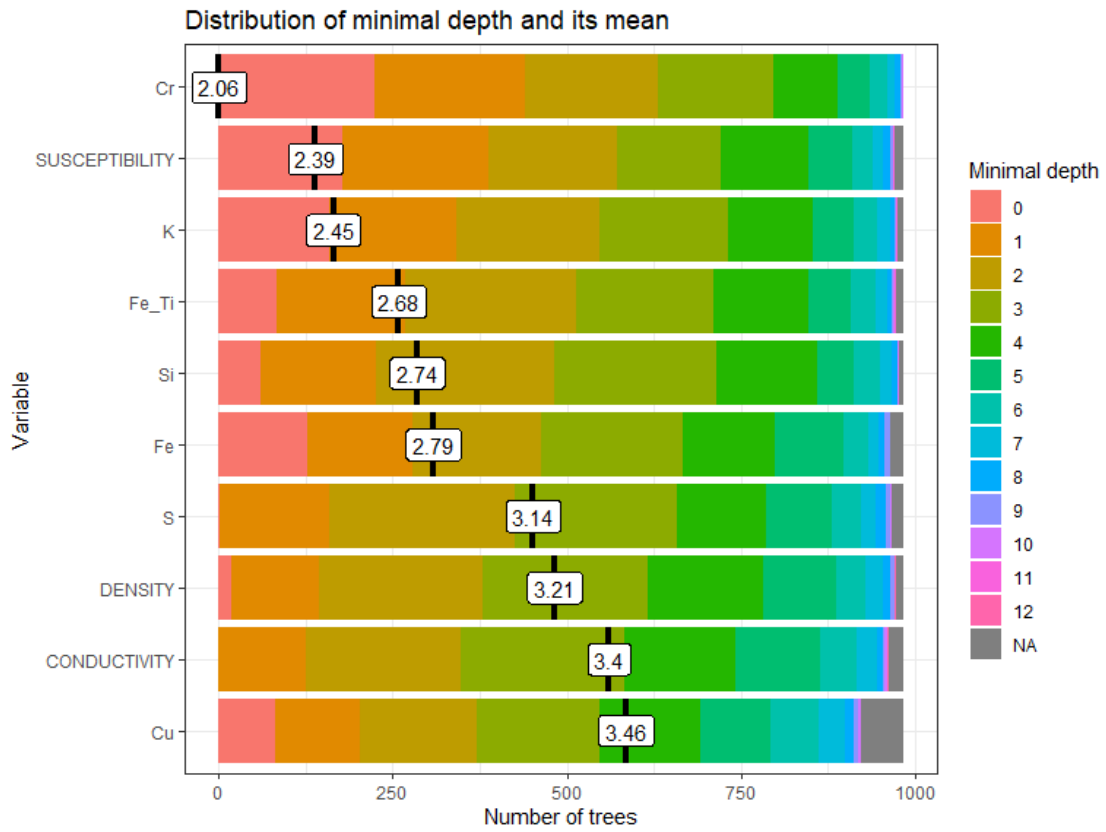
193 Fig. C. 4: AUC for balanced model

194

195

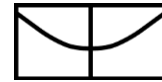


196 `plot_min_depth_distribution(minModel)`



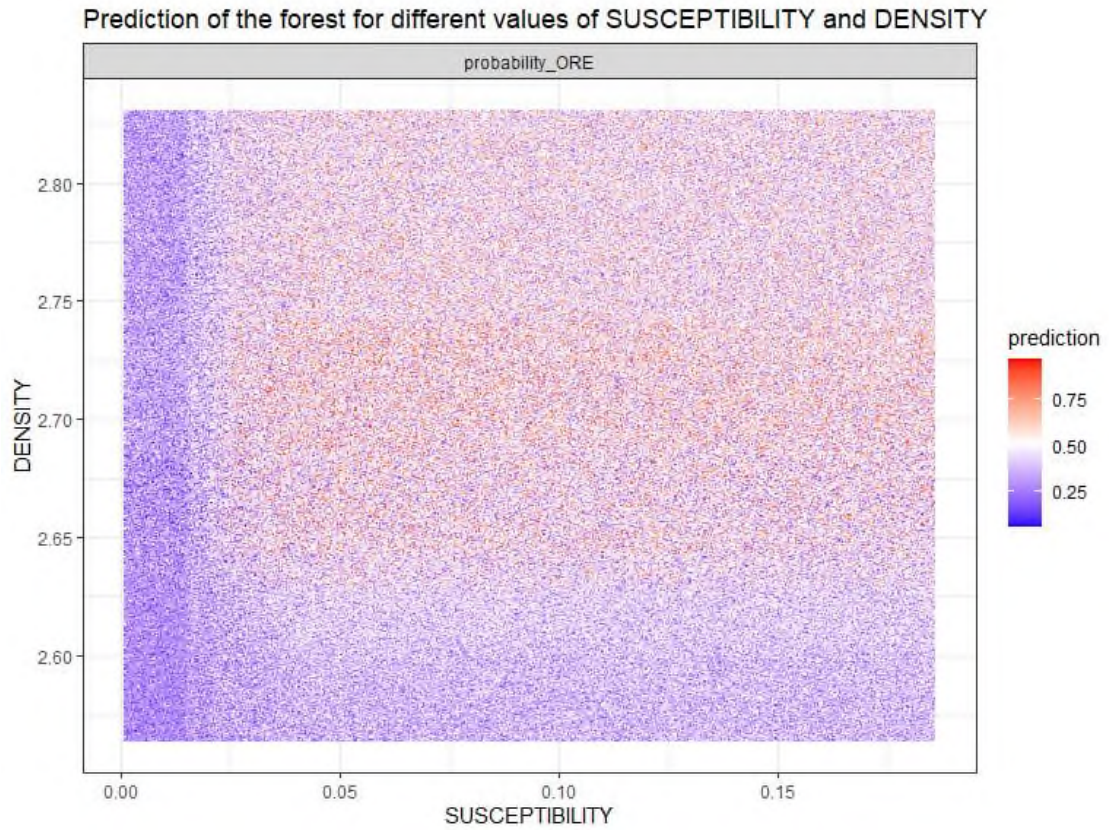
197

198 Fig. C. 5: Variable importance ranked according to the average minimal depth and number of trees
199



200

```
201 plot_predict_interaction(forest = minModel,  
202                        data = trainSplit,  
203                        variable1 = "SUSCEPTIBILITY",  
204                        variable2 = "DENSITY",  
205                        grid = 500)
```

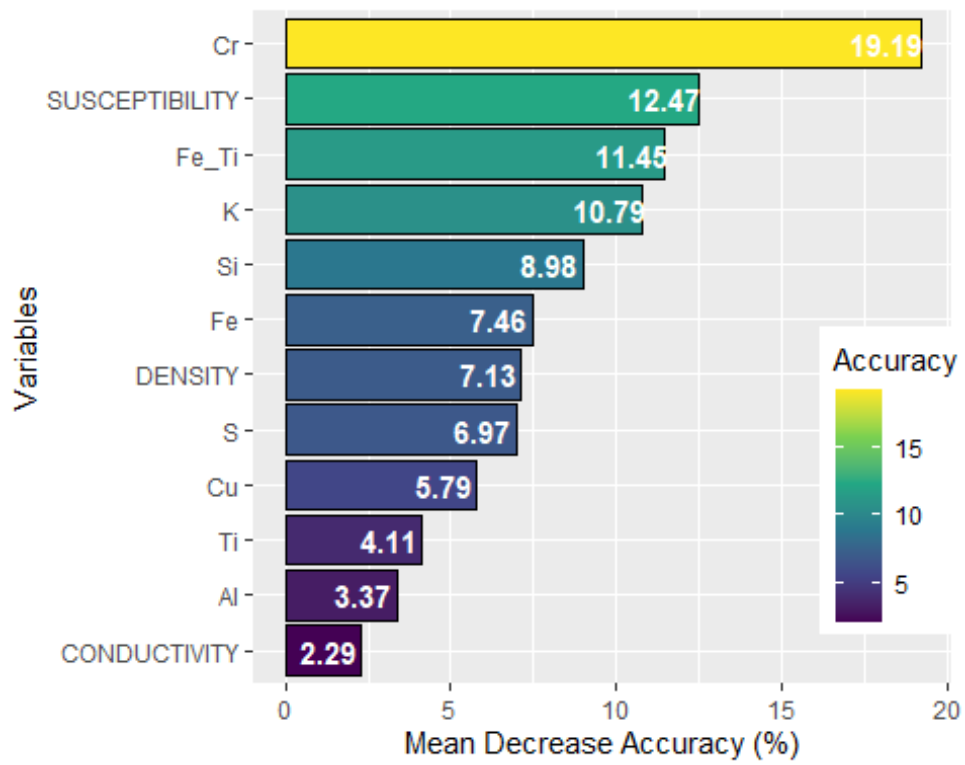
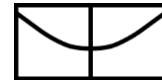


206

207 Fig. C. 6: prediction grid of the probability for been classified as ORE, according to Susceptibility and Density
208 values
209



```
210
211 confusionMatrix(as.factor(trainSplit$min), as.factor(minModel$predicted),positive = 'ORE')
212 ## Confusion Matrix and Statistics
213 ##
214 ##           Reference
215 ## Prediction BARREN ORE
216 ##   BARREN    126  19
217 ##   ORE         22 114
218 ##
219 ##           Accuracy : 0.8541
220 ##           95% CI : (0.8073, 0.8932)
221 ##   No Information Rate : 0.5267
222 ##   P-Value [Acc > NIR] : <2e-16
223 ##
224 ##           Kappa : 0.7077
225 ##
226 ##   Mcnemar's Test P-Value : 0.7548
227 ##
228 ##           Sensitivity : 0.8571
229 ##           Specificity : 0.8514
230 ##           Pos Pred Value : 0.8382
231 ##           Neg Pred Value : 0.8690
232 ##           Prevalence : 0.4733
233 ##           Detection Rate : 0.4057
234 ##   Detection Prevalence : 0.4840
235 ##   Balanced Accuracy : 0.8542
236 ##
237 ##   'Positive' Class : ORE
238 ##
239 varImp.df <- as_tibble(varImp(minModel,
240                           sort = TRUE,
241                           scale = FALSE),
242                       rownames = 'Variables')
243
244 varImp.df %>%
245   ggplot(aes(y = reorder(Variables, ORE),
246             x = round((100*ORE/sum(varImp.df$ORE)),2),
247             fill = round((100*ORE/sum(varImp.df$ORE)),2))) +
248   geom_col(col = 'black') +
249   scale_fill_viridis_c() +
250   geom_text(label = round((100*varImp.df$ORE/sum(varImp.df$ORE)),2),
251            nudge_x = -1, col = 'white',aes(fontface = c('bold')))
252             ) +
253   theme(legend.position = c(.92,.3)) +
254   labs(y = "Variables", x = "Mean Decrease Accuracy (%)", fill = 'Accuracy')
```



255

256

257

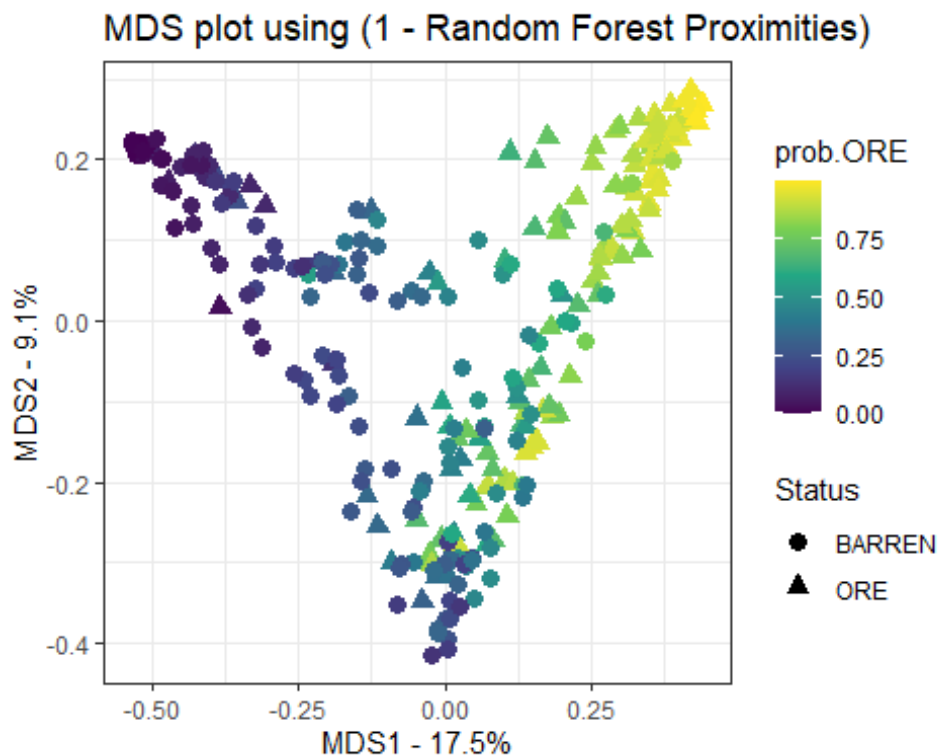
258

259

Fig. C. 7: Variable importance rank based on the Mean Decrease Accuracy parameter normalized to a percentage distribution. The bars are filled with the coded colors according to their respective importance. Abbreviations: Mag.Susceptibility – Magnetic Susceptibility; Electri. Conductivity – Electric Conductivity.



```
260 x <- df %>%
261   filter(!is.na(Cr))
262
263 y <- df %>%
264   filter(!is.na(Cr)) %>%
265   select(min)
266
267 labs <- df %>%
268   select(ID:LITHO)
269
270 pred1 <- predict(minModel,newdata = x)
271
272 pred.prob <- predict(minModel,newdata = x,type = 'prob',norm.votes = TRUE,predict.all = TR
273 UE)
274
275 prob.ORE <- pred.prob$aggregate[,2]
276
277 dc_samples <- x %>%
278   mutate(prob.ORE = all_of(prob.ORE),
279          Prediction = all_of(pred1))
280
281 ## Start by converting the proximity matrix into a distance matrix.
282 distance.matrix <- as.dist(1-minModel$proximity)
283
284 mds.stuff <- cmdscale(distance.matrix,
285                      eig=TRUE,
286                      x.ret=TRUE)
287
288 ## calculate the percentage of variation that each MDS axis accounts for...
289 mds.var.per <- round(mds.stuff$eig/sum(mds.stuff$eig)*100, 1)
290
291 ## now make a fancy looking plot that shows the MDS axes and the variation:
292 mds.values <- mds.stuff$points
293 mds.data <- data.frame(Sample=rownames(mds.values),
294                      X=mds.values[,1],
295                      Y=mds.values[,2],
296                      Status = trainSplit$min,
297                      prob = minModel$votes)
298
299 ggplot(data=mds.data, aes(x=X, y=Y, label=Sample, shape = Status)) +
300   geom_point(aes(col = prob.ORE), size = 3) +
301   theme_bw() +
302   xlab(paste("MDS1 - ", mds.var.per[1], "%", sep="")) +
303   ylab(paste("MDS2 - ", mds.var.per[2], "%", sep="")) +
304   ggtitle("MDS plot using (1 - Random Forest Proximities)") +
305   scale_color_viridis_c()
```



305

306 Fig. C. 8: Multidimensional Scaling graph based on all the Trees results. The points are classified according to
307 the Mineralization Status (i.e., Ore or Barren) and color coded by the probability of been classified as ORE.

