# Regularized Linear Regression
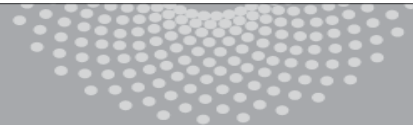
Aarti Singh

Machine Learning 10-701

Mar 22, 2023
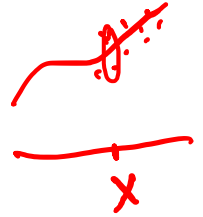
# Mean square error regression

Optimal predictor:
$$f^* = \arg\min_f \mathbb{E}[(f(X) - Y)^2]$$

$$= E[Y|X]$$

$$P(X,Y) \text{ known}$$

Empirical Minimizer:
$$\widehat{f}_n = \arg\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^{n} (f(X_i) - Y_i)^2$$

**Class of predictors**

$$\{X_i, Y_i\}_{i=1}^{n} \overset{iid}{\sim} P(X,Y)$$

$\mathcal{F}$ -   Class of Linear functions ✓
  -   Class of Polynomial functions
  -   Class of nonlinear functions

$$f(X) = X\beta$$
$$= [X^{(1)} \ldots X^{(p)}] \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

2

# Least Square solution satisfies Normal Equations

$$(\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y}$$

p x p   p x1     p x1

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1}$$

$$A = \begin{bmatrix} X_1^{(1)} \cdots X_1^{(p)} \\ \vdots \\ X_n^{(1)} \cdots X_n^{(p)} \end{bmatrix}_{n \times p}$$

If $(\mathbf{A}^T\mathbf{A})$ is invertible,

1) If dimension p not too large, analytical solution:

$$\widehat{\beta} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y} \qquad \widehat{f}_n^L(X) = X\widehat{\beta}$$

2)  If dimension p is large, computing inverse is expensive O(p³)
      Gradient descent since objective is convex (**A**ᵀ**A** $\succeq$ 0)

$$\beta^{t+1} = \beta^t - \frac{\alpha}{2}\frac{\partial J(\beta)}{\partial \beta}\bigg|_t$$
$$= \beta^t - \alpha \, \mathbf{A}^T(\mathbf{A}\beta^t - Y)$$

3

# Linear regression solution satisfies Normal Equations

$$(\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y}$$

p x p   p x1      p x1

$$A = \begin{bmatrix} X_1^{(1)} \cdots X_1^{(p)} \\ \vdots \\ X_r^{(1)} - X_n^{(p)} \end{bmatrix}_{n \times p}$$

When is $(\mathbf{A}^T\mathbf{A})$ invertible ?

Recall: Full rank matrices are invertible. What is rank of $(\mathbf{A}^T\mathbf{A})$ ?

$$A_{n\times p} = U S V^T$$
$$\underset{n\times r}{\downarrow} \quad \underset{r\times p}{\downarrow}$$
$$r\times r$$

$$\text{rank}(A) = r \leq p \qquad r \leq \min(n,p)$$

$$S = \begin{bmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_r \end{bmatrix}$$

$$A^T A = (USV^T)^T USV$$
$$= V S U^T U S V = V S^2 V^T$$
$$\underbrace{\phantom{V S^2 V^T}}$$

no. of datapoints
$$\uparrow$$
if $n < p \nearrow$ no. of features

then $A^T A$ is not invertible

$$\Rightarrow eig(A^T A) = s_1^2 \cdots s_r^2$$

# Linear regression solution satisfies Normal Equations

$$(\mathbf{A}^T\mathbf{A})\widehat{\beta} = \mathbf{A}^T\mathbf{Y} \Leftarrow$$

p x p   p x1        p x1

p equations in p unknowns ($\hat{\beta}$)

When is $(\mathbf{A}^T\mathbf{A})$ invertible ?
Recall: Full rank matrices are invertible. What is rank of $(\mathbf{A}^T\mathbf{A})$ ?

