

# Least Squares and M(C)LE

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

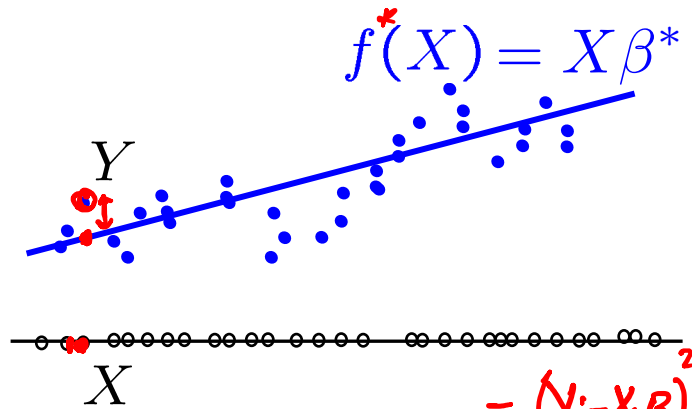
Least squares

Intuition: Signal plus (zero-mean) Noise model

$$Y = f^*(X) + \epsilon = X\beta^* + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad Y \sim \mathcal{N}(X\beta^*, \sigma^2 \mathbf{I})$$

$P(Y|X) =$



$$\hat{\beta}_{\text{MLE}} = \arg \max_{\beta} \log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n)$$

Conditional log likelihood

$$P(Y_i|X_i) \propto e^{-\frac{(Y_i - X_i \beta)^2}{2\sigma^2}}$$

$$= \arg \max_{\beta} \log \prod_{i=1}^n P(Y_i|X_i) = \arg \max_{\beta} \sum_{i=1}^n -\frac{(Y_i - X_i \beta)^2}{2\sigma^2}$$

$$= \arg \min_{\beta} \sum_{i=1}^n (X_i \beta - Y_i)^2 = \hat{\beta}$$

Least Square Estimate is same as Maximum Conditional Likelihood Estimate under a Gaussian model !

# Regularized Least Squares and M(C)AP

What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible?

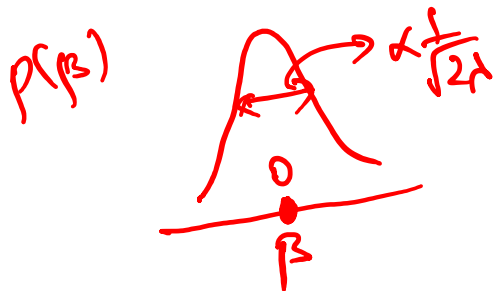
$$\text{Reg} \equiv \text{pen}(\beta) = \|\beta\|_2^2 \text{ or } \|\beta\|_1$$

$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} \underbrace{\log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n)}_{\text{Conditional log likelihood}} + \underbrace{\log p(\beta)}_{\text{log prior}}$$

$$\text{min}_{\beta} = \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

$$\begin{aligned} \log p(\beta) &\propto -\lambda \|\beta\|_2^2 \\ p(\beta) &\propto e^{-\lambda \|\beta\|_2^2} \\ &\sim N(\underline{0}, \frac{1}{2\lambda}) \end{aligned}$$

$$\begin{aligned} \beta &\sim N(0, \sigma^2) \\ e^{-\frac{\|\beta\|_2^2}{2\sigma^2}} &\equiv \frac{1}{\lambda} \end{aligned}$$



# Regularized Least Squares and M(C)AP

What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

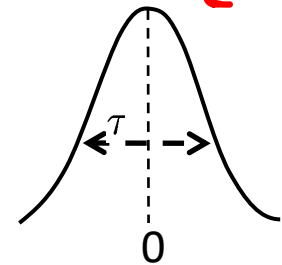
$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} \underbrace{\log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n)}_{\text{Conditional log likelihood}} + \underbrace{\log p(\beta)}_{\text{log prior}}$$