308 *# Test prob*

309

```
310 predTest <- predict(minModel,newdata = testSplit)
```

311

```
312 predTest.prob <- predict(minModel,newdata = testSplit,type = 'prob',norm.votes =  
313 TRUE,predict.all = TRUE)
```

314

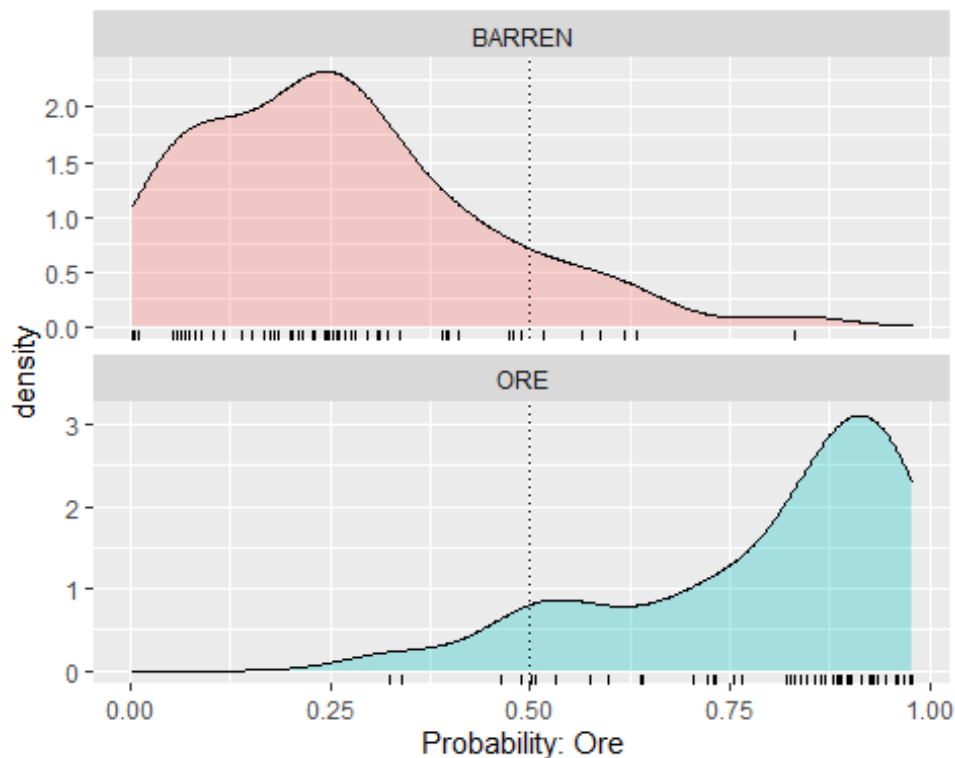
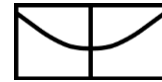
```
315 probTest.ORE <- predTest.prob$aggregate[,2]
```

316

```
317 test.prob <- testSplit %>%  
318   mutate(prob.ORE = all_of(probTest.ORE),  
319          Prediction = all_of(predTest))
```

320

```
321 test.prob %>%  
322   filter(!is.na(Fe_Ti)) %>%  
323   ggplot(aes(x = prob.ORE)) +  
324   geom_density(aes(fill = min), alpha = .3) +  
325   geom_rug() +  
326   geom_vline(xintercept = 0.5, lty = 3) +  
327   facet_wrap(min ~ .,ncol = 1,scales = 'free_y') +  
328   theme(legend.position = 'none') +  
329   labs(x = 'Probability: Ore')
```



330

331 Fig. C. 9: Ore probability analysis for all samples based on the mineralization status of the test dataset (barren
332 samples are represented in the red curve and ore samples in the blue curve). Most barren samples took a low Ore
333 probability, and the mineralized samples got the highest probabilities. The fields of True Negative (TN, i.e.,
334 barren samples predicted as non-mineralized), False Positive (FP, barren samples predicted as mineralized),
335 False Negative (FN, ore samples predicted as non-mineralized), and True Positive (TP, ore samples predicted as
336 mineralized) are indicated in the plot. The ticks at the bottom of each plot indicate the calculated probability for
337 each test dataset sample.

```
338 dc_samples <- dc_samples %>%  
339   filter(!is.na(Fe_Ti)) %>%  
340   mutate(number = 1) %>%  
341   group_by(HOLE) %>%  
342   mutate(fid = cumsum(number)) %>%  
343   ungroup()  
344  
345 h1 <-  
346   dc_samples %>%  
347   mutate(ROCK = as.factor(ROCK)) %>%  
348   filter(HOLE == 'CAN120') %>%  
349   arrange(HOLE, FROM) %>%  
350   ggplot(aes(y = fid, fill = ROCK, col = ROCK)) +  
351   geom_bar() +  
352   coord_cartesian(ylim = c(170,0)) +  
353   scale_y_reverse() +  
354   facet_wrap(. ~ HOLE, ncol = 4) +  
355   labs(y = 'Position') +  
356   theme(  
357     legend.position = 'none',  
358     axis.text.x = element_blank(),  
359     axis.title.x = element_blank(),  
360     axis.ticks.x = element_blank(),  
361     panel.grid.major.x = element_blank(),  
362     panel.grid.minor.x = element_blank())  
363  
364  
365 h2 <-  
366   dc_samples %>%
```



```
367 mutate(ROCK = as.factor(ROCK)) %>%
368 filter(HOLE == 'CAN144') %>%
369 arrange(HOLE, FROM) %>%
370 ggplot(aes(y = fid, fill = ROCK, col = ROCK)) +
371 geom_bar() +
372 coord_cartesian(ylim = c(170,0)) +
373 scale_y_reverse() +
374 facet_wrap(. ~ HOLE, ncol = 4) +
375 labs(y = 'Position') +
376 theme(
377   legend.position = 'none',
378   axis.text.x = element_blank(),
379   axis.title.x = element_blank(),
380   axis.ticks.x = element_blank(),
381   axis.title.y = element_blank(),
382   panel.grid.major.x = element_blank(),
383   panel.grid.minor.x = element_blank())
384
385
386 h3 <-
387   dc_samples %>%
388     mutate(ROCK = as.factor(ROCK)) %>%
389     filter(HOLE == 'CANIF27') %>%
390     arrange(HOLE, FROM) %>%
391     ggplot(aes(y = fid, fill = ROCK, col = ROCK)) +
392     geom_bar() +
393     coord_cartesian(ylim = c(170,0)) +
394     scale_y_reverse() +
395     facet_wrap(. ~ HOLE, ncol = 4) +
396     labs(y = 'Position') +
397     theme(
398       legend.position = 'none',
399       axis.text.x = element_blank(),
400       axis.title.x = element_blank(),
401       axis.ticks.x = element_blank(),
402       axis.title.y = element_blank(),
403       panel.grid.major.x = element_blank(),
404       panel.grid.minor.x = element_blank())
405
406
407 h4 <-
408   dc_samples %>%
409     mutate(ROCK = as.factor(ROCK)) %>%
410     filter(HOLE == 'JBA722') %>%
411     arrange(HOLE, FROM) %>%
412     ggplot(aes(y = fid, fill = ROCK, col = ROCK)) +
413     geom_bar() +
414     coord_cartesian(ylim = c(170,0)) +
415     scale_y_reverse() +
416     facet_wrap(. ~ HOLE, ncol = 4) +
417     labs(y = 'Position') +
418     theme(
419       legend.position = 'none',
420       axis.text.x = element_blank(),
421       axis.title.x = element_blank(),
422       axis.ticks.x = element_blank(),
423       axis.title.y = element_blank(),
424       panel.grid.major.x = element_blank(),
425       panel.grid.minor.x = element_blank())
426
427
428
429 prof1 <-
430   dc_samples %>%
431     filter(HOLE == 'CAN120') %>%
```



```
432 filter(!is.na(Fe_Ti)) %>%
433 mutate(Type = case_when(Prediction == min ~ 'True_PN',
434                          Prediction != min ~ 'False_PN',
435                          TRUE ~ 'ERROR!')) %>%
436 ggplot(aes(y = fid)) +
437   geom_vline(xintercept = 0.5, lty = 2, col = 'red') +
438   geom_bar(aes(fill = Type, col = NULL),
439            width = 1, alpha = .3,
440            position = 'identity') +
441   geom_segment(aes(y=fid, yend=fid, x=0, xend=prob.ORE)) +
442   geom_point(aes(x = prob.ORE, col = min), cex = 2.5) +
443   coord_cartesian(ylim = c(170,0), expand = TRUE) +
444   scale_y_reverse() +
445   facet_wrap(. ~ HOLE, ncol = 4) +
446   theme(legend.position = 'none',
447         axis.text.x = element_text(angle = -90,
448                                     hjust = .5,
449                                     vjust = .5),
450         axis.title.x = element_blank(),
451         axis.text.y = element_blank(),
452         axis.title.y = element_blank(),
453         plot.margin = unit(c(.03,.03,.03,.03), "lines")
454         ) +
455   annotate(label = 'Threshold',
456           x = 0.58, y = 160,
457           fontface = 'italic',
458           geom = 'text',
459           angle = -90,
460           colour = 'red',
461           size = 2.5) +
462   labs(y = '')
463
464 prof2 <-
465 dc_samples %>%
466 filter(HOLE == 'CAN144') %>%
467 filter(!is.na(Fe_Ti)) %>%
468 mutate(Type = case_when(Prediction == min ~ 'True_PN',
469                          Prediction != min ~ 'False_PN',
470                          TRUE ~ 'ERROR!')) %>%
471 ggplot(aes(y = fid)) +
472   geom_vline(xintercept = 0.5, lty = 2, col = 'red') +
473   geom_bar(aes(fill = Type, col = NULL),
474            width = 1, alpha = .3,
475            position = 'identity') +
476   geom_segment(aes(y=fid, yend=fid, x=0, xend=prob.ORE)) +
477   geom_point(aes(x = prob.ORE, col = min), cex = 2.5) +
478   coord_cartesian(ylim = c(170,0), expand = TRUE) +
479   scale_y_reverse() +
480   facet_wrap(. ~ HOLE, ncol = 4) +
481   theme(legend.position = 'none',
482         axis.text.x = element_text(angle = -90,
483                                     hjust = .5,
484                                     vjust = .5),
485         axis.title.x = element_blank(),
486         axis.text.y = element_blank(),
487         axis.title.y = element_blank(),
488         plot.margin = unit(c(.03,.03,.03,.03), "lines")
489         ) +
490   annotate(label = 'Threshold',
491           x = 0.58, y = 160,
492           fontface = 'italic',
493           geom = 'text',
494           angle = -90,
495           colour = 'red',
496           size = 2.5) +
```

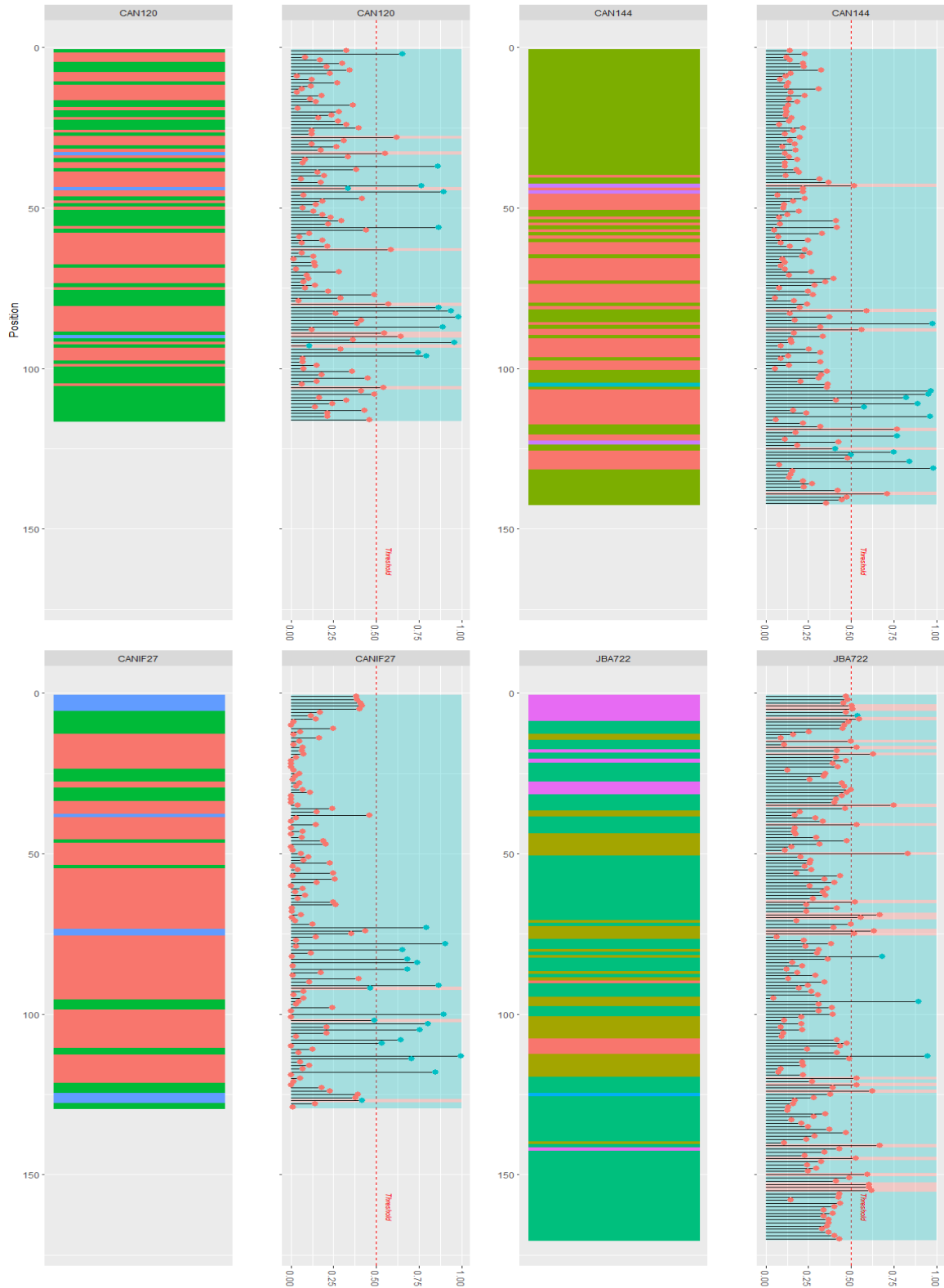
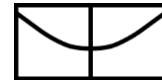


```
497   labs(y = '')
498
499 prof3 <-
500   dc_samples %>%
501   filter(HOLE == 'CANIF27') %>%
502   filter(!is.na(Fe_Ti)) %>%
503   mutate(Type = case_when(Prediction == min ~ 'True_PN',
504                           Prediction != min ~ 'False_PN',
505                           TRUE ~ 'ERROR!')) %>%
506   ggplot(aes(y = fid)) +
507   geom_vline(xintercept = 0.5, lty = 2, col = 'red') +
508   geom_bar(aes(fill = Type, col = NULL),
509           width = 1, alpha = .3,
510           position = 'identity') +
511   geom_segment(aes(y=fid, yend=fid, x=0, xend=prob.ORE)) +
512   geom_point(aes(x = prob.ORE, col = min), cex = 2.5) +
513   coord_cartesian(ylim = c(170,0), expand = TRUE) +
514   scale_y_reverse() +
515   facet_wrap(. ~ HOLE, ncol = 4) +
516   theme(legend.position = 'none',
517         axis.text.x = element_text(angle = -90,
518                                     hjust = .5,
519                                     vjust = .5),
520         axis.title.x = element_blank(),
521         axis.text.y = element_blank(),
522         axis.title.y = element_blank(),
523         plot.margin = unit(c(.03,.03,.03,.03), "lines")
524   ) +
525   annotate(label = 'Threshold',
526           x = 0.58, y = 160,
527           fontface = 'italic',
528           geom = 'text',
529           angle = -90,
530           colour = 'red',
531           size = 2.5) +
532   labs(y = '')
533
534 prof4 <-
535   dc_samples %>%
536   filter(HOLE == 'JBA722') %>%
537   filter(!is.na(Fe_Ti)) %>%
538   mutate(Type = case_when(Prediction == min ~ 'True_PN',
539                           Prediction != min ~ 'False_PN',
540                           TRUE ~ 'ERROR!')) %>%
541   ggplot(aes(y = fid)) +
542   geom_vline(xintercept = 0.5, lty = 2, col = 'red') +
543   geom_bar(aes(fill = Type, col = NULL),
544           width = 1, alpha = .3,
545           position = 'identity') +
546   geom_segment(aes(y=fid, yend=fid, x=0, xend=prob.ORE)) +
547   geom_point(aes(x = prob.ORE, col = min), cex = 2.5) +
548   coord_cartesian(ylim = c(170,0), expand = TRUE) +
549   scale_y_reverse() +
550   facet_wrap(. ~ HOLE, ncol = 4) +
551   theme(legend.position = 'none',
552         axis.text.x = element_text(angle = -90,
553                                     hjust = .5,
554                                     vjust = .5),
555         axis.title.x = element_blank(),
556         axis.text.y = element_blank(),
557         axis.title.y = element_blank(),
558         plot.margin = unit(c(.03,.03,.03,.03), "lines")
559   ) +
560   annotate(label = 'Threshold',
561           x = 0.58, y = 160,
```



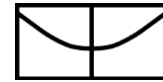
```
562     fontface = 'italic',
563     geom = 'text',
564     angle = -90,
565     colour = 'red',
566     size = 2.5) +
567   labs(y = '')

568 ggpubr::ggarrange(h1, prof1, h2, prof2, h3, prof3, h4, prof4, nrow = 2,
569                   ncol = 4,
570                   align = 'hv', widths = c(rep(c(4,4),4)),
571                   common.legend = FALSE,
572                   font.label = list(size = 16, face = 'bold'))
```



573

574 Fig. C. 10: Color-coded strip log according to the lithologies for drill cores studied in this work and respective
575 validation column. The calculated Ore probability is indicated as a bar beside the sample position for each
576 sample and drill core, and the threshold of probability is indicated as the red dashed line. The circle color shows
577 the reference values of the mineralization status at the end of the probability bar. The validation column is color-
578 coded according to the verification of the predicted and reference mineralization status.



