

Procesamiento de anomalías en métodos potenciales de prospección mediante el aprendizaje automático probabilístico

Processing total field anomalies with statistically-based machine learning

Julián L Gómez^{1,2,3}, Ana Carolina Pedraza De Marchi^{1,2,4}, Claudia L Ravazzoli^{2,5}

Resumen El aprendizaje automático es actualmente una herramienta disruptiva para el procesamiento de señales digitales y la toma de decisiones. En particular, el aprendizaje automático probabilístico permite aproximar la función de densidad de probabilidad bajo la cual se distribuyen las anomalías observadas. Utilizando el aprendizaje automático probabilístico, proponemos asistir al intérprete de datos de los métodos potenciales de prospección. Mediante la generación de alternativas estadísticamente consistentes con los datos de trabajo, el geocientista puede visualizar variaciones realistas del dato original que le permitan expandir su interpretación. Para ello, el intérprete ingresa al sistema la anomalía de interés. El sistema deduce a partir del dato suministrado una aproximación a su función de densidad de probabilidad. Luego, el sistema permite al usuario seleccionar una región del dato de entrada para generar en ella realizaciones distribuidas bajo la misma densidad de probabilidad de los datos observados. Evaluamos la novedad de los datos generados respecto del dato original en la región seleccionada, permitiendo al intérprete ponderar las propuestas obtenidas. Para inferir la función de densidad de probabilidad, utilizamos un método de adición y de remoción de ruido aleatorio sobre la cuadrícula suministrada. En la generación de datos, utilizamos un método de Monte Carlo basado en una cadena de Markov conocida como dinámica de Langevin. Algunos de los desafíos de la propuesta son el entrenamiento de un método de aprendizaje automático probabilístico con una sola base de datos y la limitación en el hardware y en el tiempo de cómputo que supone utilizar el método en una computadora personal. Presentamos una experiencia con datos sintéticos y en datos de campo de anomalía escalar de intensidad total magnética. Los resultados muestran que la propuesta puede asistir al intérprete en la delineación espacial de los cuerpos anómalos y la inversión de parámetros, tales como la dirección de magnetización.

Palabras clave Aprendizaje automático, estadística, interpretación, algoritmo, métodos potenciales de prospección.

Abstract Machine learning is currently a disruptive tool for digital signal processing and decision making. In particular, probabilistic machine learning makes it possible to approximate the probability density function under which the recorded signals are distributed. Using probabilistic machine learning, we propose to assist interpreters of potential field prospecting methods. By generating alternatives that are statistically consistent with the working data, the interpreter can visualize realistic variations of the original data that allow them to expand their insights. To do this, the interpreter enters the anomaly of interest into the system. The system deduces from the data provided an approximation to its probability density function. Then, the system allows the user to select a region of the input data to generate distributed realizations under the same probability density of the observed data. We evaluate the novelty of the data generated with respect to the original data in the selected region, allowing the interpreter to weigh the proposals obtained. To infer the probability density function, we use a

¹Facultad de Ciencias Astronómicas y Geofísicas, Universidad Nacional de La Plata, Argentina. Email: jgomez@fcaglp.unlp.edu.ar

²Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Argentina.

³YPF Tecnología S.A., Berisso, Argentina.

⁴Centro de Investigaciones Geológicas (CIG) de La Plata, Argentina.

⁵Centro de Investigaciones Geofísicas (CIGEOF), Universidad Nacional de La Plata, Argentina.

method of adding and removing random noise on the supplied grid. In generating data, we use a Monte Carlo method based on a Markov chain known as Langevin dynamics. Some of the challenges of the proposal are the training of a probabilistic machine learning method with a single database and the limitation in hardware and computing time involved in using the method on a personal computer. We present an experience on synthetic data and magnetic total intensity scalar anomaly field data. The results show that the proposal can assist the interpreter in the spatial delineation of anomalous bodies and in the inversion of parameters, such as the magnetization direction.

Keywords Machine learning, statistics, interpretation, algorithm, gravity and magnetic potential fields.

INTRODUCCIÓN

Correctamente anticipado por Flusser en 1985 (Flusser, 2023), las imágenes técnicas son hoy en día una herramienta diferencial en el desarrollo de las sociedades tecnológicas. El aprendizaje automático (ML, machine learning) y el aprendizaje profundo (deep learning) permiten resolver problemas inversos altamente no lineales a partir de imágenes o cuadrículas Burkov (2019); Chollet (2021); Murphy (2023). Actualmente, distintos métodos de ML han sido aplicados en las geociencias para resolver problemas de interpolación y clasificación. Entre ellos, para el procesamiento de datos sísmicos Liu y otros (2024), inversión de datos magnetotelúricos Zhou y otros (2024); la interpretación de reservorios a partir de imágenes de coronas (de Lima y otros, 2019); en inversión petrofísica (Li y otros, 2024) y en la inversión conjunta de datos de gravedad y magnetismo (Hu y otros, 2024).

Los métodos potenciales de prospección persiguen la interpretación de datos de anomalías de gravedad y de campo magnético para inferir propiedades del subsuelo terrestre (Blakely, 1996). Estos métodos desempeñan un papel muy importante en la exploración del petróleo y el gas, en la minería, y en estudios arqueológicos. Los registros digitales son por lo general interpolados sobre cuadrículas rectangulares para la construcción de mapas de isolíneas, constituyendo así una imagen técnica, a los fines de su posterior procesamiento e interpretación. Desde la interpretación estadística basada en el análisis de espectros de potencia propuesta por Spector y otros (1970), el método geoestadístico de Kriging (Journel, 1989) y la interpretación Bayesiana de los cuadrados mínimos, que asume que los registros son variables aleatorias debido al ruido de medición (Meju, 1995), el empleo de herramientas con bases estadísticas tiene una larga tradición en los métodos potenciales de prospección.

En el aprendizaje automático probabilístico, el proceso de entrenamiento es no supervisado y se asume que los datos observados son realizaciones de una variable aleatoria vectorial con una función de densidad de probabilidad (pdf, por sus siglas en inglés) definida Murphy (2023). El aprendizaje automático probabilístico permite aproximar, a partir de datos registrados, la pdf bajo la cual se distribuyen las observaciones. El conocimiento de la pdf permite generar muestras sintéticas de esa distribución, generando así nuevos datos sintéticos. Esta técnica tiene hoy un uso sin precedentes en tecnologías como los modelos de lenguaje grandes (p. ej., Gemini: <https://gemini.google.com/>), modelos de generación de imágenes y de videos (p. ej., Sora: <https://openai.com>).

Atentos a los avances del ML en las diferentes ramas de la ciencia y la tecnología, en este trabajo evaluamos un método de difusión probabilística (MDP) para asistir en la interpretación de datos magnéticos de los métodos potenciales de prospección. Nuestra motivación principal radica en que los métodos de MDP permiten fácilmente aproximar funciones de densidad de probabilidad asociada a datos observados. El problema de aproximar la pdf era considerado como extremadamente desafiante desde una perspectiva más clásica (Goodfellow y otros, 2016), al requerir un modelo explícito para la pdf y necesitar estimar la función de partición adecuada para normalizar la distribución (Murphy, 2023).

En este trabajo emplearemos de manera original un algoritmo de MDP, llamado *denoising score matching* (Song y otros, 2019; 2020). Por medio de este método aproximamos el gradiente respecto a la variable aleatoria del logaritmo de la pdf, (conocida también como por el nombre de función de

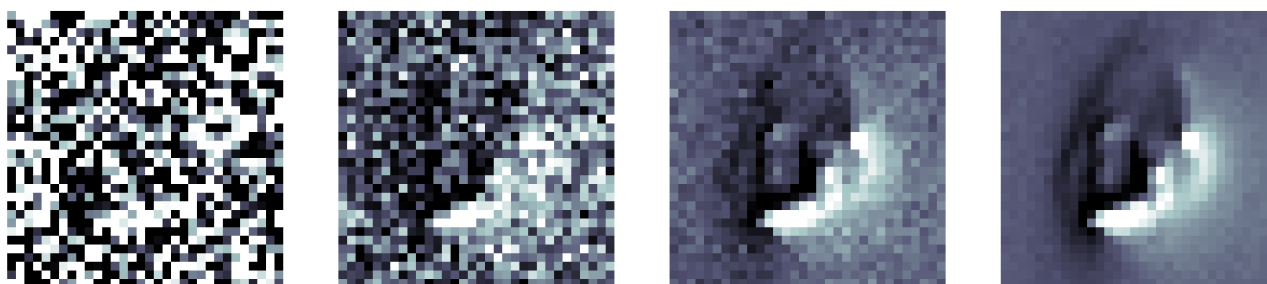


Figura 1. Ilustración del método de generación de datos estadísticamente consistente. De izquierda a derecha: una cuadrícula con valores aleatorios es transformada paulatinamente a una muestra consistente con el dato suministrado. En cada paso, el nivel de ruido inyectado disminuye. Las direcciones para navegar en el espacio de las imágenes vienen dada por el gradiente del logaritmo de la densidad de probabilidad. Se muestran 4 pasos de un total de 10.

score) y obtendremos muestras de esta distribución por medio de una cadena de Markov conocida como dinámica de Langevin (Murphy, 2023). Con este método obtendremos, al promediar muchas generaciones, una tendencia regional del dato de entrada. Esta tendencia nos permitirá generar un mapa residual que puede ser utilizado para asistir de manera complementaria en la interpretación de datos de anomalía escalar de intensidad total (TFA). Indagaremos, además, como los datos generados pueden ser utilizados para estimar rangos posibles de variación en parámetros obtenidos por medio de inversión.