If $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$, then

S - r x r

normal equations $(\mathbf{S}\mathbf{V}^\top)\hat{\beta} = (\mathbf{U}^\top\mathbf{Y})$ ⟵

r x p   p x 1        r x 1

r equations in p unknowns. Under-determined if r < p, hence no unique solution.

# Regularized Least Squares

$$Y = \begin{bmatrix} \varphi_1 \\ \vdots \\ Y_n \end{bmatrix} \quad A = \begin{bmatrix} X_1^{(1)} \cdots X_1^{(p)} \\ X_n^{(1)} \cdots X_n^{(p)} \end{bmatrix}$$

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

r equations , p unknowns – underdetermined system of linear equations
many feasible solutions

Need to constrain solution further

e.g. bias solution to "small" values of β (small changes in input don't translate to large changes in output)

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

Ridge Regression
(l2 penalty)

$$= \arg \min_\beta \ (\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) + \lambda\|\beta\|_2^2 \qquad \lambda \geq 0$$

$$2 A^T A \beta - 2 A^T Y + 2 \lambda \beta = 0$$

6

# Ridge Regression

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

Ridge Regression
(l2 penalty)

$$= \arg \min_{\beta} \ (\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) + \lambda\|\beta\|_2^2$$

$$\lambda \geq 0$$

$$= \qquad 2A^T A \beta - 2 A^T Y + \underbrace{2\lambda\beta}_{2\lambda I \beta}$$

$$= \ 2(A^T A + \lambda I)\beta - 2A^T Y \ = 0$$

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

Is $(\mathbf{A}^\top \mathbf{A} + \lambda \mathbf{I})$ invertible ?

always $\lambda > 0$ ✓

$M + \lambda I$

$\Sigma = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_p \end{bmatrix}$

$V \Sigma V^T + \lambda V V^T$

$V(\Sigma + \lambda I) V^T$

$eig(M + \lambda I) = eig(M) + \lambda \geq \lambda > 0$

7

# Understanding regularized Least Squares

$$\min_{\beta}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) + \lambda\text{pen}(\beta) = \min_{\beta} J(\beta) + \lambda\text{pen}(\beta)$$

Ridge Regression:
$$\text{pen}(\beta) = \|\beta\|_2^2$$

βs with constant $J(\beta)$
(level sets of $J(\beta)$)

Unregularized Least Squares solution

$\beta$s with constant l2 norm
(level sets of pen($\beta$))

$\beta_2$

$\beta_1$

$\hat{\beta}$

# Regularized Least Squares

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

r equations , p unknowns – underdetermined system of linear equations
many feasible solutions
Need to constrain solution further

e.g. bias solution to "small" values of β (small changes in input don't translate to large changes in output)

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

Ridge Regression
(l2 penalty)

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

$\lambda \geq 0$

Lasso
(l1 penalty)

Many β can be zero – many inputs are irrelevant to prediction in high-dimensional settings (typically intercept term not penalized)

# Regularized Least Squares

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

r equations , p unknowns – underdetermined system of linear equations
                        many feasible solutions
Need to constrain solution further

e.g. bias solution to "small" values of $\beta$ (small changes in input don't translate to large changes in output)

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

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

**Ridge Regression**
**(l2 penalty)**

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

$$\lambda \geq 0$$

**Lasso**
**(l1 penalty)**

No closed form solution, but can optimize using sub-gradient descent (packages available)

# Ridge Regression vs Lasso

$$\min_{\beta}(\mathbf{A}\beta - \mathbf{Y})^T(\mathbf{A}\beta - \mathbf{Y}) + \lambda \mathrm{pen}(\beta) = \min_{\beta} J(\beta) + \lambda \mathrm{pen}(\beta)$$

