

Introducción a la Inteligencia Artificial
Facultad de Ingeniería
Universidad de Buenos Aires



Índice

1. Terminology
2. Pipeline
3. Train-test-validation
4. Feature engineering
5. Regresión lineal



Machine Learning Terminology

- Raw vs. Tidy Data \longrightarrow raw data: los datos crudos / tidy data: data preprocesada
 - Training vs. Holdout Sets \longrightarrow crear múltiples sets para evitar overfitting (train, validación, test)
 - Baseline \longrightarrow proponer modelos sencillos para comparar. (eg: media, RL, CART)
 - Parameters vs. Hyperparameters
 - Classification vs. Regression \longrightarrow buscar una fn clasificadora / buscar una fn cont.
 - Model-Based vs. Instance-Based Learning
 - Shallow vs. Deep Learning
- Y un poco de underfitting
- entrenamiento
- checkeo del entrenamiento
- checkeo de generalización



Dataset pipeline

$\bar{\theta}$: parámetros preproc.

Acciones que generalmente se ejecutan sobre los datasets.

↑ parámetros

Obtención de datos
o synthetic dataset

Pre-procesamiento
de Missing Values

Cómputo de media,
desvío y cuantiles

Estandarización de
datos (z-score)

Ingeniería de
Features (PCA)

Data
augmentation

Split en Train,
Validation y Test

↓ parámetros

↓ $\bar{\theta}$

↓ $\bar{\theta}$

↓ $\bar{\theta}$

transform.
Escalados
Encoding
Feature selection

en total es
el paso
ante a la
extracción

Model pipeline

Pasos involucrados al entrenar un modelo de Machine Learning

Obtener el dataset
para train

Definir métricas de
evaluación y train

Calcular métricas
para modelos base

Entrenar el modelo
con el dataset train

Computar métricas
con validation

HPs
optimization

Evaluación sobre
el dataset test

D_{train}

D_{dev}

D_{test}

Dataset

train

dev

test

(80, 5, 15)

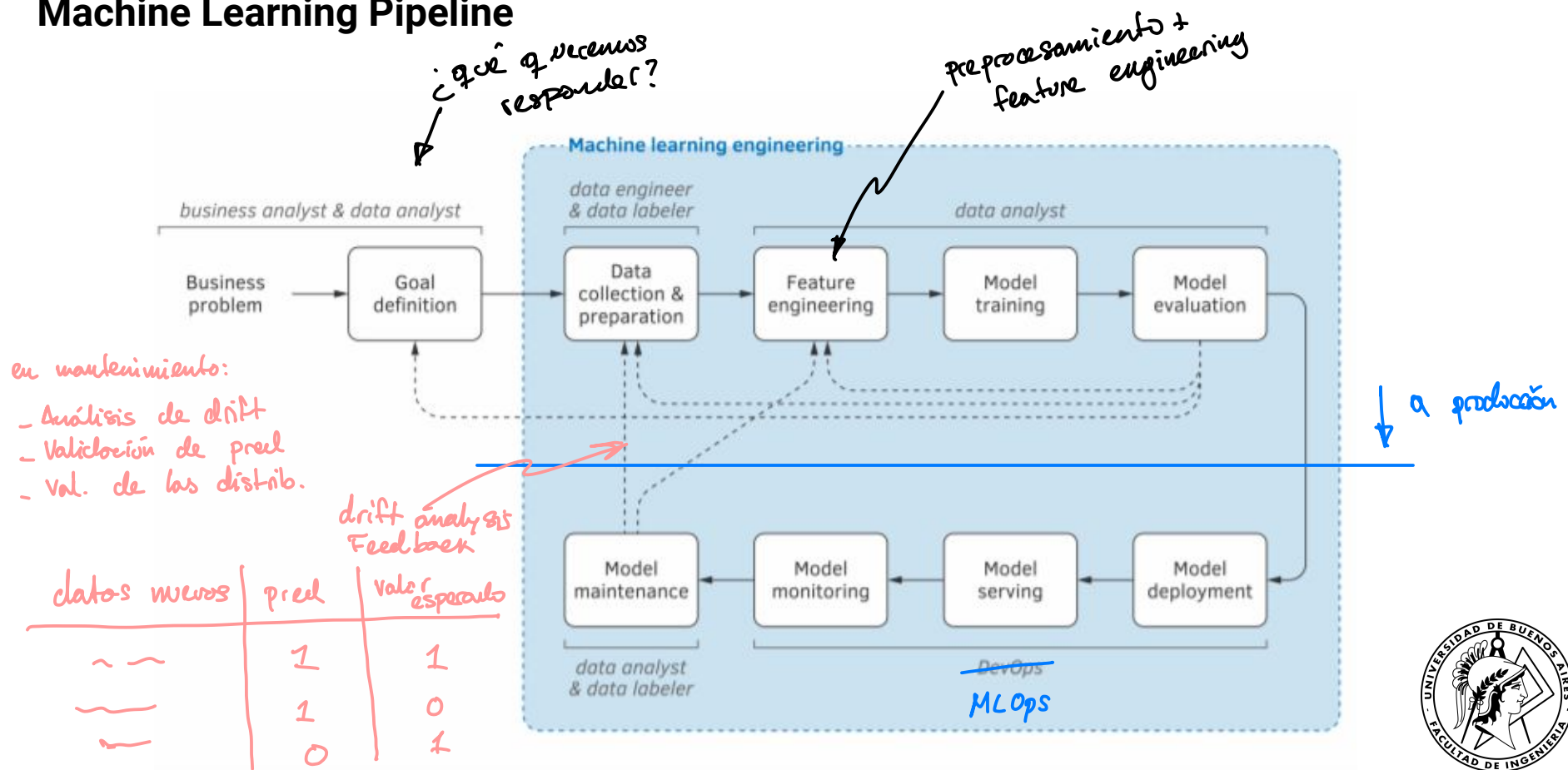
(70, 15, 15)

