

# Modelos Lineales para Clasificación

## Modelos lineales con regularización $L_1$

Juan Zamora O.



## Alta dimensionalidad en modelos lineales

Anteriormente, discutimos modelos lineales para clasificación que usaban un número moderado de variables explicativas basados en Máxima Verosimilitud (MV)

En algunos problemas, este número puede ser realmente elevado, por ejemplo en genómica pueden llegar los varios miles.

MV no funciona correctamente en estos contextos.

En problemas con alta dimensionalidad, la reducción de la cantidad de predictores es un asunto relevante.

Existen diversos métodos para la selección de subconjuntos de predictores. Por ejemplo, *Stepwise variable selection*.

Este tipo de métodos puede variar notoriamente su resultado en presencia incluso de pequeños cambios en los datos

Este tipo de métodos exhibe una alta varianza en el error de predicción.

Los métodos de regularización permiten encoger los estimadores hacia el cero e incluso realizar selección de variables vía optimización de log-verosimilitud penalizada



# Regularización de modelos lineales

Métodos de regularización derivados de estimadores de Máxima Verosimilitud (MV) se basan en la *Log-verosimilitud penalizada*  $l_P$

$$l_P(\boldsymbol{\beta}) = \sum_{i=1}^N l_i(\boldsymbol{\beta}) - \frac{\lambda}{2} J(\boldsymbol{\beta})$$

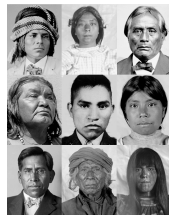
Donde  $l_i(\boldsymbol{\beta})$  representa la contribución del  $i$ -ésimo ejemplo a la log-verosimilitud,  $\lambda$  es un parámetro que indica qué tanto regularizar y  $J$  es un funcional que penaliza el tamaño de los parámetros

Asintóticamente, métodos de regularización que *encogen* los estimadores reducen la varianza<sup>1</sup>.





El conjunto de datos de diabetes en los indios Pima tiene como tarea la predicción de la aparición de diabetes en un periodo de 5 años usando detalles médicos como covariables.



Extraída desde Wikipedia.