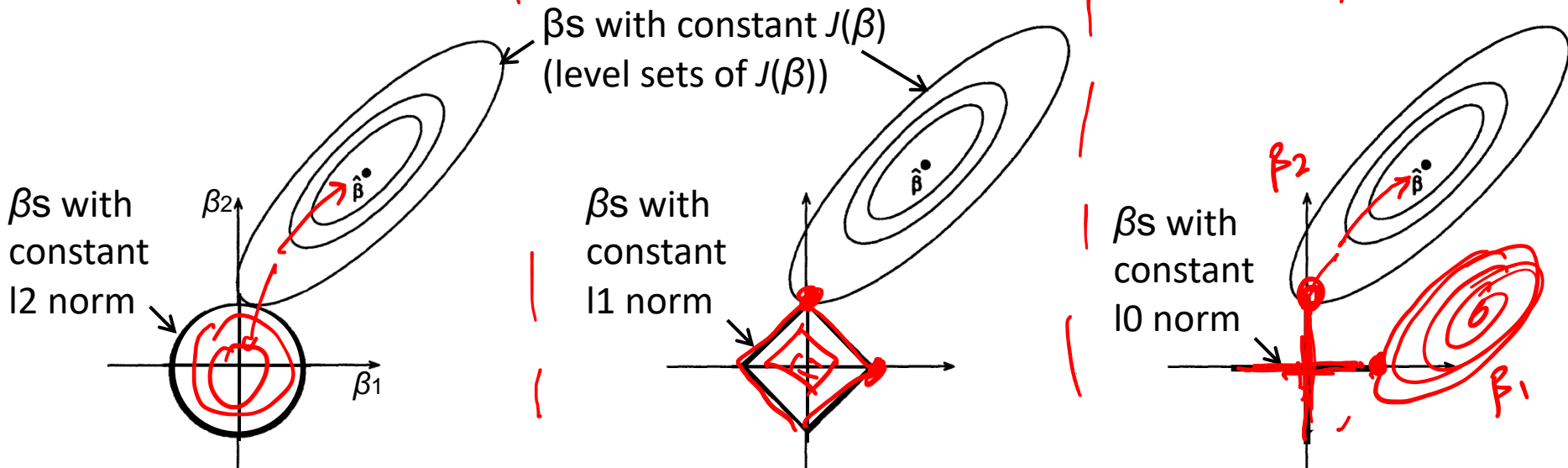$$\|\beta\|_0 = \sum_i 1_{\beta_i \neq 0}$$

Ridge Regression:
$$\mathrm{pen}(\beta) = \|\beta\|_2^2$$

Lasso:
$$\mathrm{pen}(\beta) = \|\beta\|_1$$

Ideally l0 penalty, but optimization becomes non-convex

βs with constant $J(\beta)$
(level sets of $J(\beta)$)

βs with constant l2 norm

βs with constant l1 norm

βs with constant l0 norm

**Lasso (l1 penalty) results in sparse solutions – vector with more zero coordinates**
**Good for high-dimensional problems – don't have to store all coordinates,**
**interpretable solution!**