baseline: CART

models: [knn, LDA, KSVM] } F_1 como métrica



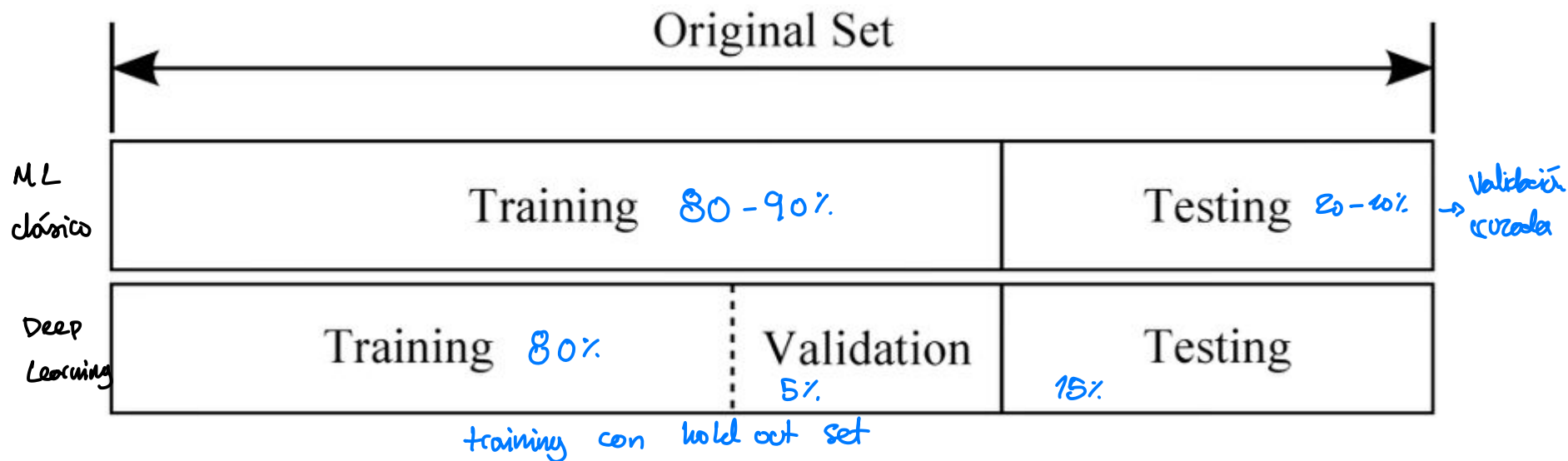
Machine Learning Pipeline



Ingeniería de Features

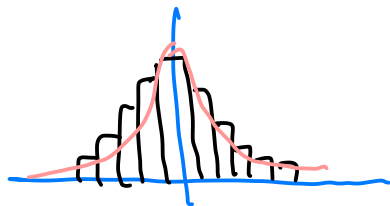
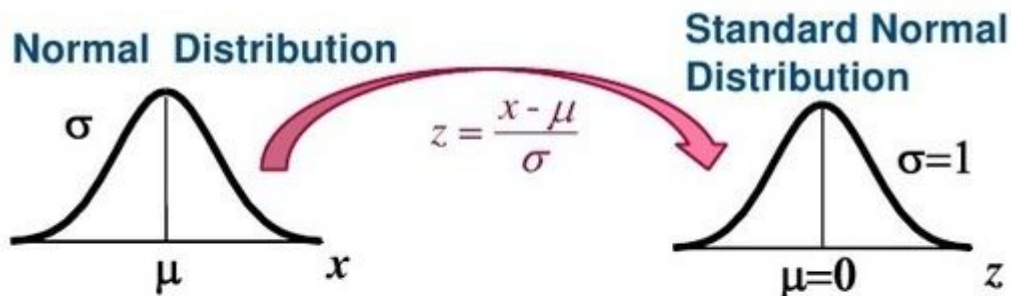
Train - test - validation

$(\bar{\theta}_{PP}, \bar{\theta}_{FE}, \bar{\phi}_M)$ → se calculan en train
y se usan en el resto.



Normalización

Muchos algoritmos de Machine Learning necesitan datos de entrada centrados y normalizados. Una normalización habitual es el z-score, que implica restarle la media y dividir por el desvío a cada feature de mi dataset.



$$\bar{x}(\hat{\mu}), s^2(\hat{\sigma}^2) \longrightarrow \hat{\mathcal{Z}} = \frac{x - \bar{x}}{s}$$

Missing Values

Es muy común en la práctica, recibir como datos de entrada, datasets que tienen información incompleta ("NaN").