1 APÊNDICE D – Código utilizado no Artigo 02

2 NOTEBOOK

3 This is the source code used in the manuscript:

4 ‘Machine learning analysis of mineral chemistry in pyrite grains from the Jacobina gold
5 deposits, São Francisco Craton, Brazil: geochemical patterns and implications to mineral
6 exploration’

7 authored by: “Guilherme Ferreira (guilherme.ferreira@cprm.gov.br)” date: “02/03/2022”

8 The following code was written in R (4.1.2).

9 The input data and complementary information can be found at:

10 <https://github.com/gferrsilva/icpms-jacobina>

11 ABSTRACT

12 We applied machine learning (ML) to process LA-ICP-MS data (45 elements) with 441 samples
13 of pyrite from gold-bearing quartz-pebble-metaconglomerate from the Serra de Jacobina
14 deposits in the São Francisco Craton, Brazil. First, the pyrite samples were described by optical
15 and scanning electron microscopy to gather information about the texture differences. Then,
16 the pyrite grains were classified according to their source and stratigraphical level: detrital and
17 epigenetic pyrite from the mineralized Jacobina Group and pyrite from the basement or
18 intrusive rocks. We used Agglomerative Clustering methods to evaluate the trace elements
19 patterns according to pyrite group, mineral source, and stratigraphic levels. Then, we
20 implemented the Uniform Manifold Approximation and Projection technique (UMAP) to
21 reduce the dimensionality of data into a two-dimensional projection to inspect the inner
22 structure of the data. This result was confirmed by the analysis of the dendrograms, which
23 show different associations of elements among detrital and epigenetic pyrites. Elements such
24 as Cu, Zn, Ag, Sb, Te, Au, Pb, and Bi are mobilized during mineral alteration and was crystallized
25 in newly formed minerals, such as chalcopyrite, pyrrhotite, and sphalerite, which are spatially
26 associated with epigenetic pyrite and free gold. These findings could explain the differences
27 in the mineral assemblage in portions of the deposits that prevail sedimentary minerals or the
28 others that were strongly modified by later alterations. In conclusion, ML is recommended in
29 the processing of mineral chemistry data because it helps to process data without discarding
30 significant variables, and the method allows to evidence the multivariate structure of data.



31 DATA WRANGLING

32 Dependencies

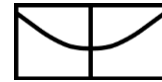
```
33 library(tidyverse) # ggplot2, tidyr, dplyr
34 library(readxl) # open XLSX data
35 library(geoquimica) # Data wrangling
36 library(umap) # Dimensionality reduction
37 library(pheatmap) # Distance matrices
38 library(dendextend) # Dendrograms
39 library(ggpubr) # Plot adjusts
```

40 Data preparation

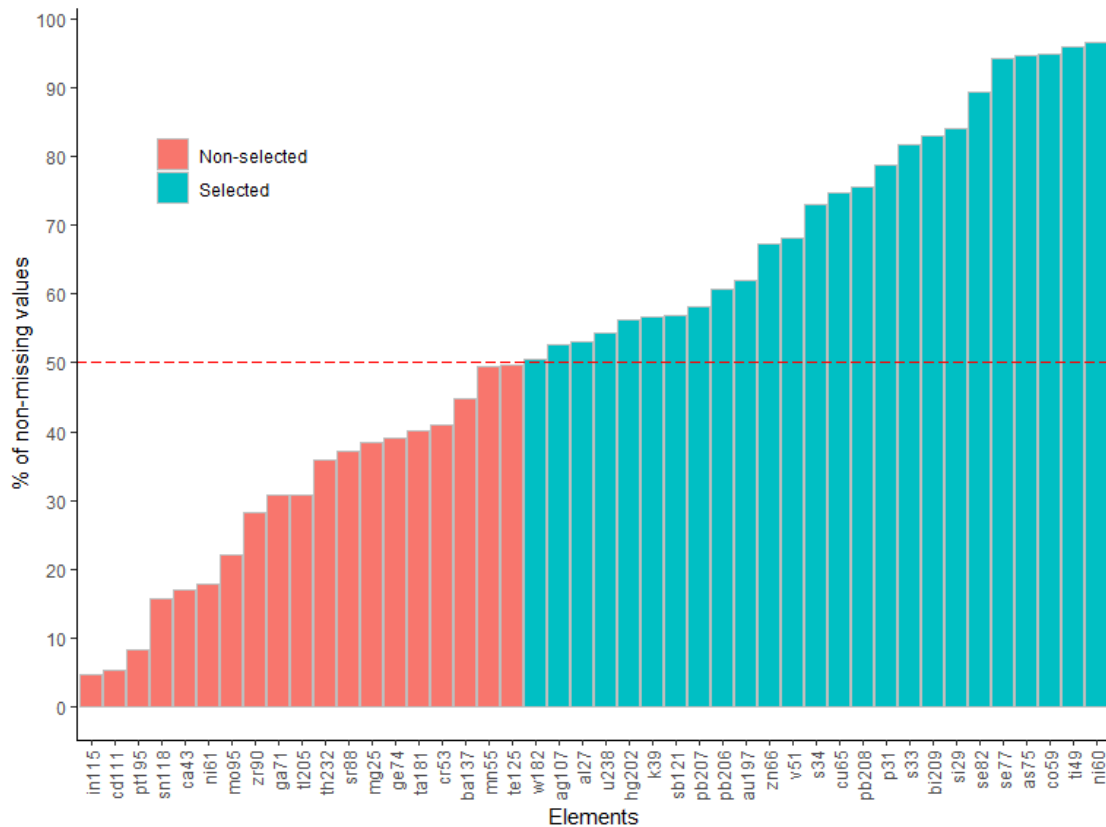
```
41 set.seed(0)
42
43 df <- data.table::fread("~/piritas_jacobina_editada_v2.csv")
44
45 # Defining variable class
46 df[,12:142] <- lapply(X = df[,12:142],FUN = as.double)
47
48 # Creating variable of imputation control
49 df1 <- df %>%
50   drop_na(`Pyrite Type`) %>%
51   mutate(impute_ni60 = ifelse(test = is.na(Ni60),yes = 'True',no = 'False'),
52          impute_co = ifelse(test = is.na(Co59),yes = 'True',no = 'False'),
53          impute_ti = ifelse(test = is.na(Ti49),yes = 'True',no = 'False'),
54          impute_v = ifelse(test = is.na(V51),yes = 'True',no = 'False'))
55
56
57 index <- df1 %>%
58   select(`Source file`:`Reef`, impute_ni60:impute_v)
59
60 statistics <- df1 %>%
61   select(-names(index)) %>%
62   select(matches('LOD$|2SE$'))
63
64 lod <- df1 %>%
65   dplyr::select(matches('LOD$'))
66
67 elems <- df1 %>%
68   dplyr::select(-names(index),-V1)
```

69 Data selection

```
70 geoquimica::elem_fillrate( data.table::fread(
71   "~/GitHub/jacobina/data/minchem/piritas_jacobina_editada_v2.csv",
72   verbose = FALSE) %>%
73   mutate_at(.vars = 12:143,.funs = as.double) %>%
74   drop_na(`Pyrite Type`) %>% # Drop wrong analysis
75   select(-(V1:`Source file`),-(DateTime:Comments)) %>%
76   janitor::clean_names()
77 ) %>%
78   filter(!str_detect(string = Column.Name, pattern = 'lod$|2se$'),
```



```
79   !Column.Name %in% c('datetime', 'generation', 'pyrite_type', 'texture', 'reef
80 ') %>%
81   arrange(Fill.Rate) %>%
82   mutate(Column.Name = fct_inorder(Column.Name)) %>%
83   # mutate()
84   ggplot(aes(x = Column.Name, y = Fill.Rate)) +
85   geom_col(aes(fill = ifelse(test = Fill.Rate < 50, 'Non-selected', 'Selected')), col
86 l = 'gray') +
87   geom_hline(yintercept = 50, lty = 5, col = 'red', size = .7) +
88   scale_y_continuous(breaks = seq(0,100,10)) +
89   theme_classic() +
90   theme(axis.text.x = element_text(angle = 90, vjust = 0.5, hjust = 1),
91         legend.position = c(.15, .8)
92   ) +
93   labs(x = 'Elements',
94        y = '% of non-missing values',
95        fill = '')
```



96

97 Fig. D. 1: list of elements ordered by the percentage of non-missing data. The 50% threshold (horizontal dashed
98 line) was used to determine if a variable could be selected for multivariate analysis. The elements Al, P, and Si
99 were not selected based on the small variability in the dataset.

100 IMPUTATION

101 LDL and Missing Value Imputation

```
102 # Imputation of LDL elements
103 half_ldl <- elems %>%
104   mutate(A127 = ifelse(test = is.na(A127), yes = `A127 LOD`/sqrt(2), no = A127),
```



```
105     Si29 = ifelse(test = is.na(Si29),yes = `Si29 LOD`/sqrt(2),no = Si29),
106     P31 = ifelse(test = is.na(P31),yes = `P31 LOD`/sqrt(2),no = P31),
107     S33 = ifelse(test = is.na(S33),yes = `S33 LOD`/sqrt(2),no = S33),
108     S34 = ifelse(test = is.na(S34),yes = `S34 LOD`/sqrt(2),no = S34),
109     K39 = ifelse(test = is.na(K39),yes = `K39 LOD`/sqrt(2),no = K39),
110     Ti49 = ifelse(test = is.na(Ti49),yes = `Ti49 LOD`/sqrt(2),no = Ti49),
111     V51 = ifelse(test = is.na(V51),yes = `V51 LOD`/sqrt(2),no = V51),
112     Co59 = ifelse(test = is.na(Co59),yes = `Co59 LOD`/sqrt(2),no = Co59),
113     Ni60 = ifelse(test = is.na(Ni60),yes = `Ni60 LOD`/sqrt(2),no = Ni60),
114     Cu65 = ifelse(test = is.na(Cu65),yes = `Cu65 LOD`/sqrt(2),no = Cu65),
115     Zn66 = ifelse(test = is.na(Zn66),yes = `Zn66 LOD`/sqrt(2),no = Zn66),
116     As75 = ifelse(test = is.na(As75),yes = `As75 LOD`/sqrt(2),no = As75),
117     Se77 = ifelse(test = is.na(Se77),yes = `Se77 LOD`/sqrt(2),no = Se77),
118     Se82 = ifelse(test = is.na(Se82),yes = `Se82 LOD`/sqrt(2),no = Se82),
119     Ag107 = ifelse(test = is.na(Ag107),yes = `Ag107 LOD`/sqrt(2),no = Ag107),
120     Sb121 = ifelse(test = is.na(Sb121),yes = `Sb121 LOD`/sqrt(2),no = Sb121),
121     W182 = ifelse(test = is.na(W182),yes = `W182 LOD`/sqrt(2),no = W182),
122     Au197 = ifelse(test = is.na(Au197),yes = `Au197 LOD`/sqrt(2),no = Au197),
123     Hg202 = ifelse(test = is.na(Hg202),yes = `Hg202 LOD`/sqrt(2),no = Hg202),
124     Pb206 = ifelse(test = is.na(Pb206),yes = `Pb206 LOD`/sqrt(2),no = Pb206),
125     Pb207 = ifelse(test = is.na(Pb207),yes = `Pb207 LOD`/sqrt(2),no = Pb207),
126     Pb208 = ifelse(test = is.na(Pb208),yes = `Pb208 LOD`/sqrt(2),no = Pb208),
127     Bi209 = ifelse(test = is.na(Bi209),yes = `Bi209 LOD`/sqrt(2),no = Bi209),
128     U238 = ifelse(test = is.na(U238),yes = `U238 LOD`/sqrt(2),no = U238)) %>%
129 mutate(impute_pb206 = ifelse(test = is.na(Pb206),yes = 'True',no = 'False'),
130        impute_pb207 = ifelse(test = is.na(Pb207),yes = 'True',no = 'False'),
131        impute_s33 = ifelse(test = is.na(S33),yes = 'True',no = 'False'),
132        impute_s34 = ifelse(test = is.na(S34),yes = 'True',no = 'False')) %>%
133 select(-matches('LOD$|2SE$')) %>%
134 geoquimica::elem_select(cut = .5)
135
136 # Imputation of missing values based on a multivariate non-parametric regression
137 imputed <- missRanger::missRanger(data = half_ldl,
138                                  pmm.k = 3,
139                                  maxiter = 10,
140                                  seed = 0,
141                                  verbose = 2)
142
143 ##
144 ## Missing value imputation by random forests
145 ##
146 ## Variables to impute:      P31, S33, S34, K39, Pb206, Pb207
147 ## Variables used to impute: Al27, Si29, P31, S33, S34, K39, Ti49, V51, Co59, N
148 ##                          i60, Cu65, Zn66, As75, Se77, Se82, Ag107, Sb121, W182, Au197, Hg202, Pb206, Pb207,
149 ##                          Pb208, Bi209, U238, impute_pb206, impute_pb207, impute_s33, impute_s34
150 ## P31 S33 K39 Pb206  Pb207  S34
151 ## iter 1:  0.9855  0.4544  0.5600  0.5095  0.2164  0.1844
152 ## iter 2:  0.5228  0.0639  0.2971  0.3116  0.1766  0.0708
153 ## iter 3:  0.5385  0.0658  0.3022  0.3387  0.1589  0.0599
```

153 Data Recode

```
154 # Recoding variables Generation, Reef, Reef_Label and Unit
155 df2 <- index %>%
156   bind_cols(imputed) %>%
157   mutate(Unit = case_when(Reef == 'Basal Reef' ~ 'Serra do Córrego',
158                          Reef == 'Main Reef' ~ 'Serra do Córrego',
159                          Reef == 'SPC' ~ 'Serra do Córrego',
```



```
160     Reef == 'LU' ~ 'Serra do Córrego',
161     Reef == 'LVLPC' ~ 'Serra do Córrego',
162     Reef == 'MSPC' ~ 'Serra do Córrego',
163     Reef == 'MPC' ~ 'Serra do Córrego',
164     Reef == 'SPC' ~ 'Serra do Córrego',
165     Reef == 'MU' ~ 'Serra do Córrego',
166     Reef == 'Holandez' ~ 'Serra do Córrego',
167     Reef == 'Maneira' ~ 'Serra do Córrego',
168     Reef == 'ITV' ~ 'Intrusive',
169     Reef == 'UMF' ~ 'Intrusive',
170     Reef == 'IQL' ~ 'Serra do Córrego',
171     Reef == 'Basement' ~ 'Basement',
172     Reef == 'CAF' ~ 'Cruz das Almas'),
173   Reef_label = case_when(Reef == 'Basal Reef' ~ 'LC',
174     Reef == 'Main Reef' ~ 'LC',
175     Reef == 'SPC' ~ 'UC1',
176     Reef == 'LU' ~ 'UC1',
177     Reef == 'LVLPC' ~ 'UC1',
178     Reef == 'MSPC' ~ 'UC1',
179     Reef == 'MPC' ~ 'UC1',
180     Reef == 'SPC' ~ 'UC1',
181     Reef == 'MU' ~ 'UC1',
182     Reef == 'Holandez' ~ 'UC2',
183     Reef == 'Maneira' ~ 'UC2',
184     Reef == 'ITV' ~ 'Intr.',
185     Reef == 'UMF' ~ 'Intr.',
186     Reef == 'IQL' ~ 'IQL',
187     Reef == 'Basement' ~ 'Base.',
188     Reef == 'CAF' ~ 'CdA'),
189   Reef = case_when(Reef == 'Basal Reef' ~ 'Lower Conglomerate',
190     Reef == 'Main Reef' ~ 'Lower Conglomerate',
191     Reef == 'SPC' ~ 'Upper Conglomerate 1',
192     Reef == 'LU' ~ 'Upper Conglomerate 1',
193     Reef == 'LVLPC' ~ 'Upper Conglomerate 1',
194     Reef == 'MSPC' ~ 'Upper Conglomerate 1',
195     Reef == 'MPC' ~ 'Upper Conglomerate 1',
196     Reef == 'SPC' ~ 'Upper Conglomerate 1',
197     Reef == 'MU' ~ 'Upper Conglomerate 1',
198     Reef == 'Holandez' ~ 'Upper Conglomerate 2',
199     Reef == 'Maneira' ~ 'Upper Conglomerate 2',
200     Reef == 'ITV' ~ 'Intrusive',
201     Reef == 'UMF' ~ 'Intrusive',
202     Reef == 'IQL' ~ 'Intermediate Quartzite',
203     Reef == 'Basement' ~ 'Basement',
204     Reef == 'CAF' ~ 'Cruz das Almas'),
205   ) %>%
206   dplyr::select(-c(`Source file`:Comments, P31,Si29,A127)) %>%
207   mutate(Generation = case_when(Reef == 'Intrusive' ~ 'Intrusive',
208     Reef == 'Basement' ~ 'Basement',
209     TRUE ~ as.character(Generation)),
210     Texture = case_when(Reef == 'Intrusive' | Reef == 'Basement' ~ 'Subhedral
211 ',
212     TRUE ~ as.character(Texture)))
```



