



# Machine learning

## regression

Mohammad-Reza A. Dehaqani

[dehqani@ut.ac.ir](mailto:dehqani@ut.ac.ir)

Slides mainly adopted from Andrew Ng and Aarti course

# Supervised Learning Tasks



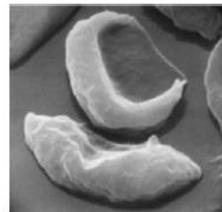
## Classification



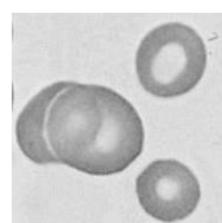
X = Document



Sports  
Science  
News



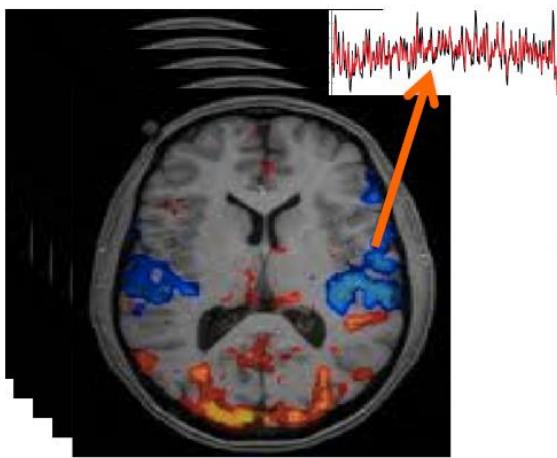
Anemic cell  
Healthy cell



X = Cell Image

Y = Diagnosis

## Regression



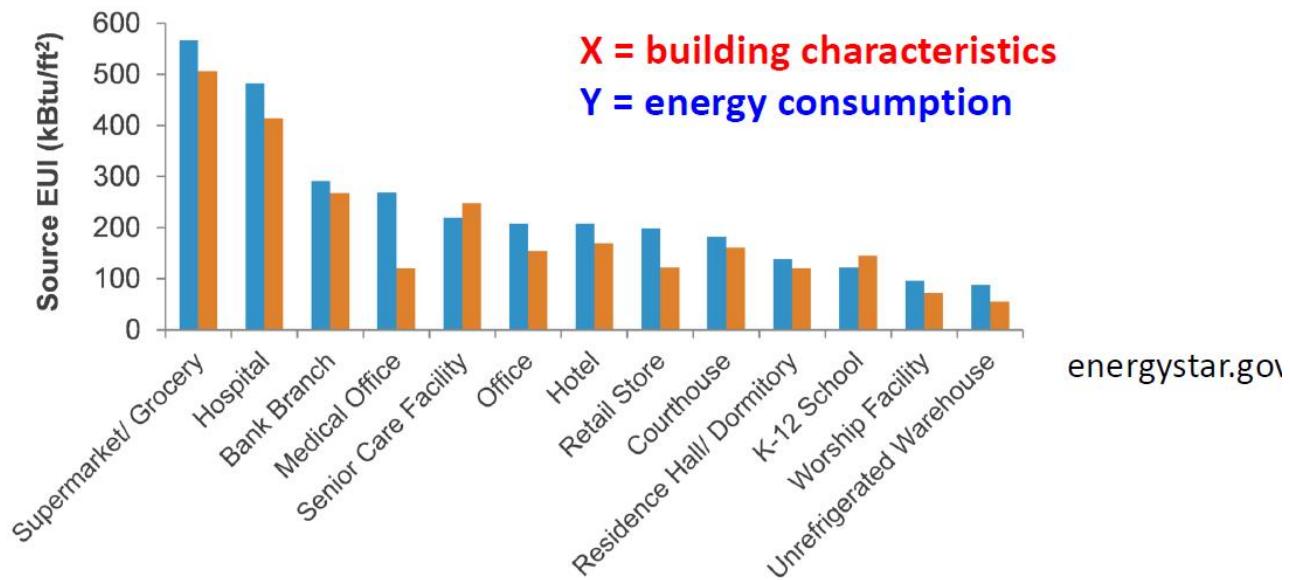
Y = Age of a subject

X = Brain Scan

# Regression Tasks



Estimating  
Energy Usage



Estimating  
Contamination



# Performance Measures



**Performance Measure:** Quantifies knowledge gained

$\text{loss}(Y, f(X))$  - Measure of closeness between true label  $Y$  and prediction  $f(X)$

Don't just want label of one test data (cell image), but any cell image  $X \in \mathcal{X}$

$$(X, Y) \sim P_{XY}$$

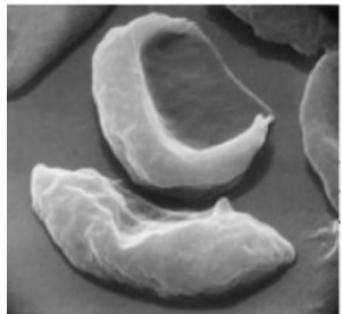
Given a cell image drawn randomly from the collection of all cell images, how well does the predictor perform on average?

$$\text{Risk } R(f) \equiv \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

# Performance Measures



**Performance Measure:** Risk  $R(f) \equiv \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$



➡ “Anemic cell”



➡ Share Price  
“\$ 24.50”

$\text{loss}(Y, f(X))$

$$\mathbf{1}_{\{f(X) \neq Y\}}$$

**0/1 loss**

Risk  $R(f)$

$$P(f(X) \neq Y)$$

**Probability of Error**

$$(f(X) - Y)^2$$

**square loss**

$$\mathbb{E}[(f(X) - Y)^2]$$

**Mean Square Error**

# Bayes Optimal Rule



Ideal goal: Construct **prediction rule**  $f^* : \mathcal{X} \rightarrow \mathcal{Y}$

$$f^* = \arg \min_f \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

*Bayes optimal rule*

Best possible performance:

*Bayes Risk*       $R(f^*) \leq R(f) \text{ for all } f$

**BUT... Optimal rule is not computable - depends on unknown  $P_{XY}$  !**

# Experience - Training Data

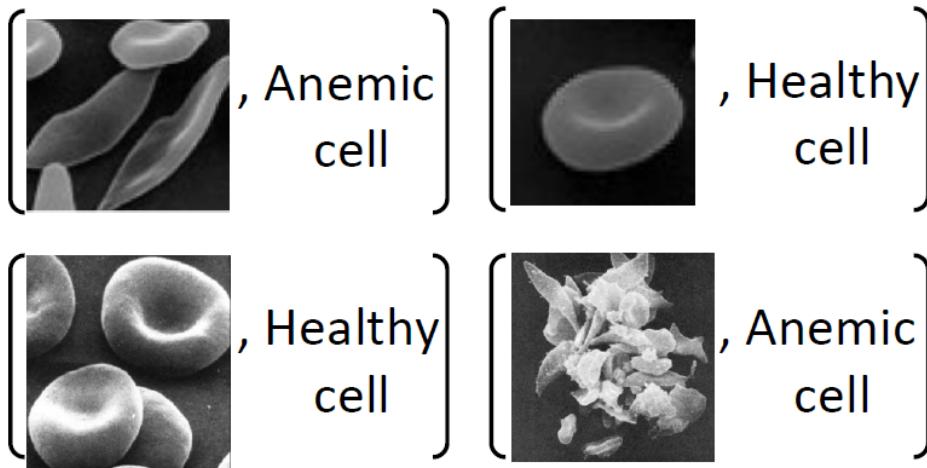


Can't minimize risk since  $P_{XY}$  unknown!

Training data (experience) provides a glimpse of  $P_{XY}$

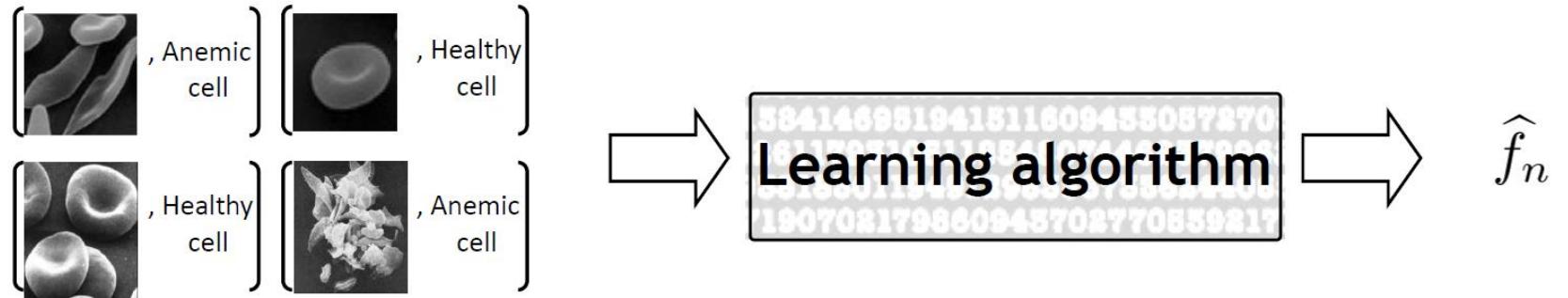
**(observed)**  $\{(X_i, Y_i)\}_{i=1}^n \stackrel{i.i.d.}{\sim} P_{XY}$  **(unknown)**

↳ independent, identically distributed



Provided by expert,  
measuring device,  
some experiment, ...

# Machine Learning Algorithm



Data  $\{(X_i, Y_i)\}_{i=1}^n$

$\hat{f}_n$  is a mapping from  $\mathcal{X} \rightarrow \mathcal{Y}$

$\hat{f}_n$  [Image of an anemic cell] = “Anemic cell”

Test data  $X$



# Empirical Risk Minimization

Optimal predictor:

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

Empirical Minimizer:

$$\hat{f}_n = \arg \min_{f \in \mathcal{F}} \left( \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \right)$$

**Empirical mean**

Law of Large Numbers:

$$\frac{1}{n} \sum_{i=1}^n [\text{loss}(Y_i, f(X_i))] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{XY} [\text{loss}(Y, f(X))]$$

# Restrict class of predictors



Optimal predictor:

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

Empirical Minimizer:

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

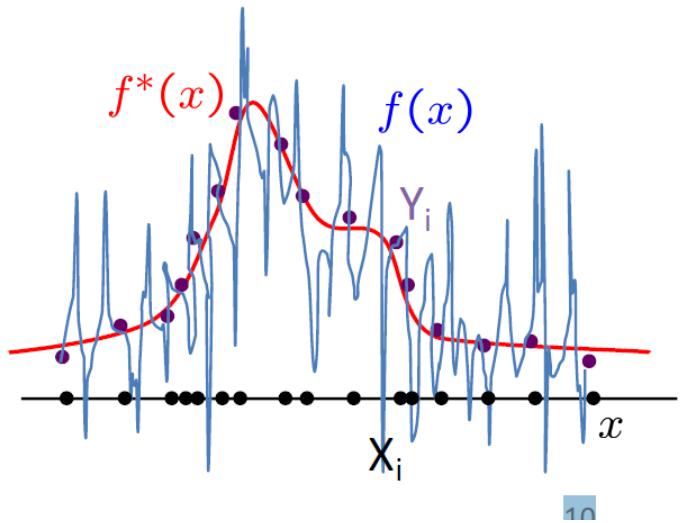
Class of predictors

Why?

Overfitting!

Empirical loss minimized by any function of the form

$$f(x) = \begin{cases} Y_i, & x = X_i \text{ for } i = 1, \dots, n \\ \text{any value,} & \text{otherwise} \end{cases}$$



# Restrict class of predictors



Optimal predictor:

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

Empirical Minimizer:

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

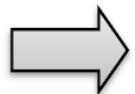
Class of predictors

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

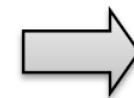
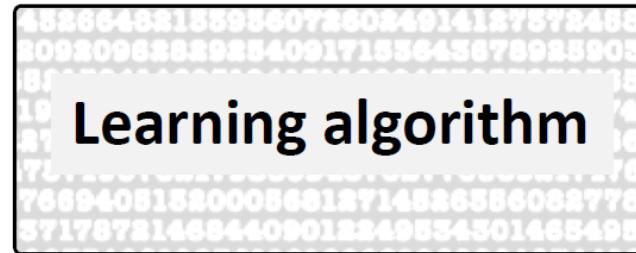
# Regression algorithms



Training data



$$\{(X_i, Y_i)\}_{i=1}^n$$



Prediction rule  
 $\hat{f}_n$

Linear Regression

Regularized Linear Regression – Ridge regression, Lasso

Polynomial Regression

Kernelized Ridge Regression

Gaussian Process Regression

Kernel regression, Regression Trees, Splines, Wavelet estimators, ...

# Linear Regression

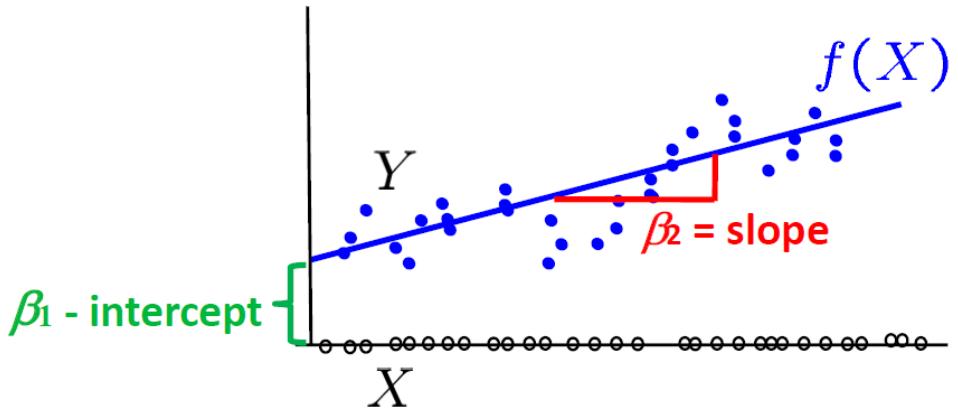


$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad \text{Least Squares Estimator}$$

$\mathcal{F}_L$  - Class of Linear functions

Uni-variate case:

$$f(X) = \beta_1 + \beta_2 X$$



Multi-variate case:

$$f(X) = f(X^{(1)}, \dots, X^{(p)}) = \beta_1 X^{(1)} + \beta_2 X^{(2)} + \dots + \beta_p X^{(p)}$$

$$= X\beta \quad \text{where} \quad X = [X^{(1)} \dots X^{(p)}], \quad \beta = [\beta_1 \dots \beta_p]^T$$

# Least Squares Estimator



$$\hat{f}_n^L = \arg \min_{f \in \mathcal{F}_L} \frac{1}{n} \sum_{i=1}^n (f(X_i) - Y_i)^2 \quad f(X_i) = X_i \beta$$



$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (X_i \beta - Y_i)^2 \quad \hat{f}_n^L(X) = X \hat{\beta}$$

$$= \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y})$$

$$\mathbf{A} = \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} X_1^{(1)} & \dots & X_1^{(p)} \\ \vdots & \ddots & \vdots \\ X_n^{(1)} & \dots & X_n^{(p)} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_n \end{bmatrix}$$

# Least Squares Estimator



$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

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

$$\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\hat{\beta}} = 0$$



# Least Square solution satisfies Normal Equations

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

p x p   p x 1      p x 1

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

$$\boxed{\hat{\boldsymbol{\beta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \quad \hat{f}_n^L(X) = X \hat{\boldsymbol{\beta}}}$$



# Least Square solution satisfies Normal Equations

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

p x p   p x 1      p x 1

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

$$\hat{\boldsymbol{\beta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \quad f_n^L(X) = X \hat{\boldsymbol{\beta}}$$

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

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

Now: What if  $(\mathbf{A}^T \mathbf{A})$  is invertible but expensive (p very large)?

# Gradient Descent

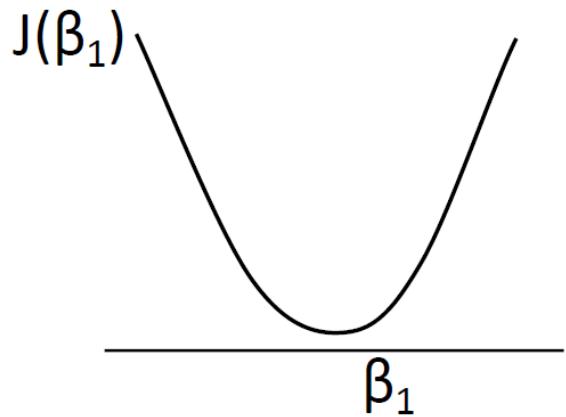


Even when  $(\mathbf{A}^T \mathbf{A})$  is invertible, might be computationally expensive if  $\mathbf{A}$  is huge.

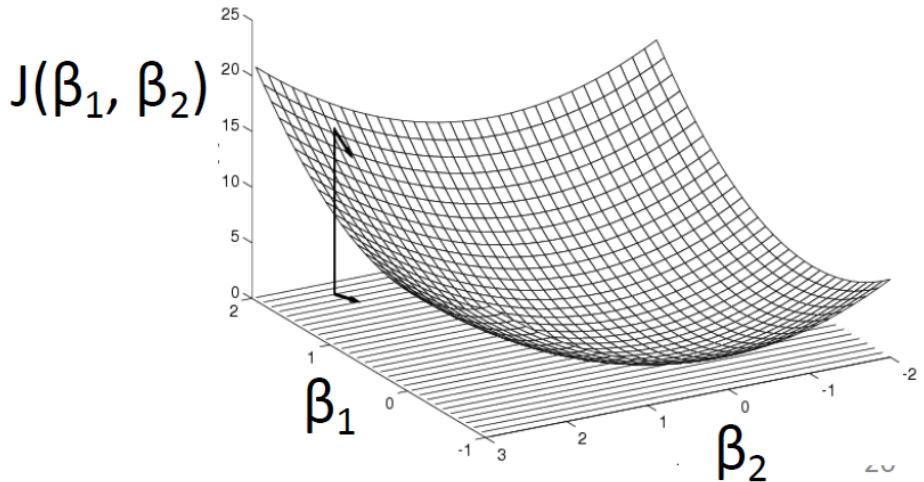
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Treat as optimization problem

Observation:  $J(\beta)$  is convex in  $\beta$ .