11

# Matlab example

```
clear all
close all
```

$n < p$

```
n = 80;      % datapoints
p = 100;   % features
k = 10;       % non-zero features

rng(20);
X = randn(n,p);
weights = zeros(p,1);
weights(1:k) = randn(k,1)+10;
noise = randn(n,1) * 0.5;
Y = X*weights +  noise;

Xtest = randn(n,p);
noise = randn(n,1) * 0.5;
Ytest = Xtest*weights + noise;
```

```
lassoWeights = lasso(X,Y,'Lambda',1,
'Alpha', 1.0);
Ylasso = Xtest*lassoWeights;
norm(Ytest-Ylasso)

ridgeWeights = lasso(X,Y,'Lambda',1,
'Alpha', 0.0001);
Yridge = Xtest*ridgeWeights;
norm(Ytest-Yridge)

stem(lassoWeights)
pause
stem(ridgeWeights)
```
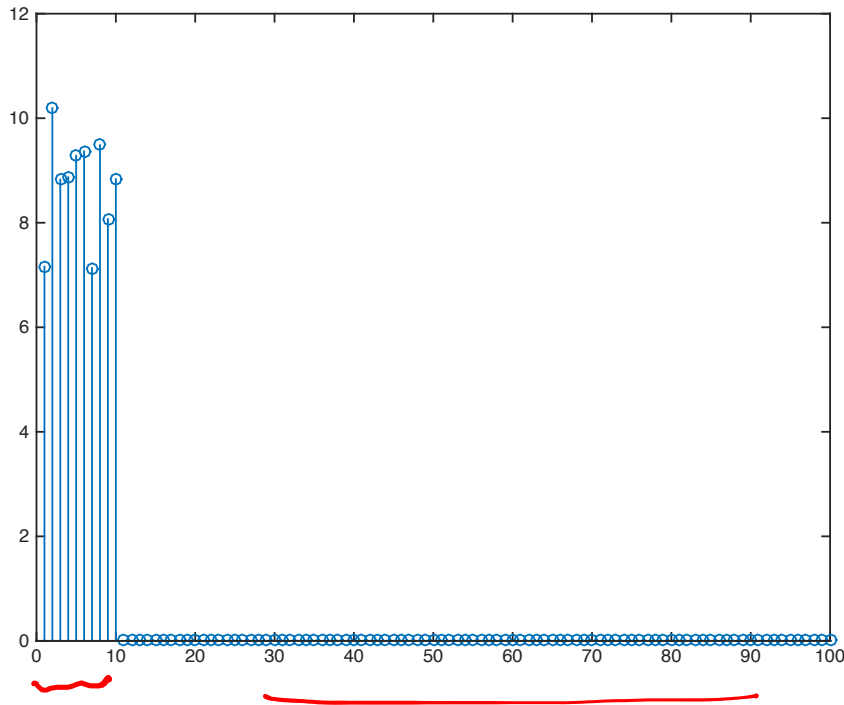
# Matlab example

Test MSE = 33.7997

Test MSE = 185.9948
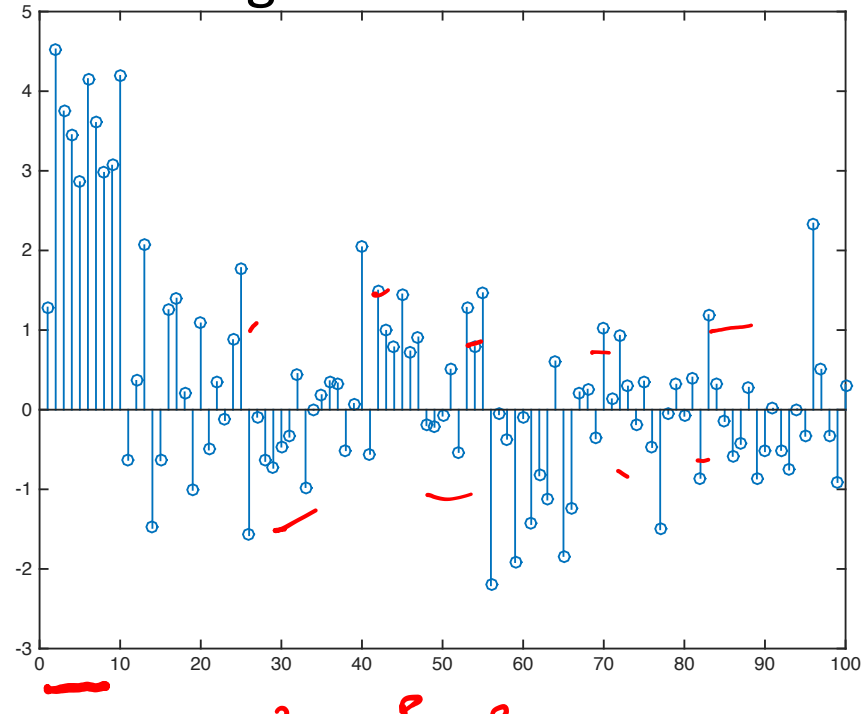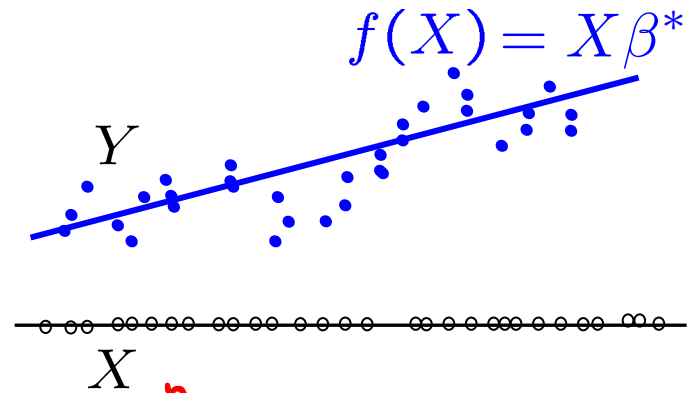