ID	City	Degree	Age	Salary	Married ?
1	Lisbon	NaN	25	45,000	0
2	Berlin	Bachelor	25	NaN	1
3	Lisbon	NaN	30	NaN	1
4	Lisbon	Bachelor	30	NaN	1
5	Berlin	Bachelor	18	NaN	0
6	Lisbon	Bachelor	NaN	NaN	0
7	Berlin	Masters	30	NaN	1
8	Berlin	No Degree	NaN	NaN	0
9	Berlin	Masters	25	NaN	1
10	Madrid	Masters	25	NaN	1



Solución 1

Una forma de solucionar el problema es remover las filas y las columnas que contienen dichos valores.

ID	City	Degree	Age	Salary	Married ?
1	Lisbon	NaN	25	45,000	0
2	Berlin	Bachelor	25	NaN	1
3	Lisbon	NaN	30	NaN	1
4	Lisbon	Bachelor	30	NaN	1
5	Berlin	Bachelor	18	NaN	0
6	Lisbon	Bachelor	NaN	NaN	0
7	Berlin	Masters	30	NaN	1
8	Berlin	No Degree	NaN	NaN	0
9	Berlin	Masters	25	NaN	1
10	Madrid	Masters	25	NaN	1

¿Filas luego columnas
ó
Columnas luego filas?



Solución 2

En columnas donde el % de NaNs es relativamente bajo, es aceptable reemplazar los NaNs por la media o mediana de la columna.

Average_Age = 26.0

ID	City	Age	Married ?
1	Lisbon	25	0
2	Berlin	25	1
3	Lisbon	30	1
4	Lisbon	30	1
5	Berlin	18	0
6	Lisbon	NaN	0
7	Berlin	30	1
8	Berlin	NaN	0
9	Berlin	25	1
10	Madrid	25	1



ID	City	Age	Married ?
1	Lisbon	25	0
2	Berlin	25	1
3	Lisbon	30	1
4	Lisbon	30	1
5	Berlin	18	0
6	Lisbon	26	0
7	Berlin	30	1
8	Berlin	26	0
9	Berlin	25	1
10	Madrid	25	1

Solución avanzada

Las técnicas mencionadas producen distorsiones en la distribución conjunta del vector aleatorio. Estas distorsiones pueden ser muy considerables y afectar en gran medida el entrenamiento del modelo. Para reducir este efecto se puede utilizar **MICE (Multivariate Imputation by Chained Equation)**

1. Se trata cada columna con missing values como la variable dependiente de un problema de regresión.
2. Se van haciendo los fits de cada columna de manera secuencial.
3. Se utiliza la regresión para completar los missing values.

Ingeniería de Features

One hot encoding \swarrow one hot encoder \searrow K-1 OHE \rightarrow top k OHE

En muchos problemas de Machine Learning, puedo tener como dato de entrada variables categóricas. Por ejemplo, una columna con información sobre el color: {rojo, amarillo, azul}

Para este tipo de información, donde no existe una relación ordinal natural entre las categorías, no sería correcto asignar números a las categorías.

Una forma más expresiva de resolver el problema es utilizar "one hot encoding" y transformar la información en binaria de la siguiente manera.

Label encoding
{ Rojo: 1,
Ama: 2,
verde: 3 }

Freq encoding:
{ Rojo: 2,
Ama: 2,
Verde: 1 }

$K^{1 \times n}$

Color
Red
Red
Yellow
Green
Yellow

$B^{K \times n}$

Red	Yellow	Green
1	0	0
1	0	0
0	1	0
0	0	1

$K=30 \quad n=300$

$\rightarrow \dim(B) = 30 \times 300 \quad \# B_{ij} \neq 0 \rightarrow$ dispersion muy alta (matrix id)

one hot encoder (OHE)
(WOE)
weight of evidence

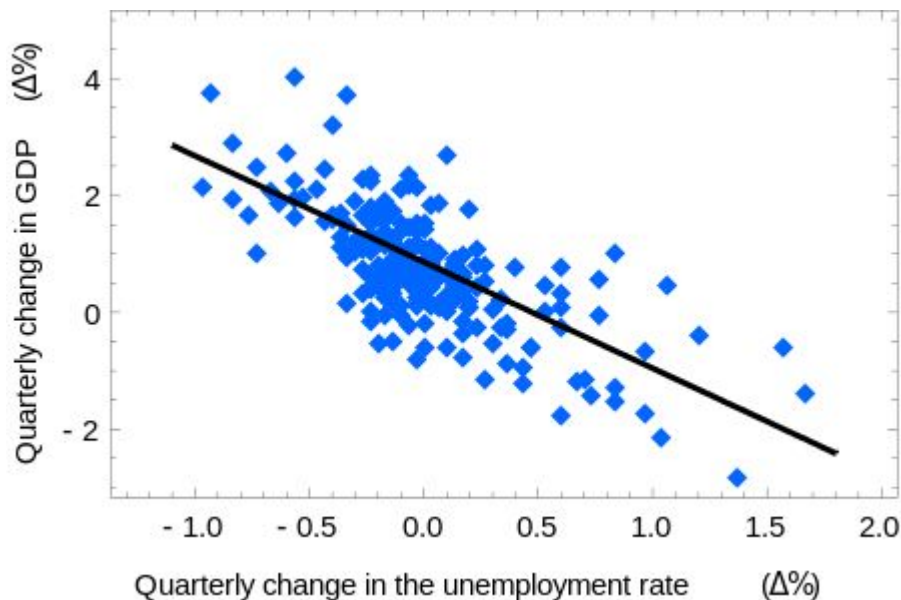
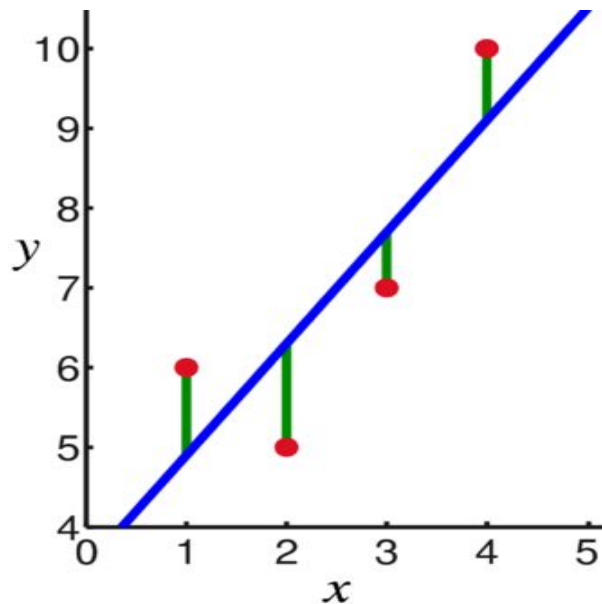
Regresión lineal

