# ENFOQUE DE APRENDIZAJE AUTOMÁTICO SUPERVISADO PARA ESTIMAR EL ACOPLE INTERPLACA: APLICACIÓN A CHILE CENTRAL.

**Sebastián Barra Cisterna**

Tesis presentada a la Facultad de Ciencias Físicas y Matemáticas para optar
al grado de Magíster en Geofísica

Junio 2023
Concepción, Chile

**Profesor Guía: Dr. Matthew Miller**
**Guía externo: Dr. Marcos Moreno**
**Comisión evaluadora**
**Dr. Roberto Benavente - Dr. Rodofo Araya - Dra. Ignacia Calisto**

UNIVERSIDAD DE CONCEPCIÓN
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS
DEPARTAMENTO DE GEOFÍSICA.

# ENFOQUE DE APRENDIZAJE AUTOMÁTICO SUPERVISADO PARA ESTIMAR EL ACOPLE INTERPLACA: APLICACIÓN A CHILE CENTRAL.

Sebastián Barra Cisterna

Tesis presentada a la Facultad de Ciencias Físicas y Matemáticas para optar
al grado de Magíster en Geofísica

Profesor Guía: Dr. Matthew Miller
Guía externo: Dr. Marcos Moreno
Comisión evaluadora
Dr. Roberto Benavente - Dr. Rodofo Araya - Dra. Ignacia Calisto

Junio 2023, Concepción, Chile.

Dedicado a mis padres

# AGRADECIMIENTOS

Quiero agradecer a todas las personas que hicieron posible el desarrollo de esta tesis, a mis amigos, mi familia, en especial a mis padres y mi pareja. Gracias por todo su apoyo durante este tiempo y darme energías cuando lo necesitaba.

También a la comisión, los profesores Roberto, Rodolfo, Ignacia y Matt, por su ayuda y buena voluntad de integrarse a este proyecto, también a Francisco Ortega y Jonathan Bedford, por su colaboración y valiosos comentarios en el artículo y sobretodo al profesor Marcos Moreno por haberme guiado durante el desarrollo de este trabajo.

# Resumen

La distribución del grado de acople es importante para estimar el tamaño y el déficit de deslizamiento de las brechas sísmicas. Las inversiones de velocidad superficial se utilizan habitualmente para estimar el acople. Aquí presentamos un enfoque innovador para inferir el grado de acople mediante algoritmos de aprendizaje supervisado. Utilizando el margen Central Chileno como área de estudio, implementamos seis algoritmos diferentes de regresión de aprendizaje supervisado: Random Forest, Decision Tree, K-nearest neighbors (KNN), Lasso, Ridge y Linear. Estos métodos se entrenaron primero con distribuciones sintéticas de acople y luego se utilizaron para inferir el acople a partir de observaciones GNSS. Probamos el rendimiento de cada algoritmo y comparamos nuestros resultados con un método de inversión estándar. Los mejores resultados se obtuvieron con la regresión de Ridge. La distribución del acople obtenida mediante aprendizaje supervisado es coherente con los mapas de acople anteriores y proporciona un RMSE de 1,88 mm/año a las observaciones GNSS. Nuestro estudio muestra que los métodos de aprendizaje automático abren nuevas vías para mejorar las estimaciones de acople y deslizamiento de fallas.

# Abstract

Imaging locking degree at faults is important for estimating the size and slip deficit of seismic gaps. GNSS velocity inversions of varying complexity are commonly used to estimate locking. Here we present an innovative approach to infer the degree of locking from surface GNSS velocities by means of supervised learning (SL) algorithms. We implemented six different SL regression methods: Random Forest, Decision Tree, K-nearest neighbors, Lasso, Ridge and Linear and apply them in the Central Chile subduction. These methods were first trained on synthetic distributions of locking and then used to infer the locking from GNSS observations. We tested the performance of each algorithm and compared our results with a least squares inversion method. Our best results were obtained using the Ridge regression. The locking distribution is consistent with previous locking maps and gives a root mean square error (RMSE) of 1.88 mm/yr to GNSS observations. Our study shows that supervised machine learning methods open new avenues for improving the locking and fault slip estimations.

# Índice general

# Índice de cuadros

# Índice de figuras

# Capítulo 1

# Introducción

## 1.1. Introducción

Las zonas de subducción son un tipo de borde de placas convergente, donde una placa de mayor densidad subduce bajo una placa menos densa. Estas regiones son de gran interés científico debido a que concentran $\sim 90\%$ de la actividad sísmica a nivel global (Bürgmann et al., 2005). En la zona de subducción N-S chilena existe un contacto entre la Placa de Nazca y la Placa Sudamericana, conocido como megathrust (Almeida et al., 2018), en esta se han registrado numerosos terremotos, incluyendo el de mayor magnitud registrado mundialmente, Mw 9.5 Valdivia 1960, y numerosos de magnitud significativa, entre los cuales se encuentran, Mw 8.2 Iquique 2014, Mw 8.3 Illapel 2015 y Mw 8.8 Maule 2010.

En la etapa intersísmica ambas placas se encuentran acopladas y se mueven en conjunto en la dirección de la subducción (hacia el este). Este contacto va acumulando esfuerzos en un periodo de décadas a siglos, hasta que eventualmente se liberan a través de grandes deslizamientos en la dirección opuesta, en un intervalo de tiempo de segundos a pocos minutos, lo que conocemos como terremoto (periodo cosísmico). Luego de esto la Placa Sudamericana tiene un deslizamiento constante sin cambios bruscos (asísmico), en la misma dirección del cosísmico, hasta que gradualmente ocurre un reacople en el *megathrust*, a esta etapa se le conoce como periodo postsísmico. Estos tres procesos en conjunto conforman el ciclo sísmico. Es importante aclarar que durante el periodo intersísmico también ocurren procesos, como sismicidad de menor magnitud y deslizamientos asísmicos, los cuales también

liberan esfuerzos, pero en menor medida.

El grado de acople es una descripción de la cinemática del *megathrust* que cuantifica la relación entre la velocidad de deslizamiento en las fallas en el periodo intersísmico y la velocidad a largo plazo de las placas (Harris and Segall, 1987; McCaffrey et al., 2000; Mazzotti et al., 2000a; Métois et al., 2014). Muchos estudios han revelado una relación espacial entre las zonas de alto acople (asperezas) y las áreas que liberan grandes deslizamientos durante los terremotos (e.g., Chlieh et al., 2008; Moreno et al., 2010; Loveless and Meade, 2011; Lin et al., 2013), por lo que es importante tener una buena estimación del grado de acople, ya que junto con tiempo de duración del periodo intersísmico, nos permiten estimar la cantidad de déficit de deslizamiento acumulado y con esto, la magnitud posible de un próximo terremoto en caso de liberar la energía acumulada.

La distribución de las zonas de acoplamiento que se presentan durante el periodo intersísmico, no pueden ser obtenidas de manera directa a partir de datos, por lo que se obtiene de manera indirecta y suele estimarse a partir de inversiones de los desplazamientos superficiales derivados de la geodesia espacial (Bürgmann et al., 2005; Moreno et al., 2010). Se puede medir la deformación superficial a través de estaciones GNSS, las cuales han aumentado en una cantidad considerable en las últimas décadas. Estas estaciones recolectan datos de posición en las componentes, Norte-Sur, Este-Oeste y Vertical. A partir de estos datos de posición es que obtenemos la tasa de deformación en superficie producto del acoplamiento de las placas (Almeida et al., 2018; Moreno et al., 2010).

Típicamente las velocidades intersísmicas se modelan como una combinación de deformación elástica de la placa cabalgante debido al acoplamiento en la interfaz de subducción. Sin embargo, obtener el acoplamiento interplaca es un problema mal condicionado con soluciones no únicas, es decir, los datos no son suficientemente informativos para determinar un modelo de acoplamiento. Además de esto, existe una falta de resolución en las partes cercanas a la fosa, debido a que solo poseemos datos en tierra firme. Hay dos principales enfoques de inversión para tratar estos problemas. Por un lado, una enfoque de inferencia Bayesiana en donde se exploran todas las soluciones posibles, mediante el muestreo de un conjunto de modelos de la distribución de probabilidad posterior de acople (e.g., Minson et al., 2013; Jolivet et al., 2020). Por otro lado, los enfoques de mínimos cuadrados buscan un modelo óptimo de déficit de slip (acople) y requieren la definición de información a

priori en forma de un término de regularización (por ejemplo, suavizado espacial) para obtener una solución estable del problema de inversión (e.g., Harris and Segall, 1987; Ortega-Culaciati et al., 2021).

Con el aumento masivo de datos procedentes de la geodesia por satélite y los instrumentos geofísicos, los algoritmos de aprendizaje automático (ML de machine learning) se utilizan cada vez más para caracterizar la física de las fallas y extraer patrones de procesos en geociencia (Kong et al., 2018; Yáñez-Cuadra et al., 2022). En sismología, se ha utilizado para mejorar la detección de terremotos y las capacidades de selección de fase (Zhao and Takano, 1999; Liao et al., 2022) . Los terremotos de laboratorio se han predicho con éxito utilizando ML, Rouet-Leduc et al. (2017) muestra que mediante el registro de las señales acústicas emitidas por una falla de laboratorio, ML puede predecir el momento del próximo terremoto de laboratorio. Estos estudios son sólo algunos ejemplos de ML en geociencias y demuestran el rápido desarrollo y rendimiento de las herramientas de ML.