**How to find the minimizer?**



# Gradient Descent



Even when  $(\mathbf{A}^T \mathbf{A})$  is invertible, might be computationally expensive if  $\mathbf{A}$  is huge.

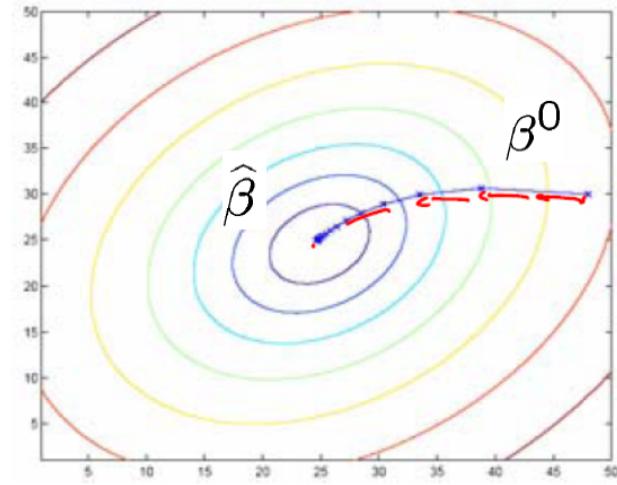
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Since  $J(\beta)$  is convex, move along negative of gradient

Initialize:  $\beta^0$

step size

$$\begin{aligned} \text{Update: } \beta^{t+1} &= \beta^t - \frac{\alpha \partial J(\beta)}{2} \Big|_t \\ &= \beta^t - \alpha \underbrace{\mathbf{A}^T (\mathbf{A}\beta^t - \mathbf{Y})}_{0 \text{ if } \hat{\beta} = \beta^t} \end{aligned}$$



Stop: when some criterion met e.g. fixed # iterations, or  $\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\beta^t} < \epsilon$ .



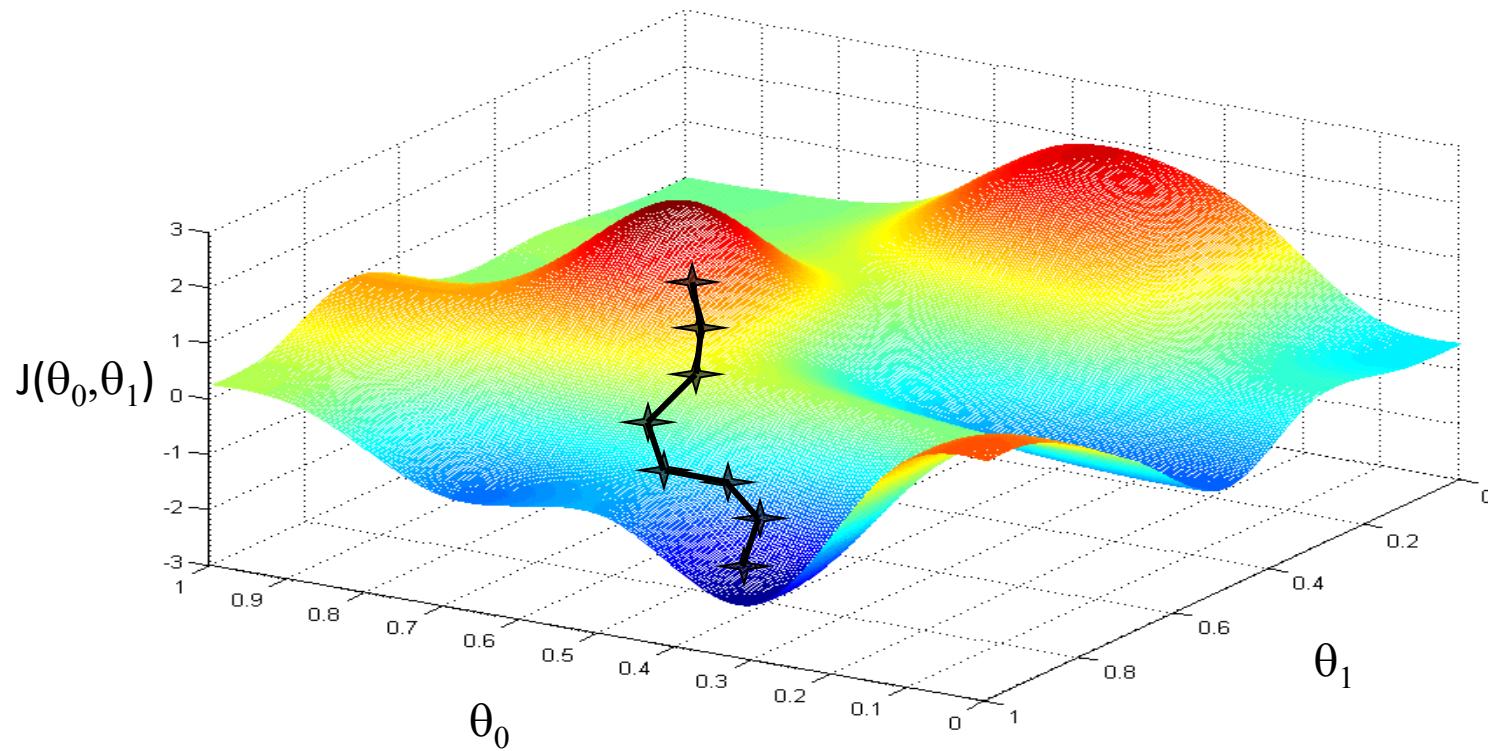
# “Batch” Gradient Descent

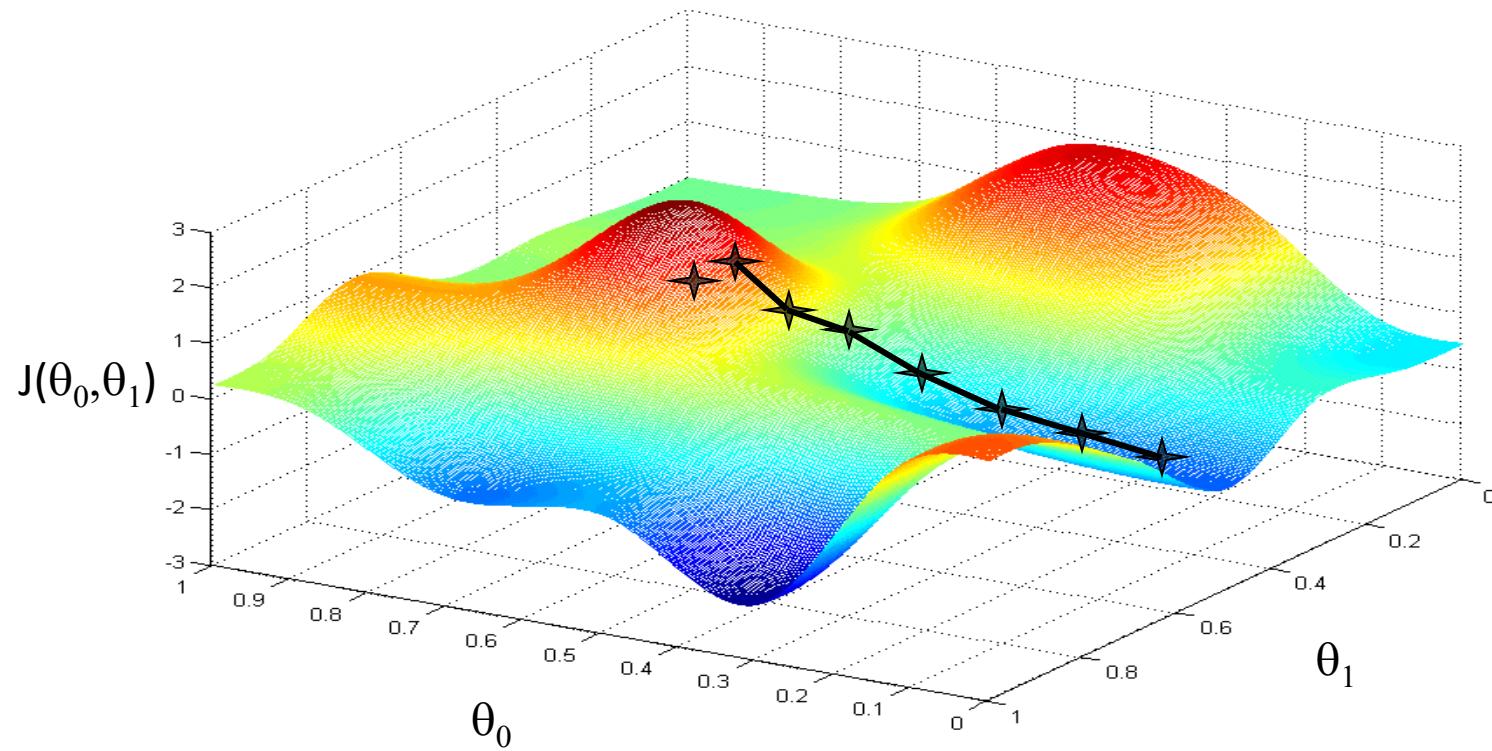
“Batch”: Each step of gradient descent uses all the training examples.



# Solution II: stochastic gradient descent

To update the parameters according to the **gradient of the error** with respect to that **single training example only**.

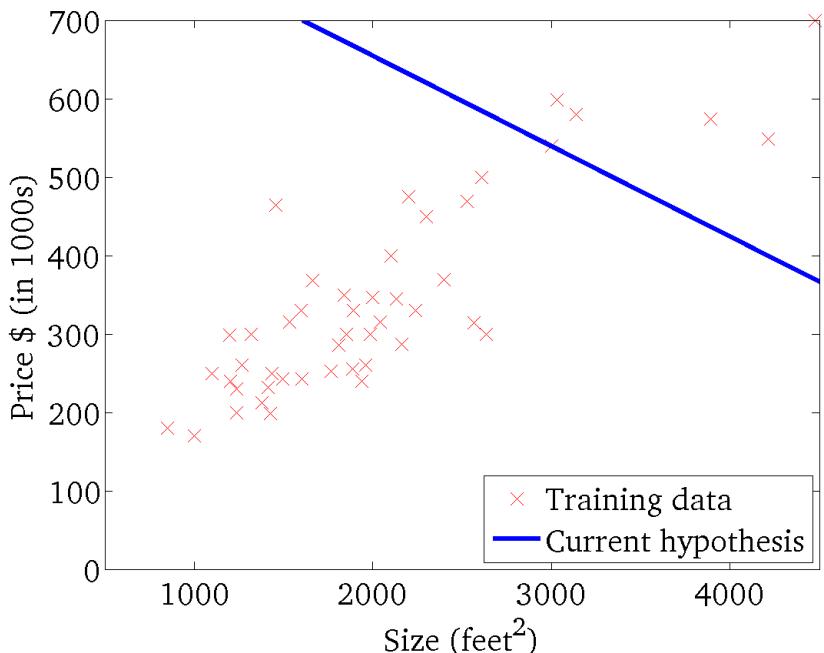






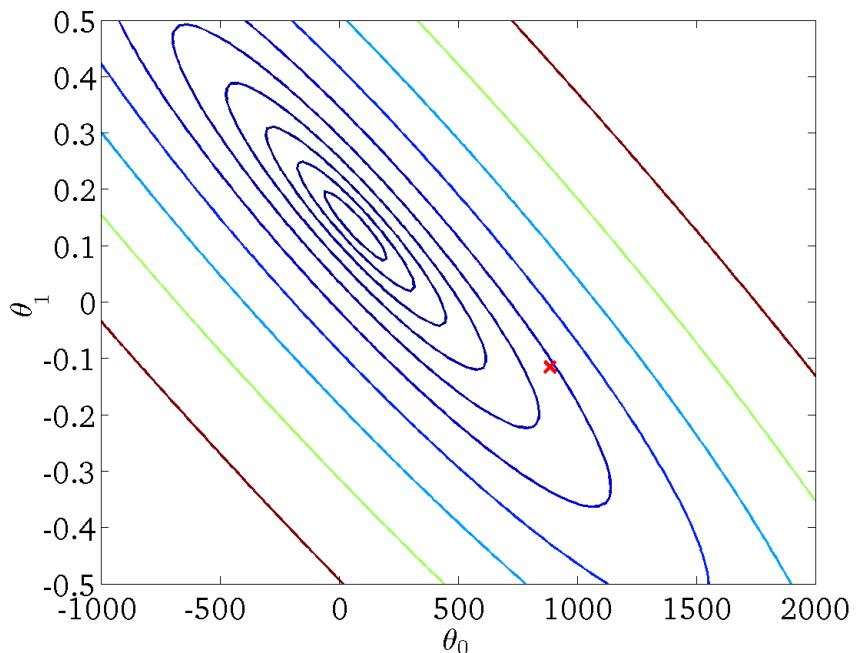
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

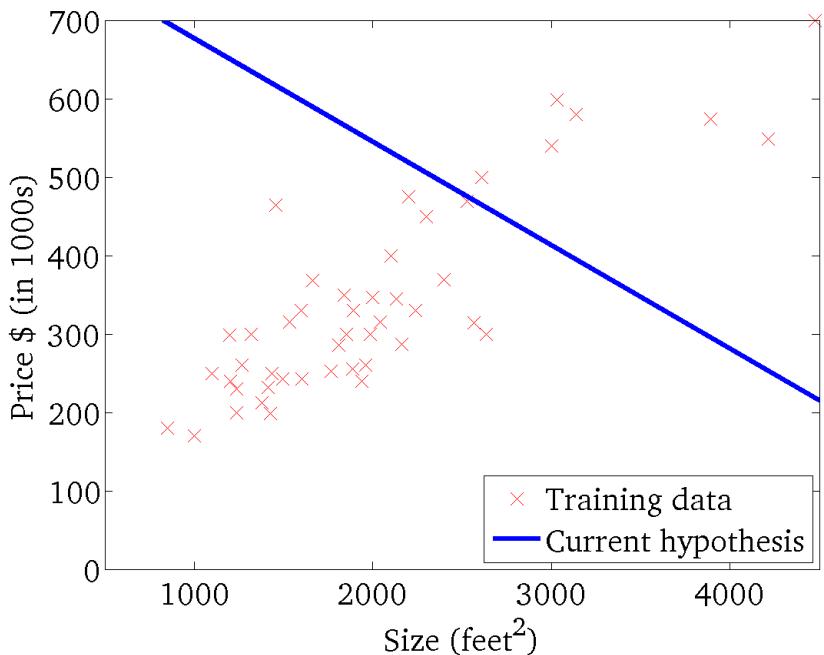
(function of the parameters  $\theta_0, \theta_1$ )





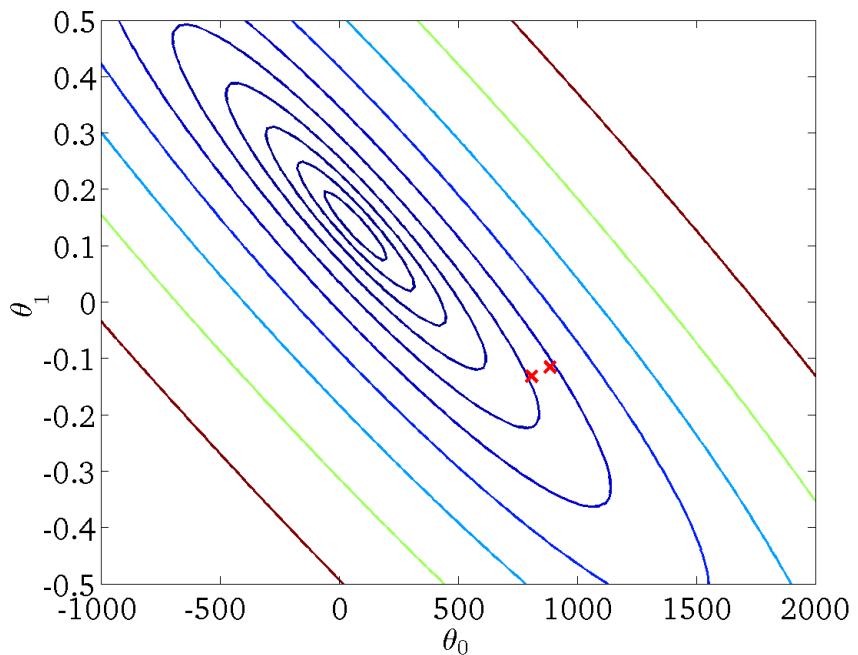
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

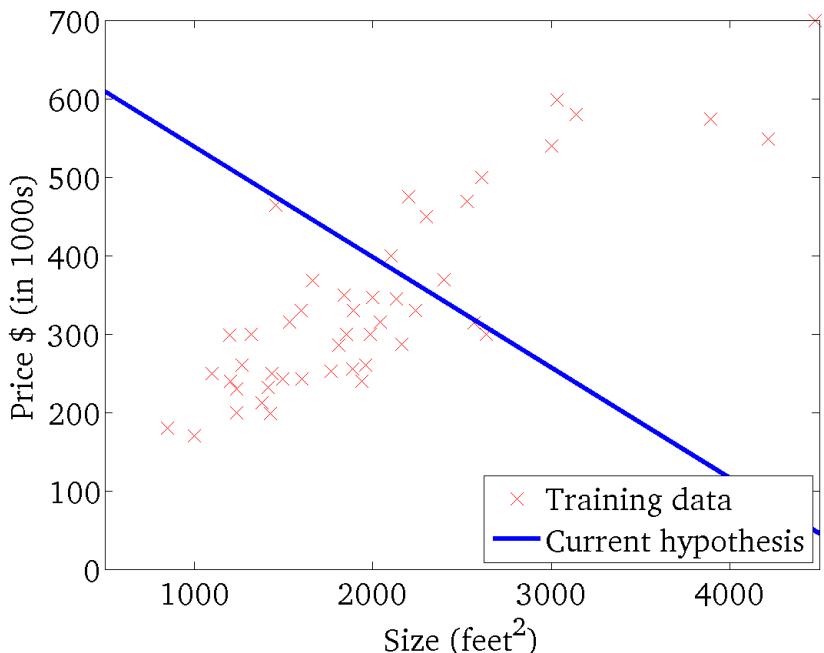
(function of the parameters  $\theta_0, \theta_1$ )