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n,$$

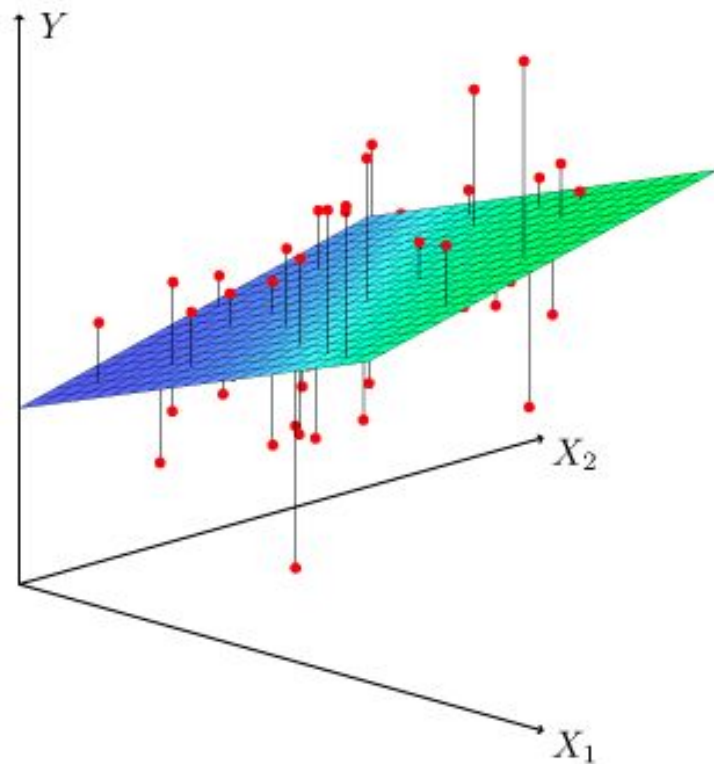


Regresión Lineal $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n,$

En ésta clase vamos a ver el framework teórico detrás de la gran mayoría de los modelos de Machine Learning: aprendizaje estadístico. Para ello, vamos a utilizar como modelo base la regresión lineal.



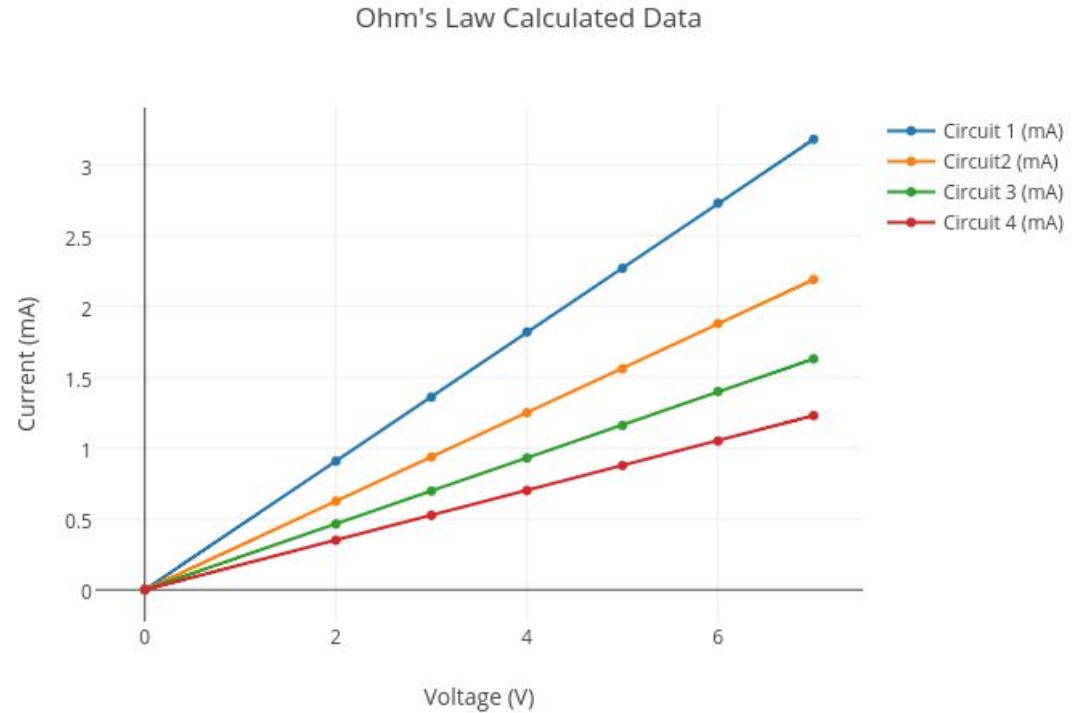
Regresión Lineal $y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} + \varepsilon_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i, \quad i = 1, \dots, n,$



