

Recap: LS Regression $y \in \mathbb{R}$, $x \in \mathbb{R}^d$

$$p(y|x, \theta) = N(\omega^\top x + b, \sigma^2)$$

$$\hat{\omega}_{MAP} = \arg \max_{\omega} \sum_{i=1}^n \log P(y_i|x_i, \theta) + \log P(\omega)$$

$$P(\omega) = N(0, t^2 I)$$

$$= \arg \min_{\omega, b} \sum (y_i - \omega^\top x_i - b)^2 + \frac{\sigma^2}{t^2} \cdot \|\omega\|^2$$

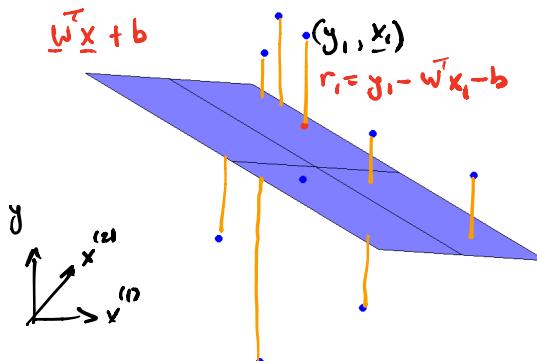
$= \lambda$

$\lambda=0$ Linear Least Squares

① Row-wise interpretation

$$\lambda=0 \quad x_1=0.41497 \quad x_2=0.24982$$

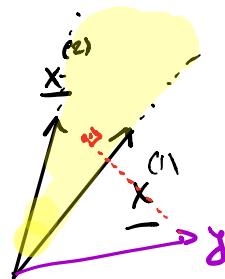
$$r_1^2 + r_2^2 + \dots + r_n^2$$



$$n \times d \quad \underline{x} = \begin{bmatrix} \underline{x}_1 \\ \vdots \\ \underline{x}_n \end{bmatrix} \quad \underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

② Column-wise

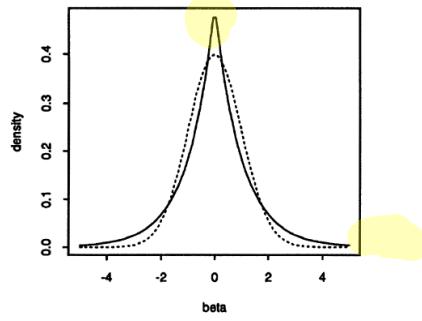
$$\underline{x} = \begin{bmatrix} \underline{x}^{(1)} & \dots & \underline{x}^{(d)} \end{bmatrix}$$



L1 Regularization (Lasso)

$$\arg \min \sum_i (y_i - \underline{w}^T \underline{x}_i - b)^2 + \lambda \|\underline{w}\|_1$$

$$\text{Laplace prior} \propto e^{-\frac{|w_1|}{t}} \cdot e^{-\frac{|w_2|}{t}} \cdots e^{-\frac{|w_d|}{t}} = e^{-\frac{1}{t} \|\underline{w}\|_1}$$

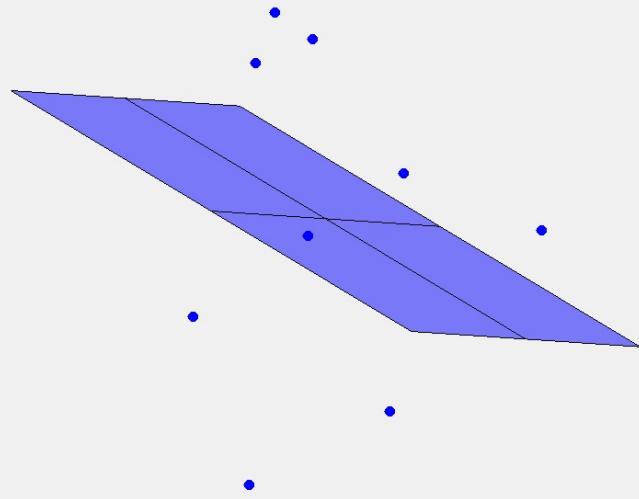


(Laplacean)

Fig. 7. Double-exponential density (—) and normal density (---); the former is the implicit prior used by the lasso; the latter by ridge regression