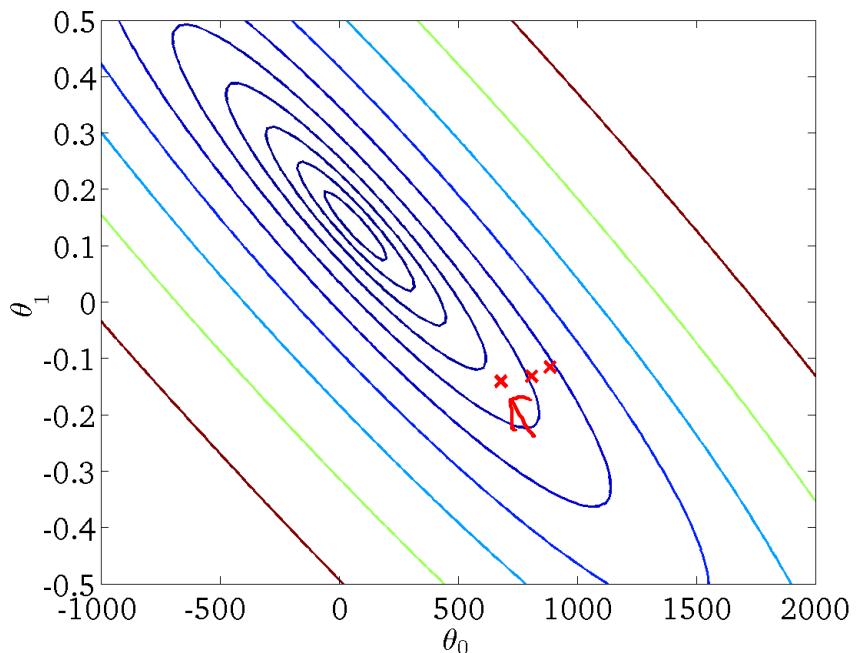
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

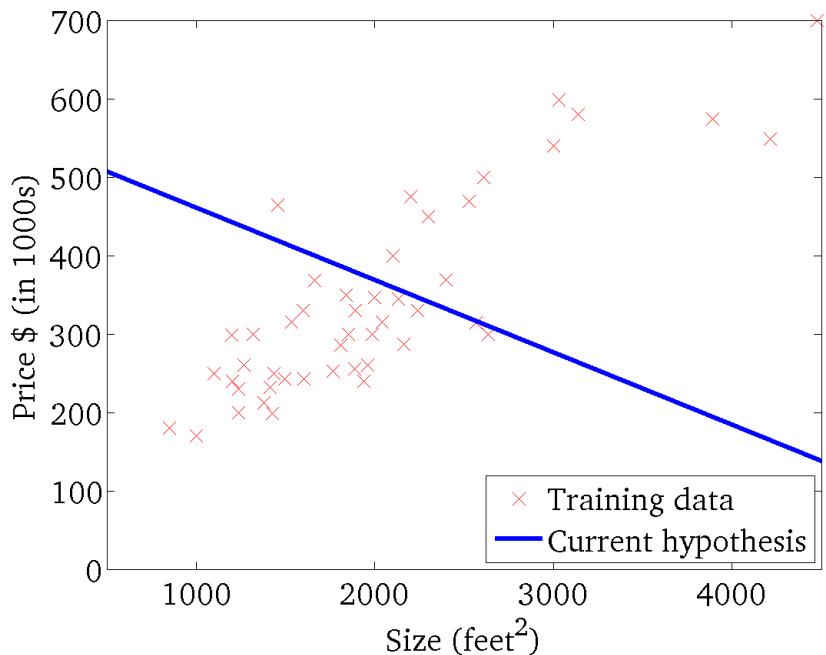
(function of the parameters  $\theta_0, \theta_1$ )





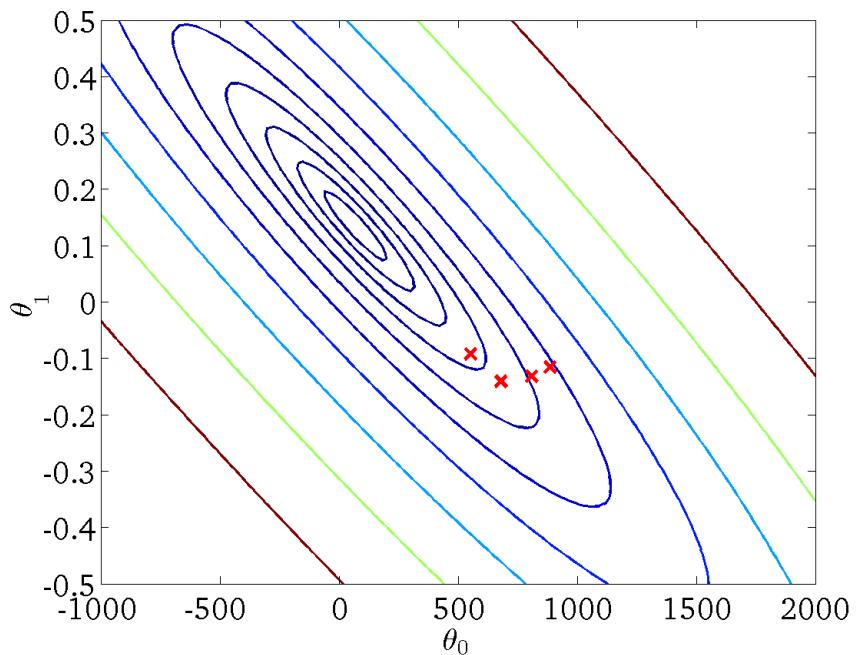
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

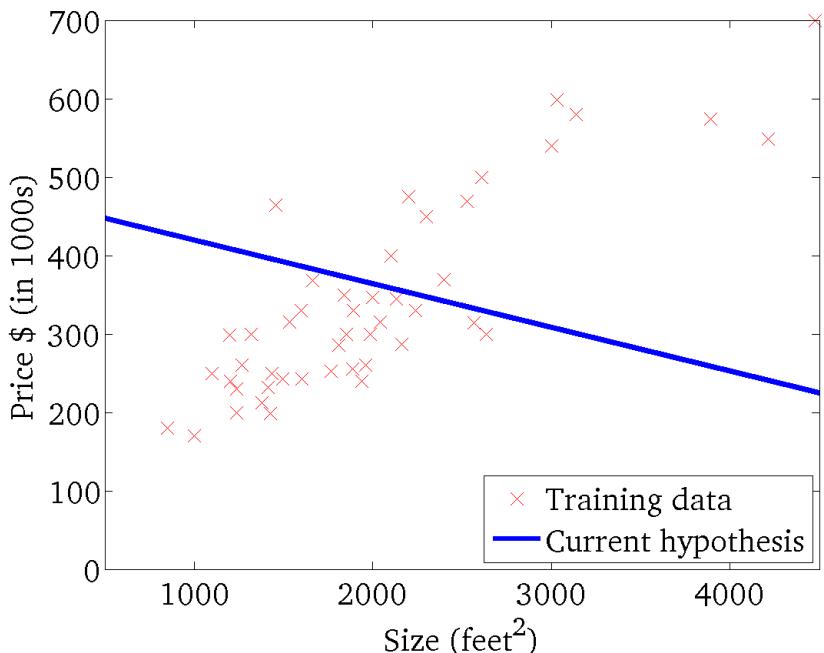
(function of the parameters  $\theta_0, \theta_1$ )





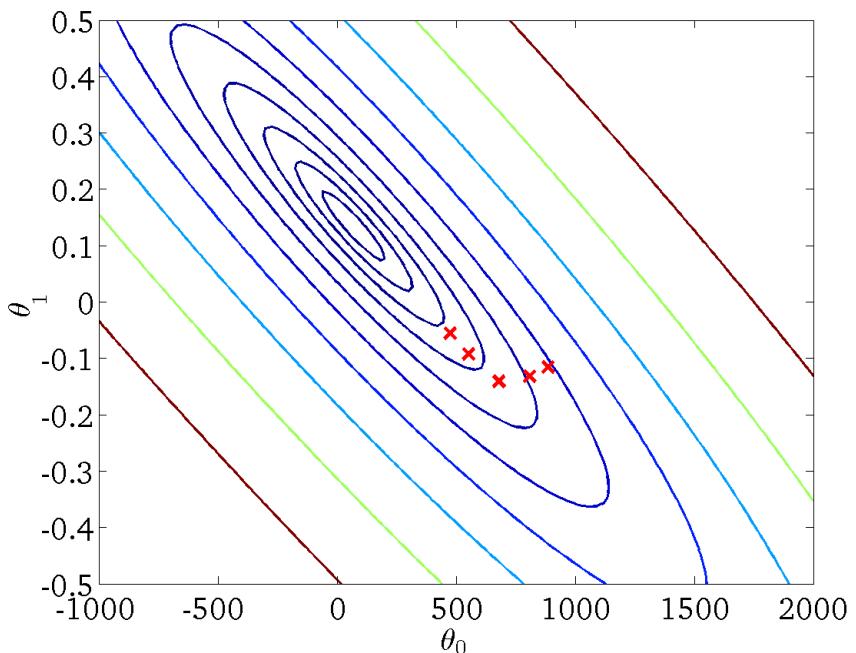
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

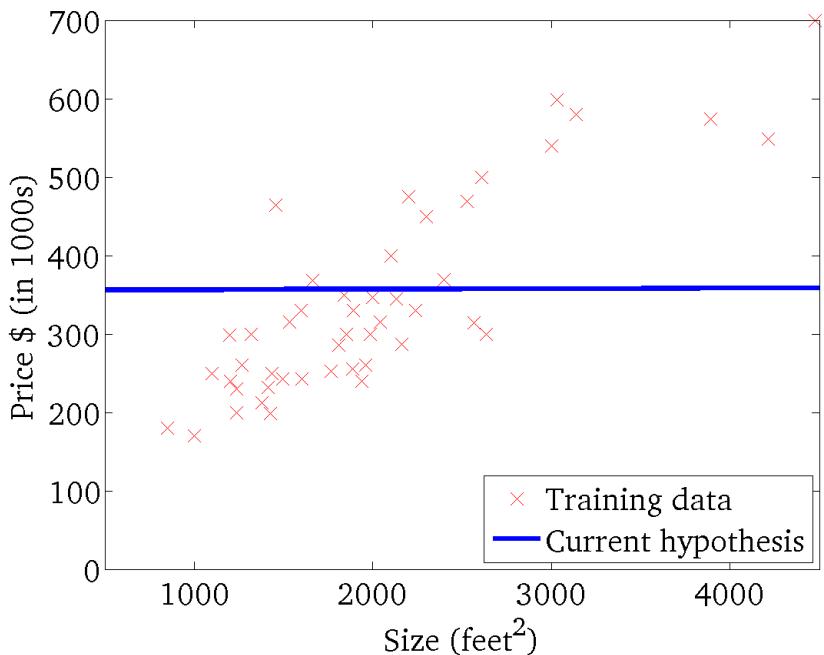
(function of the parameters  $\theta_0, \theta_1$ )





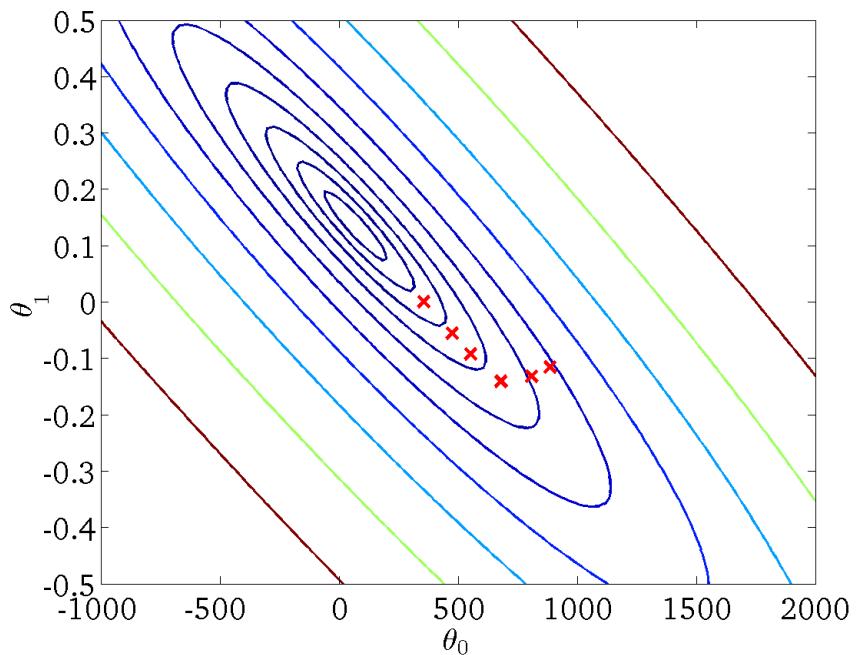
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

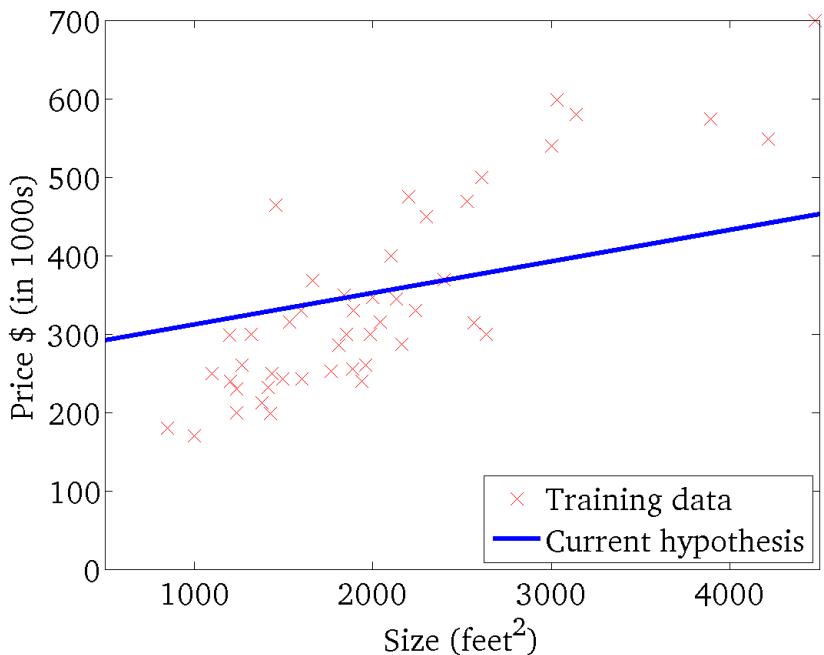
(function of the parameters  $\theta_0, \theta_1$ )





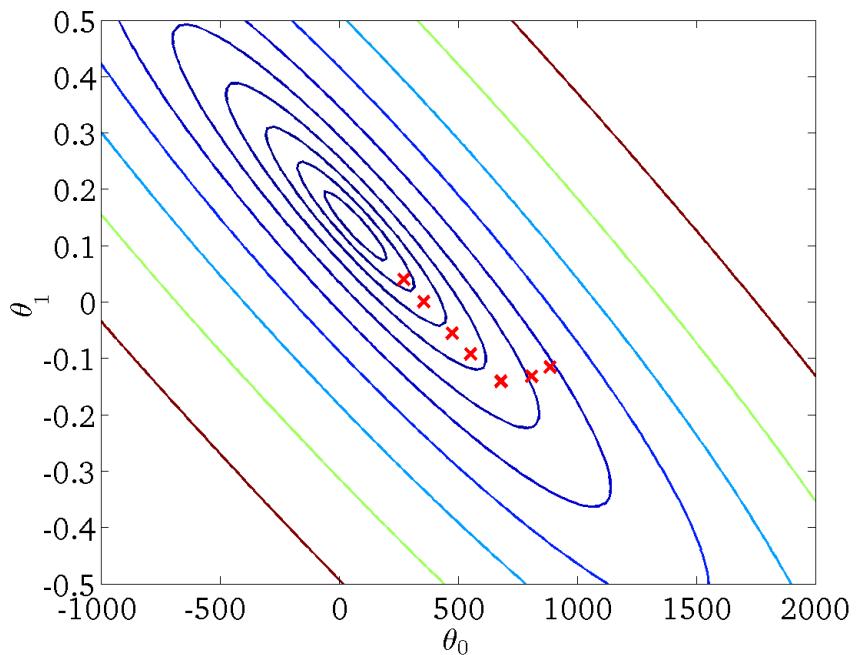
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