A continuación, describimos el método empleado, con especial atención a sus dos etapas: la aproximación de la función de score de la función de densidad de probabilidad del dato y la posterior generación de cuadrículas sintéticas. A continuación, aplicamos la metodología expuesta en un dato de campo de acceso público de TFA en la región de Montes Claros (Brasil) utilizado por Reis y otros (2020) para la determinación de la dirección de magnetización. En este ejemplo obtenemos, de las cuadrículas generadas, una tendencia regional que es removida del dato original. Además, analizamos los histogramas que resultan de la inversión de la dirección de magnetización para las distintas realizaciones y comparamos los resultados con una técnica de Monte Carlo. Finalmente, discutimos el método implementado e indicamos nuevas líneas de desarrollo e investigación.

MÉTODO

El algoritmo de *denoising score matching* (DSM) parte de una imagen aleatoria y la transforma, progresivamente, en una muestra compatible con el dato original. El método consiste de dos etapas. En la primera etapa, se aproxima la función de score asociada a la función de densidad de probabilidad del dato suministrado. La función de score se define aquí como el gradiente respecto de la variable aleatoria del logaritmo de la función de densidad de probabilidad (Elad y otros, 2023). Esta función es la que provee las direcciones en el espacio de las imágenes que conducen la imagen aleatoria inicial. En la segunda etapa se toman muestras de la función de distribución aproximada, generando así nuevas cuadrículas.

La Figura 1 muestra cómo se transforma una cuadrícula de números aleatorios en una muestra sintética compatible con el dato de campo de anomalía escalar de intensidad total magnética que será utilizada en la sección siguiente.

Aproximación de la función de score

La cuadrícula de los datos observados es considerada una variable aleatoria multidimensional con una función de densidad de probabilidad bien definida. Para inferir el score de la función de densidad de probabilidad, DSM utiliza una estrategia de adición y de remoción de ruido aleatorio a diferentes

niveles de ruido o dispersiones sobre la cuadrícula suministrada. Este proceso permite aproximar las direcciones en las cuales conducir una muestra aleatoria en el espacio de las imágenes.

Dado un modelo, $\mathbf{M}(\mathbf{x}; \mathbf{W})$, con entrada dada por la muestra \mathbf{x} y con parámetros \mathbf{W} , el DSM que implementamos busca \mathbf{W}^* que minimiza el valor esperado sobre los datos de entrenamiento del cuadrado de la diferencia entre el modelo y la función de *score* del proceso que inyecta ruido aditivo Gaussiano sobre las muestras de entrenamiento:

$$\|\mathbf{M}(\mathbf{x}_n; \mathbf{W}) - (\mathbf{x} - \mathbf{x}_n)/\sigma^2\|_2^2. \quad (1)$$

El modelo resultante, $\mathbf{M}(\mathbf{x}; \mathbf{W}^*)$, es un campo vectorial que aproxima a la función de *score*. El ruido con distribución Gaussiana de dispersión σ tiene por *score* a $(\mathbf{x} - \mathbf{x}_n)/\sigma^2$, donde \mathbf{x}_n es la muestra \mathbf{x} contaminada por el ruido aditivo. Los niveles de ruido que contaminan al dato de entrenamiento se construyen siguiendo un patrón descendente $\sigma_1 < \sigma_2 < \dots < \sigma_L$ (Song y otros, 2019; 2020). La ecuación 1 define la función de costo $\mathcal{L}_{\mathbf{W}}$ cuyo valor esperado respecto a todas las muestras \mathbf{x}_n es minimizado en el entrenamiento. En la Figura 4 (izquierda) se presenta la función de costo a distintas épocas para el caso de los datos de campo de Montes Claros. Las amplitudes utilizadas en este trabajo se presentan en la Figura 4, que van de 1 a 0.1 en diez pasos ($L = 10$). Para el entrenamiento el dato original es normalizado en el rango $[-1, +1]$. Las generaciones obtenidas serán luego reescaladas al rango original del dato.

Para el modelo que aproxima la función de *score* utilizamos una red neuronal que es la composición de una red de codificación con una red de decodificación, ambas con la arquitectura propia de una U-Net (Ronneberger y otros, 2015). Una U-net es una red convolucional que utiliza decimación, concatenación e interpolación. Los parámetros del modelo son los coeficientes de la respuesta impulsiva de los filtros de convolución bidimensional y sus respectivos sesgos. En todas las capas los filtros tienen 3×3 coeficientes. La rama de codificación consiste de 4 capas y el número de filtros por capa se incrementa en 16, 24, 32 y 64, respectivamente. La rama de decodificación consiste en tres capas, donde los filtros se reducen en número de 32 a 16 en el orden inverso al número de filtros en la capa de codificación. Las capas 1 y 7, 2 y 6, 3 y 5 concatenan sus salidas.

Para encontrar los parámetros óptimos \mathbf{W}^* , utilizamos el algoritmo de gradiente descendiente Adam (Kingma y otros, 2014) con un parámetro de aprendizaje de 10^{-5} y el resto de sus parámetros en sus valores por defecto. Estas decisiones de diseño y entrenamiento son tomadas de forma empírica mediante prueba y error, una estrategia común en el aprendizaje automático (Burkov, 2019; Chollet, 2021) y el método científico.

Generación

En la generación de datos, utilizamos un cadena de Markov conocida como dinámica de Langevin (Elad y otros, 2023). La dinámica de Langevin permite transformar una muestra aleatoria de una distribución Gaussiana, en una muestra distribuida bajo la distribución del dato original. La muestra se genera mediante la recursión

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \alpha \mathbf{M}(\mathbf{x}_{t-1}; \mathbf{W}^*)/\sigma + \sqrt{2\alpha} \mathbf{z}_t, \quad (2)$$

donde es $1 \leq t \leq T$ es el número de iteración, $\alpha > 0$ es el tamaño de paso, y \mathbf{z}_t es una matriz con números aleatorios que siguen elemento por elemento una distribución normal de media cero y varianza unitaria. El término con \mathbf{z}_t permite tomar muestras de manera estocástica de diferentes regiones de la función de densidad de probabilidad. La ecuación 2 se calcula T veces en cada nivel de ruido σ_1 a σ_L . La metodología original no ha sido reportada sobre datos geofísicos de los métodos potenciales, sino sobre imágenes naturales que poco podrían tener en común con los datos de métodos potenciales. No encontramos, o no ha sido reportado hasta la fecha, literatura específica sobre la elección de parámetros (p. ej., T , σ_1 , σ_L) en situaciones similares a nuestra aplicación.

Los detalles del algoritmo DSM pueden ser consultados en mayor profundidad en [Song y otros \(2020\)](#). Implementamos el método en el lenguaje de programación Python y utilizando la plataforma TensorFlow (<https://tensorflow.org>) para obtener \mathbf{W}^* . Otras plataformas, como Pytorch (<https://pytorch.org>) o JAX (<https://jax.readthedocs.io/>) podrían también ser utilizadas.

Caracterización de las generaciones

Para evaluar cuantitativamente cada cuadrícula o muestra generada \mathbf{x} , calculamos su coeficiente de correlación como el producto escalar

$$r(\mathbf{x}, \mathbf{d}) = \frac{\mathbf{x} \cdot \mathbf{d}}{\|\mathbf{x}\|_2 \|\mathbf{d}\|_2}, \quad (3)$$

respecto al dato original \mathbf{d} . En esta expresión, tanto la muestra generada como la anomalía original se consideran como vectores 1D al momento de hacer su producto escalar. El coeficiente de correlación de la [ecuación 3](#) provee un buen resumen del dato generado, ya que indica si la muestra generada es similar al dato de entrada y además, su signo señala si la muestra generada introduce un cambio indeseable de polaridad en las amplitudes respecto de la anomalía original.

