$\boxed{\text{Recall}}$ $Y$ data from model depending on $\beta$

$(e.g. \ Y = \beta x + \varepsilon)$. Ingredients:

- prior distribution $\pi(\beta)$
- data likelihood $f(y|\beta)$
- posterior distribution $\quad f(\beta|y) = \dfrac{f(y|\beta)\pi(\beta)}{f(y)}$

If $Y = X\beta + \varepsilon$ under A1, then $\underline{Y} = (Y_1, \dots, Y_n)^T$ iid data here

$$f(\underline{Y} \mid \beta) = \frac{1}{(2\pi v^2)^{n/2}} \exp\left( -\frac{1}{2\sigma^2}(\underline{Y} - \beta \underline{x})^T (\underline{Y} - \beta \underline{x}) \right)$$

<span style="color:red">depends on $\beta$</span>

$$\beta = \# \ , \qquad \underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \qquad \underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$$

Note: $\quad (\underline{Y} - \beta \underline{x})^T (\underline{Y} - \beta \underline{x}) = \sum_{j=1}^{n} (Y_j - \beta x_j)^2$

<u>Comment</u> The posterior mode is the value of $\beta$
that maximizes the posterior dist $f(\beta \mid \underline{Y})$

$$\text{maximize } f(\beta \mid \underline{Y})$$

$$\iff \quad \text{maximize } \log f(\beta \mid \underline{Y})$$

$$\Rightarrow \quad \text{minimize} \quad -\log f(\beta | y)$$

Thus the posterior mode minimizes

$$-\log f(\beta | y) = -\log \left\{ \frac{f(y | \beta) \pi(\beta)}{f(y)} \right\} =$$

$$-\log f(y | \beta) - \log \pi(\beta) + \log f(y) \quad \textcolor{red}{\Rightarrow \text{ does not depend on } \beta}$$

$$= \frac{1}{2\sigma^2} (y - \beta x)^T (y - \beta x) - \log \pi(\beta) \left( + \text{ constants that do not depend on } \beta \right)$$

Compare to ridge regression minimizes

$$(y - \beta x)^T (y - \beta x) + \lambda \beta^2$$

$$\Rightarrow \quad \text{set} \quad -\log \pi(\beta) = \lambda \beta^2$$

$$\Leftrightarrow \quad \pi(\beta) \propto e^{-\lambda \beta^2} \longrightarrow \text{up to normalization}$$
constants, a normal pdf

$$c \cdot e^{-\frac{1}{2(\frac{1}{2\lambda})}\left(\frac{\beta - 0}{1}\right)^2}$$

The ridge regression estimator of $\beta$ is the posterior mode in a Bayesian analysis using a normal, mean zero prior for $\beta$.
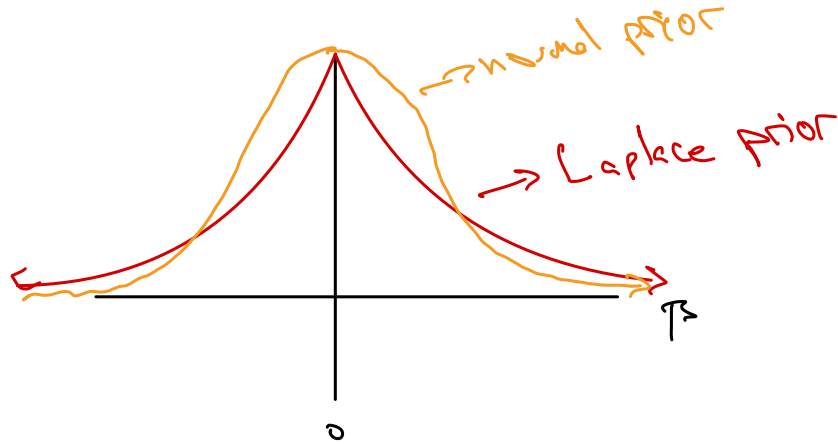
minimize

$$\frac{1}{2\sigma^2}(Y - \beta X)^T(Y - \beta X) \quad -\log \pi(\beta) + \text{const}$$

Lasso:
minimize

$$(Y - \beta X)^T(Y - \beta X) + \lambda |\beta|$$

$$\Rightarrow \quad \pi(\beta) \propto e^{-\lambda |\beta|} \rightarrow \text{double exponential or Laplace dist.}$$

$\Rightarrow$ Lasso estimator is same as posterior mode under a mean Laplace prior.



$\rightarrow$ normal prior

$\rightarrow$ Laplace prior

$\beta$

$0$

"Elastic net"

ridge          lasso

$$(\underline{y} - X\underline{\beta})^T (\underline{y} - X\underline{\beta}) + (1-\alpha) \|\beta\|_2^2 + \alpha \|\beta\|_1 \qquad \alpha \in [0,1]$$

# 7 Classification 2

$$y_i \in \{-1, +1\} \qquad \underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$$

## 7.1 Maximum margin classifier

In $p$ dimensions a <u>hyperplane</u> is flat subspace of dimension $(p-1)$

**DEF** In $p$ dim a <u>hyperplane</u> is the set of all

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \in \mathbb{R}^p \quad \text{that satisfy}$$