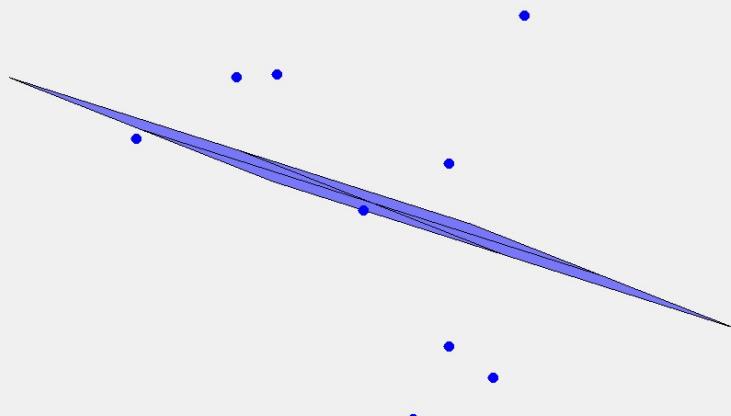
Ridge Regression ($\lambda \parallel w \parallel_2^2$)

$$\lambda=0 \quad x_1=0.41497 \quad x_2=0.24982$$



Lasso ($\lambda \parallel w \parallel_1$)

$$\lambda=40 \quad x_1=0.16088 \quad x_2=0.41715$$



Distribution of \hat{w}_{LS}

$$y_i = \bar{w}^T x_i + e_i \quad \text{e} \sim \mathcal{N}(0, \sigma^2)$$

$$\begin{aligned} \underline{y} &= \underline{x} \cdot \underline{w} + \underline{e} \quad \text{Pseudo-inverse of } \underline{x} \\ \hat{w}_{LS} &= (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y} \\ &= (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{x} \underline{w} + (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{e} \\ &= \underline{w} + (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{e} \end{aligned}$$

$$\mathbb{E} \hat{w}_{LS} = \underline{w} \quad \mathbb{E} (\hat{w}_{LS} - \underline{w}) \cdot (\hat{w}_{LS} - \underline{w})^T = \mathbb{E} (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{e} \underline{e}^T \underline{x} (\underline{x}^T \underline{x})^{-1} \\ = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{x} \sigma^2 \mathbf{I} \cdot \underline{x} (\underline{x}^T \underline{x})^{-1} \\ = \sigma^2 (\underline{x}^T \underline{x})^{-1}$$

$$\hat{w}_{LS} \sim N(\underline{w}, \sigma^2 (\underline{x}^T \underline{x})^{-1})$$

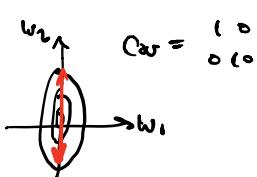
$$\underline{x} = \underline{U} \underline{\Sigma} \underline{V}^T \quad \underline{x}^T \underline{x} = \underline{V} \underline{\Sigma} \underline{U}^T \underline{\Sigma} \underline{U} \underline{V}^T = \underline{V} \underline{\Sigma}^2 \underline{V}^T$$

E.g.

$$(\underline{x}^T \underline{x}) = \begin{bmatrix} 1 & 0 \\ 0 & 0.1 \end{bmatrix}$$

$$(\underline{x}^T \underline{x})^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 10 \end{bmatrix}$$

$$(\underline{x}^T \underline{x})^{-1} = \underline{V} \underline{\Sigma}^{-2} \underline{V}^T$$

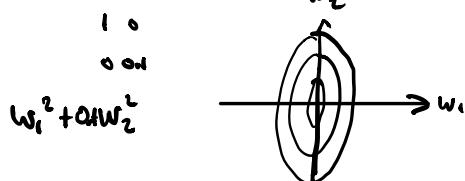


$$\|\underline{x}\underline{w} - \underline{y}\|_2^2$$

$$\text{Hessian: } \min_{\underline{w}} L(\underline{w})$$

$$\nabla^2 L(\underline{w}) = \underline{x}^T \underline{x}$$

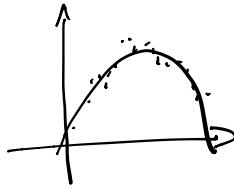
$$\text{e.g. } \underline{x}^T \underline{x} = \begin{bmatrix} 1 & 0 \\ 0 & 0.1 \end{bmatrix}$$



$$\begin{aligned} \hat{w}_{Ridge} &= (\underline{x}^T \underline{x} + \lambda \mathbf{I})^{-1} \underline{x}^T \underline{y} \\ &= \underline{V} (\underline{\Sigma}^2 + \lambda \mathbf{I})^{-1} \underline{V}^T \underline{x}^T \underline{\Sigma} \underline{y} \end{aligned}$$

$$= \nabla (\sum_i v_i \phi_i)^\top \mathbf{w} = \sum_i v_i w_i \frac{v_i}{\sigma_i^2 + \lambda}$$

Nonlinear features



$$f(x, \theta) = b + \sum_{j=0}^{m-1} w_j \phi_j(x)$$

$$\theta = \begin{bmatrix} w \\ b \end{bmatrix}$$

$\phi_j(x)$ are basis functions

$$\min \sum_{i=1}^n (y_i - f(x_i, \theta))^2$$

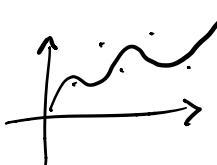
Example: $x \in \mathbb{R}$ $\phi_j(x) = x^j$ $f(x)$ is a degree $m-1$ polynomial

$$x \in \mathbb{R}^d \quad \phi(x) = (1, x^{(1)}, x^{(2)}, (x^{(1)})^2, (x^{(2)})^2, \dots)$$

$$\text{Gaussian basis } \phi_j(x) = e^{-\frac{\|x - \mu_j\|^2}{2\sigma_j^2}} \quad x \in \mathbb{R}^d$$