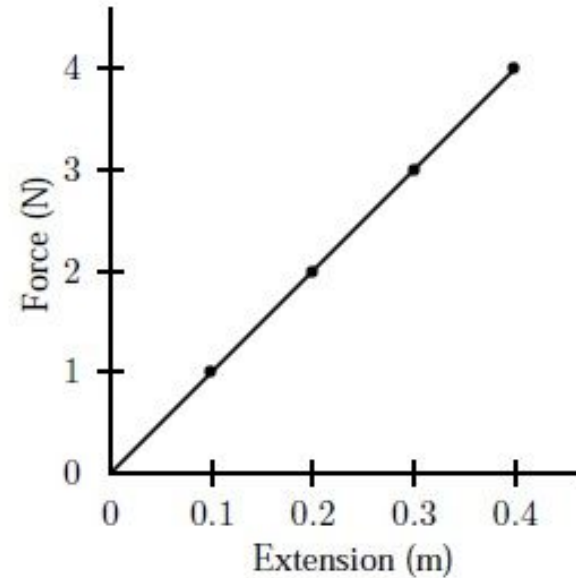
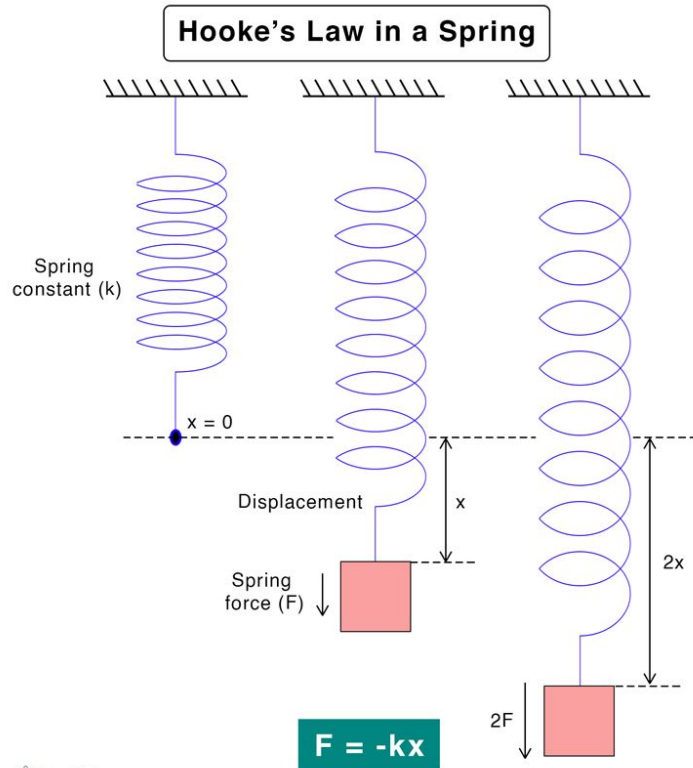
Ley de Ohm

$$I = V/R$$

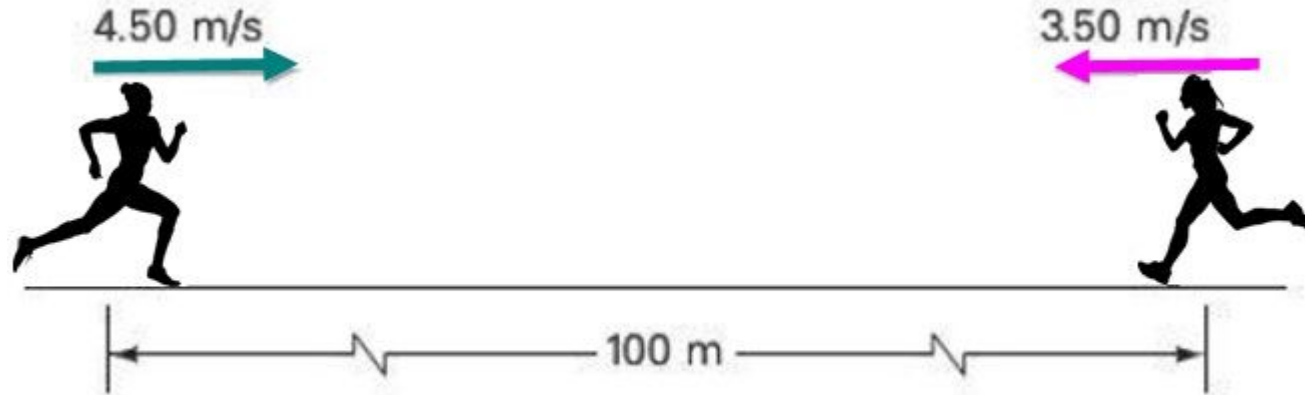
R constante



Ley de Hooke



Movimiento rectilíneo uniforme



$$x(t) = x(t_0) + V * t$$

Población de parásitos



Ejemplo: En un estudio sobre la población de un parásito se hizo un recuento de parásitos en 15 localizaciones con diversas condiciones ambientales.