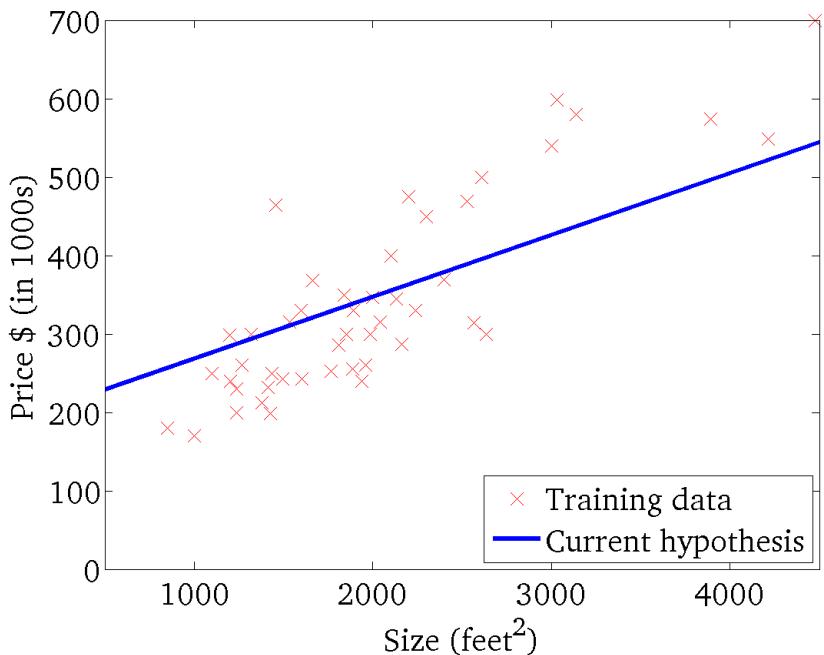
(function of the parameters  $\theta_0, \theta_1$ )





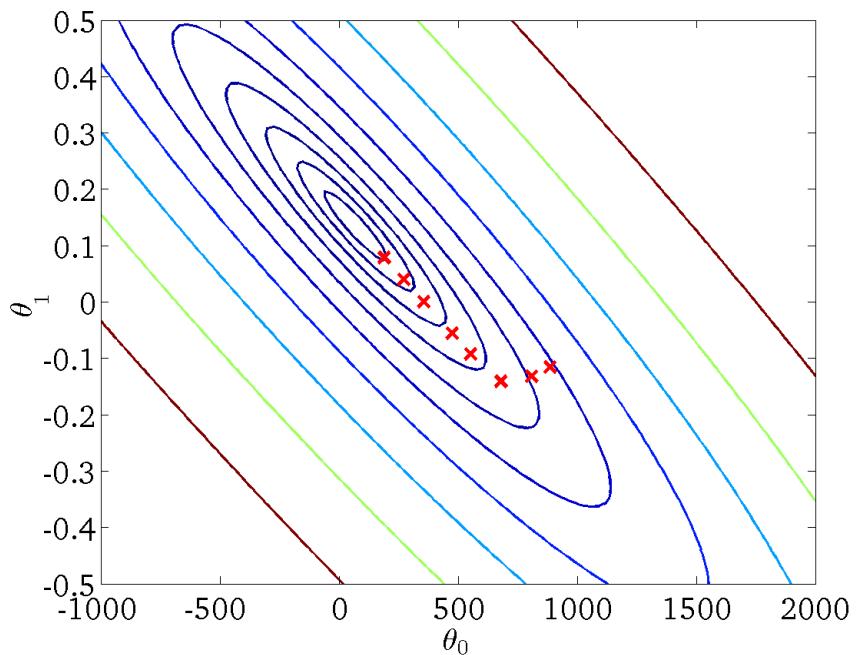
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

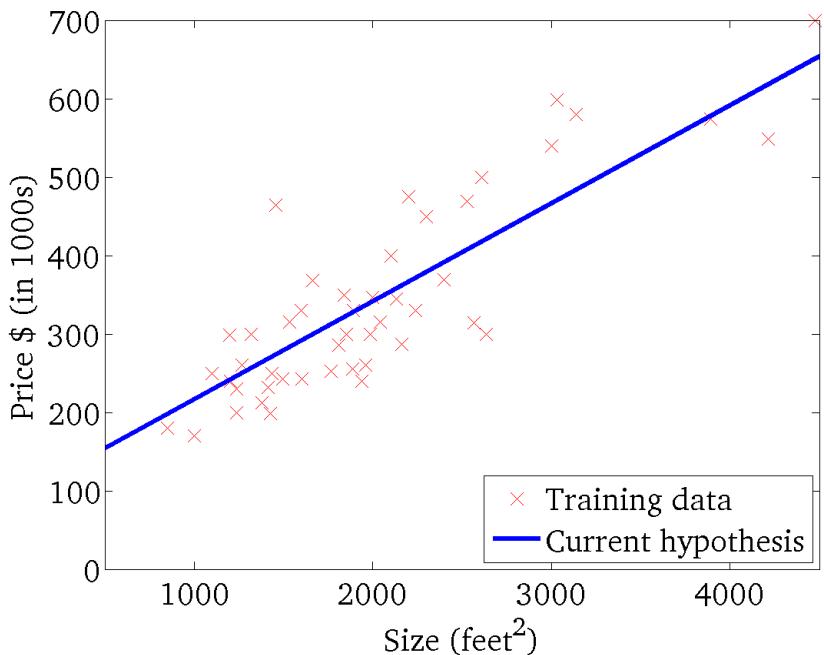
(function of the parameters  $\theta_0, \theta_1$ )





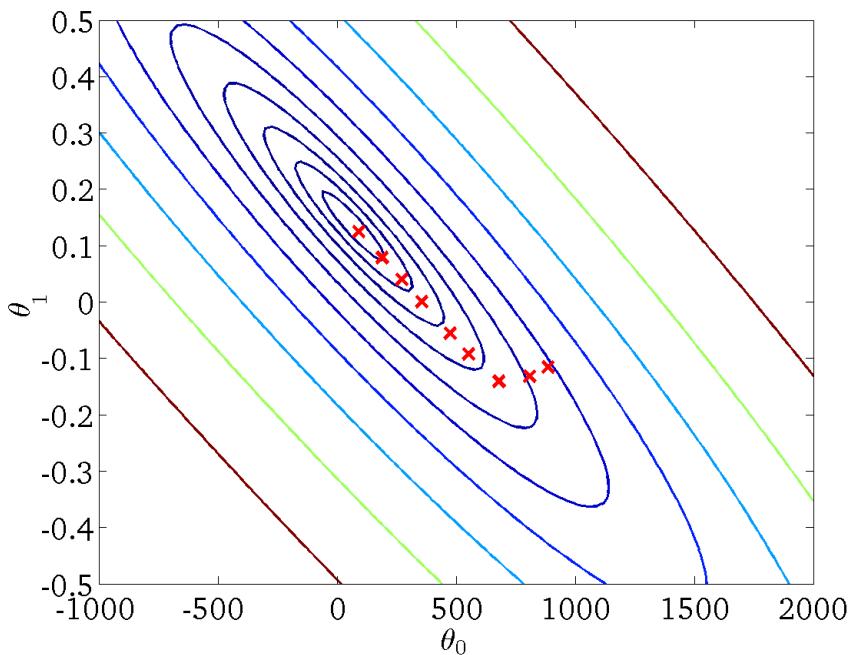
$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )





# Regularized regression

Slides adopted from Aarti course from cmu



# Geometric interpretation

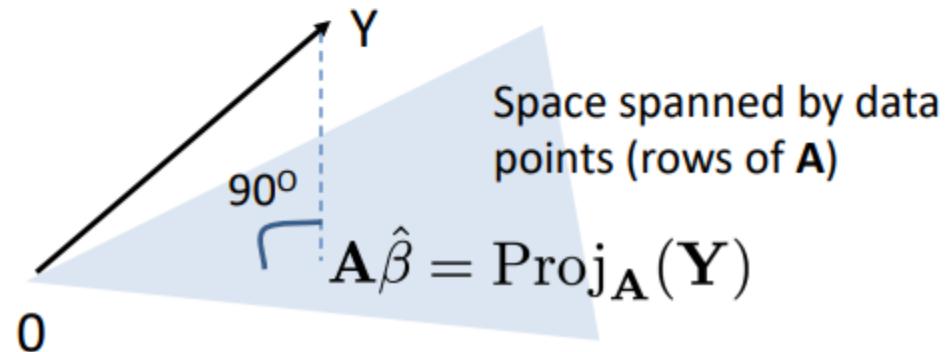
$$\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\hat{\beta}} = 0 \quad \text{gives} \quad (\mathbf{A}^T \mathbf{A}) \hat{\beta} = \mathbf{A}^T \mathbf{Y}$$

p x p   p x 1      p x 1

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

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

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





# Gradient descent in case of large dimension

$$\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\hat{\beta}} = 0 \quad \text{gives} \quad (\mathbf{A}^T \mathbf{A}) \hat{\beta} = \mathbf{A}^T \mathbf{Y}$$

p x p   p x 1      p x 1

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

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

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

- 2) If dimension p is large, computing inverse is expensive  $O(p^3)$   
Gradient descent since objective is convex  $(\mathbf{A}^T \mathbf{A} \succeq 0)$

$$\begin{aligned}\beta^{t+1} &= \beta^t - \frac{\alpha}{2} \left. \frac{\partial J(\beta)}{\partial \beta} \right|_t \\ &= \beta^t - \alpha \mathbf{A}^T (\mathbf{A} \beta^t - \mathbf{Y})\end{aligned}$$



# Rank of $\mathbf{A}^T \mathbf{A}$

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

p x p   p x 1      p x 1

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

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

Rank( $\mathbf{A}^T \mathbf{A}$ ) = number of non-zero eigenvalues of  $(\mathbf{A}^T \mathbf{A})$  = number of non-zero singular values of  $\mathbf{A}$   $\leq \min(n, p)$  since  $\mathbf{A}$  is  $n \times p$

So, rank( $\mathbf{A}^T \mathbf{A}$ ),  $r \leq \min(n, p)$       not invertible if  $r < p$  (e.g.  $n < p$   
i.e. high-dimensional setting)



# Under-determined; no unique solution

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

$p \times p$     $p \times 1$        $p \times 1$

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

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

If  $\mathbf{A} = \mathbf{U} \mathbf{S} \mathbf{V}^T$ , then normal equations  $(\mathbf{S} \mathbf{V}^T)^T \hat{\boldsymbol{\beta}} = (\mathbf{U}^T \mathbf{Y})$

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



# Regularized Least Squares; constrain the solution further

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)

$$\hat{\beta}_{\text{MAP}} = \arg \min_{\beta} \sum_{i=1}^n (Y_i - X_i \beta)^2 + \lambda \|\beta\|_2^2 \quad \text{Ridge Regression (L2 penalty)}$$

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

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

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

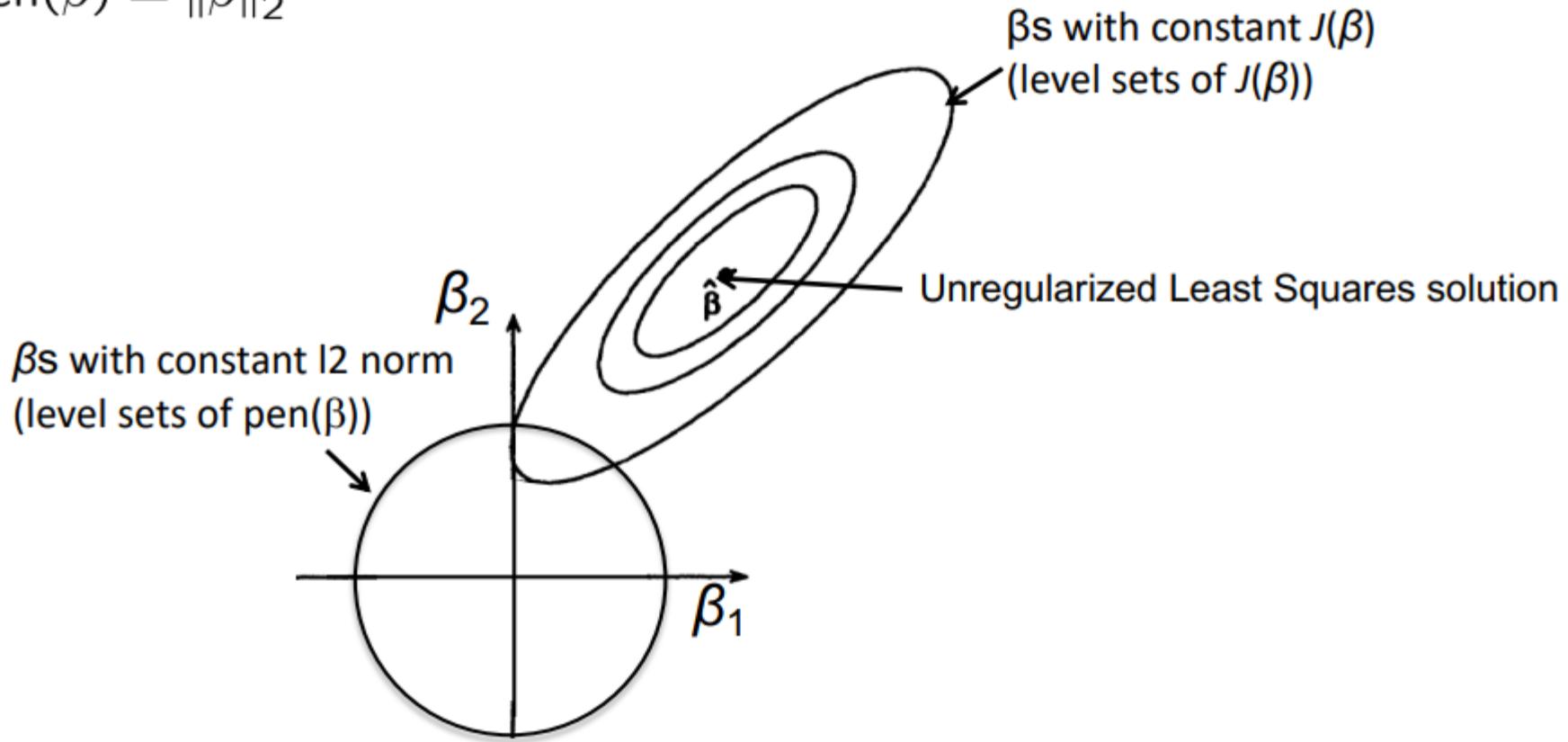
# Understanding regularized Least Squares



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

Ridge Regression:

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





# Lasso

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)

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

Ridge Regression  
(l2 penalty)

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

$\lambda \geq 0$   
Lasso  
(l1 penalty)

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

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



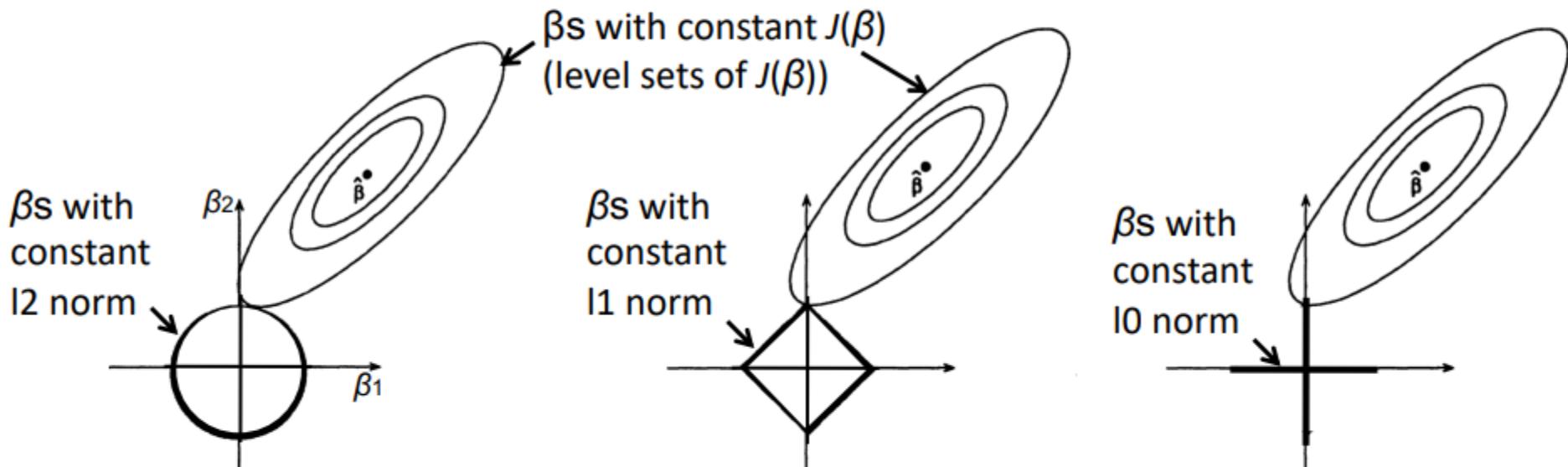
# Ridge Regression vs Lasso

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

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

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

Ideally  $\ell_0$  penalty,  
but optimization  
becomes non-convex



**Lasso ( $\ell_1$  penalty) results in sparse solutions – vector with more zero coordinates**  
**Good for high-dimensional problems – don't have to store all coordinates,**  
**interpretable solution!**



# Regularized Least Squares

– connection to MLE and  
MAP (Model-based  
approaches)



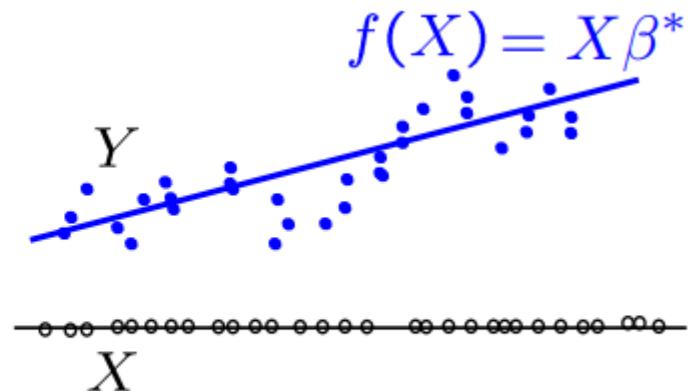
# Least Squares and MLE

Intuition: Signal plus (zero-mean) Noise model

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

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

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



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

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



# Regularized Least Squares and MAP

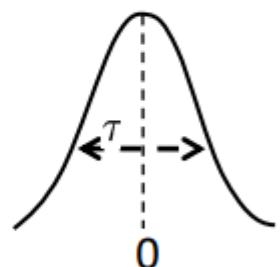
What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

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

I) 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

$\downarrow$   
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 MAP

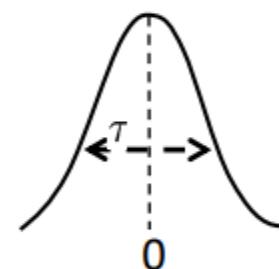
What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

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

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

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

# Regularized Least Squares and MAP

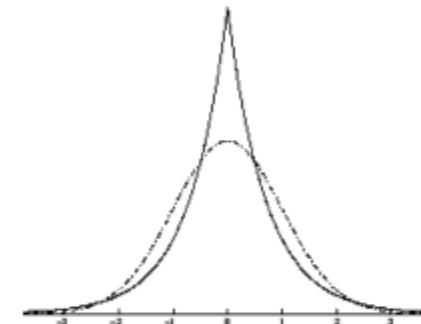


What if  $(\mathbf{A}^T \mathbf{A})$  is not invertible ?

