

6.2 Lasso

$$y = 1.2 + 2.4x_1 + 0.02x_2$$
$$1.2 \quad 2.3 \quad -1.6x_3 + \varepsilon$$

-1.7

One drawback to ridge regression is that the estimated $\hat{\beta}_j$'s always take on nonzero values, even if they are very small.

The lasso (least absolute shrinkage & selection operator) is an alternative that minimizes

$$(y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_1$$
$$= (y - X\beta)^T (y - X\beta) + \lambda \sum_{j=1}^p |\beta_j|$$

Lasso shrinks some estimated β_j 's to exactly zero, so serves as a variable selection method, & generates sparse models

G.3 Another view of ridge/lasso

An equivalent to write the ridge problem is

$$\min_{\beta} \sum_{j=1}^n (y_j - x_j^T \beta)^2$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 \leq S$$

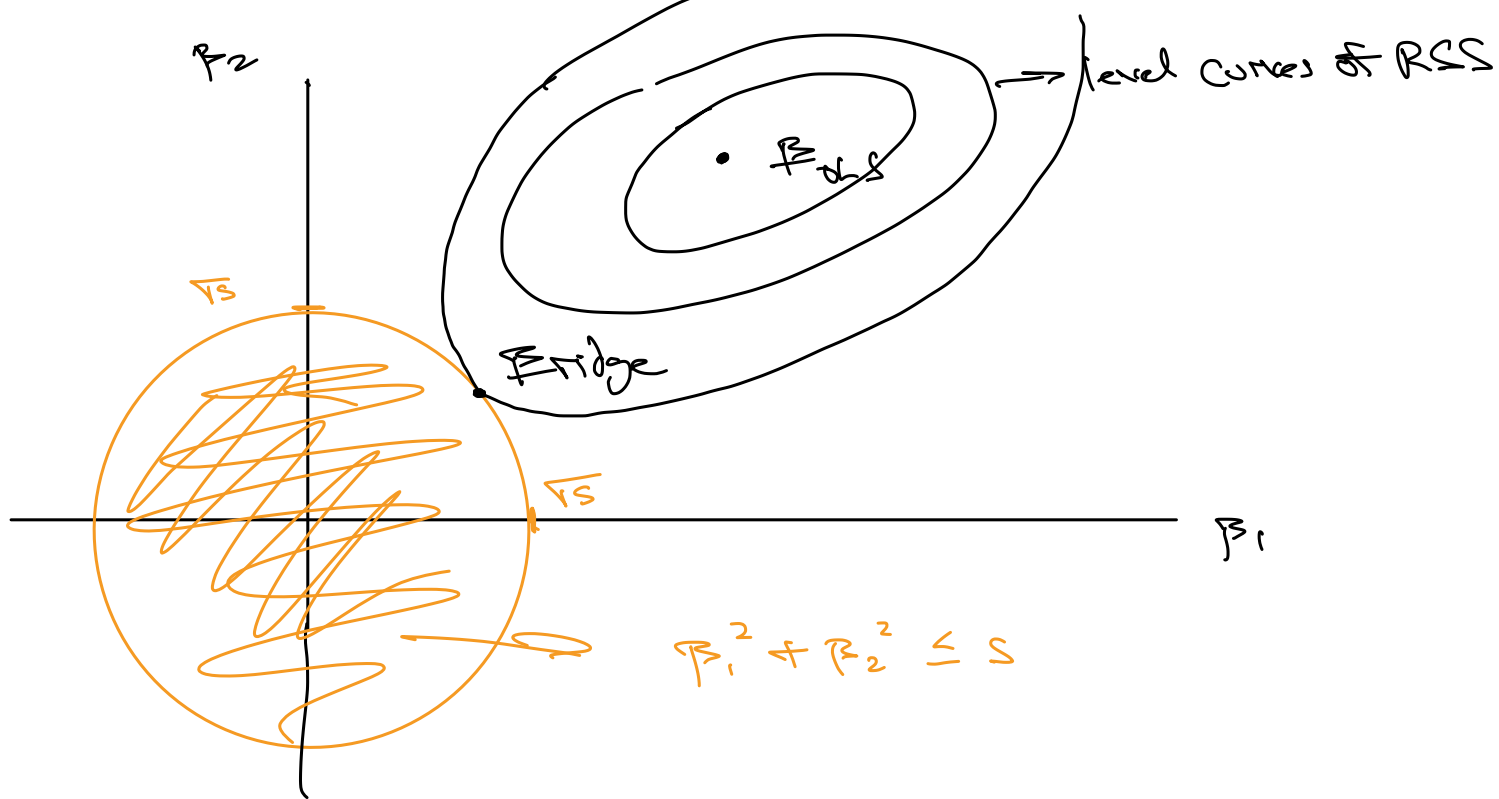
where S is a regularization parameter

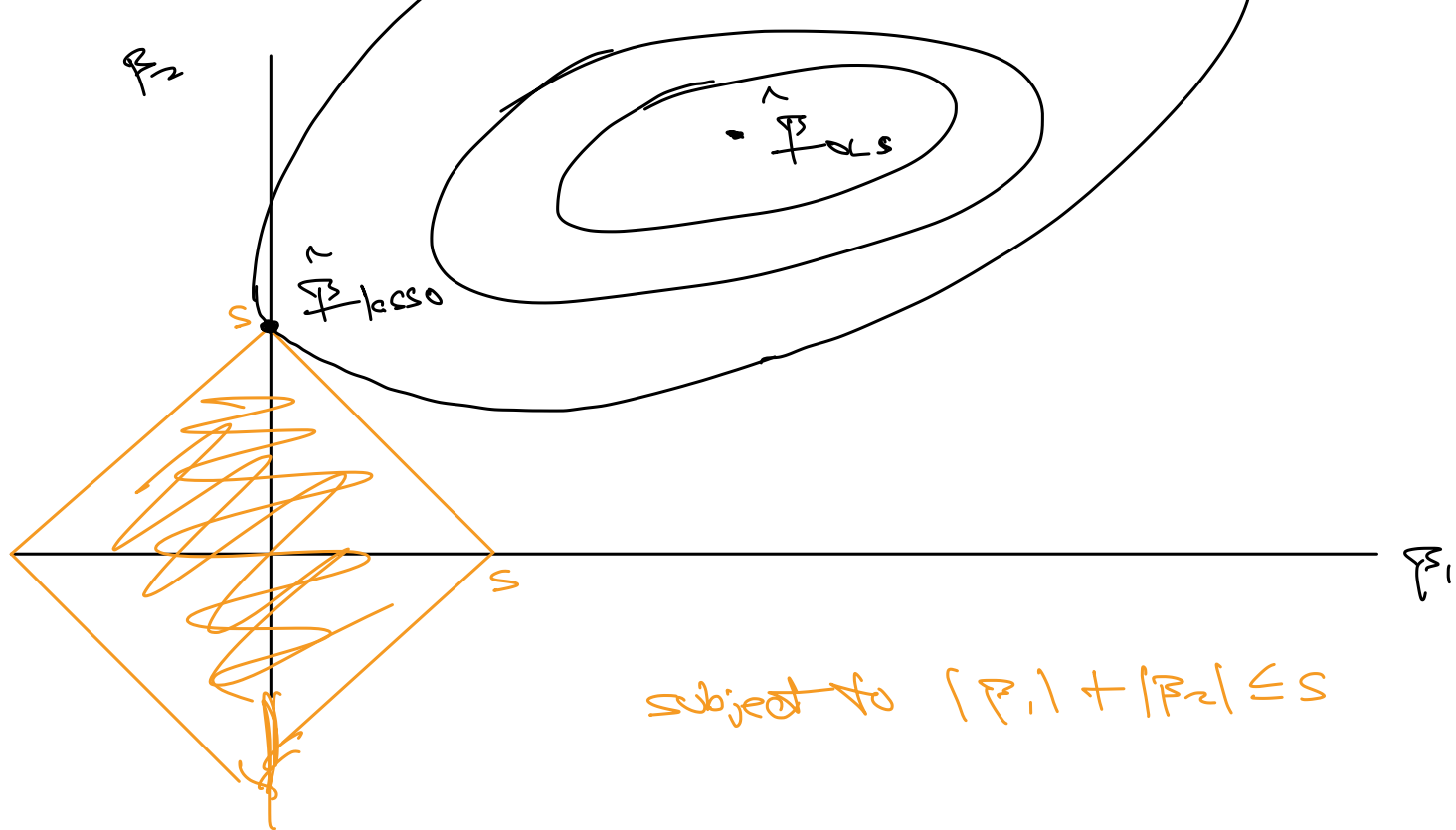
as $S \rightarrow \infty$ get the OLS estimator