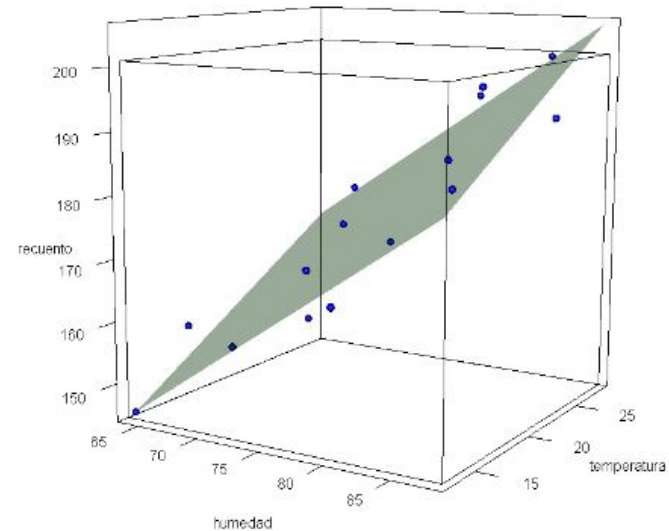
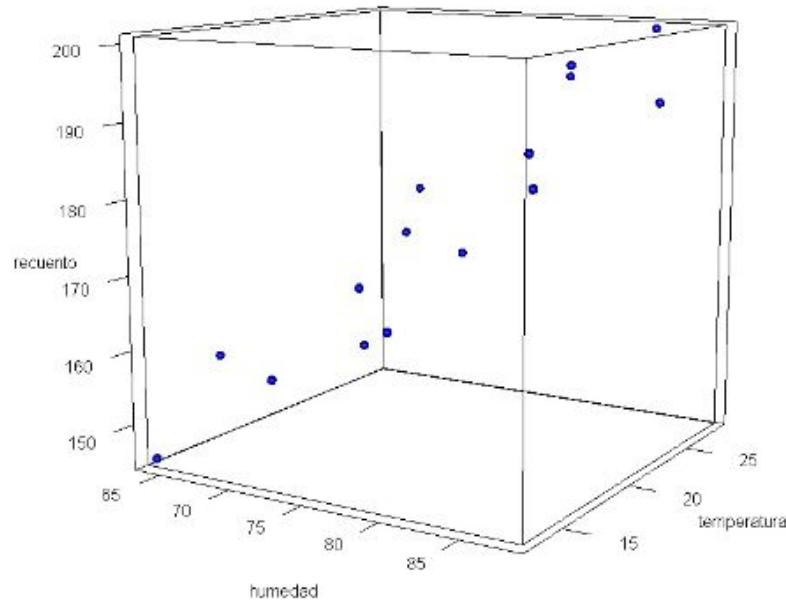
Los datos obtenidos son los siguientes:

x_1	Temperatura	15	16	24	13	21	16	22	18	20	16	28	27	13	22	23
x_2	Humedad	70	65	71	64	84	86	72	84	71	75	84	79	80	76	88
y	Recuento	156	157	177	145	197	184	172	187	157	169	200	193	167	170	192

Población de parásitos

$$\text{Recuento} = \beta_0 + \beta_1 \text{Temperatura} + \beta_2 \text{Humedad} + \epsilon$$

$$\text{Recuento} = 25.7115 + 1.5818\text{Temperatura} + 1.5424\text{Humedad}$$



Jamboard



Dado mis datos x_1, \dots, x_n con $x_i \in \mathbb{R}^D$ mediciones de mi sist y y_1, \dots, y_n con $y_i \in \mathbb{R}$ el conjunto de respuestas. llamamos a \bar{x} las variables regresivas independientes e \bar{y} variable de resp / dependiente.

En genl. buscamos encontrar la relación $y \propto \bar{x}$; $y = f(\bar{x}, \theta) + \varepsilon$

En regresión buscamos inferir $\hat{y} = \hat{f}(x)$. La prec de mi estimación de y tiene dos componentes: reducible (que depende de los datos) y otra irreducible.

Si asumimos un x fijo y \hat{f} es conocida, podemos calcular el error cuadrático medio: (entre y e \hat{y})

$$\begin{aligned} E(y - \hat{y})^2 &= E(f(x) + \varepsilon - \hat{f}(x))^2 \\ &= \underbrace{E(f(x) - \hat{f}(x))}_{\text{①}} + E(\varepsilon)^2 \end{aligned}$$

error irreducible
 $\text{Var}(\varepsilon)$

② Error reducible
(Error de estimación)

la f más sencilla es una comb. lineal \rightarrow es simple, es barata,
es explicable, \sim precisa

$$f(\bar{x}, \bar{\beta}) = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

. Supuestos del modelo lineal

- \hookrightarrow 1. los regresores son indep. 1. $P(x_1, x_2, \dots, x_n) = p(x_1) p(x_2) \dots p(x_n)$
- \hookrightarrow 2. Ausencia de colinealidad 2. no puede existir $\lambda_i / \lambda_i x_i + \lambda_j x_j = x_k$
- \hookrightarrow 3. El proceso es homocedástico 3. los E_i iid, $E_i \sim N(0, \sigma^2) \forall i$