Lasso Coefficients

Ridge Coefficients

$$\|\beta\|_2^2 = \sum_{j=1}^{p} \beta_i^2$$

# Least Squares and M(C)LE

$E[Y|x]$ ✓

Intuition: Signal plus (zero-mean) Noise model

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

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

$|x$

$f(X) = X\beta^*$

$Y$

$X$

$$\widehat{\beta}_{\mathsf{MLE}} = \arg\max_{\beta} \log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=1}^n) \sim \log \prod_{i=1}^{n} \mathcal{N}(X_i\beta^*, \sigma^2 \mathbf{I})$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\text{Conditional log likelihood}}$$

$$= \arg\max_{\beta} \sum_{i=1}^{n} \log \underbrace{\frac{1}{(\sqrt{2\pi\sigma^2})^d}}_{a} \underbrace{\exp\left(-\frac{(Y_i - X_i\beta)^2}{2\sigma^2}\right)}_{b}$$

$\log(ab) = \log a + \log b$

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

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

**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 ?

$P(\beta)$ — prior

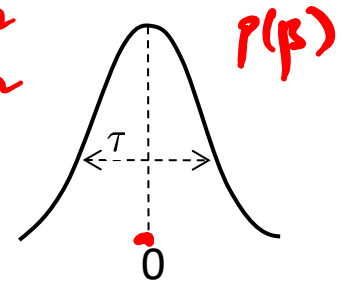$P(\beta | Data) \propto P(Data | \beta) \cdot P(\beta)$

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

I) Gaussian Prior

$\log P(\beta) \propto -\dfrac{\beta^T \beta}{2\tau^2} = -\dfrac{\|\beta\|^2}{2\tau^2}$

$P(\beta)$

$$\beta \sim \mathcal{N}(0, \tau^2 \mathbf{I}) \qquad\qquad p(\beta) \propto e^{-\beta^T \beta / 2\tau^2}$$

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

**Ridge Regression**

$\text{constant}(\sigma^2, \tau^2)$

Prior belief that β is Gaussian with zero-mean biases solution to "small" β

# Regularized Least Squares and M(C)AP

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

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\max_{\beta} \underbrace{\log p(\{Y_i\}_{i=1}^n | \beta, \sigma^2, \{X_i\}_{i=}^n}_{\text{Conditional log likelihood}} + \underbrace{\log p(\beta)}_{\text{log prior}}$$
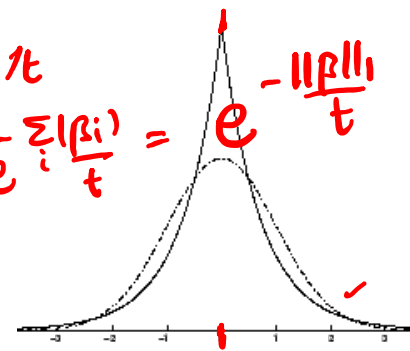
$e^{-\|\beta\|_1}$