En las primeras épocas del entrenamiento, el modelo produce por medio de la dinámica de Langevin muestras que están muy poco correlacionadas con el dato original; el coeficiente de correlación de las muestras generadas es próximo a cero. Esto significa que el dato generado es más bien ruidoso y no comparte información con el dato original. A medida que el número de épocas aumenta, el modelo genera muestras con coeficiente de correlación positivo y próximo al valor uno. Un ejemplo de ello se presenta en la [Figura 5](#), donde se observa el histograma del coeficiente de correlación de los datos generados para un mismo modelo en distintas épocas de entrenamiento. Este comportamiento nos permite monitorear el proceso de aproximación y detener las iteraciones cuando las muestras generadas presentan una alta correlación con el dato original. De continuar el entrenamiento, las muestras generadas por la dinámica de Langevin presentan un coeficiente de correlación igual a 1 y se pierde el carácter probabilístico del método; siendo las muestras generadas idénticas al dato suministrado.

APLICACIONES

En esta sección utilizamos un dato de campo de la TFA de la región de Montes Claros en el estado de Goiás, centro de Brasil. La [Figura 2](#) presenta un mapa geológico de la región basado en [Dutra y otros \(2014\)](#). En Montes Claros se ha identificado la presencia de rocas ígneas alcalinas del Cretácico, dispuestas principalmente a lo largo de un lineamiento con orientación NO-SE ([Figura 2](#)). Además de la región de Montes Claros, se destacan otros complejos alcalinos como Diorama, Córrego dos Bois, Morro do Macaco y Fazenda Buriti. Esta provincia está compuesta por rocas máficas a ultramáficas con una gran diversidad de tipos petrográficos. Las intrusiones alcalinas están rodeadas por basamento precámbrico y por rocas sedimentarias fanerozoicas de la cuenca del Paraná. Para una descripción más detallada de la geología, recomendamos consultar [2014](#).

En el año 2004, la región de Montes Claros fue recorrida por un estudio aeromagnético a una altura aproximada de 100 m desde el terreno con un intervalo de muestreo espacial de 8 metros. Las líneas de vuelo están dispuestas en dirección N-S cada 500 metros y las líneas de cruce en dirección E-O cada 5000 metros. Los datos fueron corregidos por variaciones diurnas y ajustados para el Campo Geomagnético de Referencia Internacional para la época 2004.62, con una inclinación de $-19,5^\circ$ y una declinación de $-18,5^\circ$. De estas observaciones se obtuvieron los valores de TFA empleados en este trabajo ([Figura 3](#), arriba). Las anomalías magnéticas de la TFA han sido asociadas a las rocas alcalinas en las cuales se han encontrado evidencias de magnetización remanente ([Dutra y otros, 2009](#); [Zhang y otros, 2018](#); [Dutra y otros, 2014](#)).

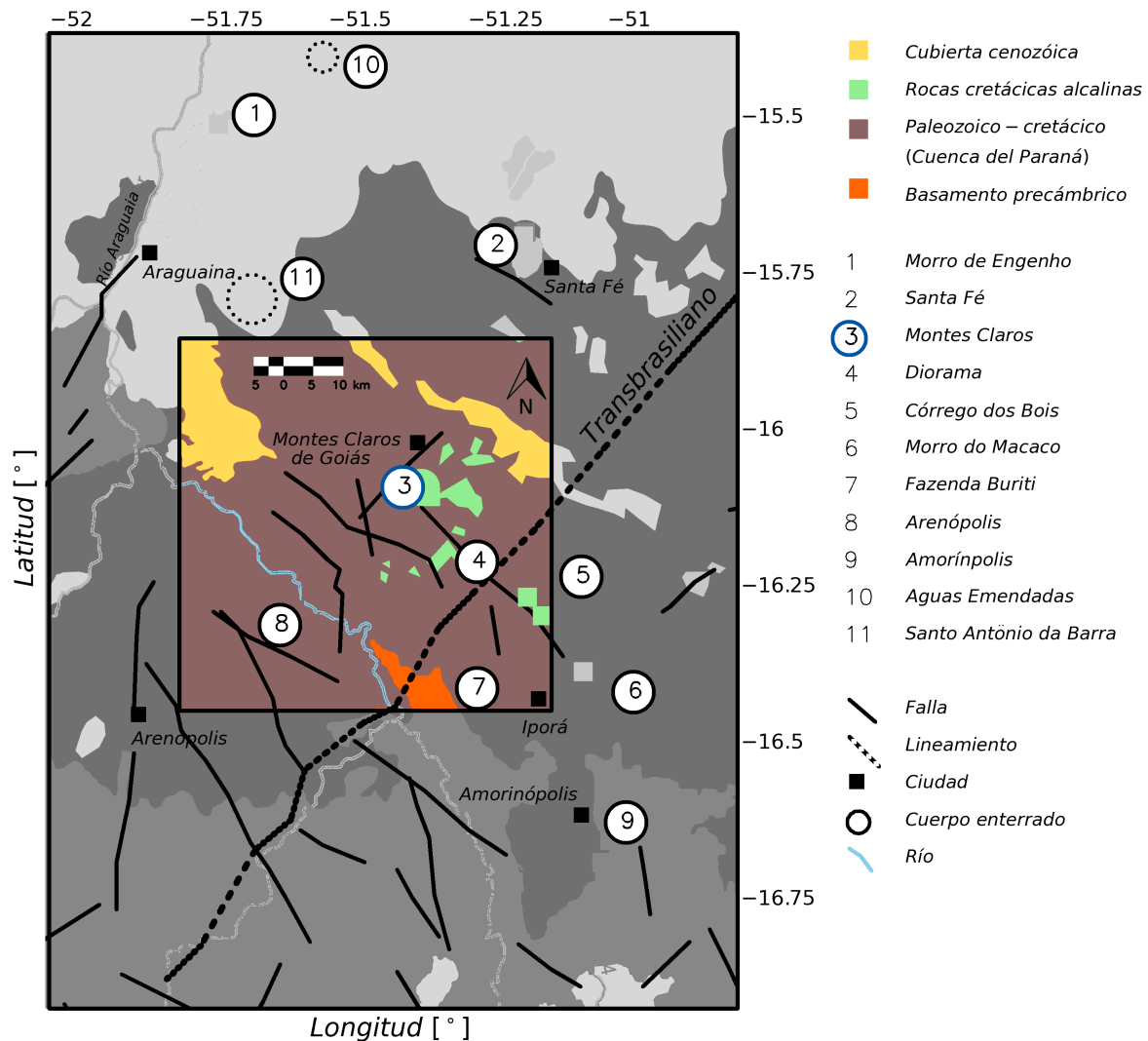


Figura 2. Mapa geológico de la región de Montes Claros. Los datos de TFA se encuentran en las inmediaciones del punto 3.

Recientemente, Reis y otros (2020) realizaron una inversión de parámetros basada en un modelo no lineal sobre la TFA observada. Los resultados obtenidos por medio del método indirecto de 2020 sugieren también la presencia de magnetización remanente.

Las observaciones originales han sido interpoladas utilizando la rutina `griddata` de la librería `interpolate` de SciPy (<https://scipy.org/>). Otro método de interpolación posible, y propio de los métodos potenciales, consistiría en utilizar el modelo de la capa equivalente (Blakely, 1996). Para conformar la entrada requerida por el algoritmo de generación, el dato interpolado es decimado resultando en una cuadrícula de 32×32 píxeles.

Entrenamiento y generación

En la Figura 4 (izquierda) presentamos la función de costo en función de las épocas de entrenamiento. Este gráfico narra el proceso de optimización que resulta en el modelo para aproximar el gradiente del logaritmo de la función de densidad de probabilidad del dato de entrada. Los niveles de ruido utilizados en el entrenamiento y luego en la generación de las cuadrículas se muestra en la Figura 4 (derecha).