El aprendizaje supervisado (SL de supervised learning) es una técnica de ML en la que se entrena un modelo para aproximar una función que relaciona datos de entrada a datos de salida basándose en un conjunto de ejemplos. Así, los métodos supervisados intentan descubrir la relación entre los atributos de entrada y un atributo objetivo. Cuando el atributo objetivo es una variable continua (por ejemplo, un número real), se llama regresión (Rokach and Maimon, 2010). Las regresiones pueden predecir un valor numérico para un nuevo punto de datos basándose en las características de entrada. Aquí, presentamos un enfoque innovador para inferir distribuciones de grados de acople utilizando métodos de regresión SL. Utilizamos seis algoritmos de regresión diferentes: Random Forest, Decision Tree, K-nearest neighbors (KNN), Lasso, Ridge y Linear, entrenados con datos sintéticos para inferir el grado de acople a partir de datos GNSS. Para esto generamos distribuciones de acoplamiento suavizadas y obtenemos su deformación superficial asociada, esperamos que al aprender de ejemplos con distribuciones estables de acople los algoritmos no requieran de una regularización o suavizado espacial adicional. Probamos el rendimiento de cada algoritmo para predecir distribuciones sintéticas del grado de acople y luego los aplicamos a observaciones GNSS reales. También comparamos nuestros resultados con un método que es comúnmente aplicado, una inversión de mínimos cuadrados ponderados regularizados.

### 1.1.1. Área de estudio

Nuestra zona de estudio es el margen Central Chileno (fig. 1.1.1), entre los 28°S y 35°S, donde la Placa de Nazca subduce bajo la Placa Sudamericana a una velocidad de convergencia de unos 66 mm/año (Angermann et al., 1999). En esta zona ocurrieron dos recientes grandes terremotos, Illapel 2015 (Mw 8,3, Tilmann et al. (2016)) y Maule 2010 (Mw 8.8, Moreno et al. (2010)). Anterior a estos dos, ocurrió el terremoto de Valparaíso 1985 (Mw 8) el cual deslizó principalmente entre los 20km y 40km de profundidad.

La zona también ha sido escenario de cuatro importantes terremotos históricos: 1647, 1730, 1822 y 1906 (Comte et al., 1986), así como varios terremotos de magnitud M $\geq$ 8 (Comte and Pardo, 1991; Lomnitz, 2004). El terremoto más significativo fue el de 1730 (Mw 9,1-9,3) (Fig. 1.1.1), que generó rupturas prácticamente en toda la zona sismogénica del área de estudio (Carvajal et al., 2017). Desde entonces, no se ha registrado ningún terremoto de magnitud > 8,5 que haya afectado la parte superficial de la megafalla en la zona comprendida entre 32°S y 34°S. Por esta razón, esta área se considera un "gap sísmico", donde la acumulación de deformación podría desencadenar un gran terremoto en el futuro (Bravo et al., 2019). Este escenario es relevante, dado que en esta zona se concentra la mayor densidad de población del país, así como importantes ciudades costeras. Por estas razones, se ha realizado una amplia investigación en esta área, y recientemente se han publicado varios modelos de acoplamiento (por ejemplo, Becerra-Carreño et al. (2022); Sippl et al. (2021)), lo que nos permite comparar nuestros resultados con estudios anteriores.

## 1.2. Hipótesis y objetivos

Según lo planteado anteriormente, el acople intersísmico siempre es abordado como un problema de inversión, por lo que proponemos la siguiente hipótesis:

Los algoritmos de aprendizaje automático de regresión supervisada permiten estimar el nivel de acoplamiento intersísmico mediante el análisis de las velocidades superficiales registradas por las estaciones GNSS con resultados similares a la estimación por mínimos cuadrados regularizados.

### 1.2.1. Objetivo general

Estimar el acople intersísmico en Chile Central (28°S-35°S) utilizando algoritmos de SL de regresión entrenados con datos GNSS sintéticos.

### 1.2.2. Objetivos específicos

- Generar datos sintéticos de slip y velocidades GNSS.

- Entrenar modelos de SL para relacionar las velocidades GNSS con el slip y por lo tanto, el acople.

- Comparar y visualizar los resultados de los diferentes modelos para el conjunto de test.

- Utilizar los modelos entrenados con los datos reales y comparar resultados cuantitativamente entre ellos y con método de mínimos cuadrados regularizados. Comparar además visualmente con resultados anteriores y con contornos de deslizamiento cosísmico.

El capítulo siguiente es de la publicación enviada, titulada "A supervised machine learning approach for estimating plate interface locking: Application to Central Chile".

**Figura 1.1.1:** Área de estudio. Las flechas azules indican velocidades derivadas de GNSS. Los contornos negros muestran los deslizamientos cosísmicos de Maule 2010 (Mw 8.8) e Illapel 2015 (Mw 8.3) reportados por Moreno et al. (2012) y Tilmann et al. (2016) respectivamente, se muestra el segmento de 1730 (Mw 9.1-9.3).

# Capítulo 2

# A supervised machine learning approach for estimating plate interface locking: Application to central Chile.

## 2.1. Abstract

Imaging locking degree at faults is important for estimating the size and slip deficit of seismic gaps. GNSS velocity inversions of varying complexity are commonly used to estimate locking. Here we present an innovative approach to infer the degree of locking from surface GNSS velocities by means of supervised learning (SL) algorithms. We implemented six different SL regression methods and apply them in the central chilean margin. These methods were first trained on synthetic distributions of locking and then used to infer the locking from GNSS observations. We tested the performance of each algorithm and compared our results with a least squares inversion method. Our best results were obtained using Ridge regression. The locking distribution is consistent with previous locking maps and gives a root mean square error (RMSE) of 1.88 mm/yr to GNSS observations. Our study shows that supervised machine learning methods open new avenues for improving the locking and fault slip estimations.

## 2.2.  Introduction

The world's largest earthquakes occur in subduction zones. They release the stresses that have accumulated due to the plate convergence in the course of the interseismic period. Many studies have revealed first-order spatial relationship between pre-seismic highly locked zones ("asperities") and areas that release large slip during earthquakes (e.g., Chlieh et al., 2008; Moreno et al., 2010; Lin et al., 2013; Loveless and Meade, 2011). Although the interpretation of plate locking in terms of the mechanical behavior of the fault is still unclear. Locking degree is a description of the kinematics in the megathrust that quantifies the relationship between the slip velocity on faults in the interseismic period and the long-term velocity of the plates (Harris and Segall, 1987; Mazzotti et al., 2000b; McCaffrey et al., 2000; Mazzotti et al., 2000a; Métois et al., 2014) and is usually estimated from inversions of surface displacements derived from space geodesy (Bürgmann et al., 2005; Moreno et al., 2010). Locked zones accumulate slip deficit at the convergence rate, while uncoupled zones creep without accumulating elastic energy.

There are several methods to model and estimate fault locking. Typically, a linear forward model is built generating Green's functions based on elastic dislocations (e.g., Okada, 1985; Aagaard et al., 2013; Nikkhoo and Walter, 2015). Using the backslip method (Savage, 1983), the forward model relates slip deficit (locking) at the fault, with the velocity predictions at observation points at the surface of the Earth. Although the forward model is linear, the estimation of fault locking is a highly ill-posed inverse problem with non-unique solutions. There are two end-member inversion approaches to deal with these problems. On one side, a Bayesian exploration of all possible solutions, through sampling an ensemble of models from the posterior probability distribution of locking (e.g., Minson et al., 2013; Jolivet et al., 2020). On the other side, least-squares approaches search for an optimal back-slip (locking) model and require the definition of prior information in the form of a regularization term (e.g., spatial smoothing) to obtain a stable solution of the ill-posed inversion problem (e.g., Harris and Segall, 1987; Ortega-Culaciati et al., 2021). The amount of regularization needs to be determined to obtain a solution of the inverse problem, using cross-validation or some other form of model class selection technique (e.g., Akaike, 1980; Hansen and O'Leary, 1993; Craven and Wahba, 1979; Sambridge et al., 2006; Becerra-Carreño et al., 2022;

Yáñez-Cuadra et al., 2022).

Moreover, prior information based on physical concepts can be used to obtain locking maps that can reduce intrinsic problems associated with loss of resolution in parts of the fault remote from GNSS land observations. For instance, Lindsey et al. (2021) use stress constraints to account for the stress shadows induced by frictional coupling of an asperity (e.g., Hetland and Simons, 2010), thereby improving the characterization of locking near the trench. Background seismicity patterns may surround highly coupled regions that have no seismicity within them (Schurr et al., 2020). Thus, seismicity geometries together with GNSS data allow for improved estimation of the degree of locking (Sippl et al., 2021).

With the massive increase in data from satellite geodesy and geophysical instruments, Machine Learning (ML) algorithms are increasingly being used to characterize fault physics and extract patterns of processes in geoscience (Kong et al., 2018). In seismology, it has been used to improve earthquake detection and phase picking capabilities (Zhao and Takano, 1999; Liao et al., 2022) . Laboratory earthquakes have been successfully predicted using ML, Rouet-Leduc et al. (2017) shows that by recording the acoustic signals emitted by a laboratory fault, ML can predict the time of the next laboratory quake. These studies are just a few examples of ML in the geosciences and demonstrate the rapid development and performance of ML tools.

Supervised Learning (SL) is a ML technique in which a model is trained to approximate a hypothetical function that maps input data to output data based on a set of examples. Thus, supervised methods attempt to discover the relationship between input attributes and a target attribute. When the target attribute is a continuous variable (e.g., a real number), it is called regression (Rokach and Maimon, 2010). Thus, regressions can predict a numeric value for a new data point based on the input features. In this work, we present a novel approach to infer images of fault locking using SL regression methods and constrained by geodetic observations.

We test our methodology in the Central Chile portion of the subduction margin where the Nazca Plate moves underneath the South American Plate at a convergence rate of about 66 mm/yr (e.g., Angermann et al., 1999). This zone has the highest of the country's population density and important coastal cities.