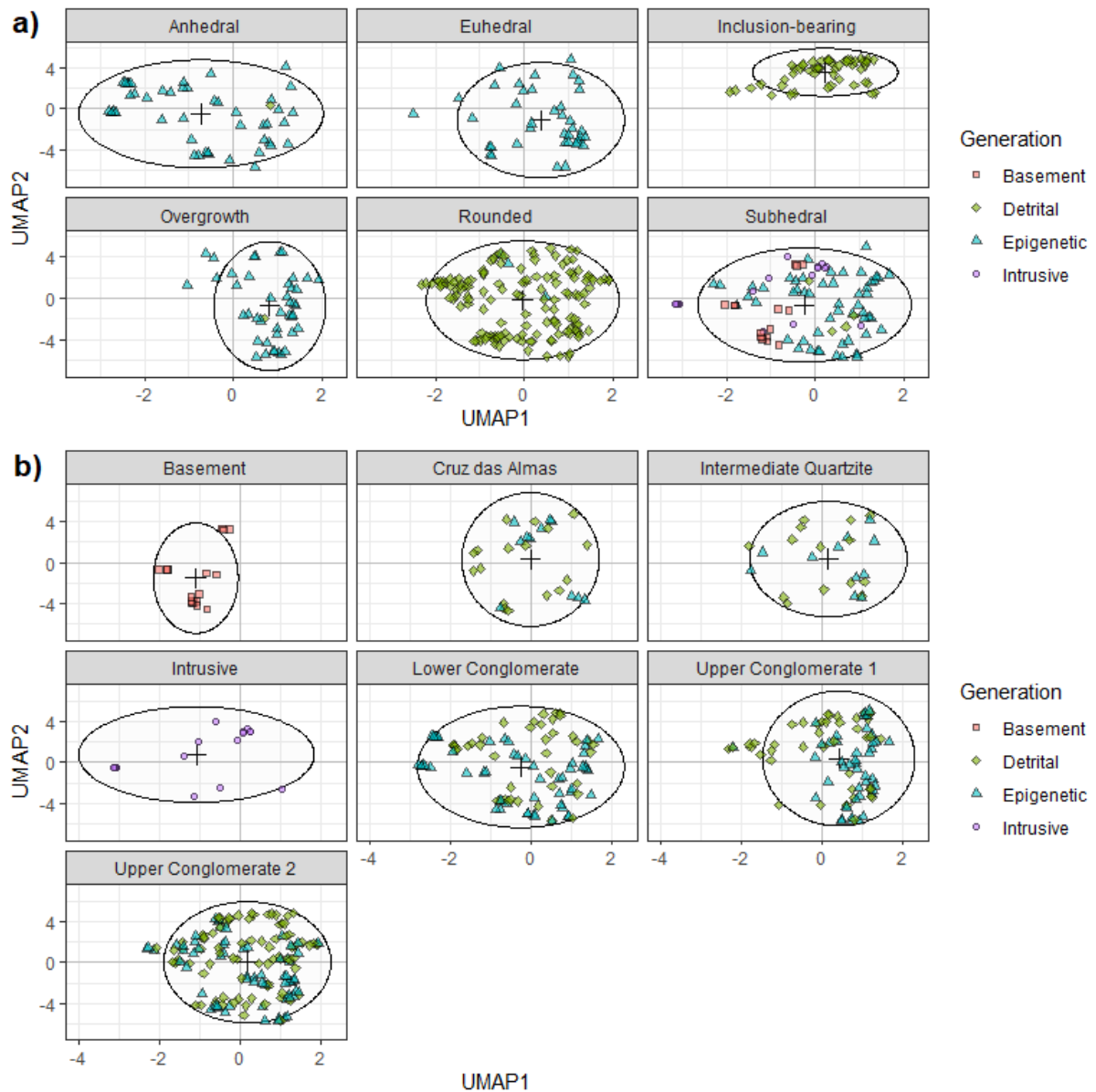
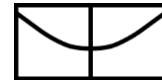
213 DIMENSIONALITY REDUCTION

214 Uniform Manifold Approximation and Projection - UMAP

```
215 # UMAP processing
216
217 dfumap <-
218   df2 %>%
219   select_if(is.numeric) %>%
220   geoquimica::elem_norm(method = 'clr') %>%
221   umap::umap()
222
223 # Data merging
224 df3 <-
225   df2 %>%
226   bind_cols(as_tibble(dfumap$layout))
227
228 # Making Plot a)
229 umap1 <-
230   df3 %>%
231   ggplot(aes(x = V1, y = V2,
232             fill = Generation,
233             shape = Generation,
234             group = Generation)) +
235   geom_vline(xintercept = 0, col = 'grey', size = .7) +
236   geom_hline(yintercept = 0, col = 'grey', size = .7) +
237   geom_point(inherit.aes = FALSE,
238             data = df3 %>%
239               group_by(Texture) %>%
240               summarize(U1mean = mean(V1),
241                         U2mean = mean(V2)),
242             mapping = aes(x = U1mean, y = U2mean),
243             cex = 3, shape = 3, col = 'black') +
244   ggforce::geom_ellipse(inherit.aes = FALSE,
245                         data = df3 %>%
246                           group_by(Texture) %>%
247                           summarize(U1mean = mean(V1),
248                                     U1sd = sd(V1),
249                                     U2mean = mean(V2),
250                                     U2sd = sd(V2)),
251                         mapping = aes(x0 = U1mean, y0 = U2mean, a = 2*U1sd, b = 2*
252 U2sd, angle = 0),
253                         fill = 'grey', alpha = .05) +
254   geom_point(alpha = .6, col = 'black') +
255   scale_shape_manual(values = c(22,23,24, 21)) +
256   facet_wrap(. ~ Texture, nrow = 2) +
257   labs(x = 'UMAP1', y = 'UMAP2') +
258   theme_bw()
259
260 # Making plot b)
261 umap2 <-
262   df3 %>%
263   ggplot(aes(x = V1, y = V2,
264             fill = Generation,
265             shape = Generation,
266             group = Generation)) +
267   geom_vline(xintercept = 0, col = 'grey', size = .7) +
```



```
268 geom_hline(yintercept = 0, col = 'grey', size = .7) +
269 geom_point(inherit.aes = FALSE,
270             data = df3 %>%
271               group_by(Reef) %>%
272               summarize(U1mean = mean(V1),
273                         U2mean = mean(V2)),
274             mapping = aes(x = U1mean, y = U2mean),
275             cex = 3, shape = 3, col = 'black') +
276 ggforce::geom_ellipse(inherit.aes = FALSE,
277                       data = df3 %>%
278                         group_by(Reef) %>%
279                         summarize(U1mean = mean(V1),
280                                   U1sd = sd(V1),
281                                   U2mean = mean(V2),
282                                   U2sd = sd(V2)),
283                       mapping = aes(x0 = U1mean, y0 = U2mean, a = 2*U1sd, b = 2*
284 U2sd, angle = 0),
285                       fill = 'grey', alpha = .05) +
286 geom_point(alpha = .6, col = 'black') +
287 scale_shape_manual(values = c(22,23,24, 21,
288                               22,23,24, 21)) +
289 facet_wrap(. ~ Reef, ncol = 3) +
290 labs(x = 'UMAP1', y = 'UMAP2') +
291 theme_bw() #+
292
293 # Arrange plots
294 ggarrange(umap1, umap2, nrow = 2,
295           labels = c('a', 'b')),
296           align = 'hv',
297           heights = c(2,3), font.label = list(size = 16, face = 'bold'))
```



298

299 Fig. D. 2: UMAP configuration for data with samples classified according to a) pyrite texture and b)
 300 stratigraphic level. The ellipses drawn in each box are centered based on the mean of the samples, and its axes
 301 are calculated according to two times the standard deviation for each UMAP coordinate.

302 DISTANCE MATRICES

```

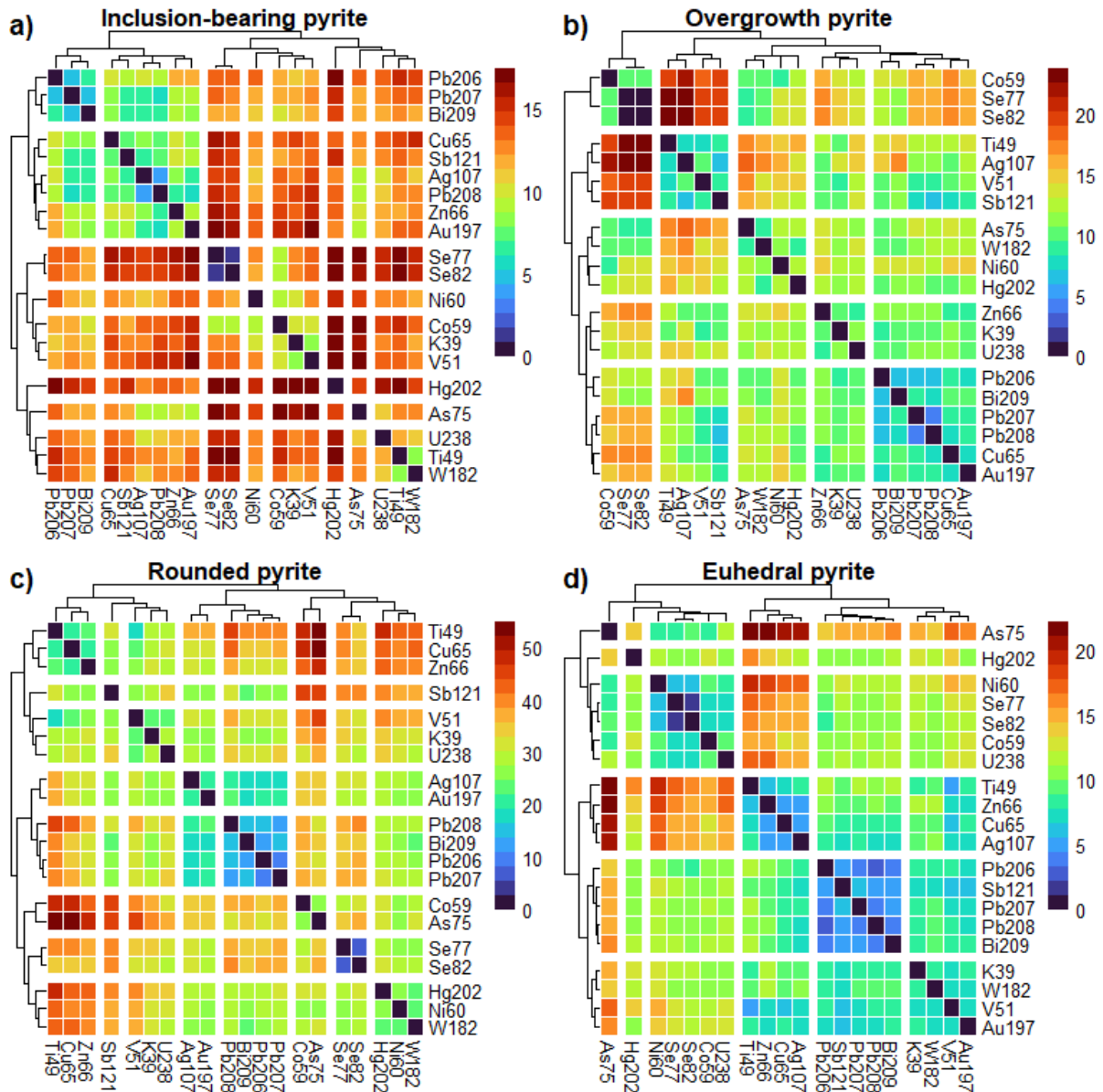
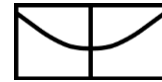
303 d <-
304   df2 %>%
305   filter(Texture == 'Inclusion-bearing') %>%
306   geoquimica::elem_norm(method = 'clr') %>%
307   geoquimica::elem_norm() %>%
308   select(K39:U238) %>%
309   t()
310
311 # Making Fig. A, Inclusion-bearing Pyrite
312 inbear <-
313   ggplotify::as.ggplot(heatmap(mat = df2 %>%
314     filter(Texture == 'Inclusion-bearing') %>%
  
```




```
315     geoquimica::elem_norm(method = 'clr') %>%
316     geoquimica::elem_norm() %>%
317     select(K39:U238) %>%
318     t() %>%
319     dist(method = 'manhattan'),
320     labels_row = rownames(d),
321     labels_col = rownames(d),
322     border_color = 'white',
323     clustering_distance_rows = 'manhattan',
324     clustering_distance_cols = 'manhattan',
325     color = pals::turbo(20),
326     cutree_rows = 8,
327     cutree_cols = 8,
328     main = "Inclusion-bearing pyrite",
329     treeheight_row = 15,
330     treeheight_col = 15,))
331 # Making Fig. B, Overgrowth Pyrite
332 overg <- ggplotify::as.ggplot(heatmap(mat = df2 %>%
333     filter(Texture == 'Overgrowth') %>%
334     geoquimica::elem_norm(method = 'clr') %>%
335     geoquimica::elem_norm() %>%
336     select(K39:U238) %>%
337     t() %>%
338     dist(method = 'manhattan'),
339     labels_row = rownames(d),
340     border_color = 'white',
341     labels_col = rownames(d),
342     color = pals::turbo(20),
343     clustering_distance_rows = 'manhattan',
344     clustering_distance_cols = 'manhattan',
345     cutree_rows = 5,
346     cutree_cols = 5,
347     main = "Overgrowth pyrite",
348     treeheight_row = 15,
349     treeheight_col = 15))
350 # Making Fig. C, Rounded Pyrite
351 rounded <- ggplotify::as.ggplot(heatmap(mat = df2 %>%
352     filter(Texture == 'Rounded') %>%
353     geoquimica::elem_norm(method = 'clr') %>%
354     geoquimica::elem_norm() %>%
355     select(K39:U238) %>%
356     t() %>%
357     dist(method = 'manhattan'),
358     labels_row = rownames(d),
359     labels_col = rownames(d),
360     border_color = 'white',
361     clustering_distance_rows = 'manhattan',
362     clustering_distance_cols = 'manhattan',
363     color = pals::turbo(20),
364     cutree_rows = 8,
365     cutree_cols = 8,
366     main = "Rounded pyrite",
367     treeheight_row = 15,
368     treeheight_col = 15))
369 # Making Fig. D, Euhedral Pyrite
370 euhedral <- ggplotify::as.ggplot(heatmap(mat = df2 %>%
```



```
371     filter(Texture == 'Euhedral') %>%
372     geoquimica::elem_norm(method = 'clr') %>%
373     geoquimica::elem_norm() %>%
374     select(K39:U238) %>%
375     t() %>%
376     dist(method = 'manhattan'),
377     labels_row = rownames(d),
378     labels_col = rownames(d),
379     border_color = 'white',
380     clustering_distance_rows = 'manhattan',
381     clustering_distance_cols = 'manhattan',
382     color = pals::turbo(20),
383     cutree_rows = 6,
384     cutree_cols = 6,
385     main = "Euhedral pyrite",
386     treeheight_row = 15,
387     treeheight_col = 15))
388 # Arrange plots
389 ggarrange(inbear, overg, rounded, euhedral,
390           ncol = 2, nrow = 2,
391           labels = c('a', 'b', 'c', 'd'), align = 'hv', hjust = 0)
```



392

393 Fig. D. 3: Dendrograms and distance matrices for pyrite grains according to grain texture: a) Inclusion-bearing,
 394 b) Overgrowth, c) Anhedral, and d) Euhedral. Rows and columns are ordered according to the agglomerative
 395 clustering, calculated by the Manhattan distance.