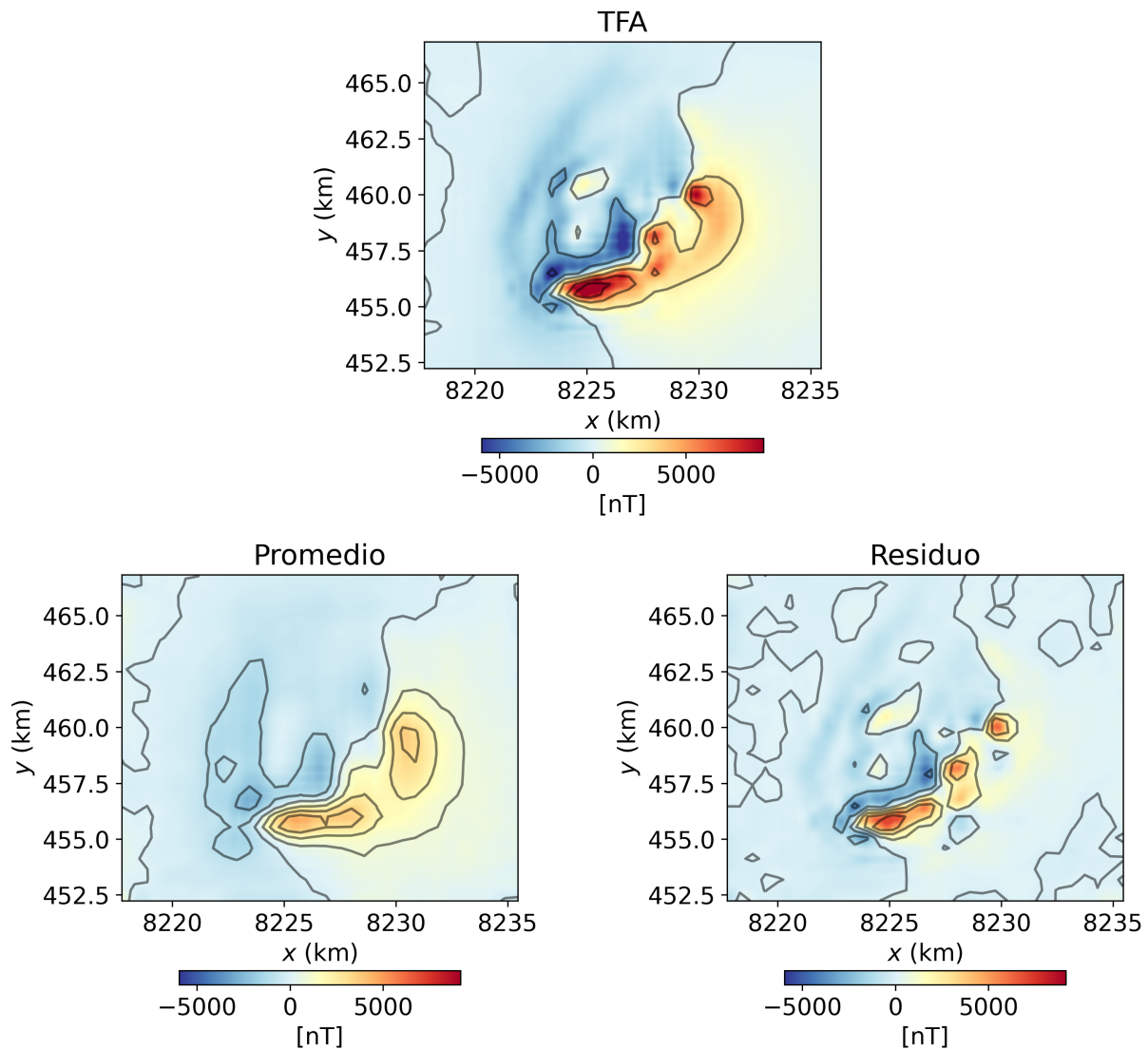


Figura 3. Arriba: cuadrícula original de la TFA en la zona de Montes Claros. Abajo: cuadrícula obtenida empleando el promedio de las generaciones (izquierda) y cuadrícula residual obtenida al remover del dato original el promedio de las generaciones obtenidas (derecha).

La [Figura 5](#) muestra la distribución de los valores del coeficiente de correlación para 500 generaciones dadas por el modelo entrenado por 10000 épocas y que es utilizado en esta demostración. La distribución de coeficientes de correlación es próxima a 0.9 lo que permite resumir que las realizaciones son similares al dato original y no presentan cambios en la polaridad de las amplitudes.

Para acelerar la generación de datos dada en la expresión (2), aplicamos a la matriz de valores aleatorios \mathbf{z}_t un filtro pasa bajos bidimensional, similar a la continuación ascendente para una altura artificial asignada por el nivel de ruido en la etapa de generación $t = 1, \dots, T$. Esta incorporación es un aporte original al algoritmo de generación dado en la [ecuación 2](#) y con ella logramos reducir notoriamente la cantidad de iteraciones. La [Figura 6](#) muestra para un modelo el efecto del parámetro T sobre los coeficientes de correlación de los datos generados. Podemos observar que utilizar un valor de T entre 5 a 10 es suficiente para obtener coeficientes de correlación próximos a 0.9. En los ejemplos de esta sección fijamos $T = 5$. La literatura por lo general recomienda valores mucho mayores para generar imágenes naturales realistas. El valor de T puede ser además ajustado observando algunas de las generaciones producidas en términos de su similitud con el dato original y nivel de ruido remanente.

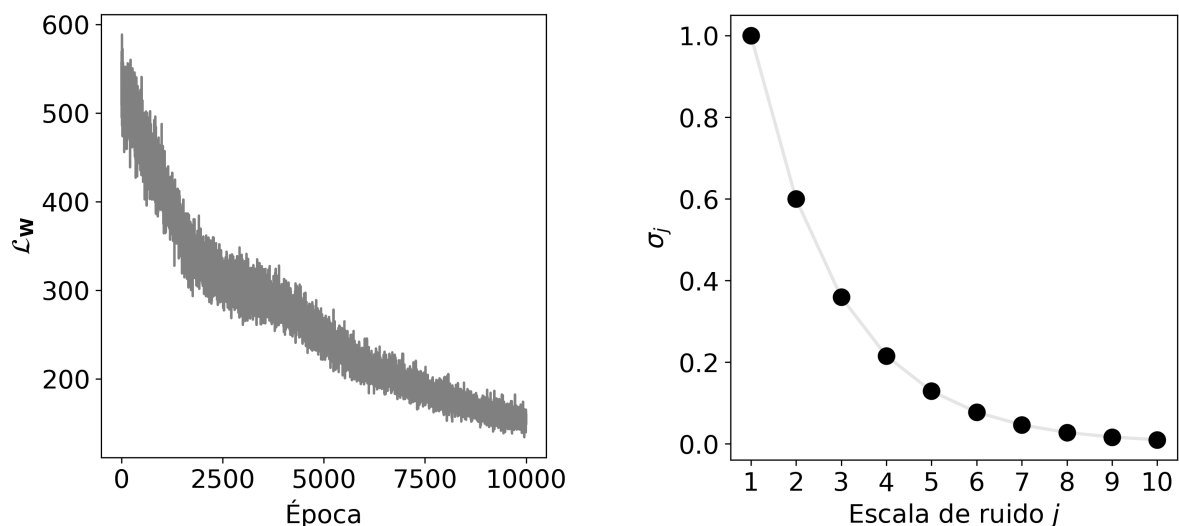


Figura 4. Izquierda: función de costo \mathcal{L}_W evaluada en cada época de entrenamiento. Derecha: niveles de ruido σ_j utilizados durante el entrenamiento.

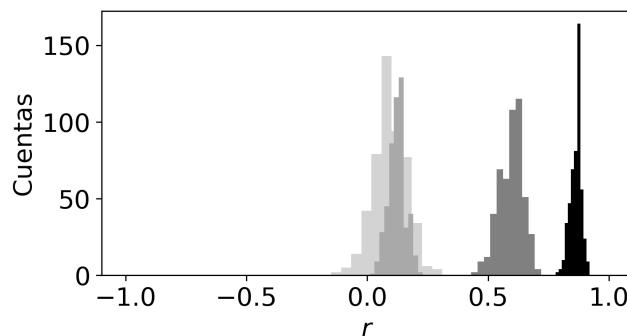


Figura 5. Los histogramas del coeficiente de correlación lineal muestran que el algoritmo parte de realizaciones poco correlacionadas con el dato original en las primeras épocas de entrenamiento (grises claros), hacia realizaciones altamente correlacionadas con el dato original (gris oscuro y negro).

Obtención de la tendencia regional

Las cuadrículas generadas por la dinámica de Langevin pueden ser utilizadas con distintos fines. En principio, ellas pueden ser empleadas individualmente como una visualización alternativa del área de estudio; de manera similar a lo indicado por [Mosegaard y otros \(1991\)](#) en base a la estrategia de Monte Carlo y, más recientemente, en el contexto del aprendizaje automático por [McAliley y otros \(2024\)](#).