Central Chile lies between the rupture zones of the Illapel 2015 (Mw 8.3 - e.g., Tilmann et al., 2016; Carrasco et al., 2019) and Maule 2010 (Mw 8.8 - e.g., Moreno et al., 2010; Lin et al., 2013) earthquakes. During the last centuries, the central chilean margin megathrust has experienced numerous M $\geq$ 8 earthquakes (Comte and Pardo, 1991; Lomnitz, 2004). The largest historical earthquake in this zone was in 1730 (Mw 9.1-9.3) (Fig. 1.1.1), which ruptured almost the entire seismogenic zone in the study area (Carvajal et al., 2017). Since then, no earthquake of magnitude > 8.5 has ruptured the shallow part of the megathrust. Therefore, the zone between Maule and Illapel ruptures is considered a seismic gap where strain accumulation could trigger a large earthquake in the future (Bravo et al., 2019). Several locking models have recently been published (e.g, Becerra-Carreño et al. (2022); Sippl et al. (2021)), allowing us to compare our results with previous studies.

## 2.3. Methods

We represent the subduction megathrust fault based on the SLAB 2.0 geometry (Hayes et al., 2018). The fault is discretized into a triangular mesh of 888 elements. We use the back-slip model (Savage, 1983) to represent inter-seismic strain accumulation due to interplate locking. We define a linear forward model relating fault slip rates (along strike and dip components) with GNSS velocities. Here, Green's functions are computed using a triangular dislocation model (Nikkhoo and Walter, 2015).

We develop a SL approach to infer images of fault locking constrained with surface GNSS velocity observations. For that purpose, we use six different regression algorithms: Random Forest(Breiman, 2001), Decision Tree (Breiman, 1984), K-Nearest Neighbors (KNN) (Cover and Hart, 1967), Linear (Hastie et al., 2009a), Ridge (Hoerl and Kennard, 1970), and Lasso (Least Absolute Shrinkage and Selection Operator) (Tibshirani, 1996). We train the regressors to predict fault locking using a set of synthetic back-slip models and their GNSS velocity predictions (As mentioned in the introduction, each back-slip estimation can be associated with fault locking by obtaining the ratio to the convergence speed). We test the performance of each algorithm for predicting synthetic back-slip distributions. For validation purposes, we compare our inferences with images

of fault locking estimated using a standard regularized weighted least squares inversion approach.

Finally, we use secular rates inferred from available GNSS observations to estimate the degree of coupling in Central Chile using both, least squares and our novel SL approach.

### 2.3.1. Locking estimation using Regularized Least-Squares Inversion

Images of fault locking degree (back-slip) can be inferred from GNSS observations by solving the General Least Squares (GLS) inversion problem

$$\min_{\mathbf{m}} \ \|\mathbf{W_x}(\mathbf{Gm} - \mathbf{d})\|_2^2 + \|\mathbf{W_h}(\mathbf{Hm} - \mathbf{h^o})\|_2^2. \tag{2.3.1}$$

The first term in the objective function of the GLS problem seeks to minimize the misfit between geodetic observations $\mathbf{d}$ and back-slip model prediction $\mathbf{Gm}$, where $\mathbf{G}$ is the Green's functions matrix relating fault back-slip $\mathbf{m}$ with the geodetic observations $\mathbf{d}$. $\mathbf{W_x}$ is a weight matrix associated with the observational and forward model prediction uncertainties. The second term is sometimes called the regularization term and introduces prior information to the inverse problem (e.g., Tarantola, 2005). Here, prior information is set in a quantity $\mathbf{h}$ linearly related to fault slip (i.e., $\mathbf{h} = \mathbf{Hm}$), where $\mathbf{H}$ is a regularization operator (e.g., discrete Laplacian), $\mathbf{h^o}$ is a reference value for $\mathbf{h}$, and $\mathbf{W_h}$ is a weight matrix related to the uncertainties in the determination of $\mathbf{h^o}$.

The solution of the GLS problem is the estimated model parameters (e.g., back-slip)

$$\widetilde{\mathbf{m}} = \widetilde{\mathbf{C}}_\mathbf{m} \left( \mathbf{G}^\intercal \mathbf{C}_\chi^{-1} \mathbf{d^{obs}} + \mathbf{H}^\intercal \mathbf{C}_\mathbf{h}^{-1} \mathbf{h^o} \right) \tag{2.3.2}$$

with covariance matrix

$$\widetilde{\mathbf{C}}_\mathbf{m} = \left( \mathbf{G}^\intercal \mathbf{C}_\chi^{-1} \mathbf{G} + \mathbf{H}^\intercal \mathbf{C}_\mathbf{h}^{-1} \mathbf{H} \right)^{-1} \tag{2.3.3}$$

representing the uncertainties of the estimated model parameters $\widetilde{\mathbf{m}}$. Equation (2.3.2) can be written in a more general form as

$$\widetilde{\mathbf{m}} = \mathbf{G}^\dagger \mathbf{d^{obs}} + \mathbf{v} \tag{2.3.4}$$

where $\mathbf{v} = \widetilde{\mathbf{C}}_{\mathbf{m}}\mathbf{H}^{\intercal}\mathbf{C}_{\mathbf{h}}^{-1}\mathbf{h^o}$ is a bias induced by prior information and $\mathbf{G}^{\dagger} = \widetilde{\mathbf{C}}_{\mathbf{m}}\mathbf{G}^{\intercal}\mathbf{C}_{\chi}^{-1}$ is the generalized inverse of $\mathbf{G}$ (e.g., Aster et al., 2013). $\mathbf{C}_{\chi}$ and $\mathbf{C}_{\mathbf{h}}$ are covariance matrix of the misfit and $\mathbf{h}$ respectively. We refer the reader to Ortega-Culaciati et al. (2021) for further details on the regularized least-squares approach for slip inversion.

In this study, we compare our SL inferences of fault locking with those obtained using a General Least Squares inversion approach, where prior information is defined using a spatially variable, sensitivity modulated Tikhonov regularization (Ortega-Culaciati et al., 2021). Here, the GLS problem becomes,

$$\min_{\mathbf{m}} \|\mathbf{W_d}(\mathbf{Gm} - \mathbf{d})\|_{\mathbf{2}}^{\mathbf{2}} + \varepsilon^{\mathbf{2}}\|\mathbf{S}^{-\frac{1}{2}}\nabla^{\mathbf{2}}\mathbf{m}\|_{\mathbf{2}}^{\mathbf{2}} \qquad (2.3.5)$$

where $\mathbf{W_d}$ the weight matrix associated with observational uncertainties, $\varepsilon^2$ is known as the regularization or damping parameter, controls the trade-off between the observational and regularization term, $\nabla^{\mathbf{2}}$ is the umbrella smoothing operator (Desbrun et al., 1999; Maerten et al., 2005), $\mathbf{S}^{-\frac{1}{2}}$ is a diagonal weight matrix whose components are $S_{ij}^{-\frac{1}{2}} = \frac{\delta_{ij}}{s_i}$, and $s_i = \frac{P_{ii}}{\max_k P_{kk}}$, with $\mathbf{P} = \mathbf{G}^{\intercal}\mathbf{C}_{\mathbf{d}}^{-1}\mathbf{G}$, is the sensitivity of fault slip. Additionally, we only allow for positive back-slip values along the up-dip direction of the fault.

## 2.3.2.   Locking estimation using ML algorithms

We trained different SL algorithms to predict the locking degree using a number of $p$ surface GNSS velocities. For that purpose, we define a training set conformed by a large number ($n$) of synthetic datasets. Each sample of the training set is formed by pairs of back-slip models (evaluated at a number of $m$ locations at the fault) and $p$ GNSS velocities predicted by the linear (physical) forward model. In a similar manner, we define a test set for evaluation purposes.

We assemble the $n$ synthetic GNSS datasets into a data matrix, commonly referred to as '$\mathbf{X}$' in ML terminology. $\mathbf{X}$ has $p$ columns, each one referring to a component of the synthetic velocity at a GNSS station, and has $n$ rows, each one referring to a different synthetic experiment of the training set. We also define a target back-slip matrix '$\mathbf{Y}$', that has $2m$ columns, each referring to a synthetic back-slip value in a given direction (dip and strike) and fault node location. In the same order as in GNSS synthetic data matrix $\mathbf{X}$, each row of $\mathbf{Y}$ refers to the corresponding

synthetic experiment. Therefore, we are solving a regression problem with a multi-output target.

We used the six regression algorithms implemented in the Scikit-learn Python library (Pedregosa et al., 2011). Each algorithm has its own hyperparameters, that control the behavior of the algorithm and must be determined externally by the user, usually using cross-validation techniques. For that purpose we use the training set.

The KNN regressor (Cover and Hart, 1967) finds a predefined number of training samples that are closest in distance to the test samples and predicts the target value of back-slip from them. The number of samples is a hyperparameter, k. The Decision Tree regressor builds a regression model in the form of a tree structure. It recursively splits the predictor variables ($\mathbf{X}$) with hyperplanes, until we get the leaf nodes. The leaf corresponds to a small region of $\mathbf{X}$ (in n-dimensional space), and we associate a real number with each leaf of the tree. Where a split is made in the logical flow of the tree, it is called a decision node (Fürnkranz, 2010). In our case, the leaf nodes represent a possible value of back-slip, and the decisions are made based on the GNSS velocity values. The trained tree makes a back-slip prediction from new GNSS velocity data by following the splits in the tree down to a leaf and returns the values in the matrix $\mathbf{Y}$. Random Forest regressor (Breiman, 2001; Geurts et al., 2006) uses an ensemble of decision trees. Here, each decision tree uses a sample from the training set with replacement, resulting in different trees. The final back-slip output is the average of all predictions.