$\beta \sim \mathcal{N}(0, \Sigma)$      $p(\beta) \propto e^{-\frac{\beta^T \Sigma^{-1} \beta}{2}}$      $\Sigma = \begin{bmatrix} \frac{1}{\tau^2} & 0 \\ 0 & \frac{1}{\tau^2} \dots \end{bmatrix}$

1) Gaussian Prior

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

$$p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

Ridge Regression

constant( $\sigma^2, \tau^2$ )

$$\hat{\beta}_{\text{MAP}} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

# Regularized Least Squares and M(C)AP

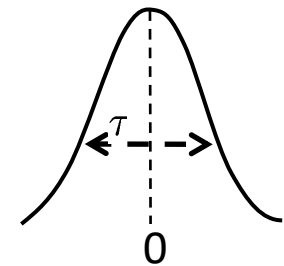
What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} \underbrace{\log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n)}_{\text{Conditional log likelihood}} + \underbrace{\log p(\beta)}_{\text{log prior}}$$

1) Gaussian Prior

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I})$$

$$p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$



$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2$$

constant( $\sigma^2, \tau^2$ )

**Ridge Regression**

Prior belief that  $\beta$  is Gaussian with zero-mean biases solution to “small”  $\beta$

# Regularized Least Squares and M(C)AP

What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible?

*Lasso*

$$\arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1$$

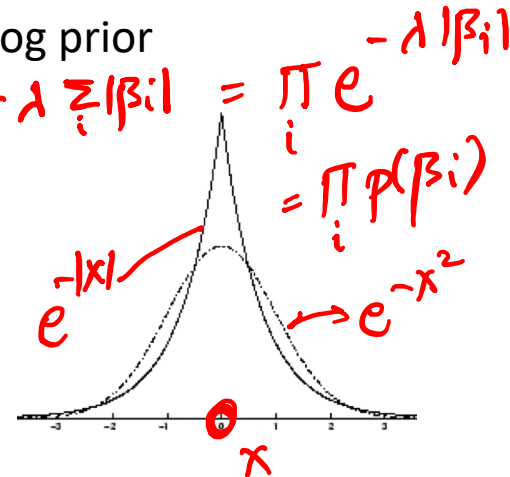
$$\hat{\beta}_{\text{MAP}} = \arg \max_{\beta} \underbrace{\log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n)}_{\text{Conditional log likelihood}} + \underbrace{\log p(\beta)}_{\text{log prior}}$$

*$-\log p(\beta) \propto \lambda \sum_i |\beta_i| \Rightarrow p(\beta) \propto e^{-\lambda \sum_i |\beta_i|} = \prod_i e^{-\lambda |\beta_i|} = \prod_i p(\beta_i)$*

II) Laplace Prior

$\beta_i \stackrel{iid}{\sim} \text{Laplace}(0, t)$

$p(\beta_i) \propto e^{-|\beta_i|/t}$



$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1$$

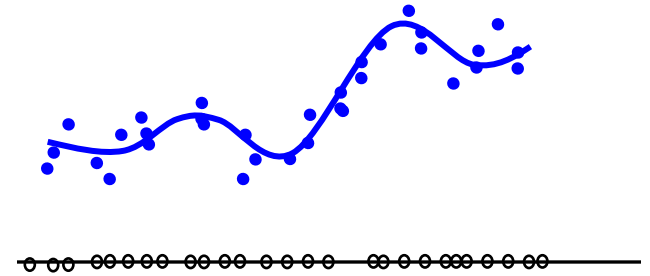
$\downarrow$   
 constant( $\sigma^2, t$ )

*Lasso*

Prior belief that  $\beta$  is Laplace with zero-mean biases solution to “sparse”  $\beta$