II) Laplace Prior

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

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

$$p(\beta) = \prod_i p(\beta_i) \propto \prod_i e^{-|\beta_i|/t} = e^{-\frac{\sum_i |\beta_i|}{t}} = e^{-\frac{\|\beta\|_1}{t}}$$

$$\widehat{\beta}_{\mathsf{MAP}} = \arg\min_{\beta} \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda\|\beta\|_1 \qquad \textcolor{red}{\text{Lasso}}$$

$$\downarrow$$
$$\text{constant}(\sigma^2, t)$$

Prior belief that β is Laplace with zero-mean biases solution to "sparse" β

# Polynomial Regression

degree m

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

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

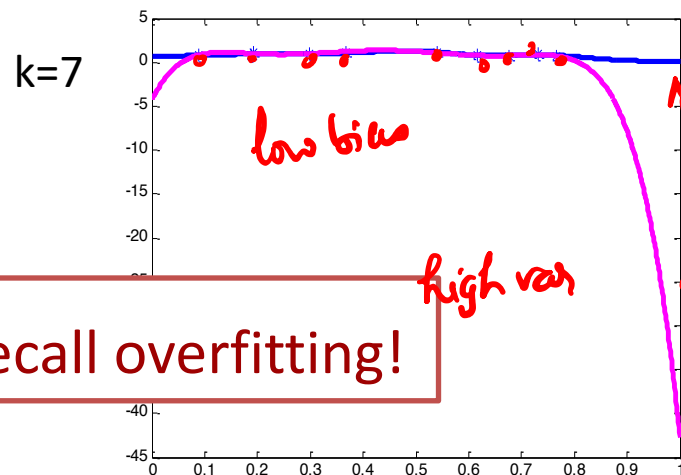$$\widehat{\beta} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{Y} \qquad\qquad \widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$$
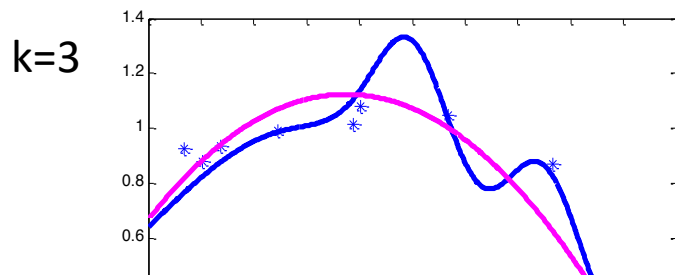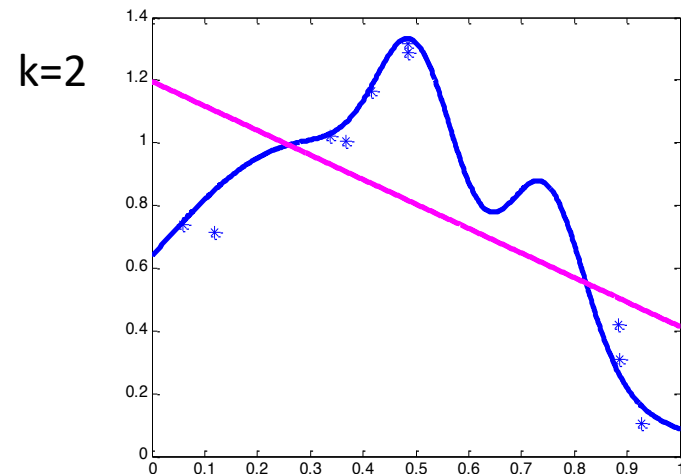
where $\mathbf{A} = \begin{bmatrix} 1 & X_1 & X_1^2 & \ldots & X_1^m \\ \vdots & & & \ddots & \vdots \\ 1 & X_n & X_n^2 & \ldots & X_n^m \end{bmatrix}$

Multivariate (p-dim) case:
$$\begin{aligned} f(X) \ = \ & \beta_0 + \beta_1 X^{(1)} + \beta_2 X^{(2)} + \cdots + \beta_p X^{(p)} \\ & + \sum_{i=1}^{p}\sum_{j=1}^{p}\beta_{ij}X^{(i)}X^{(j)} + \sum_{i=1}^{p}\sum_{j=1}^{p}\sum_{k=1}^{p}X^{(i)}X^{(j)}X^{(k)} \\ & + \ldots \text{terms up to degree m} \end{aligned}$$

17

# Polynomial Regression

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



What is the right order? Recall overfitting!