Linear regression (Hastie et al., 2009a) solves a multivariate regression model

$$\min_{\mathbf{W},\mathbf{w^0}} \|\mathbf{Y} - \hat{\mathbf{Y}}(\mathbf{W},\mathbf{w^0})\|_{\mathbf{F}}^{\mathbf{2}} \qquad (2.3.6)$$

where $\mathbf{F}$ indicates the Frobenius norm. The loss function is the linear least squares function, that seeks to minimize the residual sum of squares between the back-slip targets in the training set ($\mathbf{Y}$) and the back-slip values predicted by the linear approximation:

$$\hat{Y}_{ik} = w_k^0 + \sum_{j=1}^{p} X_{ij} W_{jk} \qquad (2.3.7)$$

where $i \in 1,...,n$ spans the $n$ samples of the training set, $k \in 1,...,2m$ spans the $m$ back-slip values at a fault for each component, and $p$ is the number of

GNSS velocities.$\hat{\mathbf{Y}}$ is the locking prediction matrix, and the unknowns $\mathbf{W}$, $\mathbf{w^0}$ are the matrix of coefficients and the intercept vector, respectively. Ridge regression (Hoerl and Kennard, 1970; Rifkin and Lippert, 2007) is similar to the linear regression, but it also adds a regularization term in the loss function that seeks to minimize the $\ell^2$-norm of the coefficients:

$$\min_{\mathbf{W},\mathbf{w^0}} \ \|\mathbf{Y} - \hat{\mathbf{Y}}(\mathbf{W}, \mathbf{w^0})\|_{\mathbf{F}}^{\mathbf{2}} + \alpha\|\mathbf{W}\|_{\mathbf{F}}^{\mathbf{2}} \tag{2.3.8}$$

where $\alpha$ is a hyperparameter that controls the strength of the regularization. Lastly, Lasso (Tibshirani, 1996; Friedman et al., 2010) solves a regression model where the loss function is the linear least squares function and the regularization is given by the $\ell^1$-norm of the coefficients:

$$\min_{\mathbf{W},\mathbf{w^0}} \ \frac{1}{2n}\|\mathbf{Y} - \hat{\mathbf{Y}}(\mathbf{W}, \mathbf{w^0})\|_{\mathbf{F}}^{\mathbf{2}} + \alpha\|\mathbf{W}\|_{\mathbf{1}} \tag{2.3.9}$$

where $\|\mathbf{W}\|_{\mathbf{1}}$ is the summation of the absolute values of all elements of the coefficient matrix $\mathbf{W}$, and $\alpha$ is a hyperparameter that controls the strength of the regularization. In both, Ridge and Lasso, such a regularization is particularly useful for mitigating the problem of multicolinearity in linear regression (Duzan and Shariff, 2015).

In Linear, Ridge and Lasso regressions, once the model is trained (i.e., $\mathbf{w^0}$ and $\mathbf{W}$ are known), a prediction of back-slip $\mathbf{y}$ given by a vector $\mathbf{d^{obs}}$ of GNSS velocity observations, can be obtained as:

$$\mathbf{y} = \mathbf{W}^{\intercal}\mathbf{d^{obs}} + \mathbf{w^0}. \tag{2.3.10}$$

Here, we note that equation (2.3.10) has the same form of the solution of the GLS inversion problem in equation (2.3.4). Therefore, the linear regressors are effectively training $\mathbf{W}$ as the transpose of the generalized inverse of a Least-Squares problem. Here, one might be tempted to write that $\mathbf{W}^{\intercal}$ is the same generalized inverse of the GLS problem (2.3.1). However, such an identity is not necessarily true, as prior information in GLS inversion and Linear/Ridge/Lasso regressions are intrinsically different. In our study, when performing GLS inversion, we define prior information as a smoothing constraint on slip to deal with the ill-posed inverse problem. Instead, for Linear, Ridge and Lasso regressions prior information

is defined directly by the estimation of the coefficients $\mathbf{W}^\intercal$ in the training process.

### 2.3.3.  Synthetic Data Generation

Synthetic GNSS horizontal velocities are calculated at the location of the GNSS stations with available secular rates observations.

We generated n = 3999 different back-slip distributions to simulate the presence of different shapes of asperities to train the ML models. We restrict the back-slip distributions to not exceed the convergence speed (66 mm/yr).

The dip-slip and strike-slip within the asperities of the forward models have variability, ranging from half to all of the convergence speed (high locking), independently in both components. Consequently, the rake angle varies, except for the minimum and maximum back slip values within the asperities, where it is an average rake angle for the study area. Outside the asperities, we imposed a back-slip range between 0 and one-fifth of the convergence speed (low locking), before applying any smoothing.

We set up distributions of 1 to 4 asperities whose locations were randomly chosen. The size of the asperities varies between 6 different possibilities, from a size of 80 km latitude x 20 km depth to large asperities of 250 km x 60 km, also chosen randomly. A Laplacian smoothing was applied to the asperity distribution, varying between 3 levels (high, medium and low smoothing). The smallest size of the range is intended to start from an asperity that has a visual impact on the surface deformation and does not fade due to smoothing, the largest size is the approximate distance from the southern part of the Illapel segment to northern part of the Maule segment. The maximum number of asperities is chosen based on locking distributions from previous work.

We generated the synthetic GNSS velocities data by running forward models and adding random Gaussian noise with different standard deviations: 0, 0.5, 1, and 2 mm/yr. In this way we generate 4 synthetic data sets, one for each error level. This methodology resulted in a comprehensive training dataset that represents a diverse interseismic slip behavior and their associated superficial velocity measurements. Figure 2.3.1 shows locking images and GNSS velocities of 8 samples of the generated dataset.

**Figure 2.3.1:** Locking degree of 8 synthetic models and their GNSS velocities (blue arrows).

We separated the dataset from the synthetic experiments into training and test data sets, with 70% and 30% data respectively. Then, to select the best hyperparameters for each algorithm, we used grid search 5-fold cross-validation on the training dataset. Next, we evaluate the model on the test data set by comparing the predicted with the target back-slip. Using the back-slip estimation, we ran a forward model to obtain the predicted velocities and compared them to the test velocities data. We tried using common preprocessing methods such as standardization and normalization, but they did not improve the quality of the test results. We believe this is because each feature is in the same unit and all represent GNSS velocities.

To ensure consistency with the physical framework, we truncated the predicted

back-slip rate to be less than the convergence rate. This affected the Linear, Ridge, and Lasso regression algorithms, which occasionally predicted higher slip values in some areas, at the cost of increasing the error in GNSS velocities predictions.

We performed locking inversion using a standard least squares method with regularization to compare the results obtained with ML algorithms.

Finally, the algorithms were tested using 135 GNSS horizontal observations. We analyzed the same compilation of GNSS velocities used by Sippl et al. (2021) and Becerra-Carreño et al. (2022). These data represent the velocity field in the years prior to 2010 (compiled by Métois et al. (2016) based on Brooks et al. (2003); Klotz et al. (2001); Vigny et al. (2009))(Fig. 1.1.1). We chose to use this data because we do not want to include the postseismic effects of the earthquakes that occurred between 2010 and 2015, and to compare our lock distribution to those previously published.

## 2.4. Results

In this section, we show the performance of the benchmarked algorithms on the test set (back-slip), synthetic GNSS velocity data, and real GNSS velocity data. The metrics used to compare the results are the Root Mean Squared Error (RMSE), the Mean Absolute Error (MAE), and the coefficient of determination (R2), the latter for the test set only.

**Figure 2.4.1:** a) Locking distribution of: a) Synthetic model of the test set, b) Random Forest regressor, c) Decision Tree regressor, d) Lasso regression, e) Inversion rwls, f) K-Nearest Neighbors regressor, g) Linear regression, h) Ridge regression. In blue, the GNSS velocities of the synthetic model. Red vectors represent the GNSS velocities due to the estimated slip distribution. In left upper corner, the root mean square of the GNSS velocity residuals.

The performance of the different algorithms and the effect of training the data with different levels of Gaussian noise are shown in Table 2.4.1, and the selected hyperparameters for each model are shown in Table 2.4.2. For back-slip estimation in the synthetic test models, the best result was obtained with the linear model trained with no additional error data, which also provided the lowest RMSE in predicting GNSS data, as expected. The RMSE of the back-slip increases for all models trained on data with higher levels of added error. Linear algorithms are more affected, especially Linear Regression, whose RMSE increases from 1.72

mm/yr to 6.26 mm/yr. Ridge has the best estimation in the higher error data sets. In general, KNN, Decision Tree, and Random Forest have higher RMSE in back-slip estimation, but these are not as sensitive to changes in data error, varying less than 1 mm/yr the RMSE of slip compared between the 0 and 2 mm/yr error level. At the higher error levels, KNN has an RMSE close to that of the linear algorithms.

The fitting time of each algorithm varies from tens of minutes (Random Forest), a few minutes (Lasso), a few seconds (Decision Tree) to less than a second (KNN, Ridge, Linear). However, the prediction time is negligible, in the order of miliseconds, for all algorithms except Random Forest, where prediction times are two orders of magnitude larger. The fitting and prediction times are shown in Table 2.4.2.

The GNSS velocities of the synthetic experiment show a similar behaviour to the back-slip, with better performance of the linear models and the models trained with no data error. However, the KNN and Random Forest estimates are still good, and the RMSEs obtained for each are less than the added error (Table 2.4.1).

Figure 2.4.1 shows the input target locking distribution of a synthetic experiment and the resulting locking from each algorithm. The inversion and the Linear, Lasso, Ridge models localize well the 3 asperities of the synthetic model, although they show some artifacts in the shallowest part. We also note that KNN, Random Forest and Decision Tree have a higher level of smoothing than the test model and can locate the larger asperities but not the smaller one to the south of the fault, this can also be better observed in Figures 2.4.4 and 2.4.5, which show the difference between the target and predicted back-slip; the nonlinear models underestimate the slip within the asperities and overestimate it around the asperities.