Con estos supuestos limitamos la forma de encontrar $\hat{\beta}$, vamos a ver 2 métodos (y un tercero como apunte):

- (A) . MSE (Mean Square error) Enfoque empírico
- (B) ML (Maximum Likelihood) Enfoque probabilístico
- (Ap) MAP (Maximum a posteriori) Enfoque bayesiano

④ MSG:

partir de un dataset $D = \{ (x_i, y_i) \quad \forall i [1, \dots, K] \quad x_i \in \mathbb{R}^{m \times 1} \}$

$$\mathcal{E}(\beta) = \sum_{n=1}^K (y_n - \hat{f}(x_n))^2 = \sum_{n=1}^K \left(y_n - \beta_0 - \sum_{i=1}^m x_i \beta_i \right)^2 \quad ①$$

Agrego a x_i un 1 para representar a $\beta_0 \Rightarrow x_i = [1, x_1, x_2, \dots, x_m]$,
entonces de ①:

$$\begin{aligned} \mathcal{E}(\beta) &= \sum_{n=1}^K \left(y_n - \sum_{i=1}^m \beta_i x_i \right)^2 \\ &= (\bar{y} - \bar{x} \bar{\beta})^t (\bar{y} - \bar{x} \bar{\beta}) \quad ② \end{aligned}$$

Vamos a minimizar ② $\Rightarrow \partial_{\bar{\beta}} \mathcal{E}(\bar{\beta}) = 0$

$$\begin{aligned} 0 = \partial_{\bar{\beta}} \mathcal{E} &= \partial_{\bar{\beta}} [(\bar{y} - \bar{x} \bar{\beta})^t (\bar{y} - \bar{x} \bar{\beta})] = -2 \bar{x}^t (\bar{y} - \bar{x} \bar{\beta}) \\ &= \bar{x}^t (\bar{y} - \bar{x} \bar{\beta}) = \bar{x}^t \bar{y} - \bar{x}^t \bar{x} \bar{\beta} \end{aligned}$$

$X^t X$: matriz de
diseño.

$$\hat{\beta} = (\bar{X}^t \bar{X})^{-1} \bar{X}^t \bar{Y}$$

$$\rightarrow \hat{y} = H y$$

$$H = X(X^t X)^{-1} X^t$$

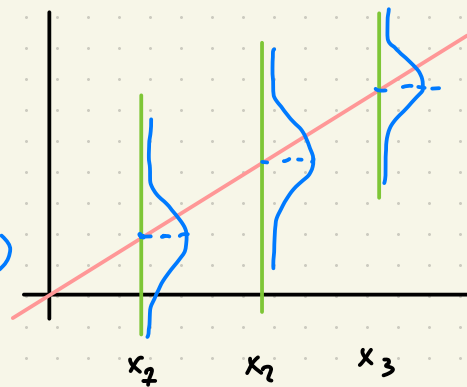
la parte más difícil de esto es calcular $(X^t X)^{-1}$. Sobre todo si $K \gg m$ (y viceversa) no para solventar esto se utiliza la *pseudo inversa* (Método Moore Pearson).

ⓑ Método de máxima verosimilitud:

$\exists \hat{\beta} / \max_{\bar{\beta}} \mathcal{L}(\beta)$. partimos de $y_i = \hat{f}(x_i, \beta) + \varepsilon$, $\varepsilon \text{ iid} \sim \mathcal{N}(0, \sigma^2)$

$$\mathbb{P}(\bar{Y} | \bar{X}, \bar{\beta}) = \mathbb{P}(y_1, \dots, y_n | x_1, \dots, x_n, \beta_0, \dots, \beta_m)$$

$$\stackrel{\text{iid}}{=} \prod_{n=1}^K \mathbb{P}(y_n | x_n, \bar{\beta}) \sim \mathcal{N}(y_n / x_n^t \bar{\beta}, \sigma^2) \textcircled{5}$$



buscamos $\hat{\beta}_{ML} = \arg \max (\mathcal{P}(\dots))$

$\mathcal{L}(\bar{\beta}) = \arg \max_{\beta} \mathcal{P}(Y|X, \beta)$ \leftarrow reemplazamos \textcircled{B}

si tomamos $\log \mathcal{L}(\beta) = \ell(\beta)$ convertimos Π en Σ :

$\ell(\beta) = - \arg \min_{\beta} (\underbrace{\log \mathcal{P}(\dots)}_{*})$

$\textcircled{*} \log \mathcal{P}(\cdot) = \frac{1}{2\sigma^2} (y_n - x_n^t \beta)^2 + c$

$\ell(\beta) = \sum_{k=1}^n \frac{1}{2\sigma^2} (y_n - x_n^t \beta)^2 = \frac{1}{2\sigma^2} \underbrace{(\bar{y} - \bar{x}\bar{\beta})^t (\bar{y} - \bar{x}\bar{\beta})}_{\|\bar{y} - \bar{x}\bar{\beta}\|^2}$