# Beyond Linear Regression

- Polynomial regression
- Regression with nonlinear features



- Kernelized Ridge Regression (Later)
- Local Kernel Regression (Later)

$$\begin{aligned}
 X &= [x^{(1)} \quad \dots \quad x^{(p)}] \\
 \rightarrow X &= \underbrace{[x^{(1)} \quad \dots \quad x^{(p)}]}_{p' \text{ features}} \underbrace{[x^{(1)2} \quad \dots \quad x^{(p)2} \quad x^{(1)}x^{(2)} \quad x^{(2)}x^{(3)} \quad \dots \quad x^{(1)}x^{(p)}]}_{(A^T A)^{-1}}
 \end{aligned}$$

# Polynomial Regression

degree m

Univariate (1-dim)  $f(X) = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_m X^m = \mathbf{X}\beta$

case:

$d=1$

where  $\mathbf{X} = [1 \ X \ X^2 \ \dots \ X^m]$ ,  $\beta = [\beta_1 \ \dots \ \beta_m]^T$

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \text{ or } (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$\hat{f}_n(X) = \mathbf{X}\hat{\beta}$$

where  $\mathbf{A} = \begin{bmatrix} 1 & X_1 & X_1^2 & \dots & X_1^m \\ \vdots & & & \ddots & \vdots \\ 1 & X_n & X_n^2 & \dots & X_n^m \end{bmatrix}$

original linear  
 $\mathbf{A} = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$

Multivariate (p-dim)  $f(X) = \beta_0 + \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$

case:

$$+ \sum_{i=1}^p \sum_{j=1}^p \beta_{ij} X^{(i)} X^{(j)} + \sum_{i=1}^p \sum_{j=1}^p \sum_{k=1}^p X^{(i)} X^{(j)} X^{(k)}$$