**Figure 2.4.2:** Estimated coupling by the different algorithms from real data. a) Random Forest regressor, b) Decision Tree regressor, c) K-nearest neighbors, d) Linear regression, e) Ridge regression, f) Lasso regression. The arrows represent horizontal GPS observations (blue) and predictions from the respectively model (red). In the right corner, the RMS of residual GPS velocities is reported.

Figure 2.4.2 shows the back-slip distribution estimated from the real GNSS data

for each algorithm. As in the case of the synthetic model estimation, the main difference is the degree of smoothing between the models. The lowest RMSE in predicting the real GNSS data is obtained by Ridge Regression trained on data with an error of 0.5 mm/yr. It has an RMSE of 1.88 mm/yr, slightly better than Linear and Lasso, and better than KNN, Decision Tree, and Random Forest (Table 2.4.1). Because of this, we will use Ridge as our best model. The back-slip distributions obtained by Ridge, Linear, and Lasso are similar, with minor differences. In contrast, the linear and Ridge models trained on data with no added noise have the highest RMSE of 13.0 and 13.44 mm/yr, respectively.

The Random Forest distribution has high smoothing and no fully locked asperities. It produces a lower mean degree of locking than that estimated by the other algorithms, but it fits the GNSS data better than the Decision Tree estimation. KNN has the best performance among the nonlinear algorithms, managing to detect an asperity between 31°S and 32°S, which coincides with the rupture zone of the Illapel earthquake. KNN also has high smoothing compared to the linear models.

In all algorithms, we can see that the data from the southernmost GNSS stations are not well fitted (predictions are smaller), probably due to their proximity to the edge of our mesh and the small amount of observations compared to the amount of GNSS sites in the northern region.

**Figure 2.4.3:** a) Ridge regression and rwls Inversion slip deficit ratio. In blue, the residuals of the GPS velocities, and the Root mean square of these residuals, white contours show coseismic slip according to Tilmann et al. (2016). b) Boxplots of the residuals of the different algorithms and the rlws inversion in both components.

**Figure 2.4.4:** a) Dip-slip distribution of a synthetic model of the test set (target), b)-h) Difference of the magnitude of Dip-slip between the target and the estimation of each algorithm. b) Random Forest regressor, c) Decision Tree regressor, d) Lasso regression, e) Inversion rwls, f) K Near Neighbors, g) Linear regression, h) Ridge regression. In blue, the GNSS velocities of the synthetic model. In left upper corner, the root mean square of the Dip-slip residuals.

**Figure 2.4.5:** a) Strike-slip distribution of a synthetic model of the test set (target), b)-h) Difference of the magnitude of strike-slip between the target and the estimation of each algorithm. b) Random Forest regressor, c) Decision Tree regressor, d) Lasso regression, e) Inversion rwls, f) K Near Neighbors, g) Linear regression, h) Ridge regression. In blue, the GNSS velocities of the synthetic model. In left upper corner, the root mean square of the strike-slip residuals.

**Table 2.4.1:** Metrics of each algorithm for slip deficit prediction and GNSS velocities.

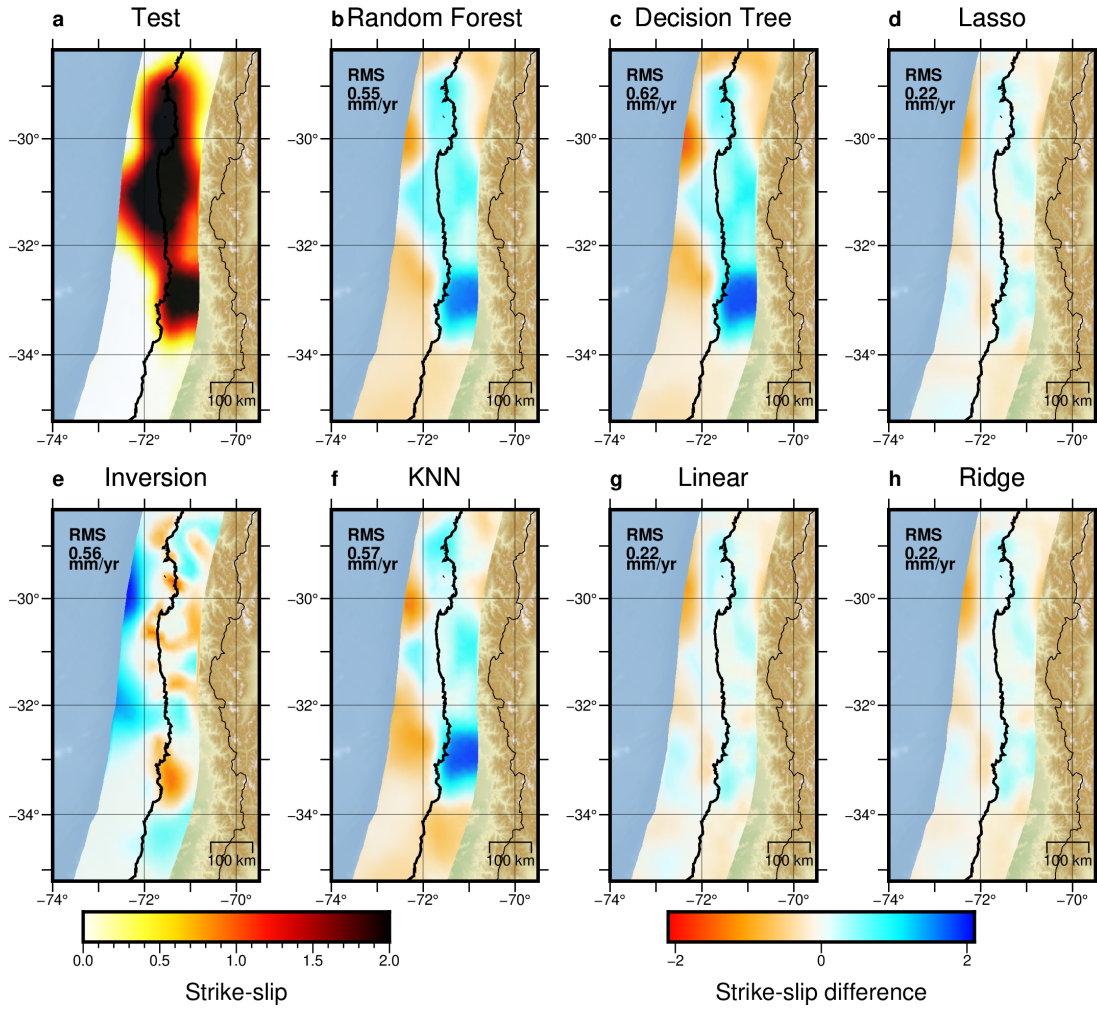| Noise level 0 mm/yr | Slip deficit rate mm/yr | | | GNSS velocities mm/yr | | Real GNSS velocities mm/yr | |
|---|---|---|---|---|---|---|---|
| | MAE | RMSE | R2 | MAE | RMSE | MAE | RMSE |
| KNN | 4.03 | 7.74 | 0.65 | 0.60 | 1.02 | 2.20 | 2.84 |
| Decision Tree | 5.71 | 9.79 | 0.45 | 1.31 | 2.13 | 3.04 | 4.75 |
| Random Forest | 4.21 | 7.35 | 0.68 | 0.75 | 1.26 | 3.01 | 4.11 |
| Linear | 0.80 | 1.72 | 0.98 | 0.01 | 0.04 | 10.58 | 13.00 |
| Ridge | 0.87 | 1.89 | 0.98 | 0.01 | 0.05 | 10.68 | 13.44 |
| Lasso | 1.33 | 2.77 | 0.95 | 0.03 | 0.07 | 2.25 | 2.91 |
| **Noise level 0.5 mm/yr** | | | | | | | |
| KNN | 4.01 | 7.76 | 0.65 | 0.78 | 1.14 | 2.20 | 2.83 |
| Decision Tree | 5.89 | 10.09 | 0.42 | 1.45 | 2.26 | 2.76 | 3.44 |
| Random Forest | 4.50 | 7.60 | 0.66 | 0.95 | 1.41 | 2.36 | 3.05 |
| Linear | 2.60 | 4.16 | 0.89 | 0.37 | 0.47 | 1.44 | 1.90 |
| Ridge | 2.45 | 4.07 | 0.89 | 0.38 | 0.48 | 1.42 | 1.88 |
| Lasso | 2.50 | 4.08 | 0.89 | 0.38 | 0.47 | 1.52 | 1.96 |
| **Noise level 1 mm/yr** | | | | | | | |
| KNN | 4.09 | 7.82 | 0.65 | 1.07 | 1.43 | 2.20 | 2.84 |
| Decision Tree | 6.02 | 10.19 | 0.41 | 1.69 | 2.44 | 3.32 | 4.39 |
| Random Forest | 4.70 | 7.83 | 0.63 | 1.24 | 1.68 | 2.35 | 3.08 |
| Linear | 3.21 | 4.99 | 0.84 | 0.75 | 0.95 | 1.44 | 1.91 |
| Ridge | 3.09 | 4.91 | 0.85 | 0.76 | 0.96 | 1.45 | 1.93 |
| Lasso | 3.07 | 4.87 | 0.85 | 0.76 | 0.95 | 1.53 | 2.00 |
| **Noise level 2 mm/yr** | | | | | | | |
| KNN | 4.35 | 8.08 | 0.62 | 1.77 | 2.25 | 2.24 | 2.91 |
| Decision Tree | 6.30 | 10.53 | 0.37 | 2.27 | 3.02 | 3.10 | 4.44 |
| Random Forest | 5.16 | 8.41 | 0.58 | 1.92 | 2.46 | 2.37 | 3.10 |
| Linear | 4.12 | 6.26 | 0.76 | 1.53 | 1.92 | 1.58 | 2.05 |
| Ridge | 3.93 | 6.13 | 0.77 | 1.55 | 1.94 | 1.61 | 2.09 |
| Lasso | 3.99 | 6.14 | 0.77 | 1.53 | 1.93 | 1.60 | 2.08 |