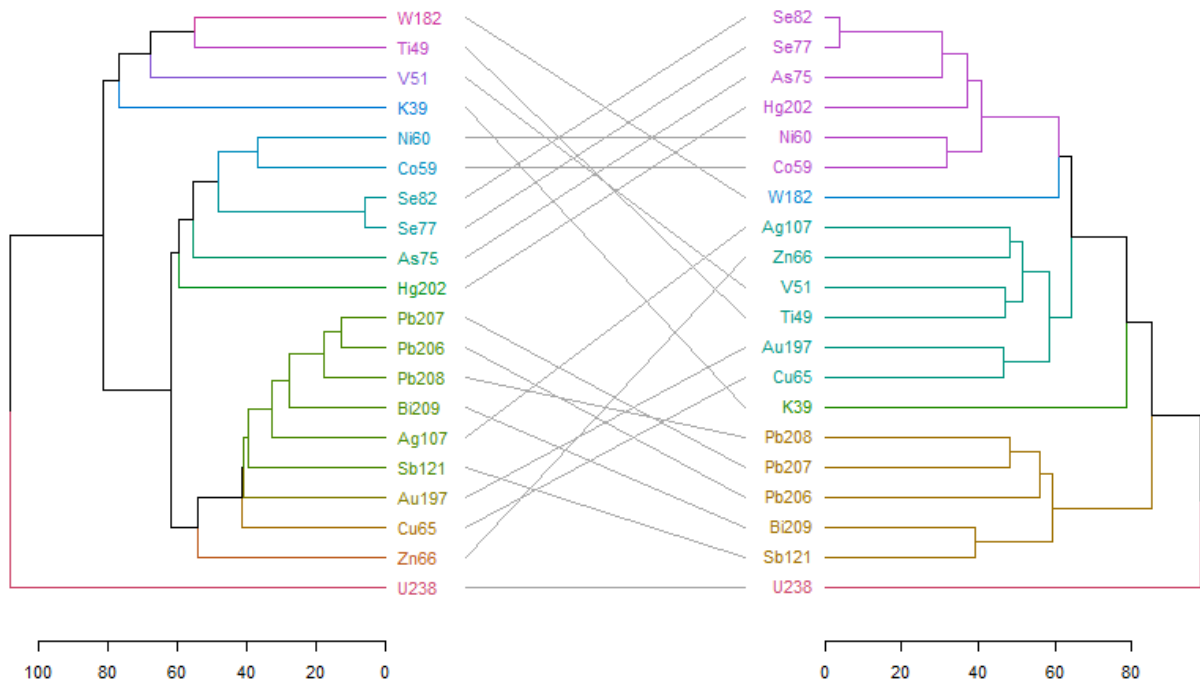
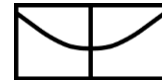
396 Compared dendrograms

```

397 d1 <-
398   df3 %>%
399   filter(Texture == 'Inclusion-bearing') %>%
400   select(K39:U238) %>%
401   geoquimica::elem_norm(method = 'clr') %>%
402   t() %>%
403   dist(method = 'manhattan') %>%
404   hclust(method = 'average') %>%
405   as.dendrogram()
406
407 d2 <-
408   df3 %>%
  
```



```
409 filter(Texture == 'Overgrowth') %>%
410 select(K39:U238) %>%
411 geoquimica::elem_norm(method = 'clr') %>%
412 t() %>%
413 dist(method = 'manhattan') %>%
414 hclust(method = 'average') %>%
415 as.dendrogram()
416
417
418 # Custom these kendo, and place them in a list
419 dl1 <-
420   dendextend::dendlist(
421     d1 %>%
422       set("branches_lty", 1) %>%
423       set("labels_col",
424         h = 40) %>%
425       set("branches_lty", 1) %>%
426       set("branches_k_color",
427         h = 40),
428     d2 %>%
429       set("branches_lty", 1) %>%
430       set("labels_col",
431         h = 60) %>%
432       set("branches_lty", 1) %>%
433       set("branches_k_color",
434         h = 60)
435   )
436
437 # Plot them together
438 dendextend::tanglegram(dl1,
439   common_subtrees_color_lines = FALSE,
440   highlight_distinct_edges = FALSE,
441   highlight_branches_lwd = FALSE,
442   margin_inner=4, intersecting = TRUE,
443   lwd=1, k_labels = NULL, k_branches = NULL)
```



444

445

Fig. D. 4: Linked dendrograms for Detrital (left) and Epigenetic pyrites (right) mapping the relation of elements.

446

IMPUTATION OF LDL AND MISSING VALUES

447

Replacement of LDL

448

```
# Function to make qq plots coded by categorical variables
```

449

```
make_qq <- function(dd, x) {
```

450

```
  dd<-dd[order(dd[[x]]), ]
```

451

```
  dd$qq <- qnorm(ppoints(nrow(dd)))
```

452

```
  dd
```

453

```
}
```

454

455

```
# Data without imputation
```

456

```
p_original <-
```

457

```
  df3 %>%
```

458

```
  filter(impute_v == 'False') %>%
```

459

```
  make_qq(dd = ., x = 'V51') %>%
```

460

```
  ggplot(aes(x = qq, y = log(V51))) +
```

461

```
  geom_qq_line(aes(sample = qq),
```

462

```
    lty = 1, col = 'gray', size = 1,
```

463

```
    alpha = .7) +
```

464

```
  geom_point(aes(col = impute_v,
```

465

```
    shape = impute_v),
```

466

```
    # col = 'white',
```

467

```
    alpha = .3) +
```

468

```
  labs(x = 'Theoretical distribution',
```

469

```
    y = 'log(V51)',
```

470

```
    col = 'Imputed?',
```

471

```
    shape = 'Imputed?') +
```

472

```
  scale_shape_manual(values = c(16,17)) +
```

473

```
  theme_classic() +
```



```
474 theme(legend.justification = 'center',
475         legend.background = element_rect(
476           fill = 'grey98'))
477
478 # Data with half of LDL imputed values
479 p_ldl <-
480   df3 %>%
481   make_qq(dd = ., x = 'V51') %>%
482   ggplot(aes(x = qq, y = log(V51))) +
483   geom_qq_line(aes(sample = qq),
484               lty = 1, col = 'gray', size = 1,
485               alpha = .7) +
486   geom_point(aes(col = impute_v,
487                 shape = impute_v),
488              # col = 'white',
489              alpha = .4) +
490   labs(x = 'Theoretical distribution',
491        y = 'log(V51)',
492        col = 'Imputed?',
493        shape = 'Imputed?') +
494   scale_shape_manual(values = c(16,17)) +
495   theme_classic() +
496   theme(legend.justification = 'center',
497         legend.background = element_rect(
498           fill = 'grey98'))
499
500 # Random Forest Regression without categorical variables
501 p_rf1 <-
502   df %>%
503   drop_na(`Pyrite Type`) %>%
504   janitor::clean_names() %>%
505   select(-matches('lod$|2se$')) %>%
506   select(-`pyrite_type`, -generation, -texture, -reef, -c(source_file:comments)) %>%
507   missRanger::missRanger(data = .,
508                           pmm.k = 5,
509                           # formula = Ni61 ~ Ni60,
510                           maxiter = 10,
511                           seed = 321321,
512                           verbose = 0,
513                           num.trees = 1000) %>%
514   bind_cols(df1$impute_ti,
515            df1$impute_v,
516            df1$impute_co,
517            df1$impute_ni60) %>%
518   rename(impute_ti = `...46`,
519          impute_v = `...47`,
520          impute_co = `...48`,
521          impute_ni60 = `...49`) %>%
522   make_qq(dd = ., x = 'v51') %>%
523   ggplot(aes(x = qq, y = log(v51))) +
524   geom_qq_line(aes(sample = qq),
525               lty = 1, col = 'gray', size = 1,
526               alpha = .7) +
527   geom_point(aes(col = impute_v,
528                 shape = impute_v),
529              alpha = .4,
530              # col = 'white'
531              ) +
532   labs(x = 'Theoretical distribution',
```



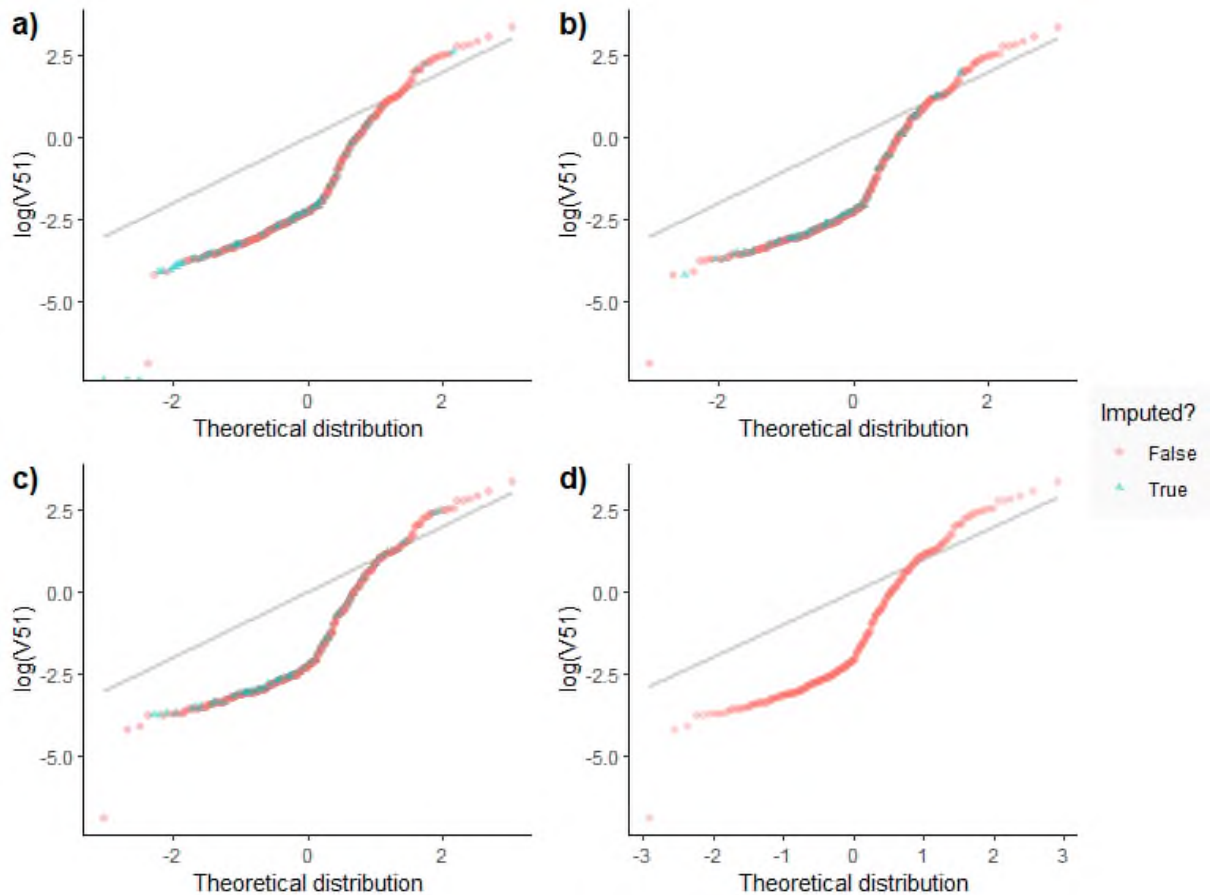
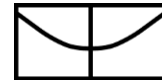
```
533     y = 'log(V51)',
534     col = 'Imputed?',
535     shape = 'Imputed?') +
536   scale_shape_manual(values = c(16,17)) +
537   theme_classic() +
538   theme(legend.justification = 'center',
539         legend.background = element_rect(
540           fill = 'grey98'))

541 ##
542 ## Missing value imputation by random forests
543 ##
544 ## Variables to impute:      mg25, al27, si29, p31, s33, s34, k39, ca43, ti49,
545 v51, cr53, mn55, co59, ni60, cu65, zn66, ga71, ge74, as75, se77, se82, sr88, zr90,
546 mo95, ag107, cd111, in115, sn118, sb121, te125, ba137, ta181, w182, pt195, au197,
547 hg202, tl205, pb206, pb207, pb208, bi209, th232, u238, ni61
548 ## Variables used to impute: v1, mg25, al27, si29, p31, s33, s34, k39, ca43, ti
549 49, v51, cr53, mn55, co59, ni60, cu65, zn66, ga71, ge74, as75, se77, se82, sr88, z
550 r90, mo95, ag107, cd111, in115, sn118, sb121, te125, ba137, ta181, w182, pt195, au
551 197, hg202, tl205, pb206, pb207, pb208, bi209, th232, u238, ni61
552 ## ni60  ti49  co59  as75  se77  se82  si29  bi209  s33 p31 pb208
553 cu65  s34 v51 zn66  au197 pb206 pb207 sb121 k39 hg202 u238 al27
554 ag107 w182 te125 mn55 ba137 cr53 ta181 ge74 mg25 sr88 th
555 232 ga71 tl205 zr90 mo95 ni61 ca43 sn118 pt195 cd111 in11
556 5
557 ## iter 1: 1.0984 0.9778 0.9965 0.8076 0.8318 0.4159 0.9323 1.0150 0.4100
558 0.9744 0.8928 0.0863 0.3518 0.6374 1.0521 0.7236 0.5422 0.2208 0.9814 0.
559 6101 0.9118 1.0325 0.8896 0.7977 0.8346 1.0043 1.0046 0.8982 0.6808 0.92
560 51 0.5294 0.7825 1.0083 0.8175 0.5663 0.5401 0.8608 0.4141 0.6516 0.7580
561 0.1158 1.0724 1.1716 0.2345
562 ## iter 2: 0.3723 0.4508 0.4851 0.4009 0.3997 0.4136 0.6449 0.6256 0.1106
563 0.4962 0.6807 0.0907 0.1298 0.2900 1.0457 0.5339 0.4181 0.2509 0.8185 0.
564 3246 0.3550 0.8766 0.8124 0.6530 0.5592 0.9023 0.6879 0.7890 0.4692 0.69
565 20 0.4680 0.5910 0.9721 0.5138 0.3880 0.3131 0.4296 0.2869 0.5159 0.4987
566 0.0920 1.0014 1.1498 0.1896
567 ## iter 3: 0.3508 0.4495 0.4654 0.4452 0.4003 0.3987 0.6406 0.6372 0.1052
568 0.5104 0.6879 0.1203 0.1243 0.2826 1.0382 0.5882 0.4123 0.2394 0.8141 0.
569 3171 0.3951 0.8912 0.8210 0.6645 0.6163 0.8923 0.6323 0.7860 0.4564 0.70
570 57 0.4843 0.5951 0.9771 0.5449 0.3799 0.2903 0.4606 0.2789 0.5045 0.4487
571 0.0907 1.0200 1.1339 0.2114

572 # Random Forests regression with caterogical variables
573 p_rf2 <-
574   df %>%
575   drop_na(`Pyrite Type`) %>%
576   janitor::clean_names() %>%
577   select(-matches('lod$|2se$')) %>%
578   select(-`pyrite_type`, -c(source_file:comments)) %>%
579   missRanger::missRanger(data = .,
580                          pmm.k = 5,
581                          maxiter = 10,
582                          seed = 321321,
583                          verbose = 2,
584                          num.trees = 1000) %>%
585   bind_cols(df1$impute_ti,
586            df1$impute_v,
587            df1$impute_co,
588            df1$impute_ni60) %>%
589   rename(impute_ti = `...49`,
```



```
590     impute_v = `...50`,
591     impute_co = `...51`,
592     impute_ni60 = `...52`) %>%
593 make_qq(dd = ., x = 'v51') %>%
594 ggplot(aes(x = qq, y = log(v51))) +
595 geom_qq_line(aes(sample = qq),
596             lty = 1, col = 'gray', size = 1,
597             alpha = .7) +
598 geom_point(aes(col = impute_v,
599              shape = impute_v),
600           alpha = .4,
601           # col = 'white'
602 ) +
603 labs(x = 'Theoretical distribution',
604      y = 'log(V51)',
605      col = 'Imputed?',
606      shape = 'Imputed?') +
607 scale_shape_manual(values = c(16,17)) +
608 theme_classic() +
609 theme(legend.justification = 'center',
610       legend.background = element_rect(
611         fill = 'grey98'))
612 # Adjusting the plots
613 ggpubr::ggarrange(p_ld1, p_rf1, p_rf2, p_original,
614                  ncol = 2, nrow = 2, align = 'hv',
615                  labels = c('a', 'b', 'c', 'd')),
616                  common.legend = TRUE, legend = 'right')
```

617

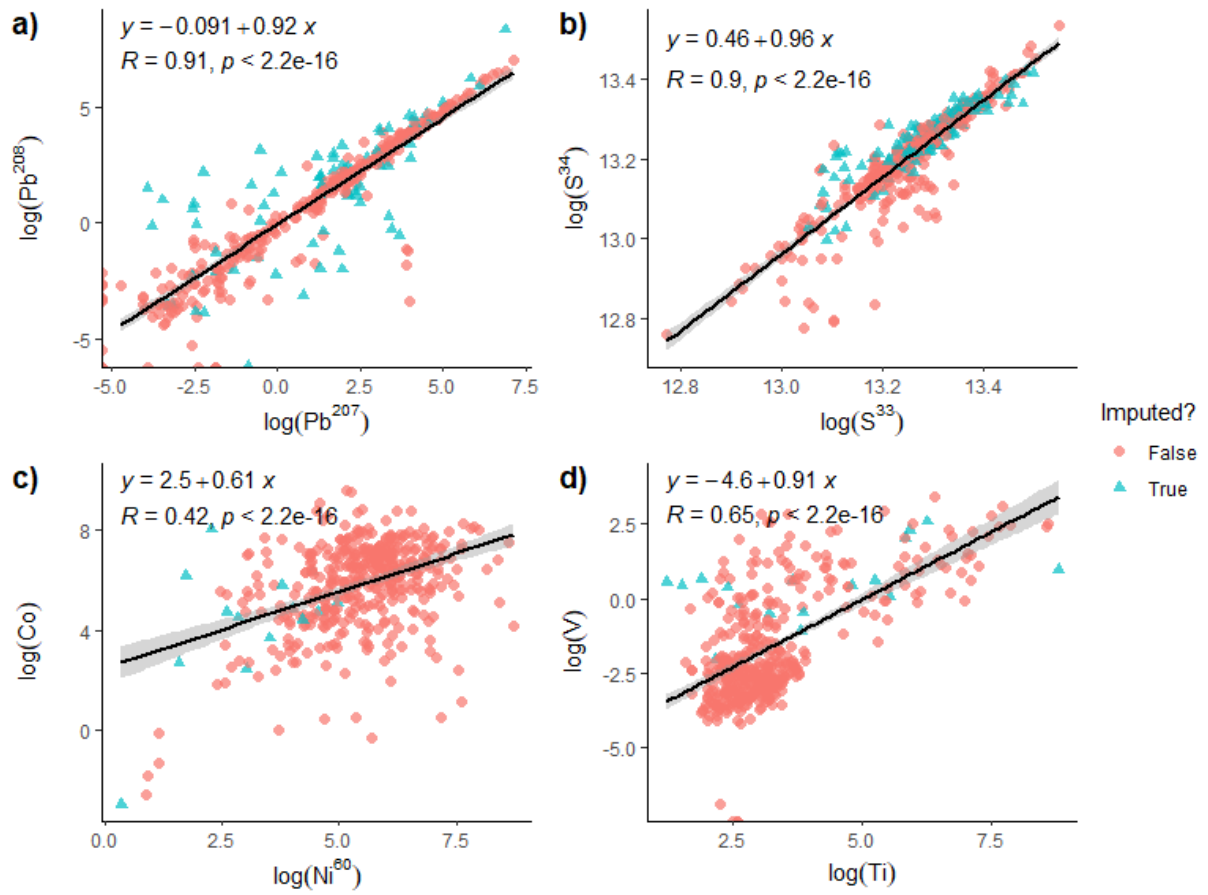
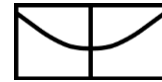
618 Fig. D. 5: Quantile-plot for the concentration of V51 on a logarithmic scale according to different distributions:
619 a) imputation based on a fraction of the detection limit, b) Random Forests imputation based only in quantitative
620 variables, c) Random Forests imputation based on quantitative and categorical variables (e.g., grain texture) and
621 d) original distribution (without imputed values).
622

