



**LUND**  
UNIVERSITY

# Fjasdf

## Subtitle

---

Johan Larsson

Department of Statistics, Lund University

April 2, 2024

# Preliminaries

---

# General Setup

- Data consists of a **fixed** matrix of features  $\mathbf{X} \in \mathbb{R}^{n \times p}$  and a response vector  $\mathbf{y} \in \mathbb{R}^n$ .
- $\mathbf{y}$  comes from a linear model, that is,

$$y_i = \beta_0^* + \mathbf{x}_i^\top \boldsymbol{\beta}^* + \varepsilon_i \quad \text{for } i \in 1, \dots, n,$$

where  $\boldsymbol{\beta}^*$  is the vector of *true* coefficients.

- $\varepsilon_i$  is the measurement noise, generated from some random variable<sup>1</sup>.

---

<sup>1</sup>No assumption on normality (yet).

# The Elastic Net

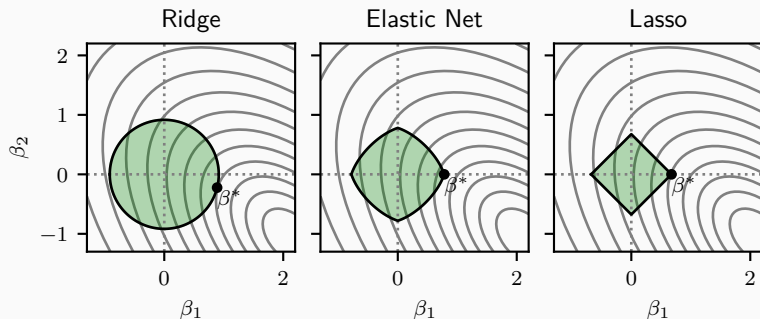
Linear regression plus a combination of the  $\ell_1$  and  $\ell_2$  penalties:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left( \frac{1}{2} \|\mathbf{y} - \beta_0 - \tilde{\mathbf{X}}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\beta}\|_2^2 \right).$$

# The Elastic Net

Linear regression plus a combination of the  $\ell_1$  and  $\ell_2$  penalties:

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left( \frac{1}{2} \|\mathbf{y} - \beta_0 - \tilde{\mathbf{X}}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2 \right).$$



**Figure 1:** The elastic net penalty is a combination of the lasso and ridge penalties

# Sensitivity to Scale

Since both the lasso and ridge penalties are norms, they are sensitive to the scale of the input features.

But what is the *optimal* scaling?

If the features are “normal”, then most people would agree that standardizing them (i.e., subtracting the mean and dividing by the standard deviation) is a good idea.

# Normalization

Let  $\mathbf{S}$  be the *scaling matrix*, which is a  $p \times p$  diagonal matrix with entries  $s_1, s_2, \dots, s_p$ . Let  $\mathbf{C}$  be the *centering matrix*, which is an  $n \times p$  matrix with each row equal to  $[c_1, c_2, c_n]^\top$ . Then the *normalized design matrix*  $\tilde{\mathbf{X}}$  is defined as  $\tilde{\mathbf{X}} = (\mathbf{X} - \mathbf{C})\mathbf{S}^{-1}$ .

**Table 1:** Common ways to normalize a matrix of features

Normalization	Centering ( $c_{1j}$ )	Scaling ( $s_j$ )
Standardization	$\frac{1}{n} \sum_{i=1}^n x_{ij}$	$\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$
Min-Max	$\min_i(x_{ij})$	$\max_i(x_{ij}) - \min_i(x_{ij})$
Unit Vector (L2)	0	$\sqrt{\sum_{i=1}^n x_{ij}^2}$
Max-Abs	0	$\max_i( x_{ij} )$
Adaptive Lasso	0	$\beta_j^{\text{OLS}}$



# Binary Features

Let's say we have a binary feature  $x_j$ , such that  $x_{ij} \in \{0, 1\}$ .

What is the “best” way to scale this feature?

# Solution for Binary Features

We assume the that normalized features are orthogonal, that is

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \text{diag}(\dots)$$

# Class Imbalance

## Mixed Data

---



## Experiments

---