Para construir un mapa de tendencia elegimos hacer un promedio simple de las generaciones obtenidas. La [Figura 3](#) (abajo, izquierda) presenta el promedio de las realizaciones generadas para el dato de Montes Claros. La tendencia obtenida es luego removida del dato original, obteniendo el residuo presentado en la [Figura 3](#) (abajo, derecha). Observamos en el mapa residual una mejoría, respecto al dato original, en la delineación de los cuerpos anómalos.

Un resultado similar al promedio de las generaciones podría ser obtenido utilizando la continuación analítica para diferentes alturas de continuación y eligiendo entre ellas la altura que produce un mapa residual adecuado ([Figura 7](#)). El método propuesto obtiene esta tendencia de manera automática, con sólo detener el entrenamiento cuando el promedio del coeficiente de correlación de las generaciones que el modelo produce se aproxima a 0.9 o un valor dado por el intérprete.

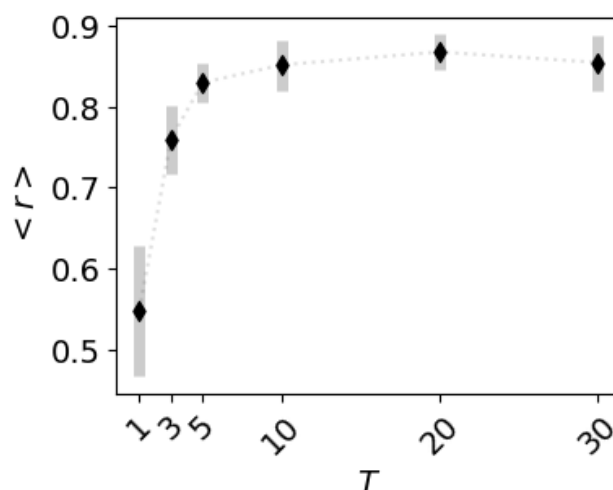


Figura 6. Elección del número de pasos generativos T . Para cada T , se calcula el promedio del coeficiente de correlación, $\langle r \rangle$, para todas las cuadrículas generadas. Alcanza con unas pocas iteraciones del proceso de generación para converger a resultados con coeficientes de correlación próximos a 0.9.

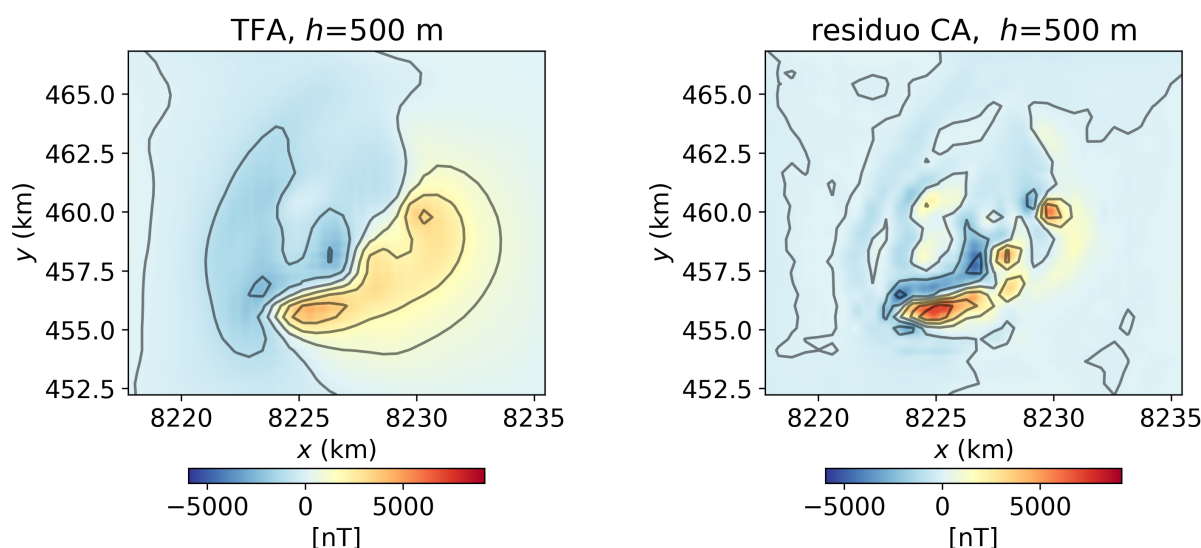


Figura 7. Izquierda: continuación analítica (CA) para $h = 500$ m. Derecha: cuadrícula residual obtenida de sustraer al dato original su continuación analítica.

Estimación de errores en la dirección de magnetización

Por lo general, se utiliza en los métodos de inversión la técnica de Monte Carlo para la estimación de errores en la determinación de los parámetros obtenidos. Esta técnica supone inyectar ruido, en lo posible “realista”, en el dato de entrada al método de inversión para estimar valores medios y desviaciones estándar de los parámetros de interés que resultan. Esta forma de proceder es utilizada para estimar la robustez de los parámetros obtenidos frente a distintos niveles de ruido en las observaciones. Ahora bien, nuestra propuesta puede ser utilizada para estimar valores medios y varianzas de contemplar el dato observado como una realización de una variable aleatoria. De esta forma, al evaluar los parámetros invertidos para cada dato generado obtenemos una medida de la distribución de los ángulos de interés en función de versiones plausibles del dato observado, y no del dato original a distintas razones de señal/ruido.

Para demostrar la técnica mencionada, construimos un dato sintético de TFA (Figura 8, izquierda) utilizando la expresión analítica de [Bhattacharyya \(1964\)](#). En esta expresión, la intensidad de la

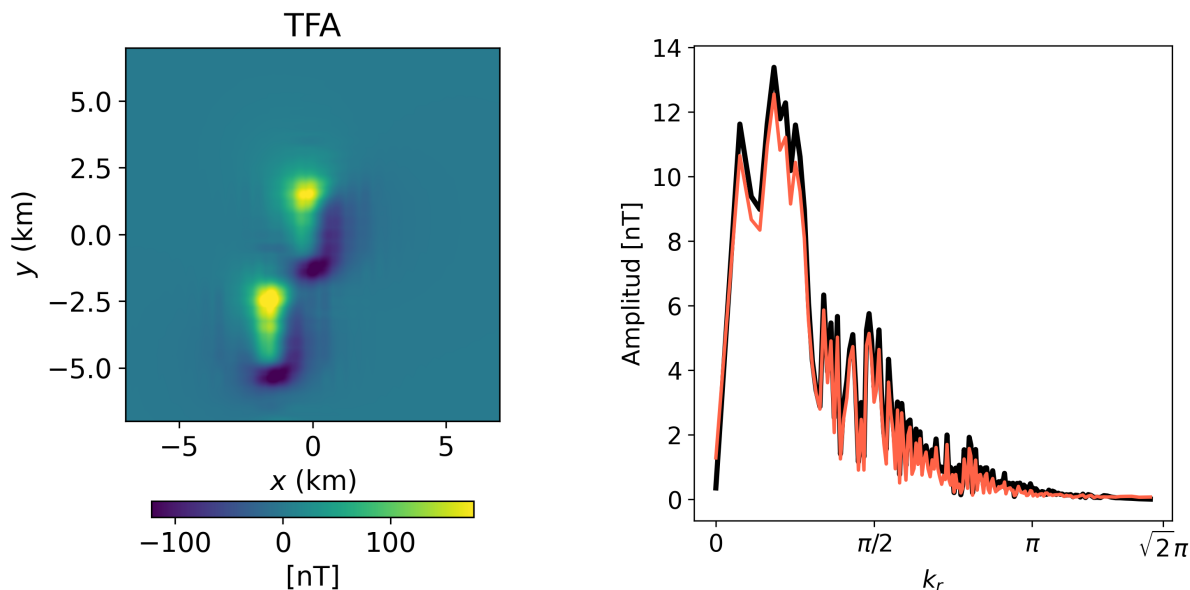


Figura 8. TFA sintética originada con dos cuerpos prismáticos. Izquierda: el cuerpo inferior tiene magnetización remanente, el cuerpo en el centro del área de trabajo tiene magnetización puramente inducida. Derecha: espectro de amplitud del dato sintético (negro) y de las generaciones (color) en función del número de onda adimensional radial.

anomalía observada viene dada por la magnitud de la magnetización en unidades de A/m. Consideramos dos prismas rectos rectangulares y con magnetización uniforme, de 1 km de espesor, 3 km de extensión en una dirección y 1 km en la otra. Un cuerpo presenta magnetización puramente inducida y otro magnetización remanente, ambos con una intensidad de polarización de 1 A/m. El campo geomagnético tiene una inclinación de -30° y una declinación de -20° . El cuerpo con magnetización remanente (en la esquina inferior derecha de la cuadrícula) presenta una inclinación de -50° y declinación de -10° . La TFA se calcula a una altura de 300 m. La Figura 8 (derecha) sintetiza la noción que los datos generados comparten las mismas características que el dato original ya que sus espectros de amplitud son muy similares en todo el rango del número de onda radial adimensional digital. Observamos que el espectro de amplitud en función del número de onda radial adimensional es similar entre la cuadrícula original y las muestras generadas.