623 Replacement of Missing Values

```
624 p1 <-  
625   df2 %>%  
626   ggplot(aes(log(Pb207), log(Pb208))) +  
627   geom_point(aes(col = impute_pb207,  
628                 shape = impute_pb206),  
629             cex = 2,  
630             alpha = .7) +  
631   geom_smooth(method = 'lm', col = 'black') +  
632   theme_classic() +  
633   ggpubr::stat_regline_equation(label.y = 8.5) +  
634   ggpubr::stat_cor(label.y = 7) +  
635   labs(col = 'Imputed?', shape = 'Imputed?',  
636        x = expression(log(Pb207)),  
637        y = expression(log(Pb208)))  
638  
639 p2 <-  
640   df2 %>%  
641   ggplot(aes(log(S33),  
642             log(S34))) +  
643   geom_point(aes(col = impute_s34,
```

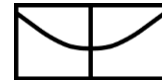


```
644         shape = impute_s34),
645         cex = 2,
646         alpha = .7) +
647 geom_smooth(method = 'lm', col = 'black') +
648 theme_classic() +
649 ggpubr::stat_regline_equation(label.y = 13.5) +
650 ggpubr::stat_cor(label.y = 13.4) +
651 labs(col = 'Imputed?', shape = 'Imputed?',
652       x = expression(log(S^33)),
653       y = expression(log(S^34)))
654
655 p3 <-
656 df2 %>%
657 ggplot(aes(log(Ni60),
658            log(Co59))) +
659 geom_point(aes(col = impute_ni60,
660               shape = impute_ni60),
661            cex = 2,
662            alpha = .7) +
663 geom_smooth(method = 'lm', col = 'black') +
664 theme_classic() +
665 ggpubr::stat_regline_equation(label.y = 10) +
666 ggpubr::stat_cor(label.y = 8.5) +
667 labs(col = 'Imputed?', shape = 'Imputed?',
668       x = expression(log(Ni^60)),
669       y = expression(log(Co)))
670
671 p4 <-
672 df2 %>%
673 ggplot(aes(log(Ti49),
674            log(V51))) +
675 geom_point(aes(col = impute_ti,
676               shape = impute_ti),
677            cex = 2,
678            alpha = .7) +
679 geom_smooth(method = 'lm', col = 'black') +
680 theme_classic() +
681 ggpubr::stat_regline_equation(label.y = 4) +
682 ggpubr::stat_cor(label.y = 2.8) +
683 labs(col = 'Imputed?', shape = 'Imputed?',
684       x = expression(log(Ti)),
685       y = expression(log(V)))
686 # Adjusting the plots
687 ggpubr::ggarrange(p1, p2, p3,p4,
688                  ncol = 2,nrow = 2,align = 'hv',
689                  labels = c('a','b','c','d')),
690                  common.legend = TRUE,legend = 'right')
```



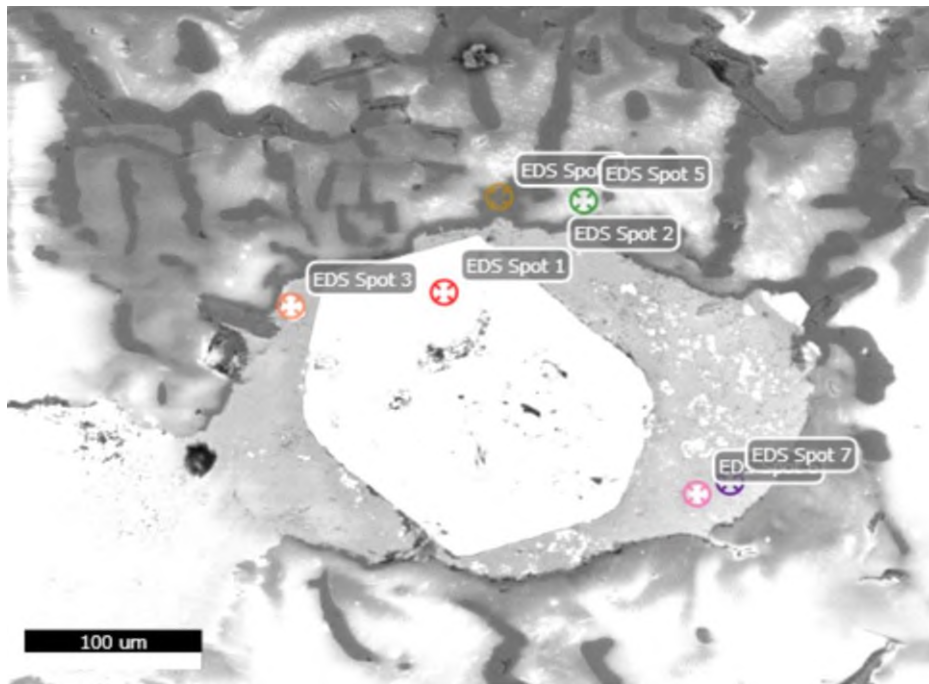
691

692 Fig. D. 6: Comparison of correlated log-transformed elements showing the relation of values according to whether
693 they originated from imputation (True, blue triangles) or not (False, orange circles). This plot was prepared based
694 on high correlated isotopes a) Pb207 - Pb208, b) S33 - S34, and based on moderated correlated elements c) Ni60-
695 Co and d) Ti - V.

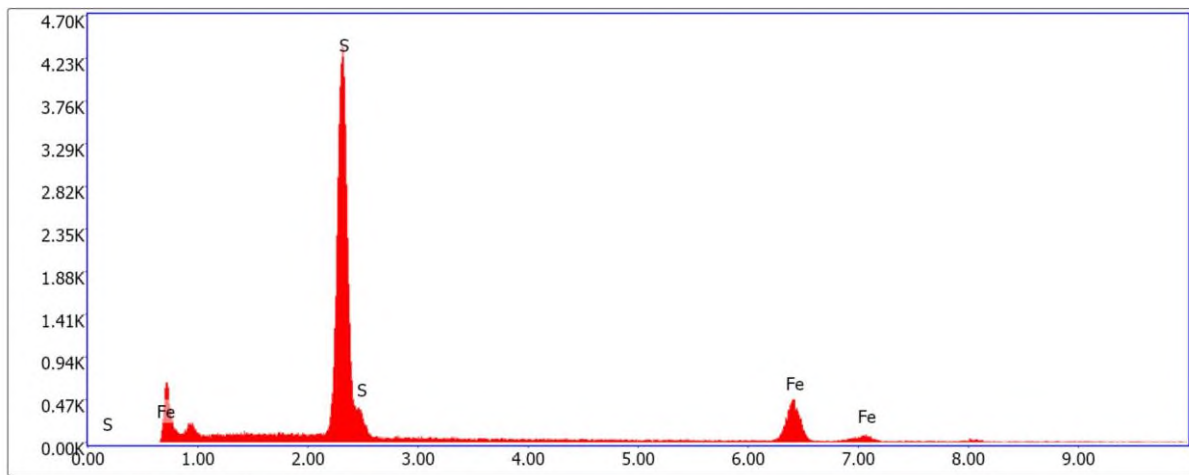


APÊNDICE F –Análises de EDS

LÂMINA HHS-443 – CAMPO 01

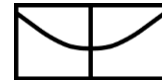


EDS Spot 1 - EDS1

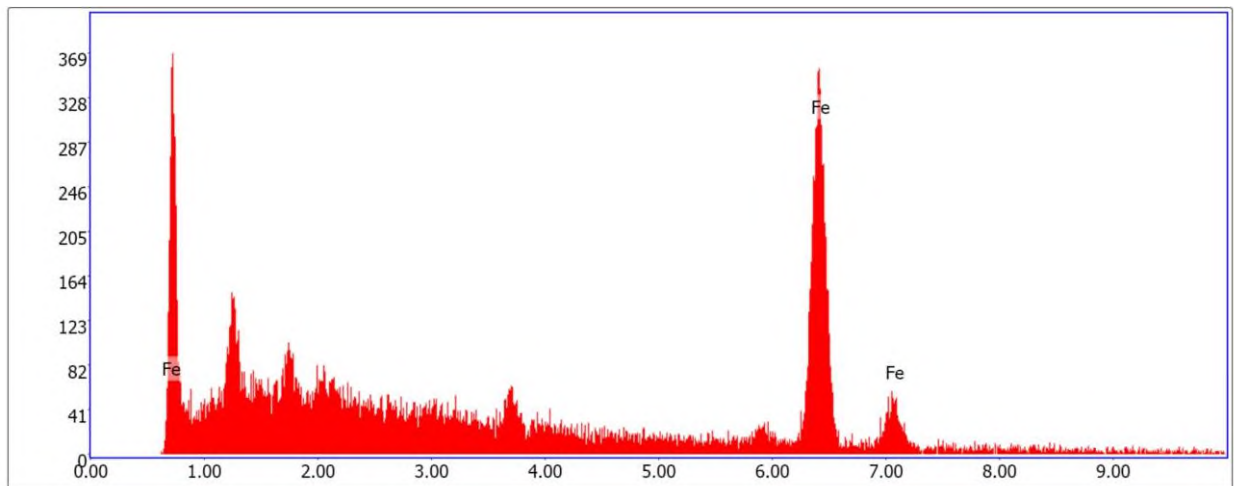


Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
S K	59.95	72.28	1616.87	2.39	0.5526	0.9625	0.9534	1.0046
FeK	40.05	27.72	222.28	4.52	0.3424	0.8511	0.9944	1.0100



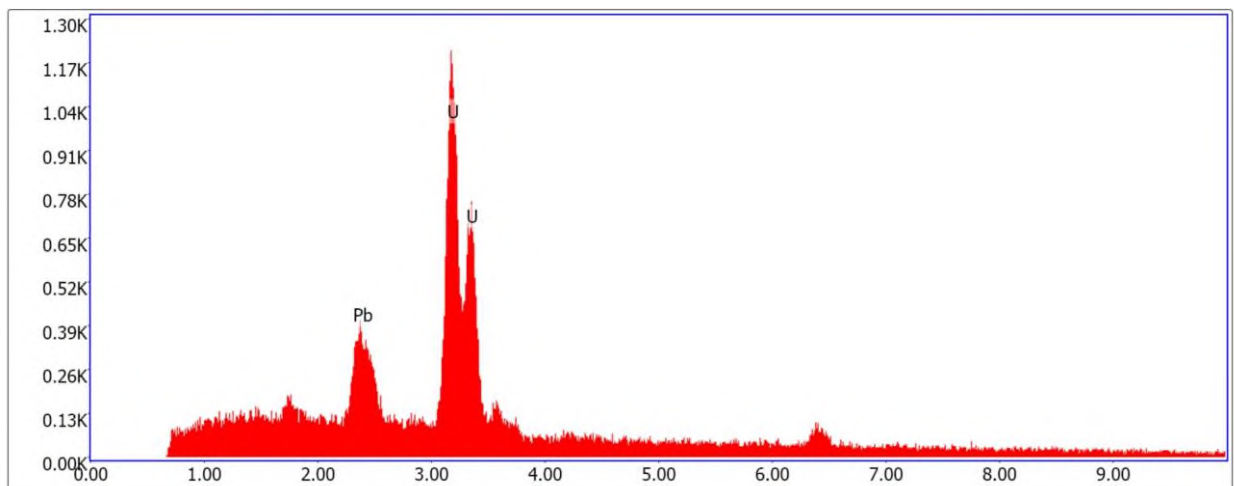
EDS Spot 2 - EDS1



Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

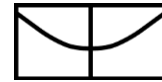
Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
FeK	100.00	100.00	158.46	5.00	0.9291	0.9266	1.0027	1.0000

EDS Spot 6 - EDS1

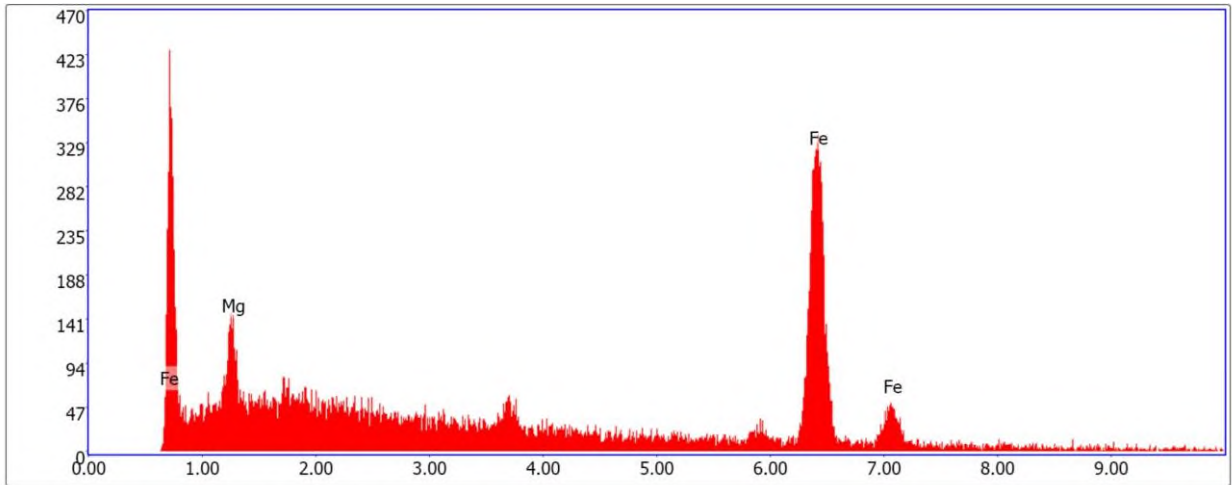


Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
PbM	16.08	18.04	109.92	6.92	0.1403	0.8525	1.0187	1.0049
U M	83.92	81.96	444.74	3.76	0.6991	0.8323	1.0020	0.9990