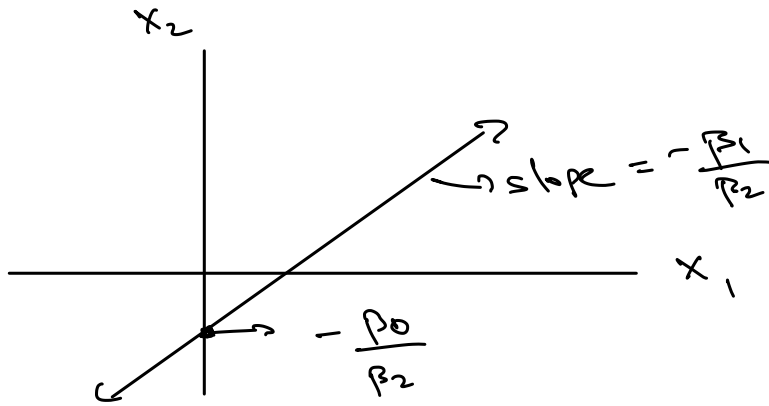$$\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p = 0 \qquad (\bigstar)$$

for some parameters $\beta_0, \beta_1, \cdots, \beta_p$

$\boxed{\text{Ex}}$  $P = 2$  $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$

$$\rightsquigarrow x_2 = \left(-\frac{\beta_0}{\beta_2}\right) + \left(-\frac{\beta_1}{\beta_2}\right) x_1$$

$$= a + b x_1$$

$x_2$

$\rightsquigarrow$ slope $= -\frac{\beta_1}{\beta_2}$

$x_1$

$-\frac{\beta_0}{\beta_2}$
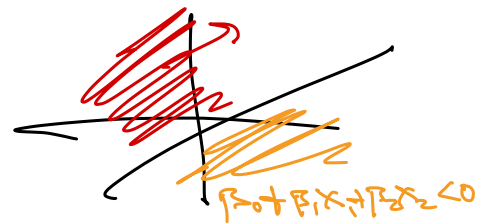
$\boxed{\text{Note}}$  If a vector $x$ does not satisfy $(\ast)$, then either

$$\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p > 0$$

$$\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p < 0$$

So a hyperplane splits space into halves $\leftarrow$

$\beta_0 + \beta_1 x_1 + \beta_2 x_2 > 0$

$\beta_0 + \beta_1 x_1 + \beta_2 x_2 < 0$

**Setup** Suppose we have training data $Y_1, \ldots, Y_n$ & features $\underline{x}_1, \ldots, \underline{x}_n$ [does not include the 1 for intercept]

$$Y_i \in \{-1, +1\} \qquad \underline{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \cdot \\ \cdot \\ x_{ip} \end{pmatrix} \rightarrow p\text{-veriate}$$

**Goal** Find a <u>separating hyperplane</u> such that

$$\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} > 0 \qquad \text{if} \quad Y_i = +1$$

$$\beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} < 0 \qquad \text{if} \quad Y_i = -1$$

$$\rightarrow \text{collapse} \quad Y_i \left( \beta_0 + \beta_1 x_{i1} + \cdots + \beta_p x_{ip} \right) > 0 \quad \text{for all } i$$

Separating hyperplane

$x_2$ axis with $+1$ points in upper-left region and $-1$ points in lower-right region, $x_1$ axis labeled.

? → do classification at a new $\underline{x}$ coordinate

Suppose $\beta_0, \ldots, \beta_p$ are known (estimated), then at a new feature $\underline{x}_* = (x_{*1}, \ldots, x_{*p})^T$, our classifier is:
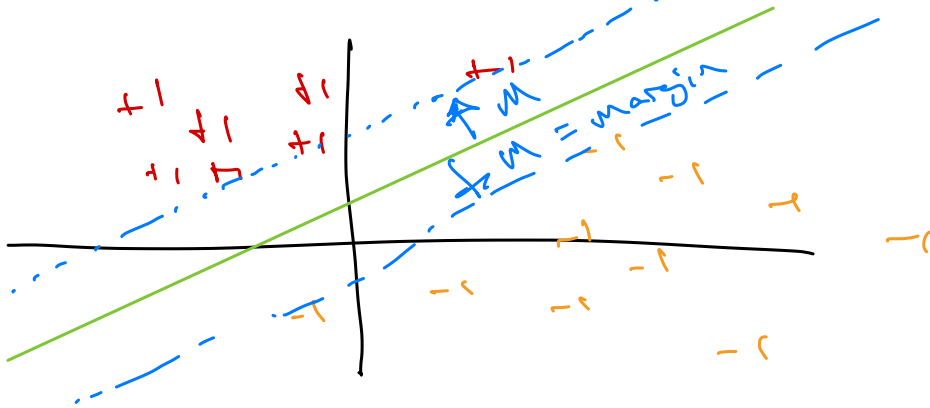
$$\hat{y} = \begin{cases} +1 & f(\underline{x}_*) > 0 \\ -1 & f(\underline{x}_*) < 0 \end{cases}$$

where $f(\underline{x}) = \beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p$

We can interpret the size/magnitude of $f$ as
confidence in our predictor.

Separable case

Suppose data are <u>linearly separable</u>. The <u>maximum
margin hyperplane</u> is the plane that separates the
classes but is as far away as possible from data.
Yields the <u>maximum margin classifier.</u>

$\boxed{\text{Note}}$ The m.m. hy. is calculated by

- maximize $M$
  $\beta_0, \beta_1, \dots, \beta_p$

- subject to $\displaystyle\sum_{j=1}^{p} \beta_j^2 = 1$

  and $y_i\left(\beta_0 + x_i^T \beta\right) \geq M$ for all $i = 1, \dots, n$