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

II) Laplace Prior

$$\beta_i \stackrel{iid}{\sim} \text{Laplace}(0, t) \quad p(\beta_i) \propto e^{-|\beta_i|/t}$$



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

↓  
constant( $\sigma^2, t$ )

**Lasso**

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



# Nonlinear Regression Bias Variance trade-off

[dehaqani@ut.ac.ir](mailto:dehaqani@ut.ac.ir)

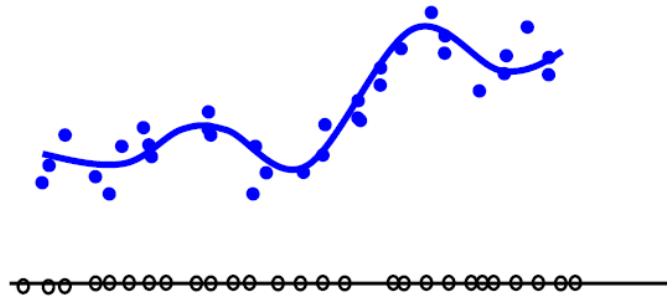
Slides are mainly adopted from cmu Aarti course and Tibshirani

# Beyond Linear Regression



Polynomial regression

Regression with nonlinear features



Kernelized Ridge Regression

Local Kernel Regression

# Polynomial Regression



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

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

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

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

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

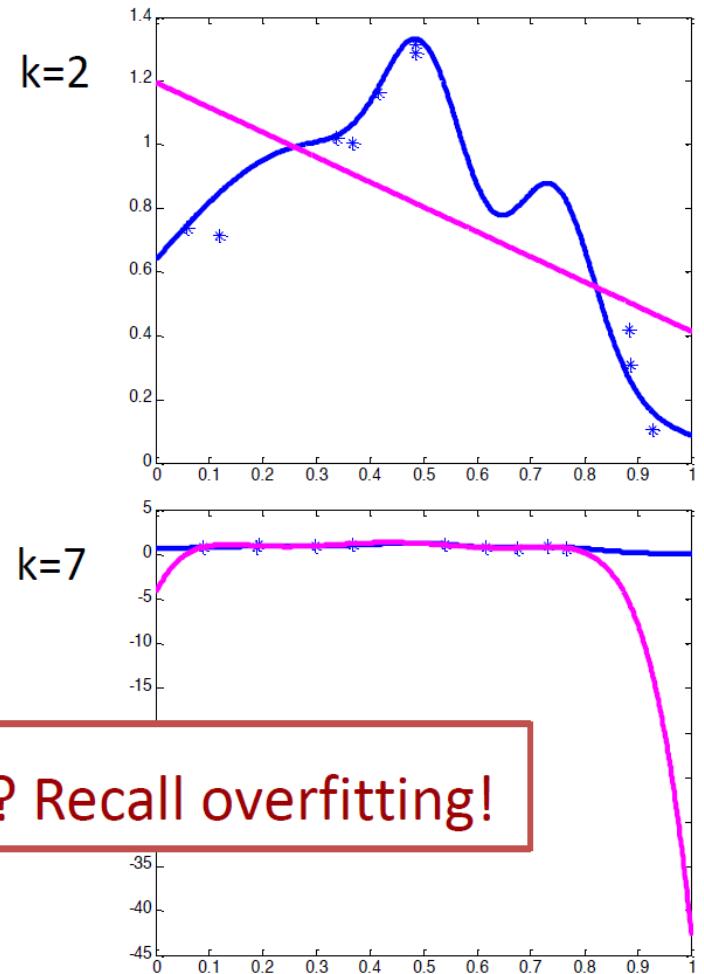
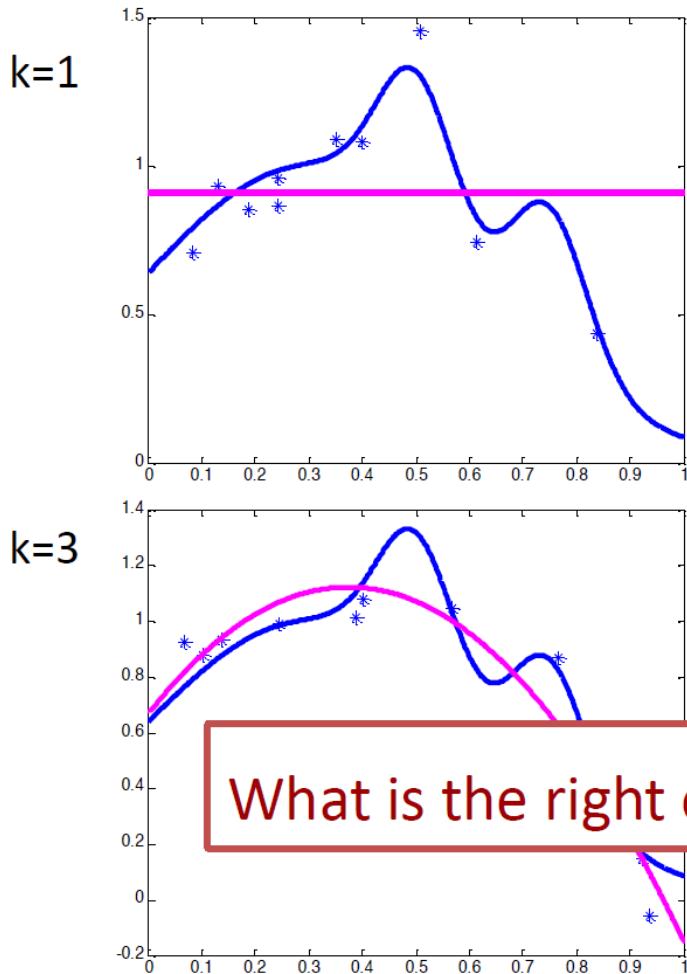
$$+ \sum_{i=1}^p \sum_{j=1}^p \beta_{ij} X^{(i)} X^{(j)} + \sum_{i=1}^p \sum_{j=1}^p \sum_{k=1}^p X^{(i)} X^{(j)} X^{(k)} + \dots \text{ terms up to degree m}$$

AA

# Polynomial Regression



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

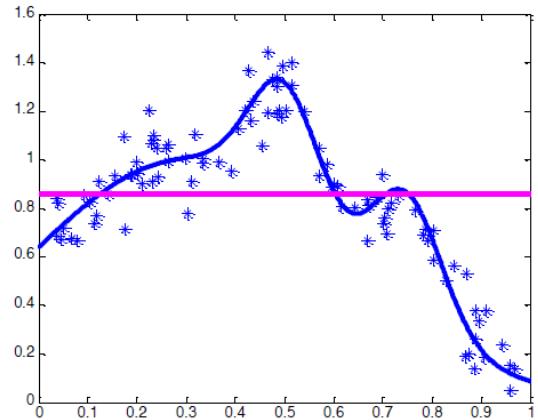
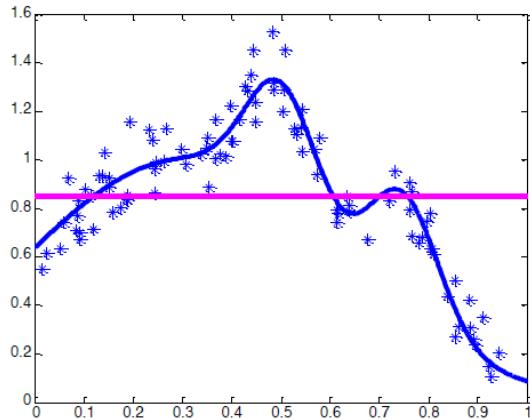
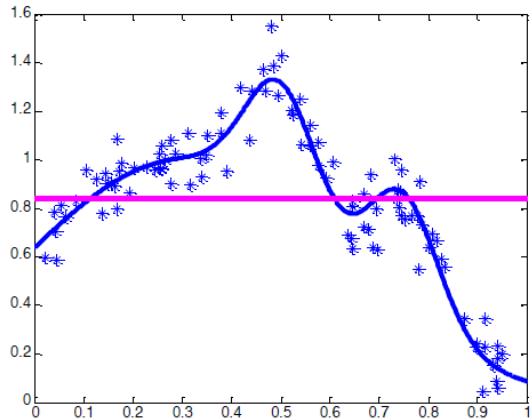


# Bias – Variance Tradeoff

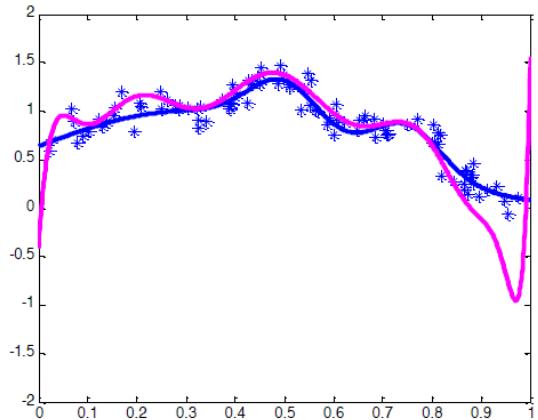
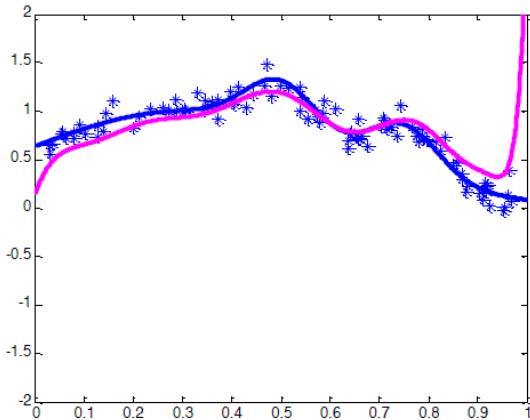
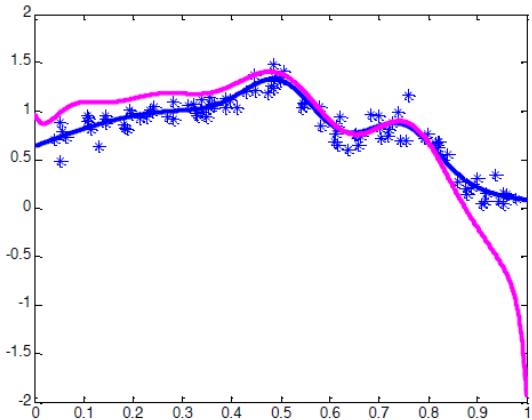


3 Independent training datasets

Large bias, Small variance – poor approximation but robust/stable



Small bias, Large variance – good approximation but unstable



# Bias – Variance Decomposition



It can be shown that

$$E[(f(X) - f^*(X))^2] = \text{Bias}^2 + \text{Variance}$$

$$\text{Bias} = E[f(X)] - f^*(X)$$

How far is the model from  
best model on average

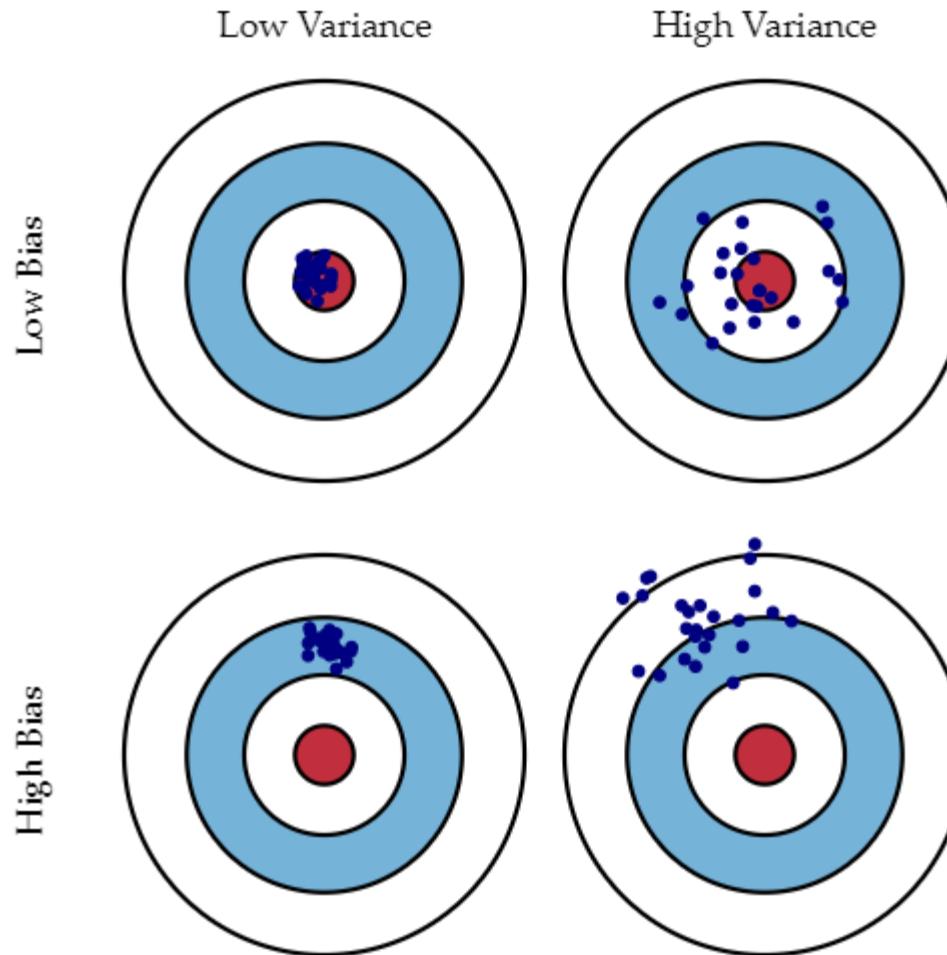
$$\text{Variance} = E[(f(X) - E[f(X)])^2]$$
 How variable is the model



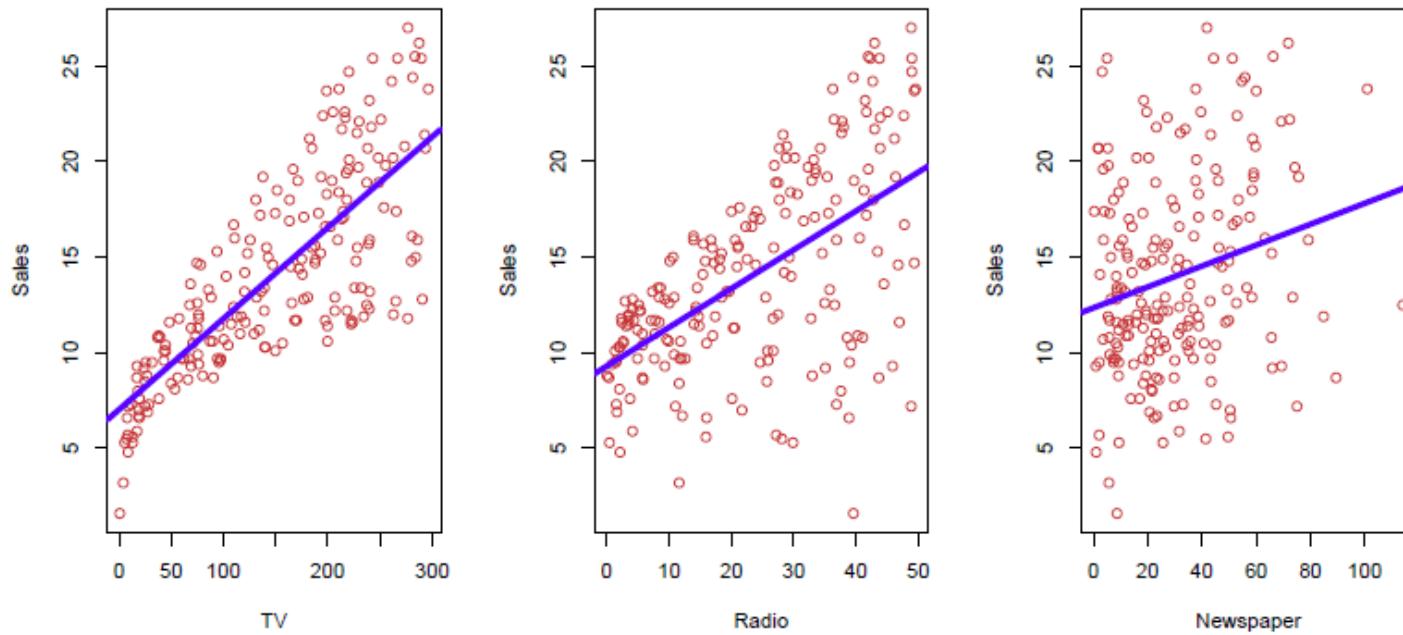
# Underfitting vs. overfitting

- The **bias** error: (**underfitting**)
  - From **erroneous assumptions** in the learning algorithm.
  - High bias can cause an algorithm to **miss the relevant relations** between features and target outputs.
- The **variance** error: (**overfitting**).
  - From **sensitivity** to **small fluctuations** in the training set.
  - High variance can cause an algorithm to **model the random noise** in the training data, rather than the intended outputs

# Graphical illustration of bias and variance



# Mathematical Definition; Statistical learning



Shown are Sales vs TV, Radio and Newspaper, with a blue linear-regression line fit separately to each. Can we predict Sales using these three? Perhaps we can do better using a model

$$\text{Sales} \approx f(\text{TV}, \text{Radio}, \text{Newspaper})$$



# Notation

Here Sales is a response or target that we wish to predict. We generically refer to the response as  $Y$ . TV is a feature, or input, or predictor; we name it  $X_1$ . Likewise name Radio as  $X_2$ , and so on. We can refer to the input vector collectively as

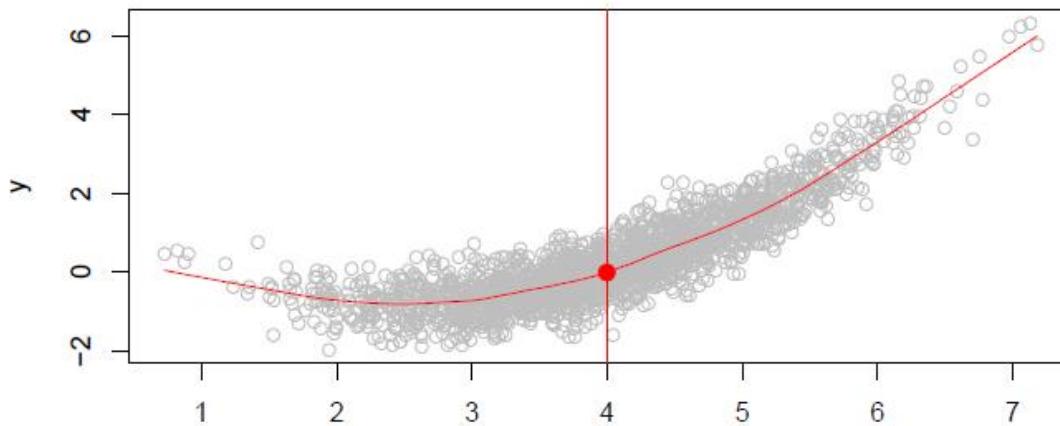
$$X = \begin{pmatrix} X_1 \\ X_2 \\ X_3 \end{pmatrix}$$

Now we write our model as

$$Y = f(X) + \varepsilon$$

where  $\varepsilon$  captures measurement errors and other discrepancies.