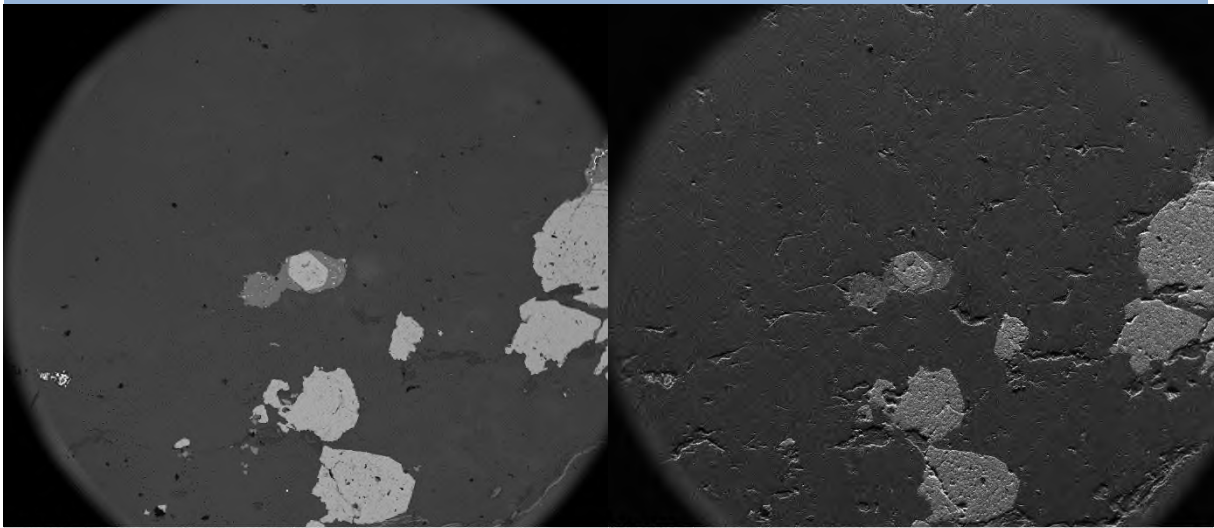


EDS Spot 7 - EDS1

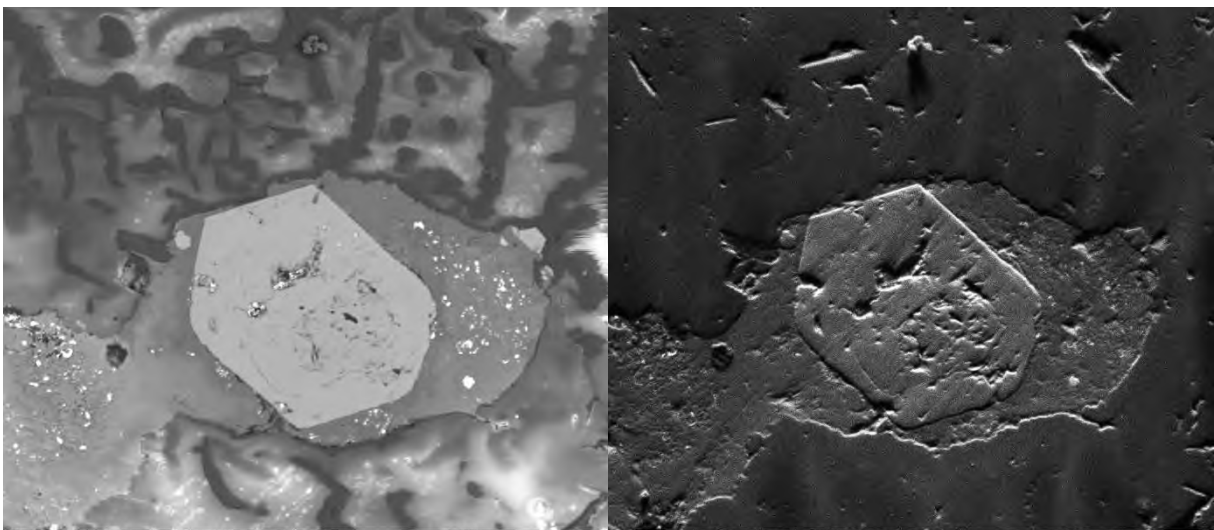


Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

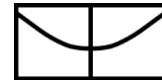
Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
MgK	3.87	8.46	25.26	21.32	0.0209	1.0570	0.5114	1.0017
FeK	96.13	91.54	163.71	4.73	0.8870	0.9199	1.0028	1.0003



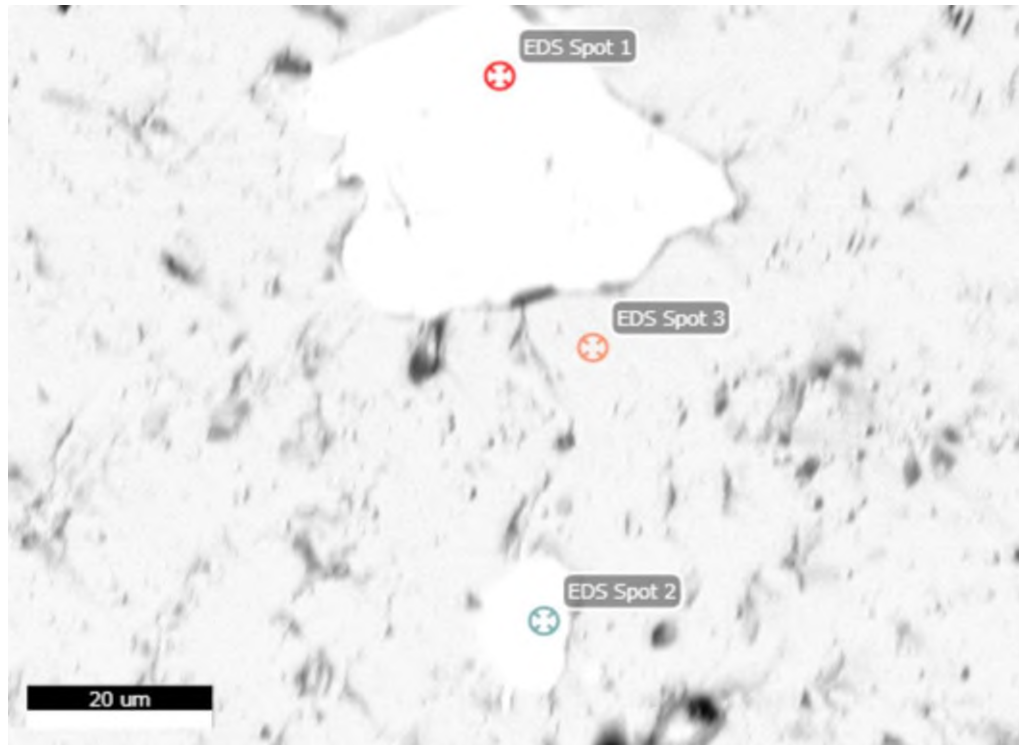
2/11/2021 4:05:31 PM HV 15.00 kV spot 6.0 WD 10.6 mm det vCD mag 50 x SEM CPRM-UnB 2/11/2021 4:05:31 PM HV 15.00 kV spot 6.0 WD 10.6 mm det ETD mag 50 x SEM CPRM-UnB



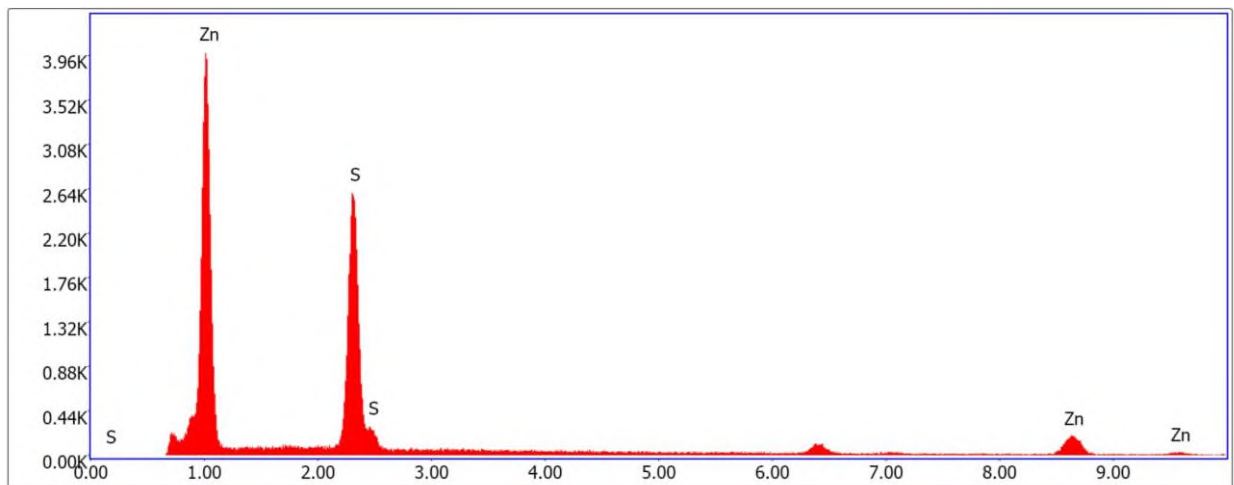
2/11/2021 4:01:31 PM HV 15.00 kV spot 6.0 WD 10.6 mm det vCD mag 312 x SEM CPRM-UnB 2/11/2021 4:01:31 PM HV 15.00 kV spot 6.0 WD 10.6 mm det ETD mag 312 x SEM CPRM-UnB



LÂMINA HHS-443 – CAMPO 02

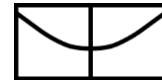


EDS Spot 1 - EDS1

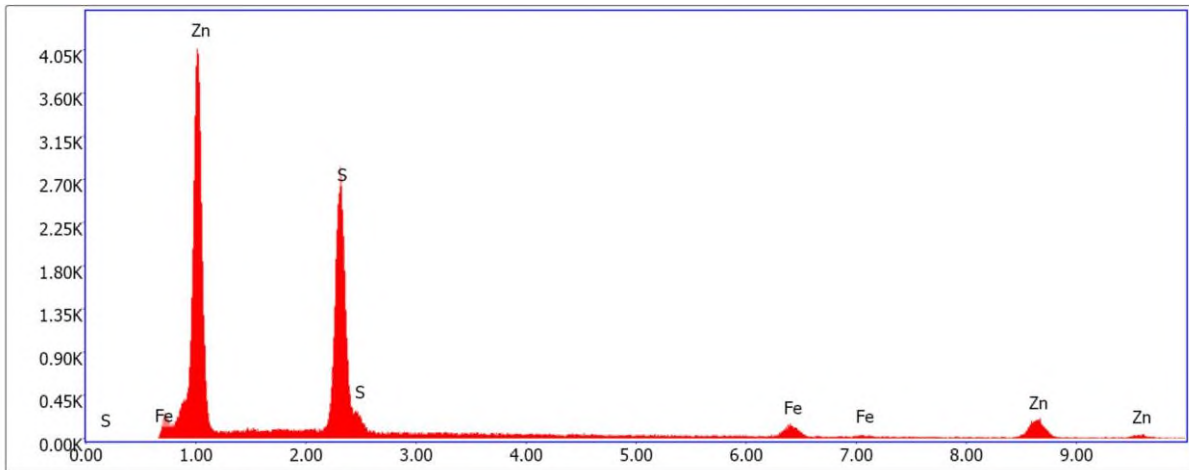


Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
S K	43.06	60.65	942.84	3.42	0.3797	0.9938	0.8843	1.0034
ZnK	56.94	39.35	103.88	8.44	0.4998	0.8677	1.0006	1.0109



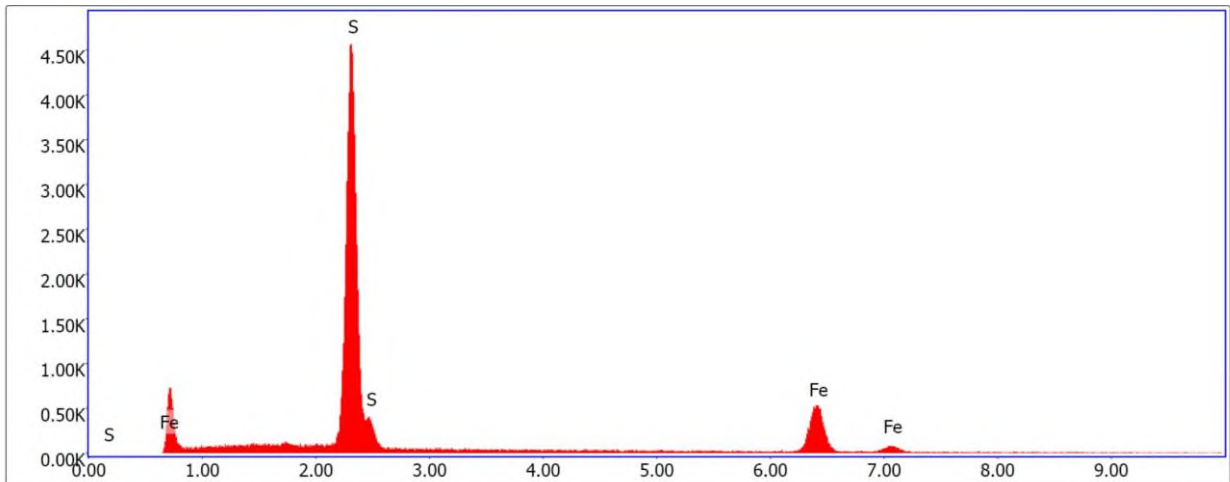
EDS Spot 2 - EDS1



Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

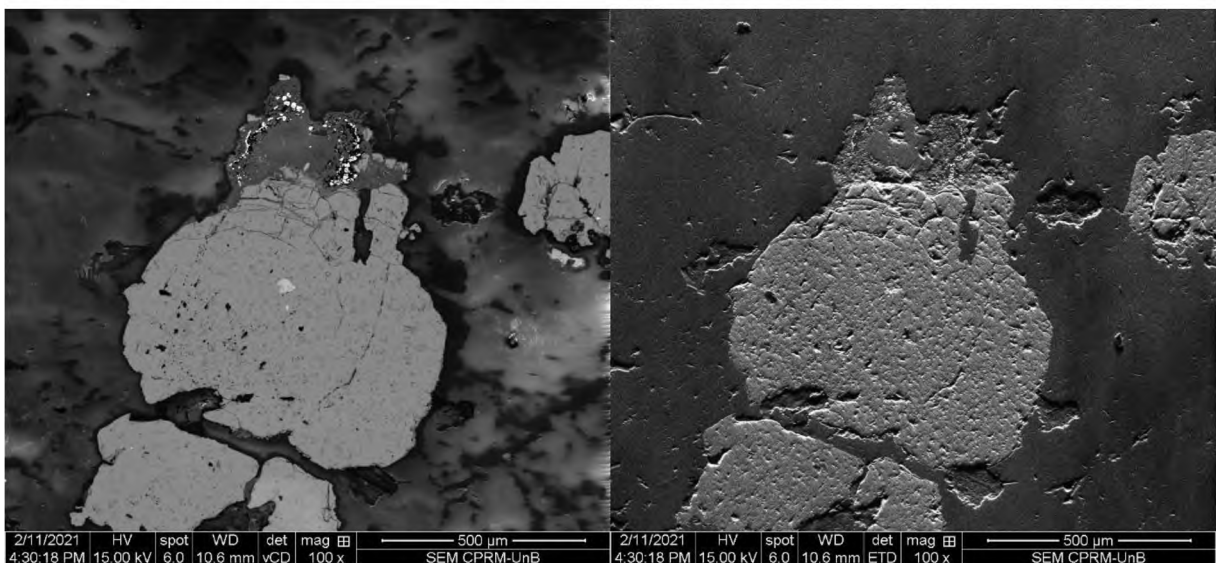
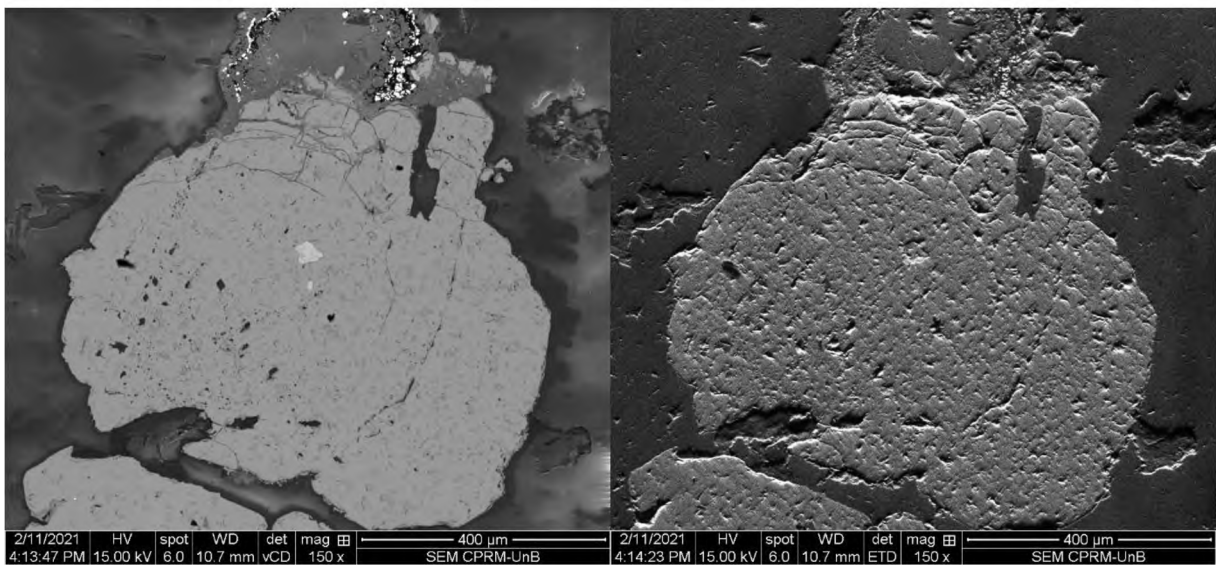
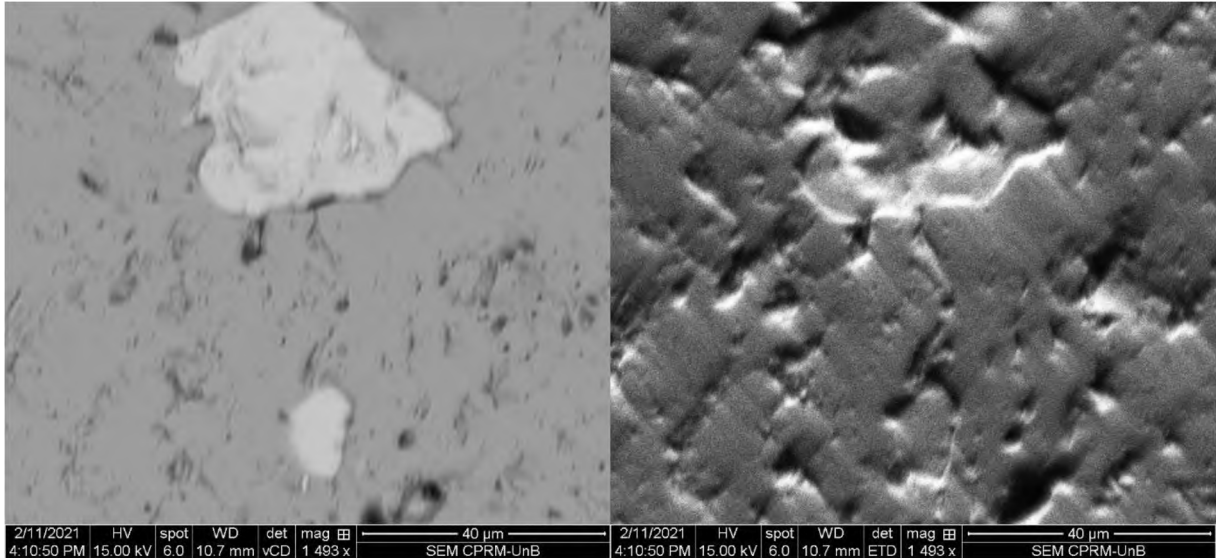
Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
S K	38.84	55.74	997.17	3.41	0.3436	0.9975	0.8834	1.0039
FeK	10.07	8.30	64.83	11.50	0.1006	0.8895	0.9944	1.1300
ZnK	51.09	35.97	109.28	7.33	0.4498	0.8719	0.9984	1.0113