## 2.5.   Discussion and conclusion

The main objective of this study is to infer interseismic locking at the Chilean subduction megathrust contrained by surface GNSS velocities. For that purpose, we develop novel approaches based on SL regression techniques. The locking

**Table 2.4.2:** Selected hyperparameters and the fitting and prediction time of each model.

| Noise level 0 mm/yr | Hyperparameters | Fit time [s] | Prediction time [s] |
|---|---|---|---|
| KNN | n_neighbors=7 | 0.02 | 5e-3 |
| Ridge | alpha=1e-07 | 0.09 | 4e-3 |
| Decision Tree | max_depth=45, max_features=200, min_samples_leaf=10 | 16.36 | 2e-3 |
| Random Forest | max_depth=35, max_features=200, n_estimators=900,n_jobs=10 | 924.8 | 0.104 |
| Linear | - | 0.41 | 4e-3 |
| Lasso | alpha=6.158e-6 | 702.79 | 7e-3 |
| **Noise level 0.5 mm/yr** | | | |
| KNN | n_neighbors=6 | 0.02 | 5e-3 |
| Ridge | alpha=500 | 0.12 | 2e-3 |
| Decision Tree | max_depth=45, max_features=200, min_samples_leaf=10 | 15.95 | 2e-3 |
| Random Forest | max_depth=35, max_features=200, n_estimators=900, n_jobs=10 | 1029.0 | 0.104 |
| Linear | - | 0.40 | 2e-3 |
| Lasso | alpha=2.636e-5 | 820.31 | 6e-3 |
| **Noise level 1 mm/yr** | | | |
| KNN | n_neighbors=7 | 0.02 | 4e-3 |
| Ridge | alpha=1000 | 0.1 | 3e-3 |
| Decision Tree | max_depth=15, max_features=200, min_samples_leaf=13 | 14.85 | 2e-3 |
| Random Forest | max_depth=35, max_features=200, n_estimators=500,n_jobs=10 | 1864.41 | 0.205 |
| Linear | - | 0.43 | 4e-3 |
| Lasso | alpha=1.128e-4 | 583.12 | 6e-3 |
| **Noise level 2 mm/yr** | | | |
| KNN | n_neighbors=8 | 0.02 | 6e-3 |
| Ridge | alpha=5000 | 0.11 | 4e-3 |
| Decision Tree | max_depth=30, max_features=100, min_samples_leaf=16 | 7.28 | 2e-3 |
| Random Forest| | max_depth=25, max_features=200, n_estimators=500,n_jobs=10 | 1842.55 | 0.205 |
| Linear | - | 0.39 | 4e-3 |
| Lasso | alpha=1.128e-4 | 757.15 | 6e-3 |

distributions for Central Chile based on the Ridge, Lasso, and Linear regression algorithms are consistent with the results of the regularized least squares inversion and have a good fit to the GNSS velocities. KNN, Decision Tree and Random Forest have a good performance on the synthetic data, but a relatively large residual predicting the real data. The box plots in Figure 2.4.3 show that the residuals have a larger variance and a more asymmetric distribution than the results of the inversion and linear algorithms. In general, the locking distributions show smoothed patterns in both the synthetic and real cases (Fig. 2.4.1 and Fig. 2.4.2).

We believe that one element that benefited the training of linear algorithms over the others is that the way we generate GNSS synthetic data is through a linear combination of smooth basis functions of back-slip. Using methods that better represent the physics of the problem to generate the synthetic data, such as a more complex nonlinear model, could change which algorithm is best suited to solve the problem, in which case linear algorithms could no longer be appropriate.

Although the locking patterns estimated by the Ridge model and the inversion are similar, there are some differences (Fig. 2.4.3). In the northwest region we see that Ridge regression predicts less coupling in the shallow zone of the subduction megathrust, compared to those obtained with the least-squares inversion. For the central asperities, Ridge regression predicts the smoothest locking distribution. For both, synthetic and real GNSS velocities, Linear, Ridge and Lasso predict smooth distributions of fault locking. Therefore, although not explicitly given in Linear Regression, the training process based on smooth synthetic back-slip scenarios provides prior information that allows for the stabilization of the ill-posed slip inversion problem. Here, we believe that $\mathbf{W}$ is trained to predict smooth models because the locking training set is smooth. Thus, not learning the anticorrelation of neighboring fault slip parameters that make inversion approaches unstable (see TextS2 in Ortega-Culaciati et al., 2021). These results make our novel approach a very practical one, as suggest that more complex prior constraints could be defined through an ad-hoc locking distributions in the training set. Thus, being an advantage over classical linear least-square inversions (e.g., GLS) where only a few forms of prior information are feasible.

First-order similarities are found between ML-estimated and published coupling model patterns (Sippl et al., 2021; Becerra-Carreño et al., 2022). Previous estimates

of the degree of locking (Tilmann et al., 2016; Sippl et al., 2021; Becerra-Carreño et al., 2022) show that the rupture zones of the Illapel and Maule earthquakes were locked, similar to our results (Fig. 2.4.3). The locking distribution of Sippl et al. (2021) is smoother than our results, showing a wider zone between 31ºS and 34ºS with a locking degree around 0.75. In this zone, the Ridge algorithm detects 3 locked asperities. Low locking is found around 32°S in both the Sippl et al. (2021) and our model (in the shallow area). Becerra-Carreño et al. (2022) detected two of the asperities found in this study; one around the 2015 Illapel earthquake and another further south, at 33°S, in the seismic gap area. Lastly, the coupling estimated by Ridge shows a small asperity below 32°S, a zone where there is low coupling in the Becerra-Carreño et al. (2022) model.

Our results demonstrate the ability of SL to estimate the locking degree and slip in the megathrust based on GNSS data. Although we used classical ML techniques in this study, further studies using neural networks could improve the estimates of locking degree and slip on faults. The versatility of SL model training allows for the inclusion of a larger number of physical parameters in slip estimation. This would allow the inclusion of variations in elastic and viscoelastic properties, and even other sources of information such as seismic or geological information. In addition, the models showed almost instantaneous speed in making fault locking predictions from the real data, which may be useful when applied to rapid coseismic slip estimation.

In conclusion, our results show that SL algorithms can be used to estimate fault locking models. Thus, contributing with a novel tool that can be used to improve our understanding of the mechanical behavior of faults and hence to enhancing earthquake hazard assessments. Although, with our approach we obtain a good fit with the real GNSS data, further investigation is needed to develop a more robust model, for instance by incorporating the uncertainties of the physical model used to train the SL regression.

# Capítulo 3

# Discusión y Conclusión

El principal objetivo de este estudio es inferir el acople intersísmico en el *megathrust* de subducción de chile central restringido por las velocidades del GNSS de superficie. Para ello, desarrollamos un enfoque novedoso basado en técnicas de regresión SL, utilizando seis diferentes algoritmos. Podemos decir que el objetivo principal fue logrado ya que las distribuciones de acople para Chile Central basadas en los algoritmos de regresión Ridge, Lasso y Lineal son consistentes con los resultados de la inversión de mínimos cuadrados regularizados (Tabla 3.0.1) y tienen un buen ajuste a las velocidades GNSS, cumpliendo la hipótesis propuesta.

Por otro lado, en el caso de Decision Tree, Random Forest y KNN, si bien logran realizar un buen ajuste en los datos sintéticos (siendo KNN el mejor entre estos) no poseen una estimación de acople similar a la obtenida por mínimos cuadrados (Tabla 3.0.1) y tienen un RMSE mayor al ajustar los datos GNSS reales (Tabla 2.4.1), por lo que estos algoritmos no cumplen la hipótesis propuesta.

Tanto para velocidades GNSS sintéticas como reales, los modelos Linear, Ridge y Lasso predicen distribuciones estables de acople de fallas. Por lo tanto, aunque no

**Cuadro 3.0.1:** Error absoluto medio (MAE) entre las estimaciones de acople de cada modelo de regresión y la obtenida a través de mínimos cuadrados regularizados.

| Model noise level 0.5 mm/yr | KNN | Decision Tree | Random Forest | Linear | Ridge | Lasso |
|---|---|---|---|---|---|---|
| MAE | 0.245 | 0.290 | 0.261 | 0.116 | 0.107 | 0.115 |

se da explícitamente en la Regresión Lineal, el proceso de entrenamiento basado en escenarios de backslip sintético suavizado proporciona información a priori que permite estabilizar el mal condicionado problema de inversión de deslizamiento. Aquí, creemos que $\mathbf{W}$, la matriz de coeficientes, está entrenado para predecir modelos suaves porque el conjunto de entrenamiento de acople es suave. Por esta razón, no aprende la anticorrelación de los parámetros de deslizamiento de fallas vecinas que hacen que los métodos de inversión sean inestables (ver TextS2 en Ortega-Culaciati et al., 2021).

Estos resultados hacen que nuestro enfoque sea muy práctico, ya que sugieren que restricciones a priori más complejas o realistas podrían definirse a través de distribuciones de acople en el conjunto de entrenamiento. Es decir, si logramos generar modelos directos que posean un suavizado basado en conceptos físicos, podemos integrar esta información de manera sencilla. Por lo tanto, es una ventaja sobre las inversiones lineales clásicas de mínimos cuadrados (por ejemplo, GLS) donde solo son factibles algunas formas de información previa. Además de integrar un mejor suavizado, la versatilidad del entrenamiento del modelo SL podría permitir la inclusión de un mayor número de parámetros físicos en la estimación del deslizamiento, como variaciones en las propiedades elásticas y viscoelásticas e incluso otras fuentes de información como información sísmica o geológica.