# What is $f(X)$ good for?



Is there an ideal  $f(X)$ ? In particular, what is a good value for  $f(X)$  at any selected value of  $X$ , say  $X = 4$ ? There can be many  $Y$  values at  $X = 4$ . A good value is

$$f(4) = E(Y|X = 4)$$

$E(Y|X = 4)$  means *expected value* (average) of  $Y$  given  $X = 4$ .

This ideal  $f(x) = E(Y|X = x)$  is called the *regression function*.



# The regression function $f(x)$

- Is also defined for vector  $X$ ; e.g.  
$$f(x) = f(x_1, x_2, x_3) = E(Y|X_1 = x_1, X_2 = x_2, X_3 = x_3)$$
- Is the *ideal* or *optimal* predictor of  $Y$  with regard to mean-squared prediction error:  $f(x) = E(Y|X = x)$  is the function that minimizes  $E[(Y - g(X))^2|X = x]$  over all functions  $g$  at all points  $X = x$ .
- $\epsilon = Y - f(x)$  is the *irreducible* error — i.e. even if we knew  $f(x)$ , we would still make errors in prediction, since at each  $X = x$  there is typically a distribution of possible  $Y$  values.
- For any estimate  $\hat{f}(x)$  of  $f(x)$ , we have

$$E[(Y - \hat{f}(X))^2|X = x] = \underbrace{[f(x) - \hat{f}(x)]^2}_{\text{Reducible}} + \underbrace{\text{Var}(\epsilon)}_{\text{Irreducible}}$$

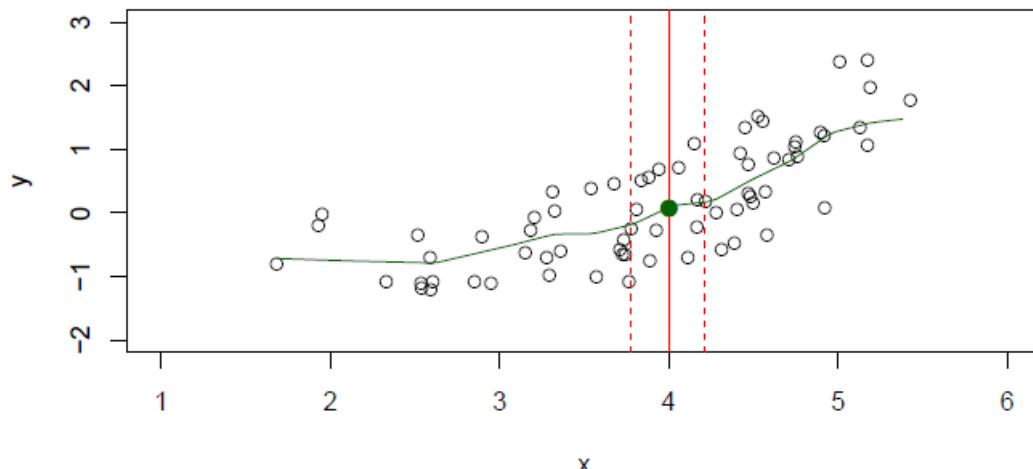


# How to estimate $f$

- Typically we have few if any data points with  $X = 4$  exactly.
- So we cannot compute  $E(Y|X = x)!$
- Relax the definition and let

$$\hat{f}(x) = \text{Ave}(Y|X \in \mathcal{N}(x))$$

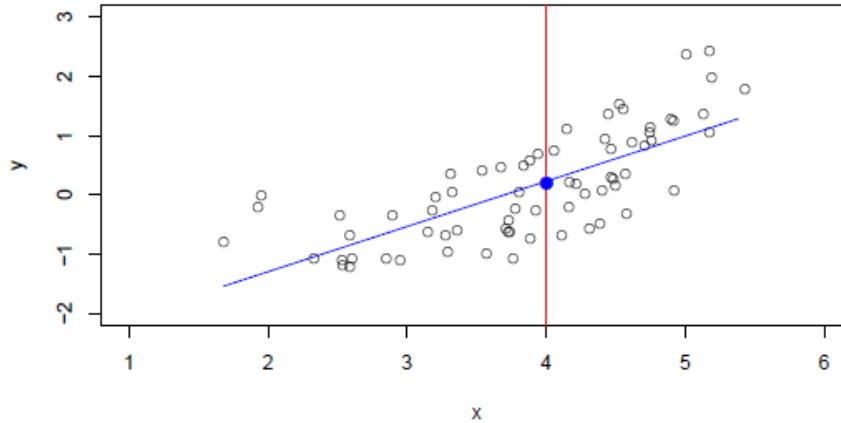
where  $\mathcal{N}(x)$  is some *neighborhood* of  $x$ .



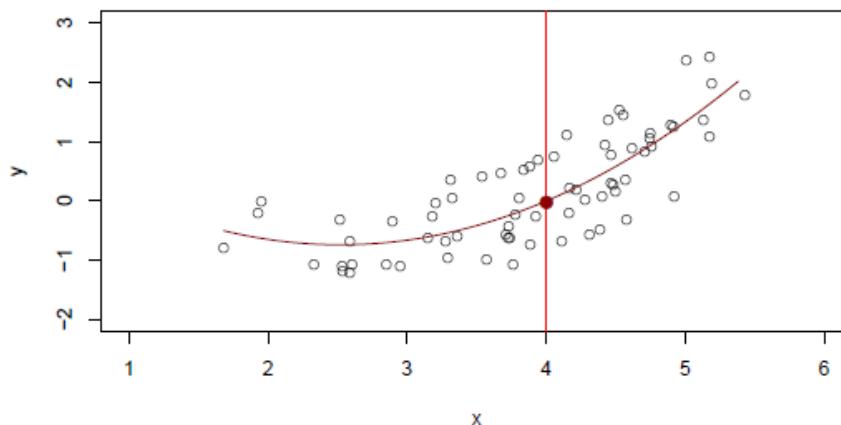
# Parametric and structured models

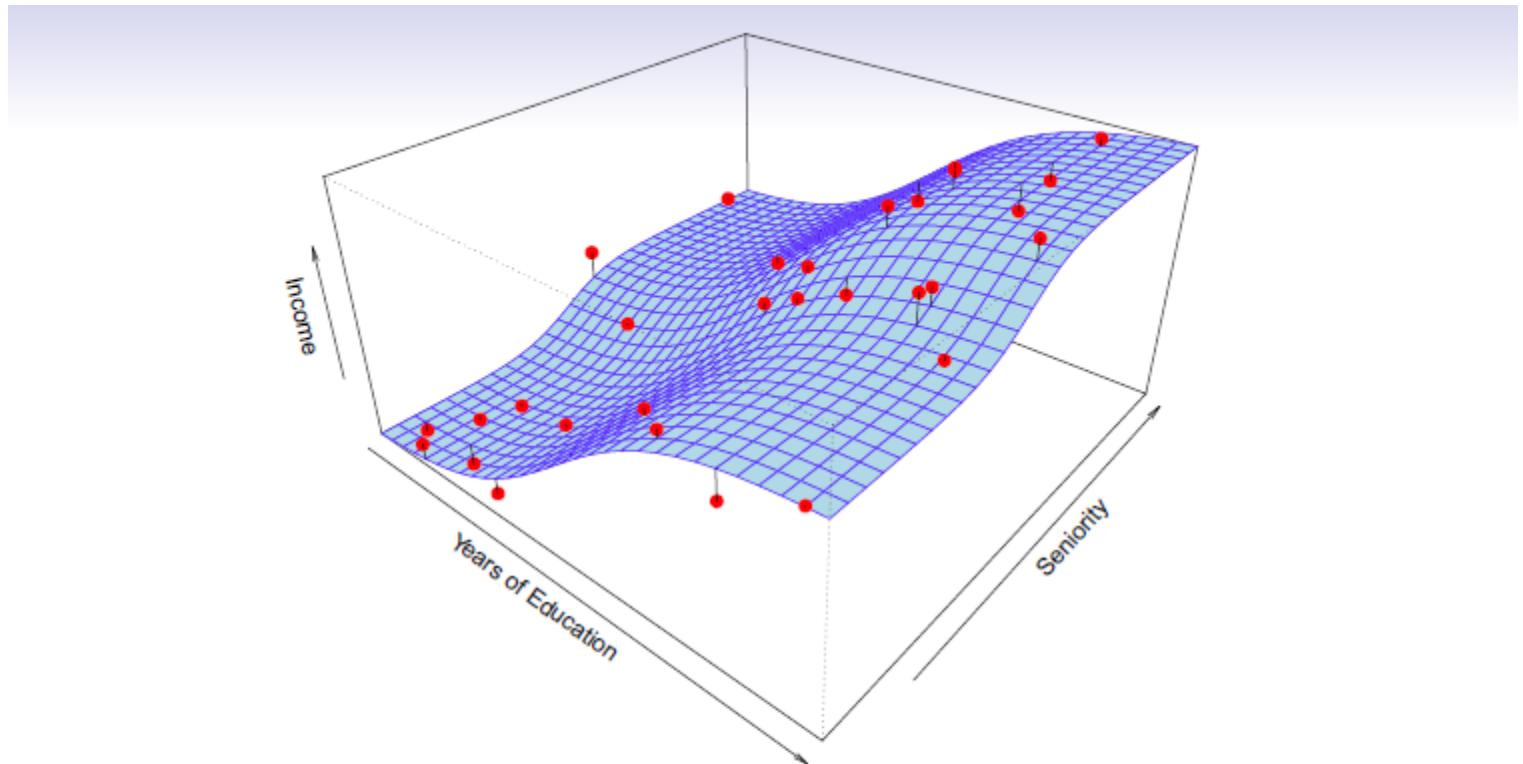


A linear model  $\hat{f}_L(X) = \hat{\beta}_0 + \hat{\beta}_1 X$  gives a reasonable fit here



A quadratic model  $\hat{f}_Q(X) = \hat{\beta}_0 + \hat{\beta}_1 X + \hat{\beta}_2 X^2$  fits slightly better.

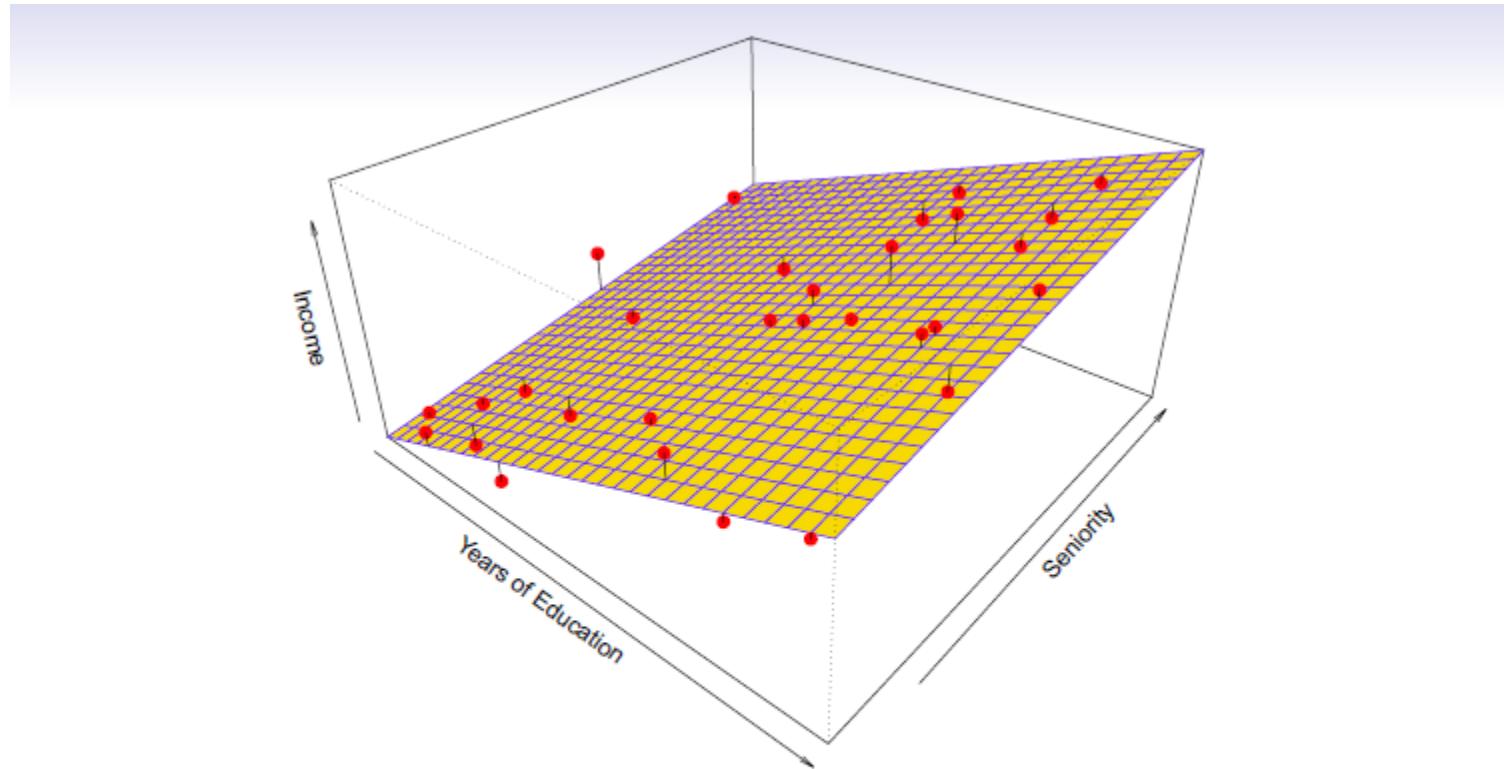




Simulated example. Red points are simulated values for **income** from the model

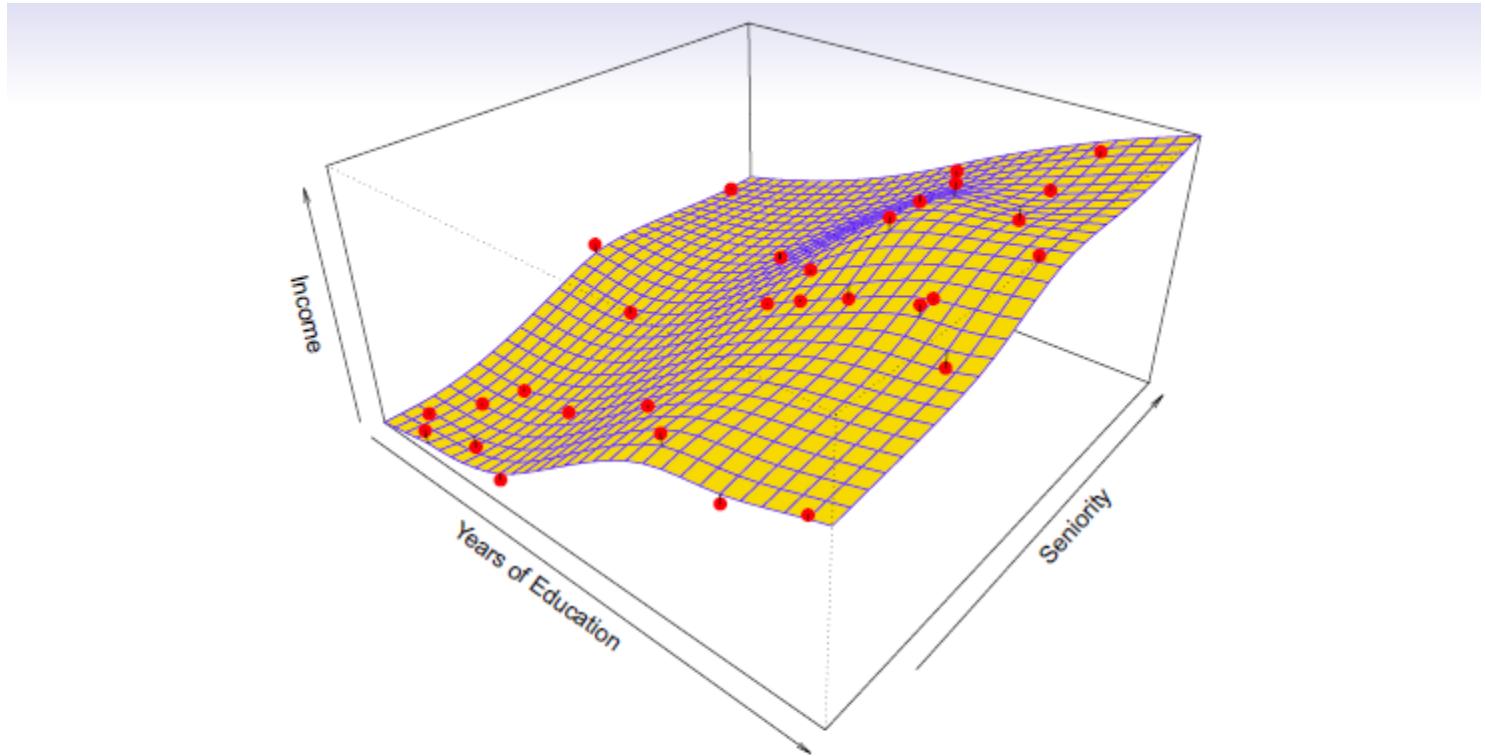
$$\text{income} = f(\text{education}, \text{seniority}) + \epsilon$$

$f$  is the blue surface.