# 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

$\phi_0(X)$

$\phi_1(X)$

$\phi_2(X)$

In general, use any nonlinear features

e.g. $e^X$, log X, 1/X, sin(X), …

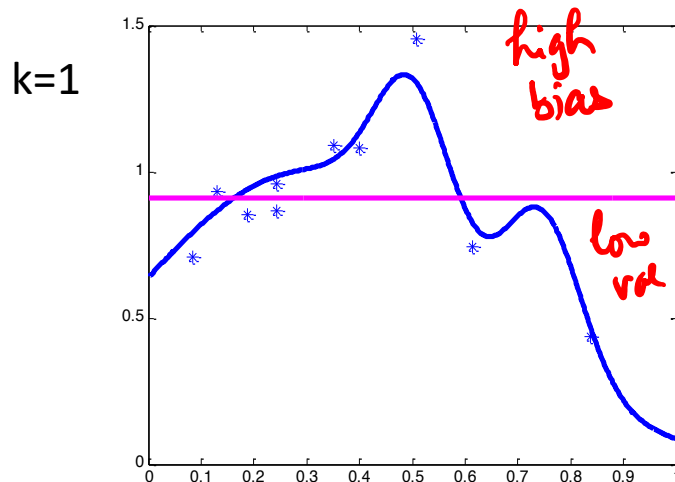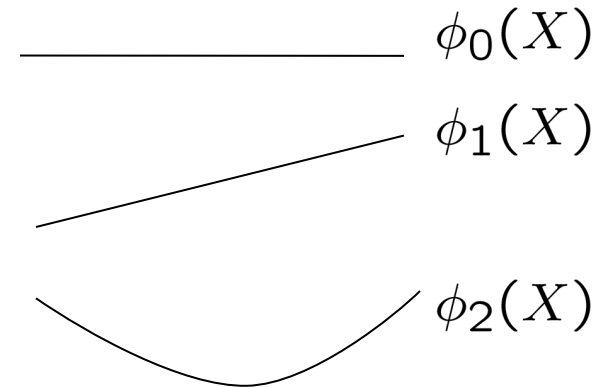$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y}$$

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

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

$$\mathbf{X} = [\phi_0(X) \; \phi_1(X) \; \ldots \; \phi_m(X)]$$

# Poll

- The maximum likelihood estimate of model parameter α for the random variable $y \sim N(\alpha\, x_1 x_2^3\,, \sigma^2)$, where $x_1$ and $x_2$ are random variables, can be learned using linear regression on n iid samples of $(x_1, x_2, y)$

    – True
    – False

$$y \sim N(\alpha\, x_1 x_2\,, \sigma^2) \checkmark$$

$$y \sim N(\alpha x_1 + x_2\,, \sigma^2) \checkmark$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \rightarrow \quad x_1 x_2^3 = z$$

$$y \sim N(\alpha z\,, \sigma^2)$$

# Can we kernelize linear regression?

# Linear (Ridge) regression

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})^{-1}\mathbf{A}^T\mathbf{Y}$$

$p \times p$

Recall

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix}$$

$$X_i \cdot X_j = X_i^T X_j$$

$$\phi(x_i) \cdot \phi(x_j)$$

$$K(x_i, x_j)$$

Hence $\mathbf{A}^T\mathbf{A}$ is a p x p matrix whose entries denote the (sample) correlation between the features

NOT inner products between the data points – the inner product matrix would be $\mathbf{A}\mathbf{A}^T$ which is n x n (also known as Gram matrix)

Using dual formulation, we can write the solution in terms of $\mathbf{A}\mathbf{A}^T$

22

# Ridge regression

$$\min_\beta \sum_{i=1}^n (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})^{-1}\mathbf{A}^T\mathbf{Y}$$

Similarity with SVMs

Primal problem:

$$\min_{\beta, z_i} \sum_{i=1}^n z_i^2 + \lambda\|\beta\|_2^2$$

$$\text{s.t. } z_i = Y_i - X_i\beta$$

$\alpha_i$

SVM Primal problem:

$$\min_{w, \xi_i} C \sum_{i=1}^n \xi_i + \frac{1}{2}\|w\|_2^2$$

$$\text{s.t. } \xi_i = \max(1 - Y_i X_i w, 0)$$

Lagrangian:

