→ Make the size of the coefficients small ⟹ "a continuous version of sparsity".

② SHRINKAGE METHODS:

LINEAR REGRESSION (LS): GIVEN $X, \vec{y}$ ONE FIND $\vec{\beta}^{*}$

$$\min_{\vec{\beta}} \| \vec{y} - X\vec{\beta} \|^2 = \min_{\vec{\beta}} L_{LS}(\vec{\beta})$$

$$(X^T X)^{-1} X^T$$

THE SOLUTION IS:

$$\vec{\beta}^{*} = \text{ARGMIN} \| \vec{y} - X\vec{\beta} \|^2 = X^T \vec{y}$$

$$\hat{y} = \hat{f}(\vec{x}_{in}) = \vec{x}_{in} \vec{\beta}^{*}$$

2.1 RIDGE REGRESSION    $>0$

$$L_{RIDGE}(\vec{\beta}) \quad \| \vec{y} - X\vec{\beta} \|^2 + \boxed{\lambda} \| \vec{\beta} \|^2$$

( THE IDEA IS TO "ELIMINATE" SOME DESCRIPTORS. IN PRACTICE WE WEIGHT THEM LOW )

$$\min_{\vec{\beta}} L_{RIDGE}(\vec{\beta}) = \min_{\vec{\beta}} \| \vec{y} - X\vec{\beta} \|^2 + \underline{\lambda \| \beta \|^2}$$

||| EQUIVALENT    ← Penalty term

$$\min_{\vec{\beta}} \| \vec{y} - X\vec{\beta} \|^2$$

SUBJECT TO    $\| \vec{\beta} \|^2 \leq \underline{t}$

$\lambda$ and $t$ are inversely related

THIS PROBLEM CAN BE SOLVED EXPLICITLY:

$$L_{RIDGE}(\vec{\beta}) = \vec{y}^T \vec{y} - \vec{\beta}^T X^T \vec{y} - \vec{y}^T X \vec{\beta} + \vec{\beta}^T (X^T X + \lambda I) \vec{\beta}$$

$$\nabla_{\vec{\beta}} L_{RIDGE} = -2 X^T \vec{y} + 2(X^T X + \lambda I) \beta$$

$$X^T \vec{y} = (X^T X + \lambda I) \vec{\beta}^{*}$$

$$\vec{\beta}_{RIDGE}^{*} = (X^T X + \lambda I)^{-1} X^T \vec{y}$$

WE CAN CHECK THAT THE HESSIAN IS POSITIVE DEFINITE

BIAS: $\vec{y} = X\vec{\beta} + \vec{\varepsilon}$    $\mathbb{E}[\vec{\varepsilon}] = 0_v$

$$\mathbb{E}[\vec{\beta}_{RIDGE}^{*}] = \mathbb{E}[(X^T X + \lambda I)^{-1}(X^T X)\vec{\beta} + (X^T X + \lambda I)^{-1} X^T \varepsilon] =$$

$$= (X^T X + \lambda I)^{-1}(X^T X)\vec{\beta}$$

$$(X^T X) = V D V^T \quad V V^T = V^T V = I \quad (X^T X)^{-1} = V D^{-1} V^T$$

$\underset{\text{EIGENDECOMPOSITION}}{\uparrow}$

$D$ DIAGONAL    $\sigma(D) \equiv \sigma(X^T X)$

• $V$ contains the eigenvectors as columns

• $D = \text{diag}(d_i)$ has the eigenvalues on the diagonal.

$$\text{BIAS} = \mathbb{E}\left[\vec{\beta}_{RIDGE}^*\right] - (X^TX + \lambda I)^{-1}X^TX\vec{\beta} = V\left[(D+\lambda I)^{-1}D - I\right]V^T\vec{\beta}$$

$$\mathcal{D} = (D+\lambda I)^{-1}D - I \quad \longleftarrow \quad \text{EVERYTHING IS DIAGONAL}$$

$$\mathcal{D}_{ii} = \left[\frac{d_i}{d_i+\lambda} - 1\right] = -\frac{\lambda}{d_i+\lambda}$$

$$\Rightarrow \text{BIAS} = -\lambda V\left[(D+\lambda I)^{-1}\right]V^T\vec{\beta} = -\lambda(X^TX+\lambda I)^{-1}\vec{\beta}$$

As EXPECTED AS $\lambda \to 0$ BIAS $\to 0$ BECAUSE WE RECOVER LEAST SQUARES

As $\lambda \to \infty$ BIAS $\to -\vec{\beta}$ ($\approx 100\%$ ERROR)