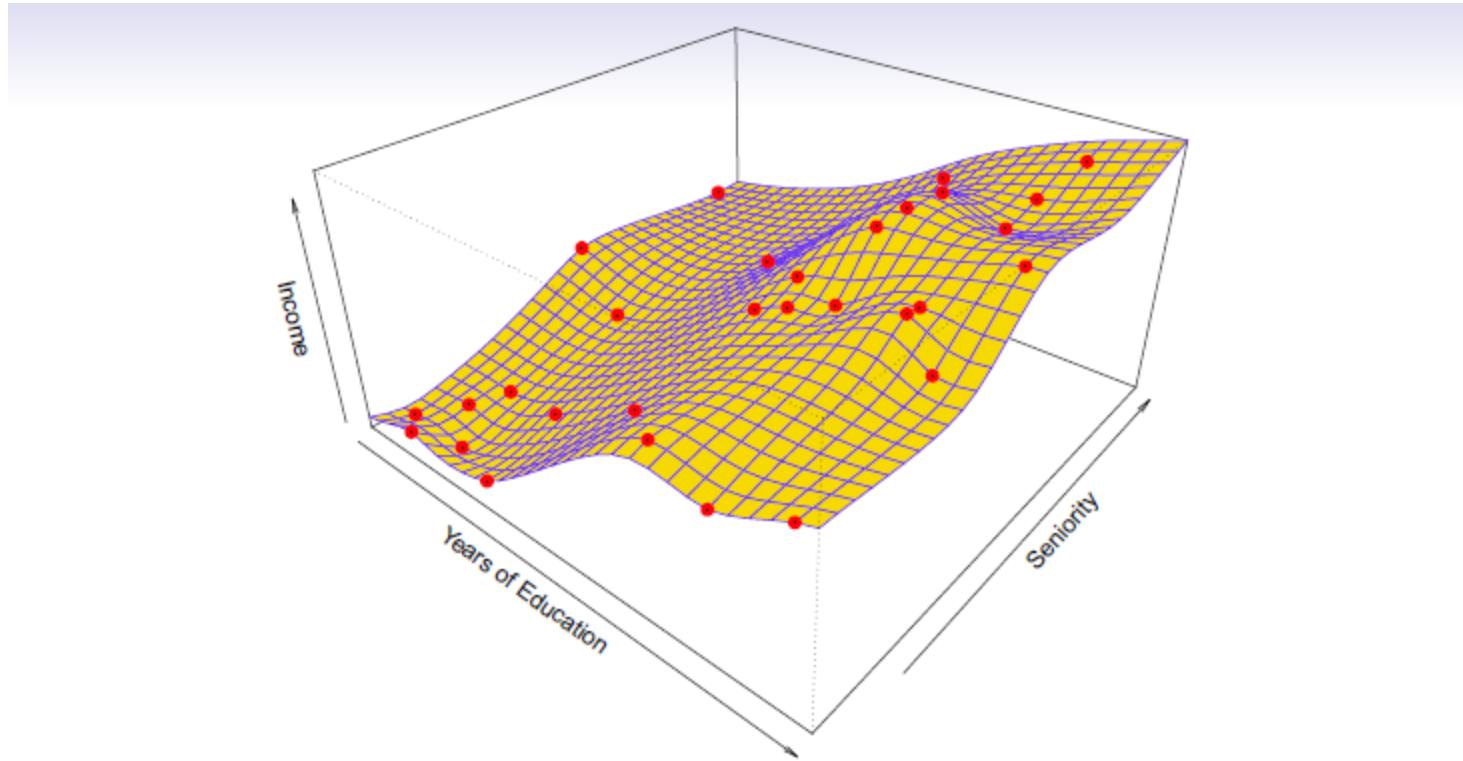


Linear regression model fit to the simulated data.

$$\hat{f}_L(\text{education}, \text{seniority}) = \hat{\beta}_0 + \hat{\beta}_1 \times \text{education} + \hat{\beta}_2 \times \text{seniority}$$



More flexible regression model  $\hat{f}_S(\text{education}, \text{seniority})$  fit to the simulated data. Here we use a technique called a *thin-plate spline* to fit a flexible surface. We control the roughness of the fit



Even more flexible spline regression model  
 $\hat{f}_S(\text{education}, \text{seniority})$  fit to the simulated data. Here the fitted model makes no errors on the training data! Also known as *overfitting*.

# Some trade-offs



- Prediction accuracy versus interpretability.
  - Linear models are easy to interpret; thin-plate splines are not.
- Good fit versus over-fit or under-fit.
  - How do we know when the fit is just right?
- Parsimony versus black-box.
  - We often prefer a simpler model involving fewer variables over a black-box predictor involving them all.

# Assessing Model Accuracy



Suppose we fit a model  $\hat{f}(x)$  to some training data  $\mathbf{Tr} = \{x_i, y_i\}_1^N$ , and we wish to see how well it performs.

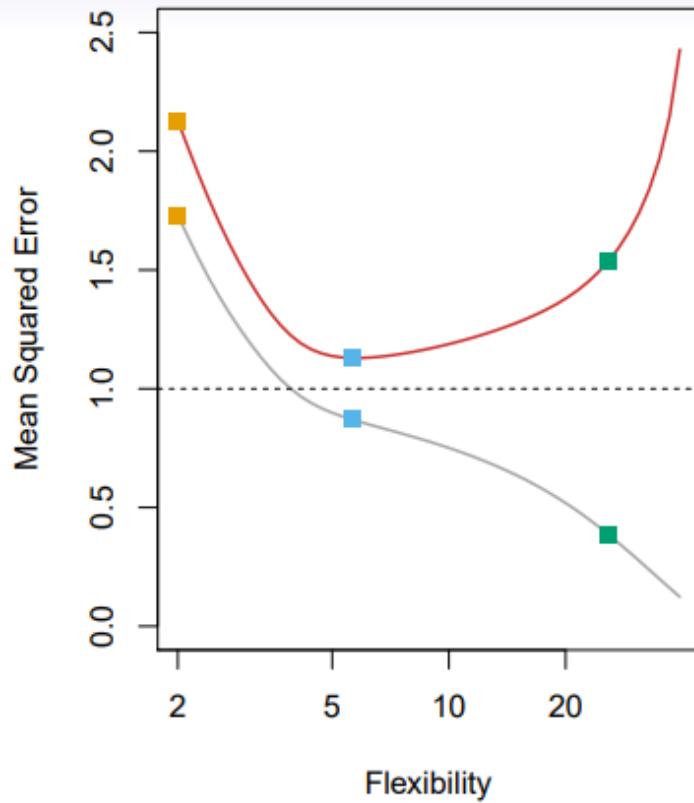
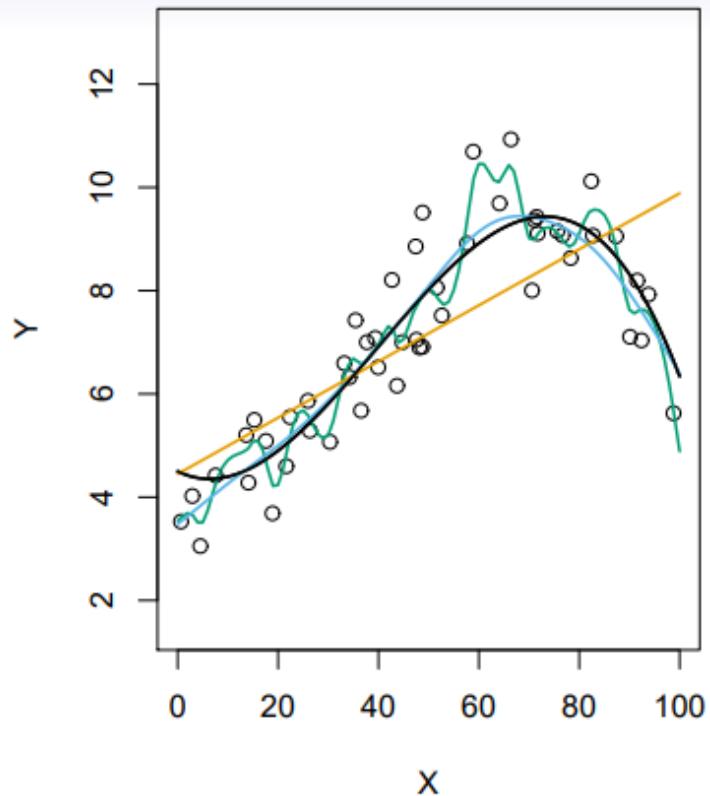
- We could compute the average squared prediction error over  $\mathbf{Tr}$ :

$$\text{MSE}_{\mathbf{Tr}} = \text{Ave}_{i \in \mathbf{Tr}} [y_i - \hat{f}(x_i)]^2$$

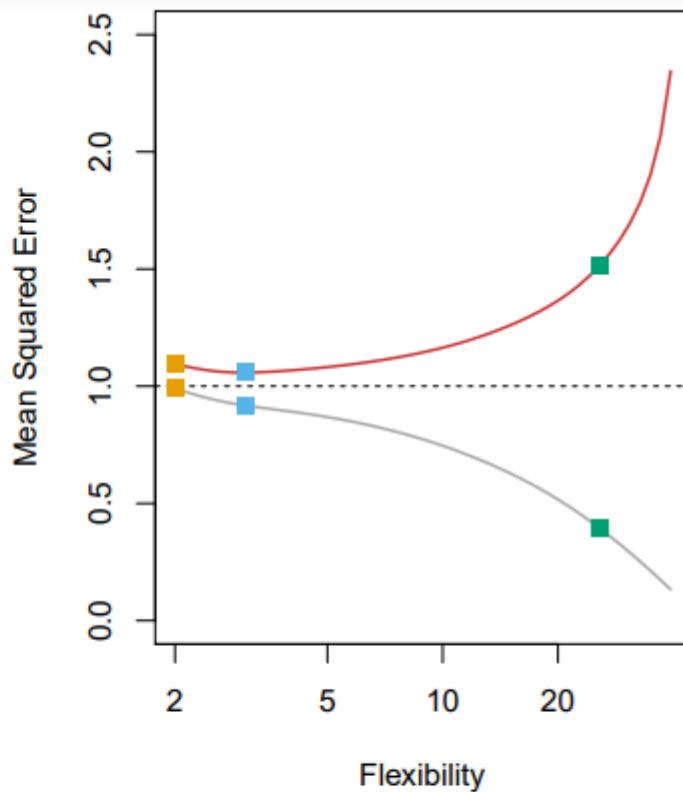
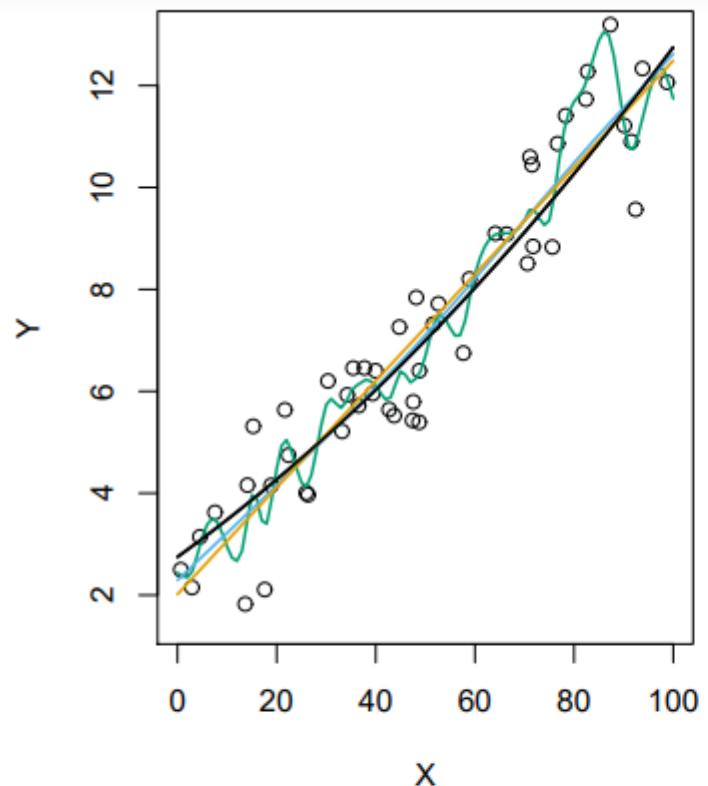
This may be biased toward more overfit models.

- Instead we should, if possible, compute it using fresh *test* data  $\mathbf{Te} = \{x_i, y_i\}_1^M$ :

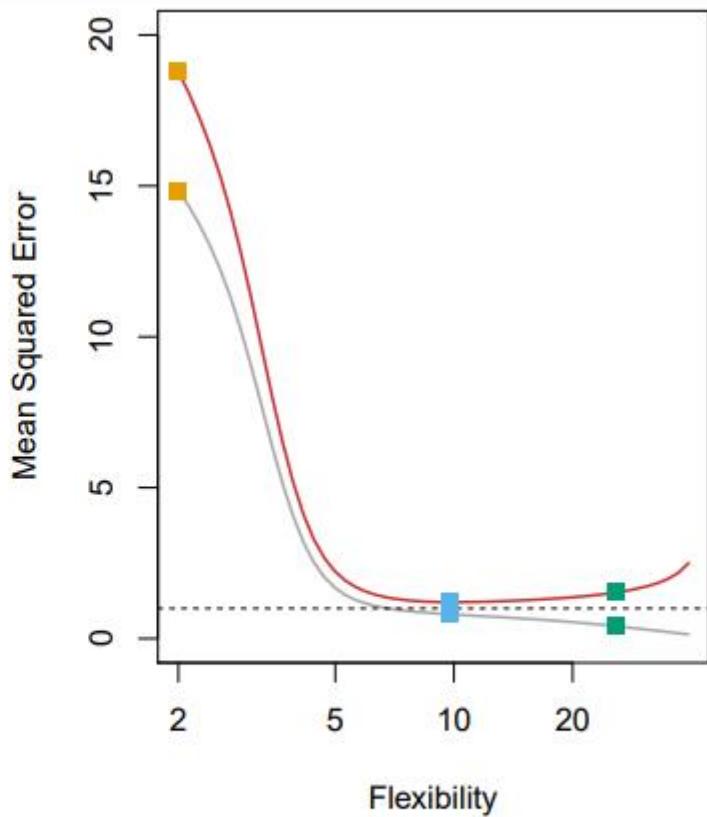
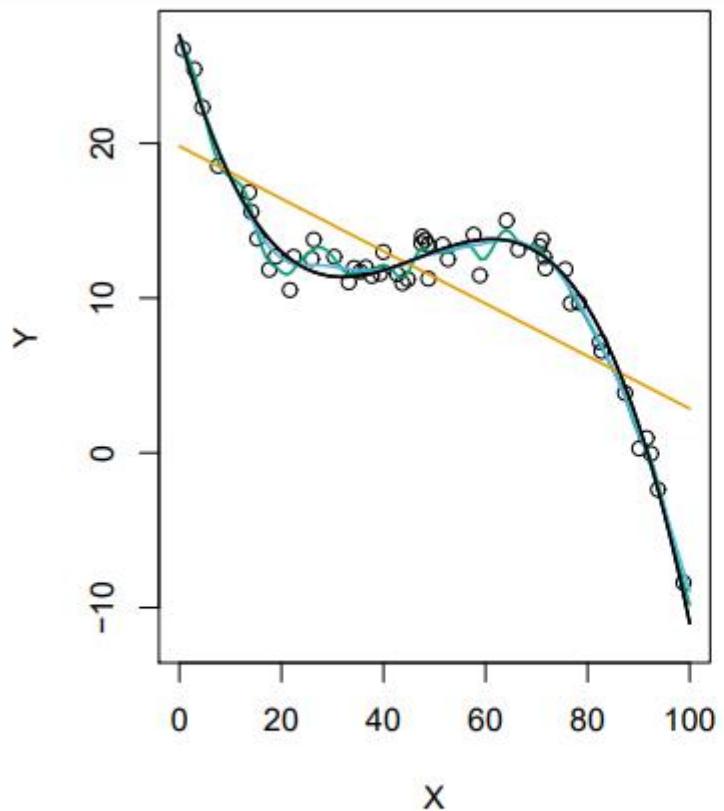
$$\text{MSE}_{\mathbf{Te}} = \text{Ave}_{i \in \mathbf{Te}} [y_i - \hat{f}(x_i)]^2$$



Black curve is truth. Red curve on right is  $MSE_{Te}$ , grey curve is  $MSE_{Tr}$ . Orange, blue and green curves/squares correspond to fits of different flexibility.