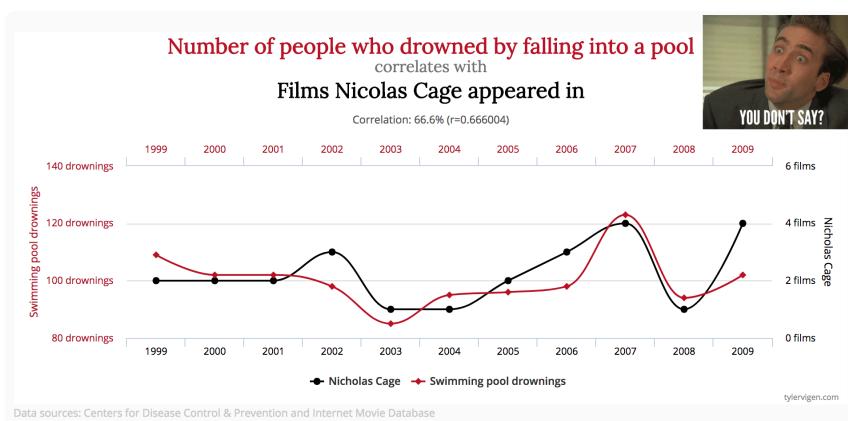
$$\min \|y - \underline{\underline{\theta}} \underline{\underline{w}}\|_2^2$$

→ overfit



$$\underline{\underline{\theta}} = \begin{bmatrix} \phi_0(x_1) & \phi_{m-1}(x_1) \\ \vdots & \vdots \\ \phi_0(x_n) & \phi_{m-1}(x_n) \end{bmatrix}$$

→ spurious correlation

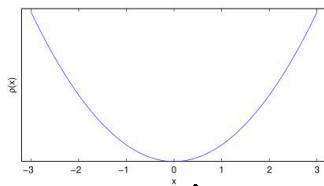
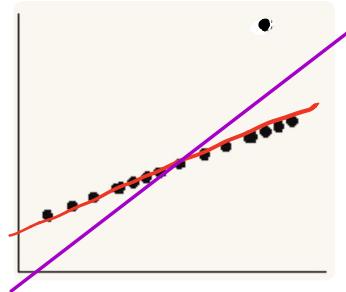


Robust Regression

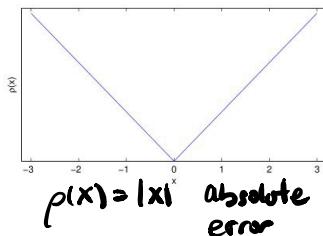
LS is not robust to outliers

$$\min \sum_{i=1}^n \rho(y_i - \mathbf{w}^\top \mathbf{x}_i - b)$$

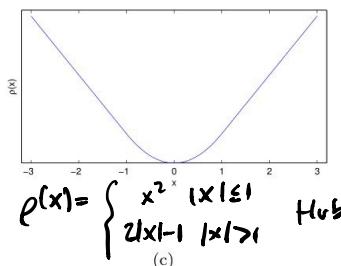
where ρ is a robust loss function



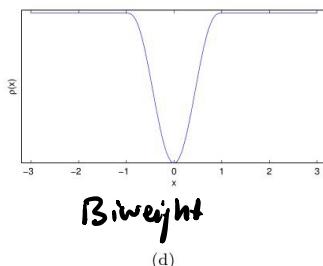
(a)



(b)



(c)



(d)

Kernels

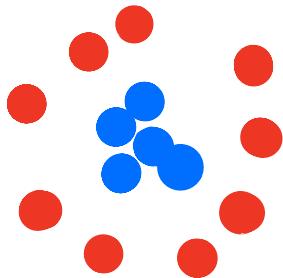
Non linear feature maps: transform feature vectors via

$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$$

and apply a linear method (e.g. least-squares, SVM) to the transformed data $\Phi(\mathbf{x}_1), \dots, \Phi(\mathbf{x}_n)$.

Classification: $u \rightarrow \text{sign}(\mathbf{w}^\top \Phi(\mathbf{x}) + b)$

Example $x = \begin{bmatrix} x^{(1)} \\ x^{(2)} \end{bmatrix} \in \mathbb{R}^2$ Consider $\Phi(x) = \begin{bmatrix} x^{(1)} \\ x^{(2)} \\ x^{(1)}, x^{(2)} \\ (x^{(1)})^2 \\ (x^{(2)})^2 \end{bmatrix}$



Can be separated by a circular classifier

$$y \rightarrow \text{sign}\left((x^{(1)} - c_1)^2 + (x^{(2)} - c_2)^2 - r^2\right)$$

This is a linear classifier in the transformed space $\Phi(x)$