VARIANCE:

$$\text{VAR}(\vec{\beta}_{RIDGE}) = \mathbb{E}\left[\left(\vec{\beta}_{RIDGE} - \mathbb{E}[\vec{\beta}_{RIDGE}^*]\right)^2\right]$$

$$= \mathbb{E}\left[(X^TX+\lambda I)^{-1}X^T\underbrace{\vec{\xi}\vec{\xi}^T}_{\sigma^2 I}X(XX^T+\lambda I)^{-1}\right] =$$

$$= \sigma^2(X^TX+\lambda I)^{-1}(X^TX)(X^TX+\lambda I)^{-2} =$$

$$= \sigma^2 V\underbrace{\left[(D+\lambda I)^{-1}D(D+\lambda I)^{-1}\right]}_{\mathcal{P}}V^T \qquad \mathcal{P}_{ii} = \frac{d_i}{(d_i+\lambda)^2}$$

As $\lambda \to \infty$ $\mathcal{P}_{ii} \to 0$ (QUADRATICALLY)



## 2.2 LASSO (TIBSHIRANI)

$$L_{LASSO}(\vec{\beta}) = \|\vec{y} - X\vec{\beta}\|^2 + \lambda\|\vec{\beta}\|_1 \qquad \boxed{\|\vec{\beta}\|_1 = \sum_{i=1}^{P}|\beta_i|}$$

$$\boxed{\min_{\vec{\beta}} L_{LASSO}(\vec{\beta})} \iff \boxed{\begin{array}{c}\min_{\vec{\beta}} \|\vec{y} - X\vec{\beta}\|^2 \\ \text{SUBJECT TO} \quad \|\vec{\beta}\|_1 \le t\end{array}}$$

~~WE CAN'T FIND AN ANALYTICAL SOLUTION BUT THIS IS A CONVEX CAN BE FOUND~~

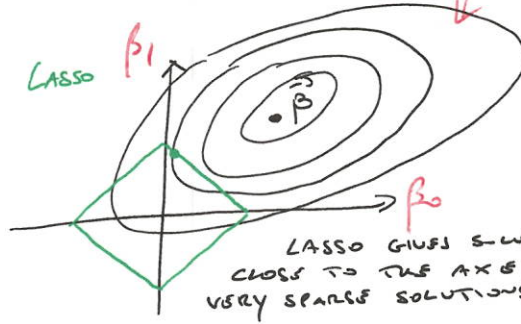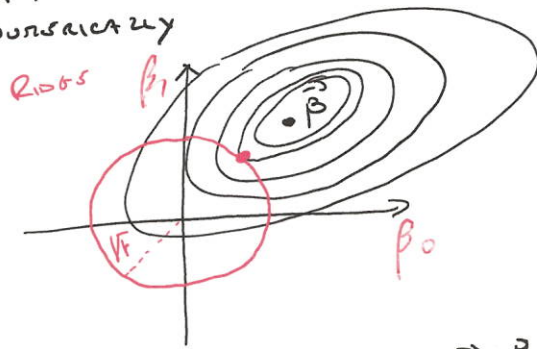No analytical solution for LASSO (contrary to Ridge)

but it is a convex problem $\left[\begin{array}{c}\text{optimize a convex function} \\ \text{over a convex set}\end{array}\right]$

So we can use convex optimisation techniques (quadratic programming) to optimise globally.

.... optimisation viewpoint:

WE CAN OPTIMISATION PROBLEM THAT HAS A SOLUTION THAT CAN B NUMERICALLY

RIDGE

$\beta_1$ ↑

$\beta_0$

LASSO $\beta_1$ ↑

$\beta_0$

level of LS: curves

$\| \vec{y} - X\hat{\beta} \|^2 = k$

LASSO GIVES SOLUTIONS CLOSE TO THE AXES, i.e. VERY SPARSE SOLUTIONS.

OTHER VARIATIONS: $\lambda \| \vec{\beta} \|^q$ FOR DIFFERENT $q$

$\beta_1$ ↑ $q=\infty$

$\beta_0$

↑ $q=4$

↑ $q=2$

Ridge

↑ $q=1$

LASSO

↑ $q<1$

↑ $q=0$

convex sets

Non-convex sets

optimisation is doable

Optimisation is difficult

$q=0$ is the "$l_0$"-pseudonorm case

Where we only have solutions where some of the parameters $\beta_i$ are zero ≡ Subset selection

⇓

i.e., "Sparse" models

↳ Difficult to optimise for sparsity

this whole area is called regularisation

or

controlling for model complexity

Another direction is to mix the penalty terms in objective function:

Example:

Elastic net :

$$L_{EN}(\vec{\beta}) = \| \vec{y} - X\vec{\beta} \|_2^2 + \lambda \left[ \alpha \| \vec{\beta} \|_1 + (1-\alpha) \| \vec{\beta} \|_2 \right]$$