$= \frac{1}{2\sigma^2} \|\bar{y} - \bar{x}\bar{\beta}\|^2 \textcircled{4}$

optimizamos ④:

$$\partial_{\beta} l(\beta) = 0 \rightarrow \partial_{\beta} \left(\frac{1}{2\sigma^2} (\bar{y} - \bar{x}\bar{\beta})^t (\bar{y} - \bar{x}\bar{\beta}) \right) = 0$$

$$\partial_{\beta} \left(y^t y - 2 y^t x \beta + \beta^t x^t x \beta \right) = 0$$

$$0 - 2 y^t x + 2 \beta^t x^t x = 0$$

$$\rightarrow -y^t x + \beta^t x^t x = 0 \rightsquigarrow \beta^t x^t x = y^t x$$

$$\beta^t = y^t x (x^t x)^{-1}$$

$$\hat{\beta}_{ML} = (x^t x)^{-1} x^t y$$

MAP (Maximum a posteriori) Enfoque Bayesiano

En los métodos que vimos anteriormente no ponemos suposiciones sobre los parámetros θ . El método MAP propone asumir la distribución 'a priori' $p(\theta)$. Esto, restringe los valores que pueden tomar. Vamos a considerar $p(\theta) \sim \mathcal{N}(0, 1)$, esto va a limitar el valor de $\theta \in [-2, 2]$ con alta probabilidad (esto es $\pm 2 \sigma_\theta$). Teniendo el dataset (X, Y) , en vez de maximizar la fn. de verosimilitud, vamos a buscar los parámetros θ que maximizan la distribución a posteriori $p(\theta/x, y)$. Si aplicamos el teorema de Bayes:

Teorema de Bayes:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(\theta/x, y) = \frac{P(y/x, \theta) P(\theta)}{P(y/x)}$$

M1

En la ec. M1 vamos a buscar θ_{MAP} que maximice la distrib. a posteriori.

Vamos a utilizar un truco similar al log usado en ML.

$$\log(P(\theta/x, y)) = \log(P(y/x, \theta)) + \log(P(\theta)) + \underline{\text{cte.}}$$

no depende de θ

M2

Para encontrar θ_{MAP} , planteamos:

$$\theta_{\text{MAP}} \in \operatorname{argmin} \{-\log P(y/x, \theta) - \log P(\theta)\}$$

Para esto vamos a considerar:

$$-\partial_{\theta} \log p(\theta|x, y) = -\partial_{\theta} \log p(y|x, \theta) - \partial_{\theta} \log p(\theta)$$

sabiendo que $p(\theta) \sim \mathcal{N}(\phi, b^2 \mathbb{I})$, $\phi = [0, \dots, 0] \in \mathbb{R}^D$; $b^2 \mathbb{I} = \begin{bmatrix} b & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & b \end{bmatrix}$ podemos obtener:

$$-\partial_{\theta} \log p(\theta|x, y) = \partial_{\theta} \left(\frac{1}{2\sigma^2} (y - \Phi \theta)^t (y - \Phi \theta) + \frac{1}{2b^2} \theta^t \theta + \text{cte} \right) \quad (M3)$$

donde Φ es la matriz de features $[\mathbb{1}^t, \bar{x}] = \begin{bmatrix} 1 & x_{n1} & \dots & x_{n1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{D1} & \dots & x_{Dn} \end{bmatrix}$

A partir de (M3):

$$-\partial_{\theta} \log p(\theta|x, y) = \frac{1}{\sigma^2} (\theta^t \Phi^t \Phi - y^t \Phi) + \frac{1}{b^2} \theta^t$$

tomando $-\partial_{\theta} \log p(\theta|x, y) = 0$

$$\frac{1}{\sigma^2} (\theta^t \Phi^t \Phi - y^t \Phi) + \frac{1}{b^2} \theta^t = 0 \quad \leadsto \quad \theta^t \left(\frac{1}{2\sigma^2} \Phi^t \Phi + \frac{1}{b^2} \mathbb{I} \right) - \frac{1}{\sigma^2} y^t \Phi = 0$$

Continuando:

$$\Theta^t \left(\Phi^t \Phi + \frac{\sigma^2}{b^2} \mathbb{I} \right) = y^t \Phi \leadsto \Theta^t = y^t \Phi \left(\Phi^t \Phi + \frac{\sigma^2}{b^2} \mathbb{I} \right)^{-1}$$

Con esto obtenemos el estimador MAP

$$\Theta_{\text{MAP}} = \left(\Phi^t \Phi + \frac{\sigma^2}{b^2} \mathbb{I} \right)^{-1} \Phi^t y$$

Si vemos el resultado obtenido es muy similar al obtenido previamente salvo por el término $\sigma^2/b^2 \mathbb{I}$. Este término nos asegura que el término a invertir sea simétrico y definido estricto positivo. Esto asegura la existencia de la inversa $\Rightarrow \Theta_{\text{MAP}}$ tiene solución única.

Finalmente, Θ_{MAP} tiene un efecto regularizador sobre los parámetros que luego aprovecharemos.

Bibliografía

- The Elements of Statistical Learning | Trevor Hastie | Springer ← MAP
- An Introduction to Statistical Learning | Gareth James | Springer
- Deep Learning | Ian Goodfellow | <https://www.deeplearningbook.org/>
- Stanford | CS229T/STATS231: Statistical Learning Theory | <http://web.stanford.edu/class/cs229t/>
- Mathematics for Machine Learning | Deisenroth, Faisal, Ong ← MAP
- Artificial Intelligence, A Modern Approach | Stuart J. Russell, Peter Norvig