En nuestro estudio utilizamos diferentes algoritmos de ML, notamos que algoritmos relativamente sencillos y rápidos de entrenar (Linear y Ridge) obtuvieron resultados incluso mejores que Random Forest, el cual es considerablemente más caro computacionalmente. Creemos que esto se debe a la forma en que generamos datos sintéticos, ya que los datos de velocidad superficial fueron generados a través de una combinación lineal de las funciones de Green y el *backslip*. El uso de modelos que representen mejor la dinámica del problema para generar los datos sintéticos, como un modelo mecánico friccional, podría cambiar qué algoritmo es el más adecuado para resolver el problema, ya que relaciones no lineales entre la fricción y la deformación superficial no podrían ser representadas con los algoritmos lineales. En este caso el uso de algoritmos más complejos, como redes neuronales podrían ser más adecuados.

En conclusión, nuestros resultados muestran que los algoritmos SL de regresión se pueden usar para estimar modelos de acople de fallas, cumpliendo la hipótesis propuesta. Por lo tanto, contribuye con una herramienta novedosa que se puede

utilizar para mejorar nuestra comprensión del comportamiento mecánico de las fallas y, por lo tanto, para mejorar las evaluaciones de riesgo de terremotos.

# Bibliografía

Aagaard, B. T., Knepley, M. G., and Williams, C. A. (2013). A domain decomposition approach to implementing fault slip in finite-element models of quasi-static and dynamic crustal deformation. *Journal of Geophysical Research: Solid Earth*, 118(6):3059–3079.

Akaike, H. (1980). Likelihood and the Bayes procedure. *Trabajos de Estadistica Y de Investigacion Operativa*, 31:143–166.

Almeida, R., Lindsey, E. O., Bradley, K., Hubbard, J., Mallick, R., and Hill, E. M. (2018). Can the updip limit of frictional locking on megathrusts be detected geodetically? quantifying the effect of stress shadows on near-trench coupling. *Geophysical Research Letters*, 45(10):4754–4763.

Angermann, D., Klotz, J., and Reigber, C. (1999). Space-geodetic estimation of the nazca-south america euler vector. *Earth and Planetary Science Letters*, 171(3):329–334.

Aster, R. C., Borchers, B., and Thurber, C. H. (2013). *Parameter estimation and inverse problems*. Elsevier.

Becerra-Carreño, V., Crempien, J. G. F., Benavente, R., and Moreno, M. (2022). Plate-locking, uncertainty estimation and spatial correlations revealed with a bayesian model selection method: Application to the central Chile subduction zone. *Journal of Geophysical Research: Solid Earth*, 127(10):e2021JB023939.

Bravo, F., Koch, P., Riquelme, S., Fuentes, M., and Campos, J. (2019). Slip distribution of the 1985 valparaíso earthquake constrained with seismic and deformation data. *Seismological Research Letters*.

Breiman, L. (1984). *Decision Tree*, page 368. Routledge, Boston, MA.

Breiman, L. (2001). Random forests. *Machine Learning*, 45(1):5–32.

Brooks, B. A., Bevis, M., Smalley Jr., R., Kendrick, E., Manceda, R., Lauría, E., Maturana, R., and Araujo, M. (2003). Crustal motion in the southern andes (26°–36°s): Do the andes behave like a microplate? *Geochemistry, Geophysics, Geosystems*, 4(10):1085.

Bürgmann, R., Kogan, M. G., Steblov, G. M., Hilley, G., Levin, V. E., and Apel,

E. (2005). Interseismic coupling and asperity distribution along the kamchatka subduction zone. *Journal of Geophysical Research: Solid Earth*, 110(B7).

Carrasco, S., Ruiz, J. A., Contreras-Reyes, E., and Ortega-Culaciati, F. (2019). Shallow intraplate seismicity related to the illapel 2015 mw 8.4 earthquake: Implications from the seismic source. *Tectonophysics*, 766:205–218.

Carvajal, M., Cisternas, M., and Catalán, P. A. (2017). Source of the 1730 chilean earthquake from historical records: Implications for the future tsunami hazard on the coast of metropolitan chile. *Journal of Geophysical Research: Solid Earth*, 122(5):3648–3660.

Chlieh, M., Avouac, J. P., Sieh, K., Natawidjaja, D. H., and Galetzka, J. (2008). Heterogeneous coupling of the sumatran megathrust constrained by geodetic and paleogeodetic measurements. *Journal of Geophysical Research: Solid Earth*, 113(B5).

Comte, D., Eisenberg, A., Lorca, E., Pardo, M., Ponce, L., Saragoni, R., Singh, S. K., and Suárez, G. (1986). The 1985 central chile earthquake: A repeat of previous great earthquakes in the region? *Science*, 233(4762):449–453.

Comte, D. and Pardo, M. (1991). Reappraisal of great historical earthquakes in the northern chile and southern peru seismic gaps. *Nat Hazards*, 4:23–44. Physics of Earthquake Rupture Propagation.

Cover, T. and Hart, P. (1967). Nearest neighbor pattern classification. *IEEE transactions on information theory*, 13(1):21–27.

Craven, P. and Wahba, G. (1979). Smoothing noisy data with spline functions - estimating the correct degree of smoothing by the method of generalized cross-validation. *Numerische Mathematik*, 31(4):377–403.

Desbrun, M., Meyer, M., Schröder, P., and Barr, A. H. (1999). Implicit fairing of irregular meshes using diffusion and curvature flow. pages 317–324.

Duzan, H. and Shariff, N. S. B. M. (2015). Ridge regression for solving the multicollinearity problem: review of methods and models. *Journal of Applied Science*, 15(3):392–404.

Friedman, J. H., Hastie, T., and Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software*, 33(1):1–22.

Fürnkranz, J. (2010). *Decision Tree*, pages 263–267. Springer US, Boston, MA.

Geurts, P., Ernst, D., and Wehenkel, L. (2006). Extremely randomized trees. *Machine Learning*, 63:3–42.

Hansen, P. C. and O'Leary, D. P. (1993). The use of the L-curve in the regularization of discrete ill-posed problems. *SIAM Journal on Scientific Computing*, 14(6):1487–1503.

Harris, R. A. and Segall, P. (1987). Detection of a locked zone at depth on the parkfield, california, segment of the san andreas fault. *Journal of Geophysical Research*, 92(B8):7945–7962.

Hastie, T., Tibshirani, R., and Friedman, J. (2009a). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, chapter 3, pages 43–56. Springer series in statistics. Springer.

Hastie, T., Tibshirani, R., Friedman, J. H., and Friedman, J. H. (2009b). *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer.

Hayes, G. P., Moore, G. L., Portner, D. E., Hearne, M., Flamme, H., Furtney, M., and Smoczyk, G. M. (2018). Slab2, a comprehensive subduction zone geometry model. *Science*, 362(6410):58–61.

Hetland, E. A. and Simons, M. (2010). Post-seismic and interseismic fault creep II: transient creep and interseismic stress shadows on megathrusts. *Geophysical Journal International*, 181(1):99 – 112.

Hoerl, A. E. and Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67.

Jolivet, R., Simons, M., Duputel, Z., Olive, J.-A., Bhat, H. S., and Bletery, Q. (2020). Interseismic loading of subduction megathrust drives long-term uplift in northern chile. *Geophysical Research Letters*, 47(8):e2019GL085377.

Klotz, J., Khazaradze, G., Angermann, D., Reigber, C., Perdomo, R., and Cifuentes, O. (2001). Earthquake cycle dominates contemporary crustal deformation in central and southern andes. *Earth and Planetary Science Letters*, 193(3):437–446.

Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., and Gerstoft, P. (2018). Machine learning in seismology: Turning data into insights. *Seismological Research Letters*, 1(90):3–14.

Liao, W.-Y., Lee, E.-J., Chen, D.-Y., Chen, P., Mu, D., and Wu, Y.-M. (2022). Red-pan: Real-time earthquake detection and phase-picking with multitask attention network. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–11.

Lin, Y.-n. N., Sladen, A., Ortega-Culaciati, F., Simons, M., Avouac, J.-P., Fielding, E. J., Brooks, B. A., Bevis, M., Genrich, J., Rietbrock, A., Vigny, C., Smalley, R., and Socquet, A. (2013). Coseismic and postseismic slip associated with the 2010 maule earthquake, chile: Characterizing the arauco peninsula barrier effect. *Journal of Geophysical Research: Solid Earth*, 118(6):3142–3159.

Lindsey, E., Mallick, R., and Hubbard, J. e. a. (2021). Slip rate deficit and earthquake potential on shallow megathrusts. *Nat. Geosci.*, 14(5):321–326.

Lomnitz, C. (2004). Major earthquakes of chile: a historical survey 1535-1960.

*Seismological Research Letters*, 75(3):368–378. Physics of Earthquake Rupture Propagation.

Loveless, J. P. and Meade, B. J. (2011). Spatial correlation of interseismic coupling and coseismic rupture extent of the 2011 mw = 9.0 tohoku-oki earthquake. *Geophysical Research Letters*, 38(17):190–192.

Maerten, F., Resor, P., Pollard, D., and Maerten, L. (2005). Inverting for slip on three-dimensional fault surfaces using angular dislocations. *Bulletin of the Seismological Society of America*, 95(5):1654–1665.