La etapa de entrenamiento y de generación de datos sigue los mismos lineamientos que los detallados para el dato de campo. Al generar realizaciones plausibles del dato de entrada y computar los ángulos según el método Max-Min (Fedi y otros, 1994) para cada realización, se obtienen los histogramas de la Figura 9. En ellos puede observarse que los parámetros invertidos pueden mejorar su precisión respecto a los valores obtenidos del dato original. Para el ángulo de inclinación, el valor medio y su desviación estándar contiene a los valores de los parámetros utilizados en el modelo directo.

A modo de comparación, la Figura 10 presenta los histogramas para la dirección de magnetización que surgen de realizar la inversión de parámetros para distintas realizaciones de ruido aditivo Gaussiano y ruido aditivo Gaussiano coloreado, ambos con una dispersión del 10 % de la amplitud máxima de la anomalía. Para obtener coloreado, los números aleatorios son procesados con un filtro pasa bajos de Butterworth de orden 8 y frecuencia de corte en el número de onda angular adimensional radial de $\pi/2$. Estos histogramas evidencian la sensibilidad del método Max-Min frente a distintas realizaciones de ruido. Los valores obtenidos para los ángulos de los histogramas de Monte Carlo con ruido blanco distan de los valores reales del modelo directo y de los obtenidos de los histogramas que utilizan las muestras generadas. La similitud entre los histogramas de Monte Carlo con ruido coloreado y los histogramas que surgen de las muestras generadas por el método implementado fue posible gracias a la elección acertada de la amplitud del ruido y de la frecuencia de corte del filtro de paso bajo.

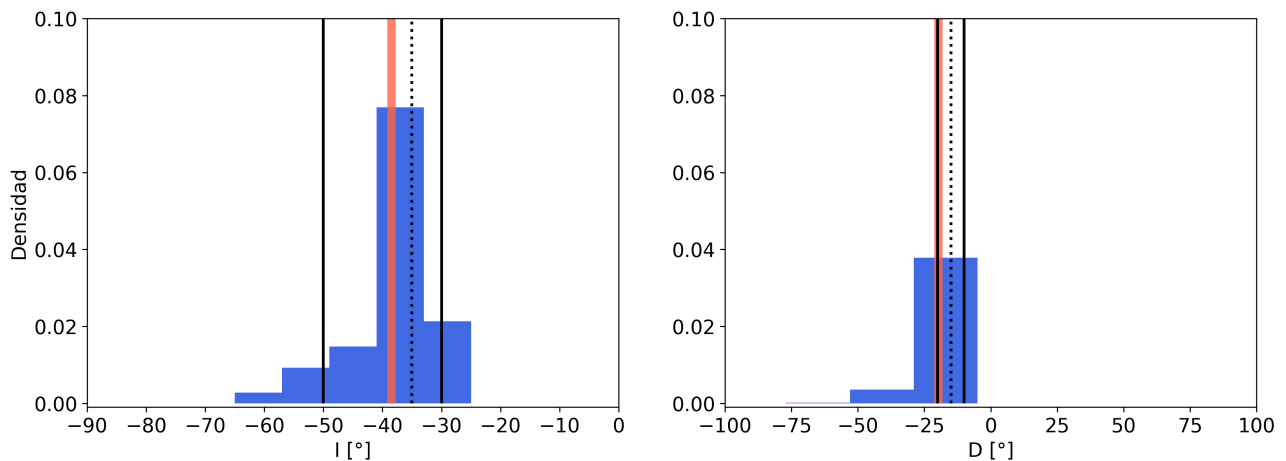


Figura 9. Histogramas de la dirección de magnetización utilizando los datos generados: Izquierda: ángulo de inclinación. Derecha: ángulo de declinación. La línea llena de color indica el valor promedio de los ángulos obtenidos respecto de los datos sintéticos generados. La línea de puntos indica el valor obtenido respecto del dato observado. La línea llena en negro indica los ángulos reales para cada prisma.

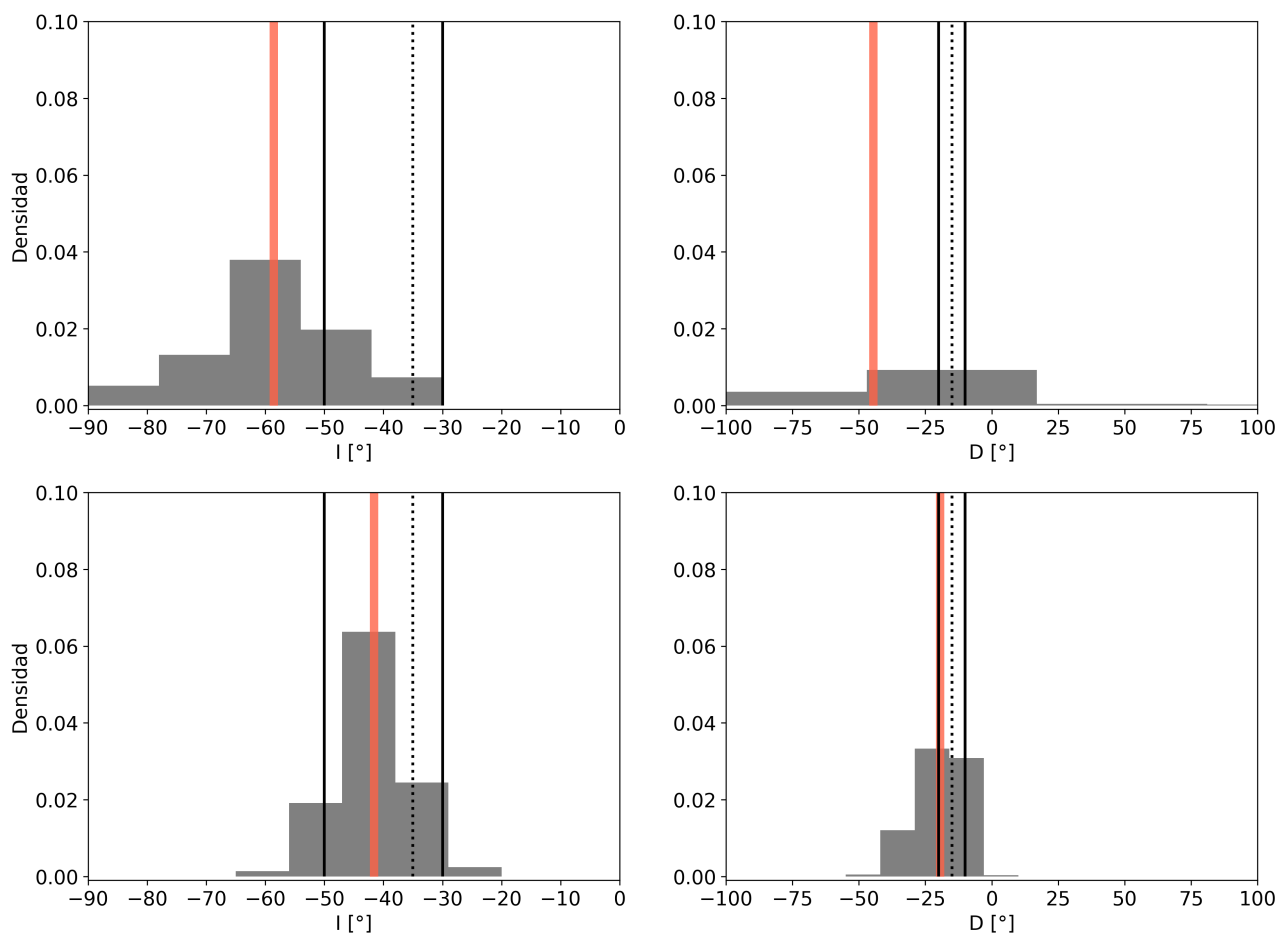


Figura 10. Histogramas de la dirección de magnetización utilizando el dato original y sumando ruido distribuido de manera normal con desviación del 10 % de la amplitud del dato. Primera fila: ruido blanco. Segunda fila: ruido coloreado por filtro pasa bajos de Butterworth. La línea llena de color indica el valor promedio de los ángulos obtenidos respecto de los datos sintéticos generados. La línea de puntos indica el valor obtenido respecto del dato observado. La línea llena en negro indica los ángulos reales para cada prisma.