as $S \rightarrow 0$ the model becomes the null model $Y = \beta_0 + \epsilon$

$$\text{For lasso, } \min_{\beta} \sum_{i=1}^n (y_i - x_i^T \beta)^2$$

$$\text{subject to } \sum_{j=1}^p |\beta_j| \leq S$$





subject to $|\beta_1| + |\beta_2| \leq s$

6.4 A Bayesian interpretation

In Bayesian statistics, we represent beliefs / uncertainty / knowledge about a parameter β using a prior distribution

then, given data, update beliefs / uncertainty / confidence estimates in a posterior distribution.

Think of $\gamma \sim U(0, \beta)$ or $\text{Exp}(\beta)$ or $N(\beta, 1)$ or $N(0, \beta)$

We get data from a data likelihood, for y_1, \dots, y_n

$$f(y | \beta) = \text{data likelihood}$$

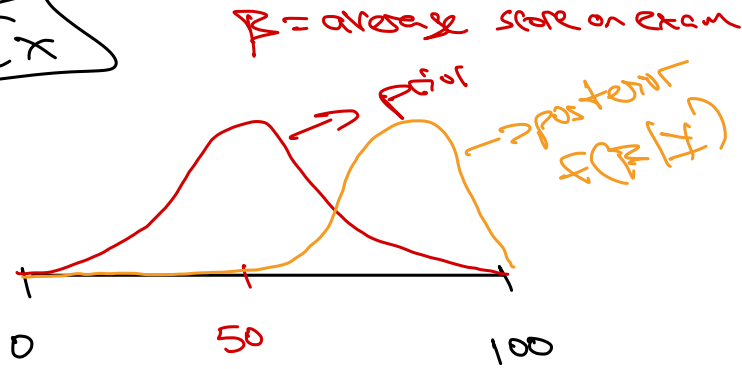
we seek

$$f(\beta | y) = \frac{f(y | \beta) \pi(\beta)}{f(y)}$$

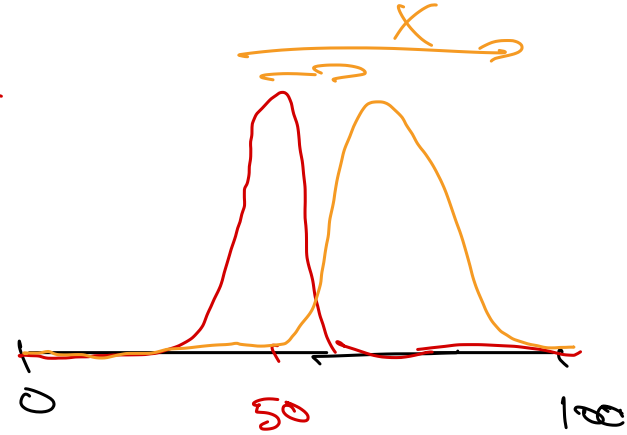
where $\pi(\beta)$ is a prior distribution, $\pi(\beta|y)$ is
posterior distribution.

$\pi(\beta)$ represents our beliefs about β before seeing any data.

Ex



get data \Rightarrow everybody
> 90%



same data