Mazzotti, S., Le Pichon, X., Henry, P., and Miyazaki, S.-I. (2000a). Full interseismic locking of the nankai and japan-west kurile subduction zones: An analysis of uniform elastic strain accumulation in japan constrained by permanent gps. *Journal of Geophysical Research: Solid Earth*, 105(B6):13159–13177.

Mazzotti, S., Le Pichon, X., Henry, P., and Miyazaki, S.-I. (2000b). Full interseismic locking of the nankai and japan-west kurile subduction zones: An analysis of uniform elastic strain accumulation in japan constrained by permanent gps. *Journal of Geophysical Research: Solid Earth*, 105(B6):13159–13177.

McCaffrey, R., Long, M. D., Goldfinger, C., Zwick, P. C., Nabelek, J. L., Johnson, C. K., and Smith, C. (2000). Rotation and plate locking at the southern cascadia subduction zone. *Geophysical Research Letters*, 27(19):3117–3120.

Minson, S. E., Simons, M., and Beck, J. L. (2013). Bayesian inversion for finite fault earthquake source models I—theory and algorithm. *Geophysical Journal International*, 194(3):1701–1726.

Moreno, M., Rosenau, M., and Oncken, O. (2010). 2010 maule earthquake slip correlates with pre-seismic locking of andean subduction zone. *Nature*, 467(7312):198–202.

Métois, M., Vigny, C., and Socquet, A. (2016). Interseismic coupling, megathrust earthquakes and seismic swarms along the chilean subduction zone (38°–18°s) jo - pure and applied geophysics. *Geophysical Journal International*, 173(1431–1449).

Métois, M., Vigny, C., Socquet, A., Delorme, A., Morvan, S., Ortega, I., and Valderas-Bermejo, C.-M. (2014). Gps-derived interseismic coupling on the subduction and seismic hazards in the atacama region, chile. *Geophysical Journal International*, 196(644–655).

Nikkhoo, M. and Walter, T. R. (2015). Triangular dislocation: an analytical, artefact-free solution. *Geophysical Journal International*, 201(2):1119–1141.

Okada, Y. (1985). Surface deformation due to shear and tensile faults in a half-space. *Bulletin of the Seismological Society of America*, 75(4):1135–1154.

Ortega-Culaciati, F., Simons, M., Ruiz, J., Rivera, L., and Díaz-Salazar, N. (2021). An epic tikhonov regularization: Application to quasi-static fault slip inversion. *Journal of Geophysical Research: Solid Earth*, 126(7):e2020JB021141.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.

Rifkin, R. and Lippert, R. (2007). Notes on regularized least squares. Massachusetts Institute of Technology, Tech. Rep. MIT-CSAIL-TR-2007-025.

Rokach, L. and Maimon, O. (2010). *Supervised Learning*, pages 133–147. Springer US, Boston, MA.

Rouet-Leduc, B., Hulbert, C., Lubbers, N., Barros, K., Humphreys, C. J., and Johnson, P. A. (2017). Machine learning predicts laboratory earthquakes. *Geophysical Research Letters*, 44(18):9276–9282.

Sambridge, M., Gallagher, K., Jackson, A., and Rickwood, P. (2006). Trans-dimensional inverse problems, model comparison and the evidence. *Geophysical Journal International*, 167(2):528–542.

Savage, J. C. (1983). A dislocation model of strain accumulation and release at a subduction zone, journal of geophysical research: Solid earth. *Journal of Geophysical Research: Solid Earth*, 88(4984-4996):1654–1660.

Schurr, B., Moreno, M., Tréhu, A. M., Bedford, J., Kummerow, J., Li, S., and Oncken, O. (2020). Forming a mogi doughnut in the years prior to and immediately before the 2014 m8.1 iquique, northern chile, earthquake. *Geophysical Research Letters*, 47(16):e2020GL088351.

Sippl, C., Moreno, M., and Benavente, R. (2021). Microseismicity appears to outline highly coupled regions on the central chile megathrust. *Journal of Geophysical Research: Solid Earth*, 126(11):e2021JB022252.

Tarantola, A. (2005). *Inverse problem theory : and methods for model parameter estimation*. Society for industrial and applied mathematics, Philadelphia (PA). 1987 version of the book completely rewritten.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288.

Tilmann, F., Zhang, Y., Moreno, M., Saul, J., Eckelmann, F., Palo, M., Deng, Z., Babeyko, A., Chen, K., Baez, J. C., Schurr, B., Wang, R., and Dahm, T. (2016). The 2015 illapel earthquake, central chile: A type case for a characteristic earthquake? *Geophysical Research Letters*, 43(2):574–583.

Vigny, C., Rudloff, A., Ruegg, J.-C., Madariaga, R., Campos, J., and Alvarez, M. (2009). Upper plate deformation measured by GPS in the Coquimbo Gap, Chile. *Physics of the Earth and Planetary Interiors*, 175(1):86–95. Earthquakes in subduction zones: A multidisciplinary approach.

Yáñez-Cuadra, V., Ortega-Culaciati, F., Moreno, M., Tassara, A., Krumm-Nualart,

N., Ruiz, J., Maksymowicz, A., Manea, M., Manea, V. C., Geng, J., and Benavente, R. (2022). Interplate Coupling and Seismic Potential in the Atacama Seismic Gap (Chile): Dismissing a Rigid Andean Sliver. *Geophysical Research Letters*, pages 1–26.

Zhao, Y. and Takano, K. (1999). An artificial neural network approach for broadband seismic phase picking. *Bulletin of the Seismological Society of America*, 89(3):670–680.

# Apéndice A

# Apéndice

### A0.1.   Hiperparámetros y validación cruzada

Como se mencionó anteriormente, los hiperparámetros son parámetros de un algoritmo que no se van aprendiendo según se entrena el modelo, sino que son predefinidos antes de empezar el entrenamiento, por lo que su elección es importante en la performance del modelo entrenado. Para hacer esto se suele hacer una búsqueda por grilla de los hiperparámetros, probando diferentes combinaciones y escoger los que dan un mejor resultado. Sin embargo, si hacemos esta prueba utilizando siempre el mismo conjunto de entrenamiento y validación, puede que sesgue el resultado, ya que el modelo entrenado depende también de los datos que tenga el conjunto de entrenamiento. Para lidiar con este problema de selección utilizamos validación cruzada Hastie et al. (2009b).

La validación cruzada es uno de los métodos de remuestreo de datos más utilizados para estimar el error de predicción de los modelos y ajustar los hiperparámetros del modelo.

La validación cruzada por k grupos (k-fold cross-validation) es un método en el cual se separa el set de datos de entrenamiento en un número k de subconjuntos, en el caso de este estudio usamos 5. Se utilizan 4 de esos subconjuntos como datos para entrenar los algoritmos y el quinto es utilizado para validación. Luego se entrena un modelo, se realiza una predicción de los datos de validación y se obtiene una métrica para estimar que tan bueno es el modelo, en este caso escogimos el RMSE, la cual es habitual para algoritmos de regresión. Este procedimiento se

repite un número k de veces, cambiando en cada iteración el subconjunto que se utiliza para la validación, de forma que todos los subconjuntos pasan como set de entrenamiento y de validación. A continuación, se promedia el RMSE obtenido en cada iteración y esta queda como la medida para el conjunto de hiperparámetros testeados. Este método se puede aplicar conjunto a una búsqueda de grilla, por lo que luego se prosigue con el siguiente set de hiperparámetros y finalmente se escogen los hiperparámetros con menor RMSE.

## A0.2.  Hiperparámetros de Decision Tree y Random Forest

Como se mencionó, el árbol de decisión regresor (Decision Tree) construye un modelo de regresión en forma de estructura de árbol. Divide recursivamente los datos de entrenamiento en nuevos conjuntos hasta obtener los nodos hoja, que son un pequeño subconjunto de datos y tienen asociado un número real dependiendo del promedio del atributo objetivo respectivo a este subconjunto. Cuando se realiza una división en el flujo lógico del árbol, se denomina nodo de decisión (Fürnkranz, 2010). En nuestro caso, los nodos de hoja representan un posible valor de back-slip y las decisiones se toman en función de los valores de velocidad GNSS. El árbol entrenado realiza una predicción del deslizamiento hacia atrás a partir de los nuevos datos de velocidad del GNSS siguiendo las divisiones del árbol hasta una hoja y devuelve los valores en la matriz $\mathbf{Y}$.

Los árboles de decisión tienen varios hiperparámetros, en este trabajo hemos calibrado 3, max_depth, max_features y min_samples_leaf. Max_depth se refiere a la profundidad máxima del árbol, es decir llegado a ese número de decisiones los nodos (un subconjunto de datos) no se pueden seguir separando y se considera un nodo hoja. Max_features es el número de atributos que consideramos para tomar las decisiones, en este caso es el número de datos de velocidades GNSS que se van a considerar (antiintuitivamente el número máximo no siempre es mejor). Min_samples_leaf se refiere al número mínimo de muestras requeridas para estar en un nodo hoja. Un punto de división a cualquier profundidad solo se considerará si deja al menos min_samples_leaf muestras de entrenamiento en cada una de las ramas izquierda y derecha.

El regresor Random forest (Breiman, 2001; Geurts et al., 2006) utiliza un conjunto de árboles de decisión. Aquí, cada árbol de decisión utiliza un muestreo del conjunto de entrenamiento elegido aleatoriamente con reemplazo (bootstraping).

El resultado final de la predicción de Random Forest es la media de todas las predicciones.

En este trabajo se calibraron 3 hiperparámetros, max_depth, max_features y n_estimators. Los primeros 2 representan lo mismo que en los árboles de decisión, el tercero, n_estimators, controla la cantidad de árboles de decisiones en Random Forest, por lo que es uno de los hiperparámetros más importantes. Al tener un mayor número de árboles, Random forest también tiene un mayor tiempo de entrenamiento y es más caro computacionalmente.