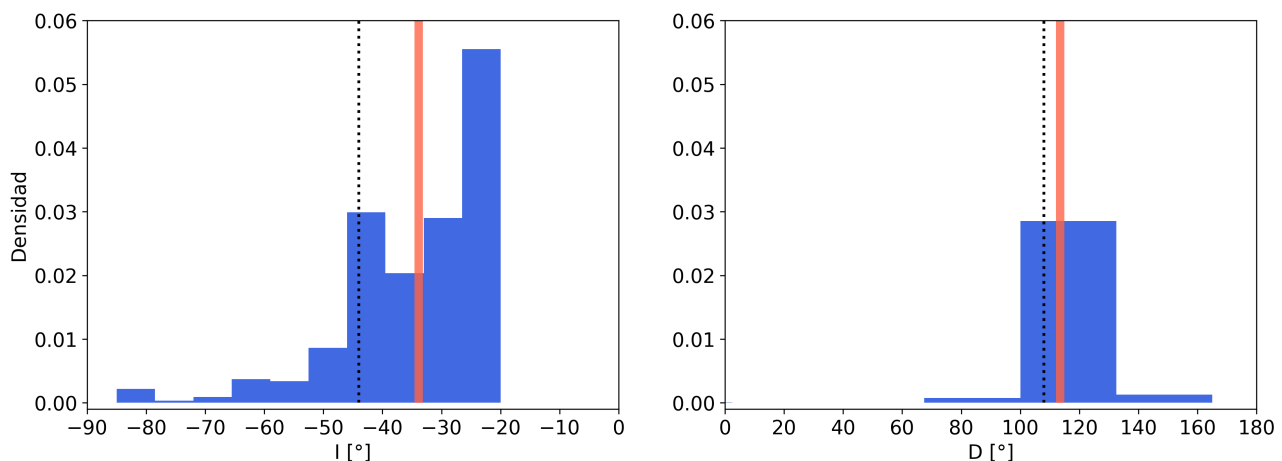


Figura 11. Histogramas de la dirección de magnetización para Montes Claros: Izquierda: ángulo de inclinación. Derecha: ángulo de declinación. La línea llena indica el valor promedio de los ángulos obtenidos respecto de los datos sintéticos generados. La línea de puntos indica el valor obtenido respecto del dato observado. El método propuesto genera una distribución de valores posibles de la dirección de magnetización y podría dar indicios de la presencia de magnetización remanente.

Por último, analizamos nuevamente el dato de campo de Montes Claros. La dirección de la magnetización para el dato real luego de ser dispuesto en una cuadrícula mediante interpolación cúbica y de ser decimado a 32×32 tiene una inclinación de -45° y una declinación de 110° . La Figura 11 presenta histogramas de los valores invertidos de los ángulos de inclinación y declinación para los distintos mapas residuales obtenidos de remover al dato original decimado el dato generado. Los ángulos obtenidos para cada inversión son promediados y sus desviaciones estándar son calculadas. Utilizando el método Max-Min obtenemos un rango de valores para la inclinación de $-40^\circ \pm 12^\circ$ y para la declinación de $113^\circ \pm 20^\circ$. A partir del dato de entrada, podemos obtener un rango posible de variación de la dirección de magnetización calculada por el método indirecto. Esta estimación no está vinculada a consideraciones respecto al nivel de ruido en los datos observados. Para el ángulo de inclinación, el carácter bimodal del histograma en la Figura 11 (izquierda) sugiere que soluciones próximas a -45° o -25° podrían también ser contempladas.

DISCUSIÓN

El promedio de las generaciones puede ser analizado como una tendencia regional ya que presenta un contenido de longitudes de onda medias. La decisión de interrumpir el entrenamiento en base al promedio del coeficiente de correlación entre los datos generados y el dato original implica una suerte de diseño de un filtro. El proceso de generación de muestras dado por la dinámica de Langevin comienza desde las longitudes de onda largas en los primeros niveles de ruido, hasta las longitudes de onda más cortas en los últimos niveles de ruido. Para un modelo que no fue entrenado hasta llegar al *overfitting*, esto implica que las generaciones obtenidas se detienen antes de poder reproducir detalles de longitudes de onda corta. Esto es utilizado para aislar, por medio del mapa residual, los contenidos de longitudes de onda corta del dato original y realzar los cuerpos anómalos.

Para este trabajo utilizamos la métrica del coeficiente de correlación para monitorear de manera práctica las imágenes generadas. Dado que el coeficiente de correlación puede ser negativo o positivo, nos permite inferir si la cuadrícula generada está invirtiendo las amplitudes de interés. Naturalmente, otras métricas pueden ser posibles, resultando aún más informativas que nuestra elección. El campo de la visión por computador (González y otros, 2002) cuenta con una multitud de métricas que pueden ser analizadas a los fines de profundizar nuestra propuesta.

Una forma de obtener un resultado similar al mapa residual es por medio de la continuación ascendente

(Blakely, 1996) con una elección apropiada de la altura de la prolongación analítica. La Figura 7 presenta la continuación ascendente (CA) y el residuo respecto al dato original para una altura de 500 m, los cuales pueden ser comparados cualitativamente con la tendencia regional y el mapa residual obtenidos por el método de aprendizaje automático. De esta manera es posible interpretar la propuesta de este trabajo en términos de una CA para una altura adecuada que debe ser definida por el intérprete. La CA es una opción con un nivel de cómputo despreciable respecto a nuestra propuesta, ya que sólo utiliza la transformada rápida de Fourier. Además de esta interpretación del promedio de las generaciones en términos de una CA, el método propuesto ofrece al intérprete la posibilidad de analizar las generaciones individualmente. Muchas de estas generaciones presentan rasgos de alta y baja frecuencia que pueden complementar el análisis de la cuadrícula original. La interpretación de los datos generados de manera individual representa una línea de trabajo a desarrollar.

Una aplicación alternativa de nuestro procedimiento podría ser la de considerar cada dato generado de manera individual y en función de su coeficiente de correlación. Por ejemplo, las cuadrículas generadas pueden ser ordenadas en función del coeficiente de correlación y luego el intérprete las utiliza para ver qué estructuras del dato de campo se destacan en ellas. Las cuadrículas generadas con un coeficiente de correlación próximo a 1 introducen una menor novedad, mientras que las cuadrículas generadas con coeficiente de correlación cercanos a cero pueden no contener información de valor, siendo más ruidosas respecto al dato de interés.

En el caso estudiado, los cuerpos anómalos de interés comparten un contenido similar de longitudes de onda. El promedio de las generaciones remueve adecuadamente el contenido de longitudes de onda largas si la decisión sobre el corte en el coeficiente de correlación es acertada. Una forma práctica de monitorear el coeficiente de correlación consiste en evaluar su histograma respecto a las generaciones producidas por un modelo a medida que transcurren las épocas. El usuario puede detener el entrenamiento del modelo cuando se observa que el coeficiente de correlación de las muestras que son generadas satisfacen un cierto valor prefijado; por ejemplo, poseer en promedio un valor cercano a 0.9.

La necesidad de decimar la cuadrícula obtenida de los registros de campo puede ser relajada con el acceso a *hardware* que permita un nivel de cómputo mayor en un tiempo adecuado. Los resultados obtenidos fueron procesados en una *notebook* personal de trabajo con una CPU CORE-i7 de Intel, demorando la etapa de entrenamiento de 10000 épocas aproximadamente 2 minutos. La generación de los datos no reviste un tiempo de espera mucho menor que el entrenamiento y depende linealmente del número de muestras generadas que el intérprete desea producir. En nuestro caso, generar 500 muestras de 32×32 píxeles demandó un tiempo de 30 segundos.

A diferencia del resto de las aplicaciones de los métodos generativos difusivos (Murphy, 2023), en este trabajo aplicamos la estrategia de aproximación de la función de *score* considerando solamente un dato de entrenamiento, que es la cuadrícula aportada por el usuario. Esta es una situación no contemplada en el método propuesto originalmente por Song y otros (2019). En la etapa de generación, aplicamos un filtro pasa bajos para acelerar la convergencia de los datos generados. Esta estrategia es un uso original de los métodos generativos para su empleo en las geociencias.

En este trabajo no realizamos generación condicionada (Saharia y otros, 2023) la cual otorga cierto control inicial sobre las cuadrículas a ser generadas. Por ejemplo, Liu y otros (2024) utilizan la generación condicionada para la interpolación de trazas en datos sísmicos 2D. Analizamos a futuro incorporar la opción de generación condicionada, utilizando para ello diferentes cuadrículas de entrenamiento obtenidas del mismo dato de campo.