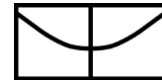
EDS Spot 3 - EDS1



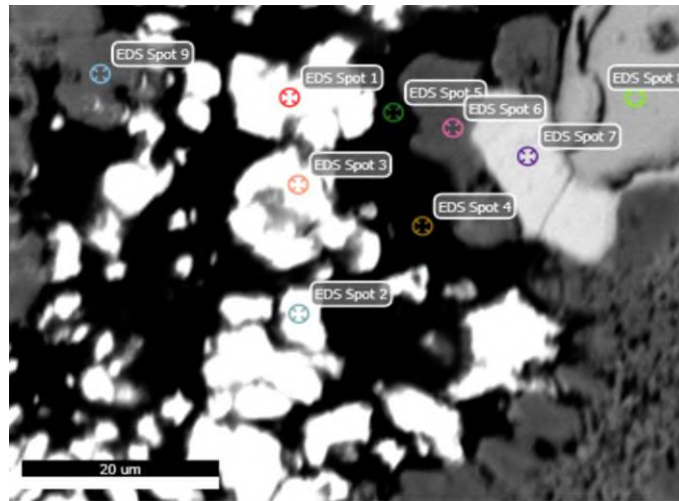
Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
S K	56.21	69.10	1667.77	2.44	0.5178	0.9666	0.9481	1.0051
FeK	43.79	30.90	268.71	4.23	0.3760	0.8554	0.9949	1.0089

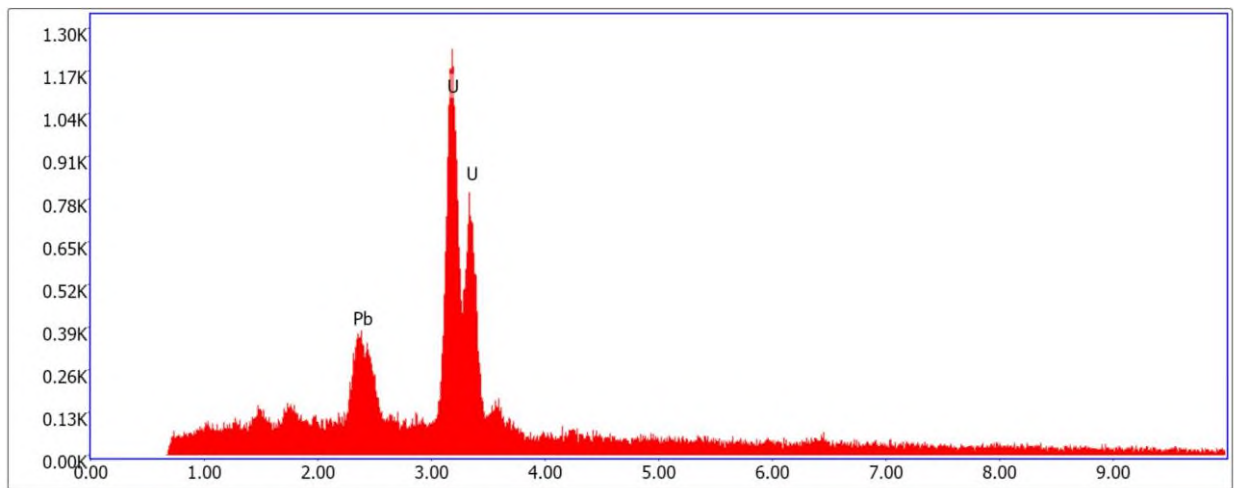




LÂMINA HHS-443 – CAMPO 03



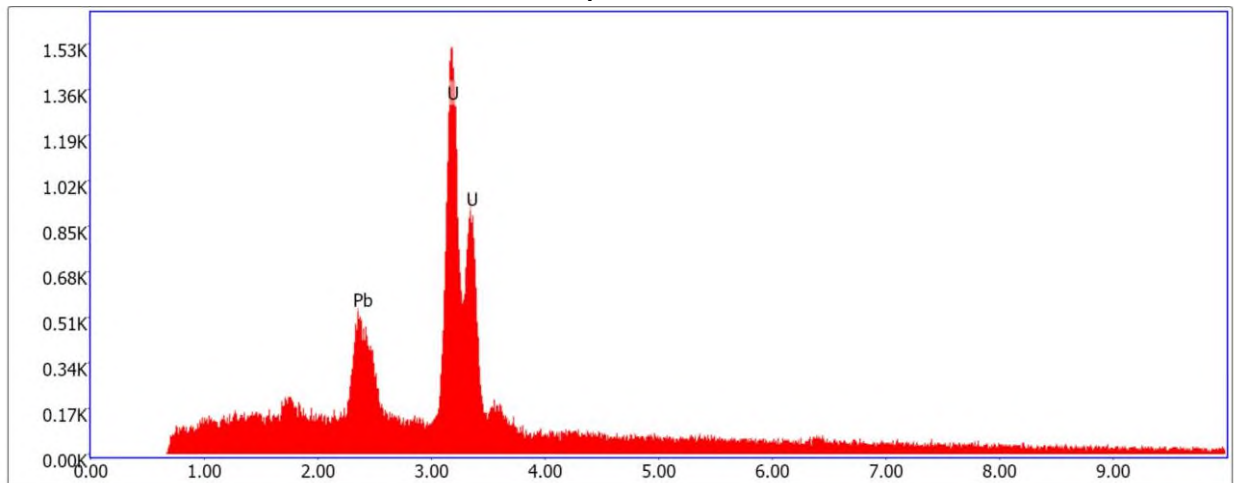
EDS Spot 1 - EDS1



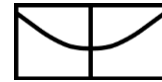
Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
PbM	14.25	16.03	97.50	9.47	0.1243	0.8529	1.0181	1.0050
U M	85.75	83.97	456.98	3.53	0.7176	0.8327	1.0060	0.9991

EDS Spot 2 - EDS1

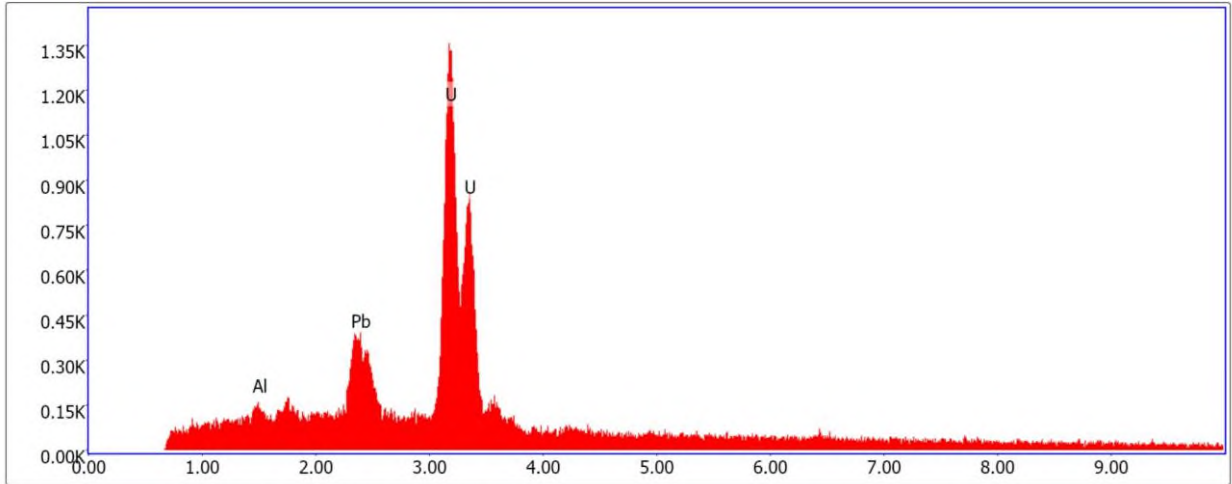


Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD



Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
PbM	16.28	18.26	141.21	6.68	0.1421	0.8525	1.0188	1.0049
U M	83.72	81.74	562.59	3.28	0.6971	0.8322	1.0016	0.9990

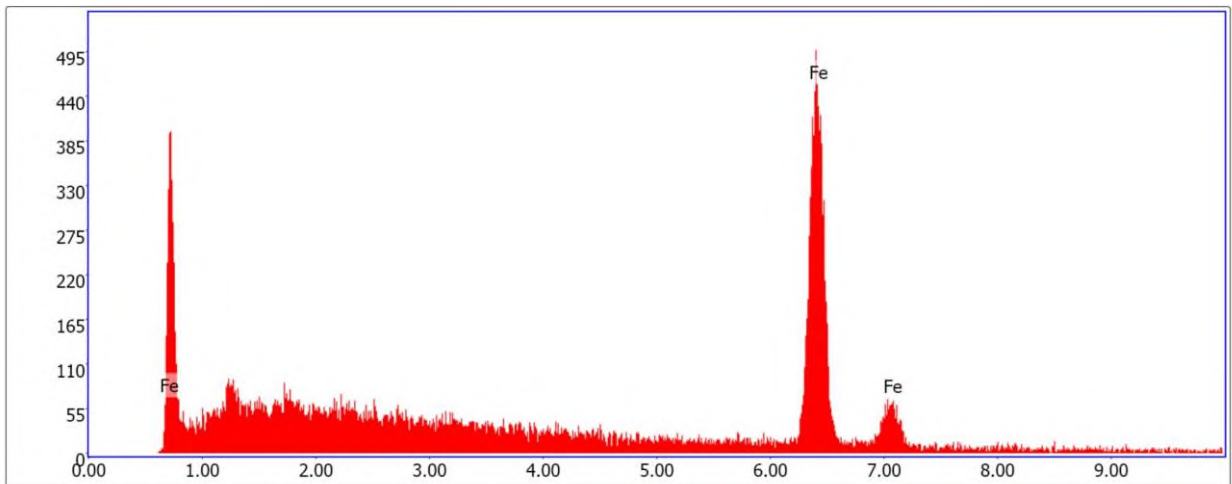
EDS Spot 3 - EDS1



Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

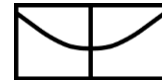
Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
AlK	0.87	7.02	17.67	40.27	0.0069	1.3337	0.5984	0.9998
PbM	13.89	14.66	109.20	8.39	0.1208	0.8497	1.0184	1.0050
U M	85.25	78.32	522.51	3.41	0.7117	0.8294	1.0075	0.9991

EDS Spot 6 - EDS1

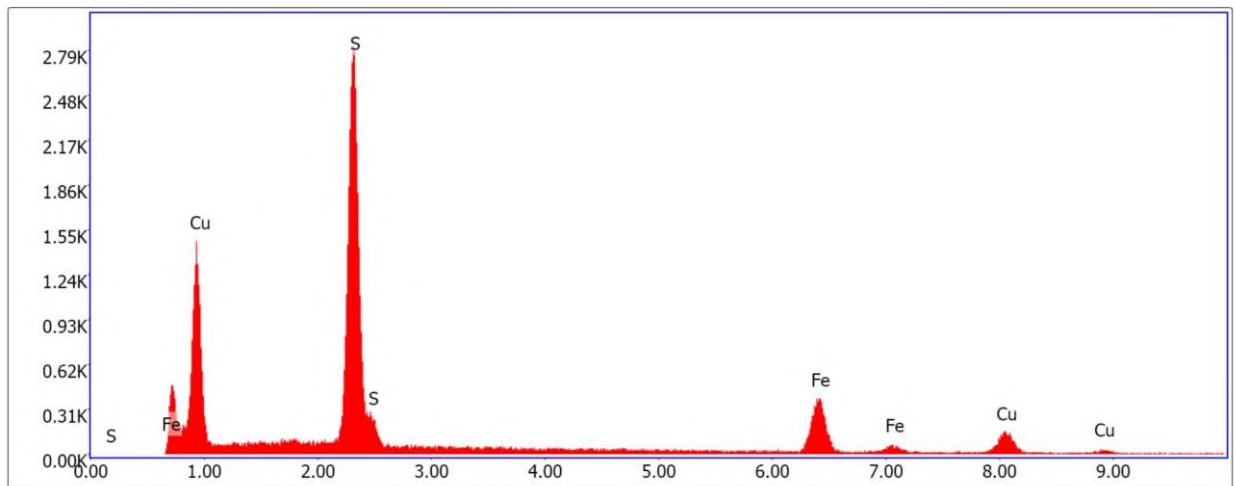


Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
FeK	100.00	100.00	232.26	4.21	0.9292	0.9266	1.0028	1.0000

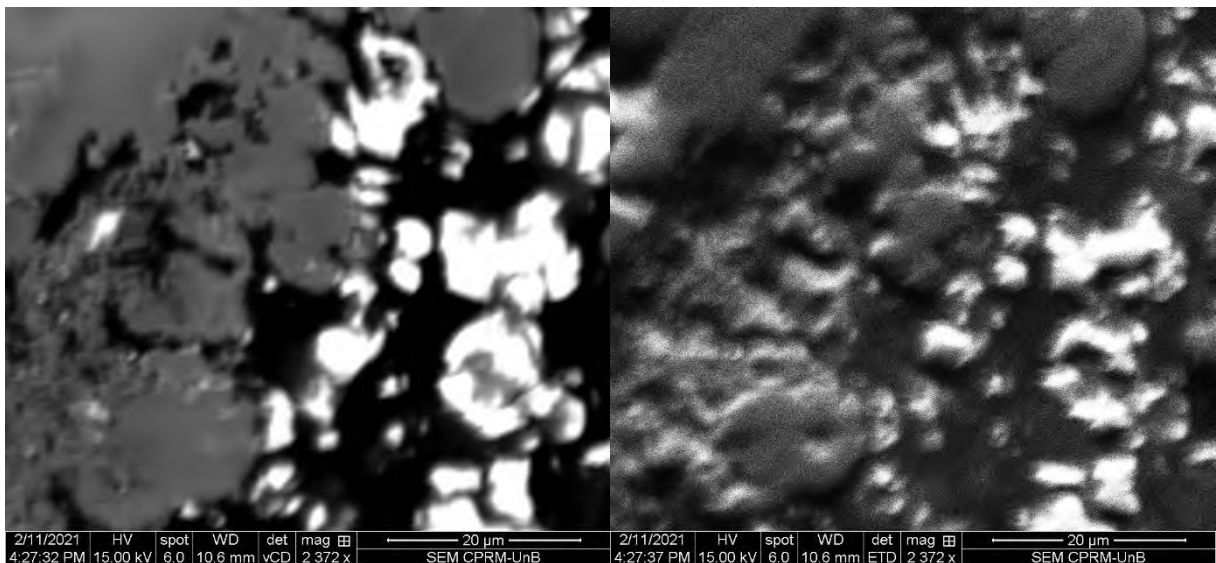


EDS Spot 7 - EDS1



Lsec: 30.0 0 Cnts 0.000 keV Det: Apollo X-SDD

Element	Weight %	Atomic %	Net Int.	Error %	Kratio	Z	A	F
S K	40.13	55.34	1032.16	3.01	0.3638	0.9918	0.9091	1.0054
FeK	31.35	24.81	183.50	5.37	0.2915	0.8822	0.9960	1.0583
CuK	28.52	19.85	78.79	9.48	0.2461	0.8596	0.9904	1.0134



2/11/2021 4:27:32 PM HV 15.00 kV spot 6.0 WD 10.6 mm det vCD mag 2.372 x SEM CPRM-UnB 2/11/2021 4:27:37 PM HV 15.00 kV spot 6.0 WD 10.6 mm det ETD mag 2.372 x SEM CPRM-UnB