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

$\alpha_i$ – Lagrange parameter, one per training point

# Ridge regression (dual)

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})^{-1}\mathbf{A}^T\mathbf{Y}$$

Dual problem:

$$\max_{\alpha} \min_{\beta, z_i} \sum_{i=1}^{n} z_i^2 + \lambda\|\beta\|^2 + \sum_{i=1}^{n} \alpha_i(z_i - Y_i + X_i\beta)$$

$\alpha = \{\alpha_i\}$ for i = 1,..., n

Taking derivatives of Lagrangian wrt $\beta$ and $z_i$ we get:

$$\beta = -\frac{1}{2\lambda}\mathbf{A}^\top\alpha \qquad z_i = -\frac{\alpha_i}{2}$$

Dual problem: $\quad \max_{\alpha} \quad -\frac{\alpha^\top\alpha}{4} - \frac{1}{4\lambda}\alpha^\top\mathbf{A}\mathbf{A}^\top\alpha - \alpha^\top\mathbf{Y}$

n-dimensional optimization problem

# Ridge regression (dual)

$$\min_{\beta} \sum_{i=1}^{n} (Y_i - X_i\beta)^2 + \lambda\|\beta\|_2^2 \qquad \widehat{\beta} = (\mathbf{A}^T\mathbf{A} + \lambda\mathbf{I})^{-1}\mathbf{A}^T\mathbf{Y}$$

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

Dual problem:

$$\max_{\alpha} \; -\frac{\alpha^\top\alpha}{4} - \frac{1}{4\lambda}\alpha^\top\mathbf{A}\mathbf{A}^\top\alpha - \alpha^\top\mathbf{Y} \qquad \Rightarrow \widehat{\alpha} = -\left(\frac{\mathbf{A}\mathbf{A}^\top}{\lambda} + \mathbf{I}\right)^{-1} 2\,\mathbf{Y}$$

can get back $\quad \hat{\beta} = -\dfrac{1}{2\lambda}\mathbf{A}^\top\hat{\alpha} \quad = \mathbf{A}^\top(\mathbf{A}\mathbf{A}^\top + \lambda\mathbf{I})^{-1}\mathbf{Y}$

Weighted average of
training points

Weight of each training point (but typically not sparse)

25

# Kernelized ridge regression

$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \quad \checkmark$$

Using dual, can re-write solution as:

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

How does this help?

- Only need to invert n x n matrix (instead of p x p or m x m)
- More importantly, kernel trick!

> $\mathbf{A}\mathbf{A}^T$ involves only inner products between the training points
> BUT still have an extra $\mathbf{A}^T$

Recall the predicted label is $\widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$

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

$$\underbrace{\phantom{\mathbf{X}\mathbf{A}^T}}_{K_{x,x_i}} \qquad \underbrace{\phantom{\mathbf{A}\mathbf{A}^T}}_{K_{x_i,x_i}}$$

$\mathbf{X}\mathbf{A}^T$ contains inner products between test point $\mathbf{X}$ and training points!

# Kernelized ridge regression

$$\widehat{\beta} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{Y} \qquad \widehat{f}_n(X) = \mathbf{X}\widehat{\beta}$$

Using dual, can re-write solution as:

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

How does this help?
- Only need to invert n x n matrix (instead of p x p or m x m)
- More importantly, kernel trick!

$$\widehat{f}_n(X) = \mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y} \quad \text{where} \quad \begin{array}{l} \mathbf{K}_X(i) = \boldsymbol{\phi}(X) \cdot \boldsymbol{\phi}(X_i) \\ \mathbf{K}(i,j) = \boldsymbol{\phi}(X_i) \cdot \boldsymbol{\phi}(X_j) \end{array}$$

Work with kernels, never need to write out the high-dim vectors

Ridge Regression with (implicit) nonlinear features $\boldsymbol{\phi}(X)$! $\quad f(X) = \phi(X)\beta$