CONCLUSIONES

Implementamos un método de aprendizaje automático probabilístico para obtener datos sintéticos de anomalías magnéticas. Las muestras generadas pueden considerarse versiones plausibles del dato

original, ya que comparten la misma función de densidad de probabilidad que los datos observados. En particular, el método fue aplicado para observaciones de anomalía magnética escalar de intensidad total en dos aplicaciones: la determinación de una tendencia regional y en la estimación de valores plausibles en la inversión de parámetros.

Para la estimación de la tendencia regional, se procede a promediar las realizaciones generadas. En este trabajo asignamos una medida cuantitativa del parentesco de cada generación con el dato original mediante el coeficiente de correlación. En el ejemplo de campo de Montes Claros, la tendencia regional construida a partir de los datos generados permite construir de manera automática un mapa residual que facilita la delineación de cuerpos anómalos.

En una segunda aplicación, se estima la dirección de magnetización mediante el método Max-Min. El método permite construir histogramas para los ángulos de inclinación y declinación magnética de los cuerpos anómalos. Estos histogramas resultan en valores que pueden aportar valor en la interpretación, conduciendo a valores más robustos que los obtenidos por el método tradicional de Monte Carlo. En particular, probamos con un ejemplo sintético que los valores promediados de distintas realizaciones conducen a una estimación de mayor precisión que por medio de la estrategia Monte Carlo; salvo que se elijan de manera acertada la amplitud del ruido y la frecuencia de corte para simular ruido de banda limitada. En el ejemplo de Montes Claros, el método puede proveer al intérprete de una distribución de valores posibles de la dirección de magnetización de la zona de trabajo e indicar la presencia de magnetización remanente. La presencia de magnetización remanente se deduciría simplemente de analizar si la dirección de magnetización invertida por el método Max-Min difiere de la dirección de magnetización inducida por el campo principal actual.

Como describimos en las aplicaciones, el método implementado admite una interpretación sencilla en términos de una continuación ascendente para una altura de continuación acertada en la construcción de la tendencia regional, y de una apropiada elección de ruido realista para definir rangos de variación en la inversión de parámetros. El mayor beneficio del mapa residual obtenido por el método propuesto reside en el hecho que se obtiene el mapa residual de manera automática, utilizando en el promedio mapas que son estadísticamente consistentes con el dato observado. Al utilizar la continuación analítica para el mismo fin, se requiere del intérprete la elección acertada de una altura de continuación.

Otras variantes de la metodología expuesta de aprendizaje automático probabilístico aguardan por ser desarrolladas. Entre ellas, la generación condicionada para independizar al usuario del método de interpolación asignado para disponer las observaciones en una cuadrícula rectangular y la generación en un espacio latente para la generación de cuadrículas con un mayor número de muestras.

Agradecimientos Agradecemos a los revisores anónimos y editores, a YPF Tecnología S. A., al Centro de Investigaciones Geológicas de La Plata y a la Facultad de Ciencias Astronómicas y Geofísicas de La Plata. Este trabajo fue financiado por la Universidad Nacional de La Plata y el Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET).

REFERENCIAS

- Bhattacharyya, B. K. (1964). Magnetic anomalies due to prism-shaped bodies with arbitrary polarization. *Geophysics*, 29(4), 517-531.
- Blakely, R. J. (1996). *Potential theory in gravity and magnetic applications*. Cambridge University Press.
- Burkov, A. (2019). *The hundred-page machine learning book*. Andriy Burkov.
- Chollet, F. (2021). *Deep learning with Python*. Simon and Schuster.
- de Lima, R. P., Suriamin, F., Marfurt, K. J. y Pranter, M. J. (2019). Convolutional neural networks as aid in core lithofacies classification. *Interpretation*, 7(3), SF27-SF40. <https://doi.org/10.1190/INT-2018-0245.1>
- Dutra, A. C. y Marangoni, Y. R. (2009). Gravity and magnetic 3-D inversion of Morro do Engenho complex, central Brazil. *Journal of South American Earth Sciences*, 28(2), 193-203.

- Dutra, A. C., Marangoni, Y. y Trindade, R. I. F. (2014). Aeromagnetic and physical-chemical properties of some complexes from Goiás Alkaline Province. *Brazilian Journal of Geology*, 44, 361-373. <https://doi.org/10.5327/Z2317-4889201400030003>
- Elad, M., Kwar, B. y Vaksman, G. (2023). Image denoising: The deep learning revolution and beyond – a survey paper. *SIAM Journal on Imaging Sciences*, 16(3), 1594-1654. <https://doi.org/10.1137/23M1545859>
- Fedi, M., Florio, G. y Rapolla, A. (1994). A method to estimate the total magnetization direction from a distortion analysis of magnetic anomalies. *Geophysical Prospecting*, 42(3), 261-274.
- Flusser, V. (2023). *El universo de las imágenes técnicas: elogio de la superficialidad*. Caja Negra Editora.
- González, R. C. y Woods, R. E. (2002). *Digital image processing* (2ª ed.). Prentice Hall.
- Goodfellow, I., Bengio, Y. y Courville, A. (2016). *Deep learning*. MIT press.
- Hu, Y., Wei, X., Wu, X., Sun, J., Huang, Y. y Chen, J. (2024). Three-dimensional cooperative inversion of airborne magnetic and gravity gradient data using deep-learning techniques. *Geophysics*, 89(1), WB67-WB79. <https://doi.org/10.1190/geo2023-0225.1>
- Journel, A. G. (1989). *Fundamentals of geostatistics in five lessons* (Vol. 8). American Geophysical Union.
- Kingma, D. P. y Ba, J. (2014). Adam: A Method for stochastic optimization (v. 1). <https://doi.org/10.48550/arXiv.1412.6980>
- Li, P., Liu, M., Alfarraj, M., Tahmasebi, P. y Grana, D. (2024). Probabilistic physics-informed neural network for seismic petrophysical inversion. *Geophysics*, 8(2), M17-M32. <https://doi.org/10.1190/geo2023-0214.1>
- Liu, Q. y Ma, J. (2024). Generative interpolation via a diffusion probabilistic model. *Geophysics*, 89(1), V65-V85. <https://doi.org/10.1190/geo2023-0182.1>
- McAliley, W. A. y Li, Y. (2024). Stochastic inversion of geophysical data by a conditional variational autoencoder. *Geophysics*, 89(1), WA219-WA232. <https://doi.org/10.1190/geo2023-0147.1>
- Meju, M. A. (1995). *Geophysical data analysis: Understanding inverse problem Theory and Practice*. Society of Exploration Geophysicists. <https://doi.org/10.1190/1.9781560802570>
- Mosegaard, K. y Tarantola, A. (1991). Monte Carlo analysis of geophysical inverse problems. *SEG Technical Program Expanded Abstracts*, 940-940. <https://doi.org/10.1190/1.1888770>
- Murphy, K. P. (2023). *Probabilistic machine learning: Advanced topics*. MIT press.
- Reis, A. L. A., Jr., V. C. O. y Barbosa, V. C. F. (2020). Generalized positivity constraint on magnetic equivalent layers. *Geophysics*, 85(6), J99-J110. <https://doi.org/10.1190/geo2019-0706.1>
- Ronneberger, O., Fischer, P. y Brox, T. (2015). U-Net: Convolutional networks for biomedical image segmentation. *Medical Image Computing and Computer-Assisted Intervention*, 234-241. <https://doi.org/10.48550/arXiv.1505.04597>
- Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D. J. y Norouzi, M. (2023). Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4), 4713-4726. <https://doi.org/10.1109/TPAMI.2022.3204461>
- Song, Y. y Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32. <https://doi.org/10.48550/arXiv.1907.05600>
- Song, Y. y Ermon, S. (2020). Improved techniques for training score-based generative models. *Advances in Neural Information Processing Systems*, 33, 12438-12448. <https://doi.org/10.48550/arXiv.2006.09011>
- Spector, A. y Grant, F. (1970). Statistical models for interpreting aeromagnetic data. *Geophysics*, 35(2), 293-302. <https://doi.org/10.1190/1.1440092>
- Zhang, H., Ravat, D., Marangoni, Y. R., Chen, G. y Hu, X. (2018). Improved total magnetization direction determination by correlation of the normalized source strength derivative and the reduced-to-pole fields. *Geophysics*, 83(6), J75-J85. <https://doi.org/10.1190/geo2017-0178.1>
- Zhou, H., Guo, R., Li, M., Yang, F., Xu, S. y Abubakar, A. (2024). Feature-based magnetotelluric inversion by variational autoencoder using a subdomain encoding scheme. *Geophysics*, 89(1), WA67-WA83. <https://doi.org/10.1190/geo2022-0774.1>