+ ... terms up to degree m

# Polynomial Regression

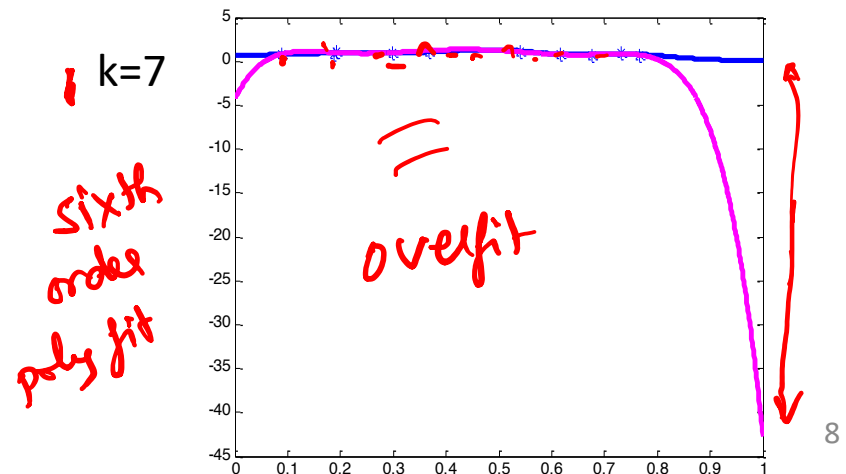
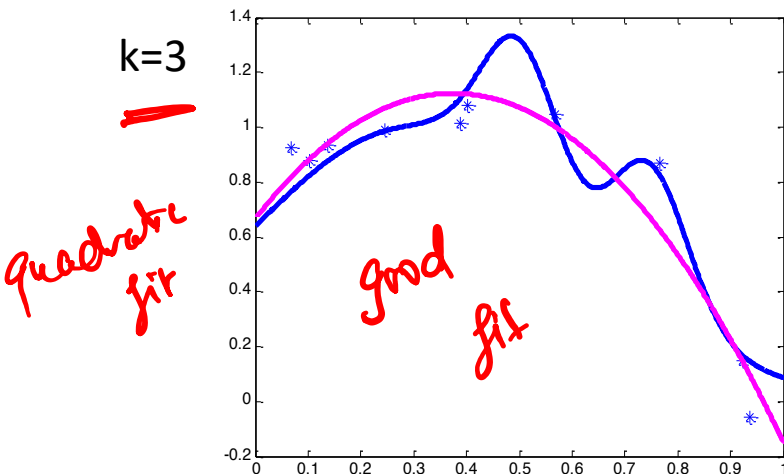
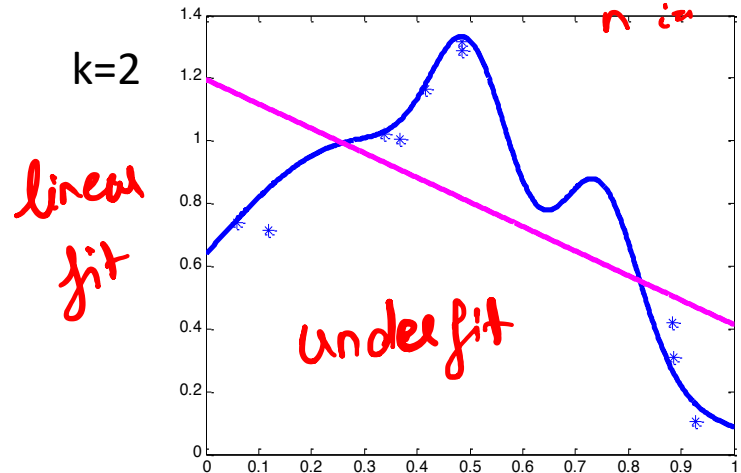
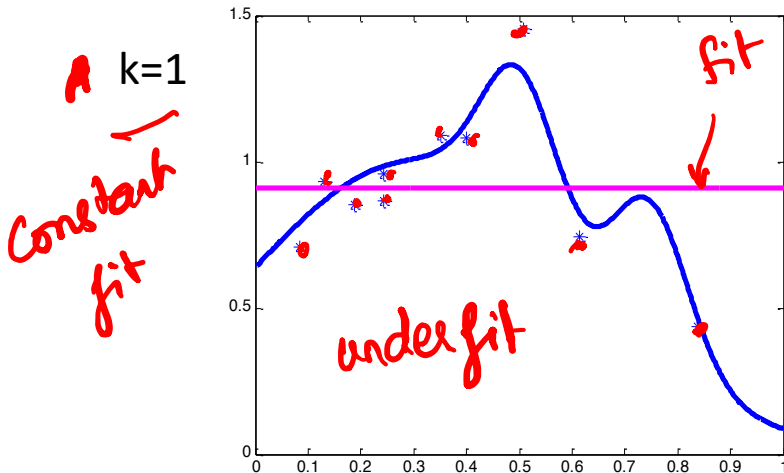
$y_i = f^*(x_i) + \epsilon_i$

blue function

Polynomial of order  $k$ , equivalently of degree up to  $k-1$

$E[(y - f(x))^2]$

$\rightarrow \frac{1}{n} \sum_{i=1}^n$





# Regression with nonlinear features

$$f(X) = \sum_{j=0}^m \beta_j X^j = \sum_{j=0}^m \beta_j \phi_j(X)$$

Weight of each feature ← Nonlinear features

In general, use any nonlinear features

e.g.  $e^X$ ,  $\log X$ ,  $1/X$ ,  $\sin(X)$ , ...

$$\hat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$$

or

$$(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y}$$

$$\mathbf{A} = \begin{bmatrix} \phi_0(X_1) & \phi_1(X_1) & \dots & \phi_m(X_1) \\ \vdots & & \ddots & \vdots \\ \phi_0(X_n) & \phi_1(X_n) & \dots & \phi_m(X_n) \end{bmatrix}$$

$$\hat{f}_n(X) = \mathbf{X} \hat{\beta}$$

$$\mathbf{X} = [\phi_0(X) \quad \phi_1(X) \quad \dots \quad \phi_m(X)]$$