Here the truth is smoother, so the smoother fit and linear model do really well.



Here the truth is wiggly and the noise is low, so the more flexible fits do the best.



# Bias-Variance Trade-off

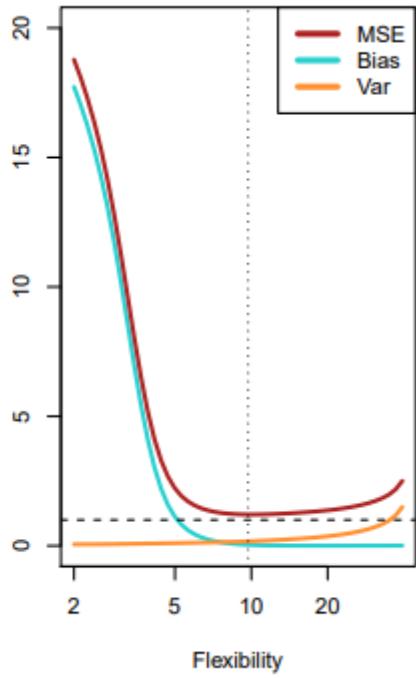
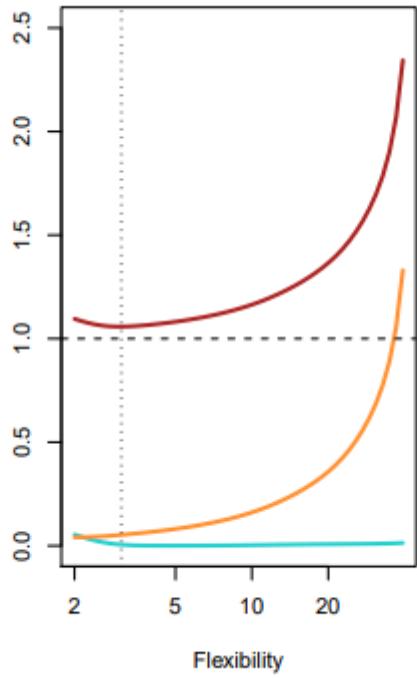
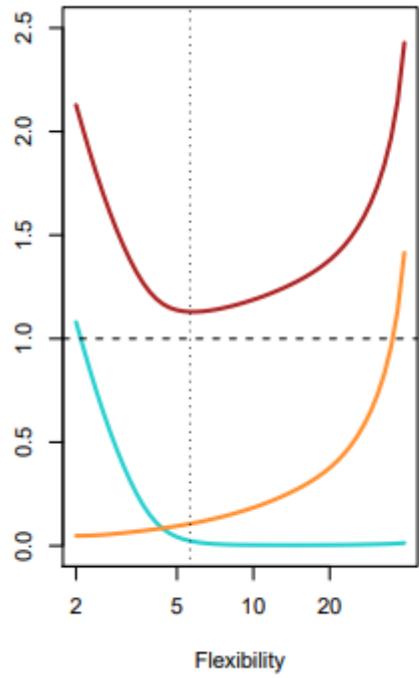
Suppose we have fit a model  $\hat{f}(x)$  to some training data  $\text{Tr}$ , and let  $(x_0, y_0)$  be a test observation drawn from the population. If the true model is  $Y = f(X) + \epsilon$  (with  $f(x) = E(Y|X = x)$ ), then

$$E \left( y_0 - \hat{f}(x_0) \right)^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon).$$

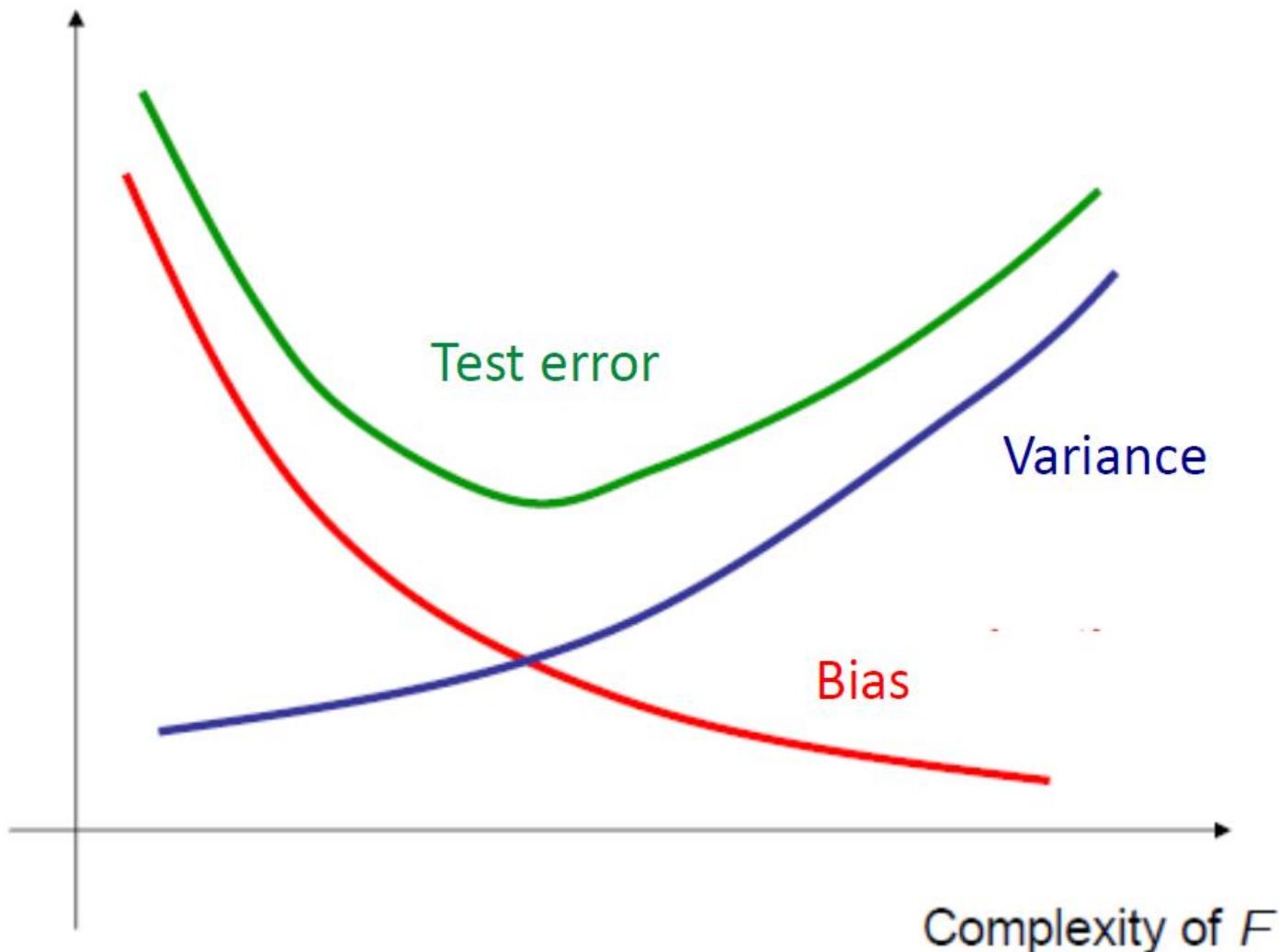
The expectation averages over the variability of  $y_0$  as well as the variability in  $\text{Tr}$ . Note that  $\text{Bias}(\hat{f}(x_0)) = E[\hat{f}(x_0)] - f(x_0)$ .

Typically as the *flexibility* of  $\hat{f}$  increases, its variance increases, and its bias decreases. So choosing the flexibility based on average test error amounts to a *bias-variance trade-off*.

# Bias-variance trade-off for the three examples



# Effect of Model Complexity



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

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

or

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

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

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

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



# Ridge regression (dual)

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

Primal problem:

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

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

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)

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

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  $i = 1, \dots, n$

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

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

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

n-dimensional optimization problem



# Ridge regression (dual)

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

Dual problem:

$$\max_{\alpha} -\frac{\alpha^T \alpha}{4} - \frac{1}{4\lambda} \alpha^T \mathbf{A} \mathbf{A}^T \alpha - \alpha^T \mathbf{Y} \quad \Rightarrow \hat{\alpha} = - \left( \frac{\mathbf{A} \mathbf{A}^T}{\lambda} + \mathbf{I} \right)^{-1} 2 \mathbf{Y}$$

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

Weighted average of training points

Weight of each training point (but typically not sparse)

# Kernelized ridge regression



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

Using dual, can re-write solution as:

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

How does this help?

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

We will come back  
to it later !!

$\hat{f}_n(X) = \mathbf{K}_X (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{Y}$  where

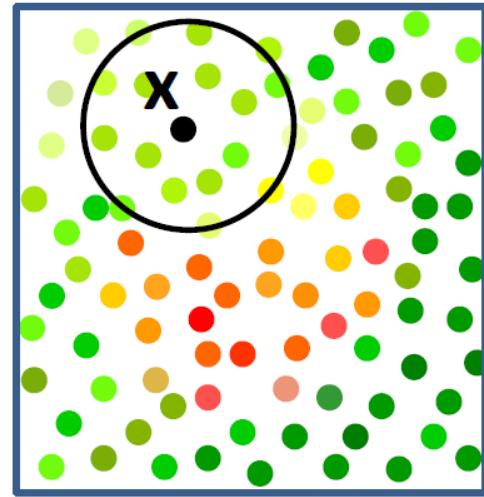
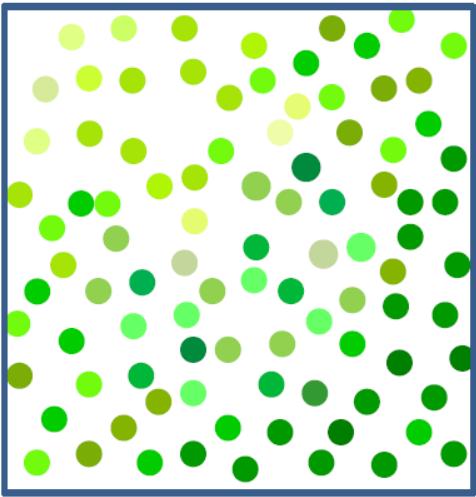
$\mathbf{K}_X(i, j) = \phi(X_i) \cdot \phi(X_j)$

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

# Local Kernel Regression



- What is the temperature in the room? at location  $x$ ?



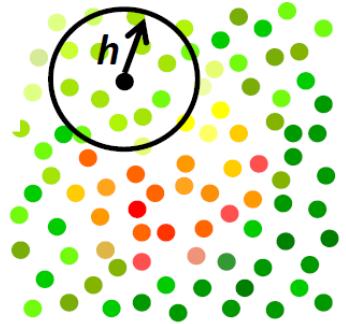
$$\hat{T} = \frac{1}{n} \sum_{i=1}^n Y_i$$

Average

$$\hat{T}(x) = \frac{\sum_{i=1}^n Y_i \mathbf{1}_{||X_i - x|| \leq h}}{\sum_{i=1}^n \mathbf{1}_{||X_i - x|| \leq h}}$$

"Local" Average

# Local Kernel Regression



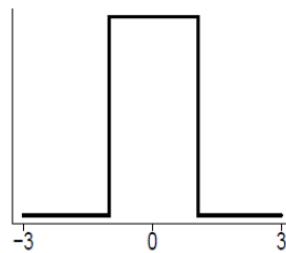
- Nonparametric estimator akin to kNN
- Nadaraya-Watson Kernel Estimator

$$\hat{f}_n(X) = \sum_{i=1}^n w_i Y_i \text{ Where } w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

- Weight each training point based on distance to test point
- Boxcar kernel yields local average

boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$



$$K(x) \geq 0,$$

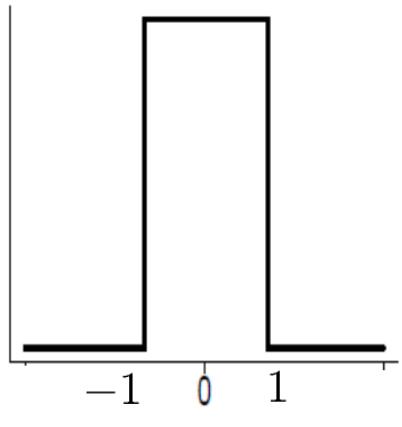
$$\int K(x)dx = 1$$

# Kernels

$$K\left(\frac{X_j - x}{\Delta}\right)$$

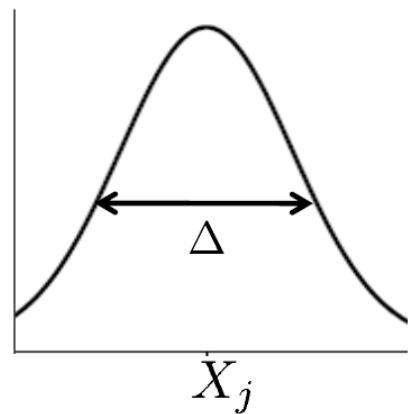
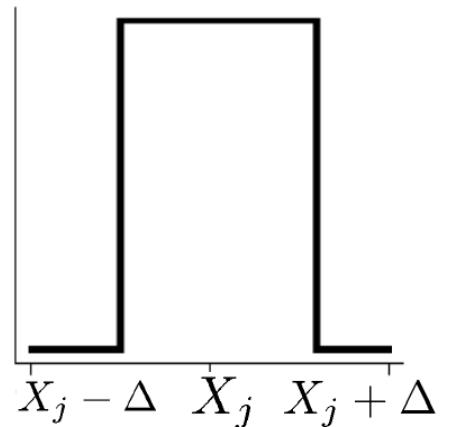
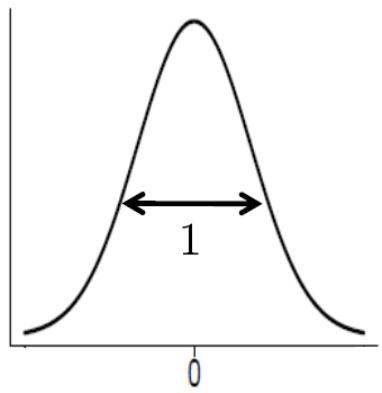
boxcar kernel :

$$K(x) = \frac{1}{2}I(x),$$



Gaussian kernel :

$$K(x) = \frac{1}{\sqrt{2\pi}}e^{-x^2/2}$$



# Choice of kernel bandwidth $h$

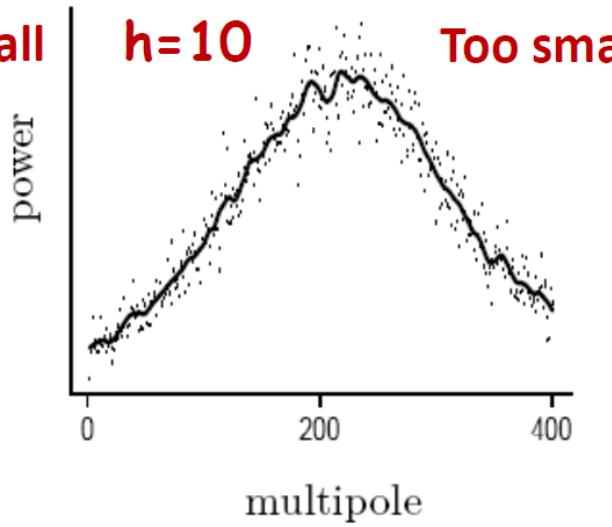
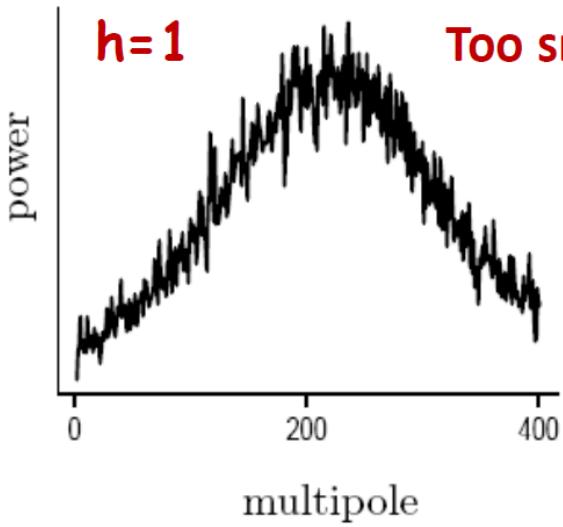
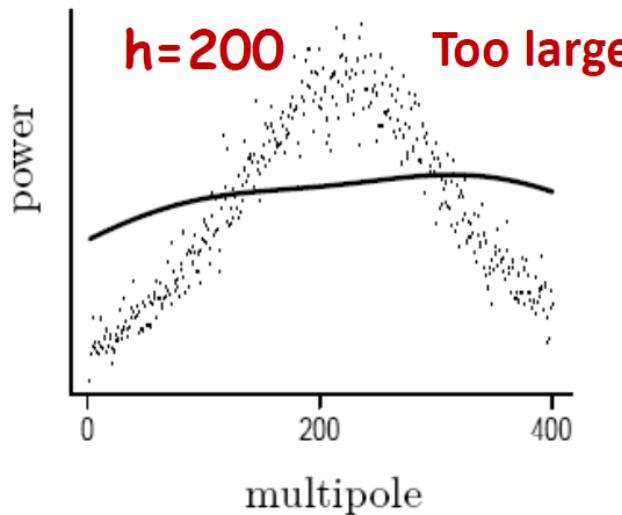
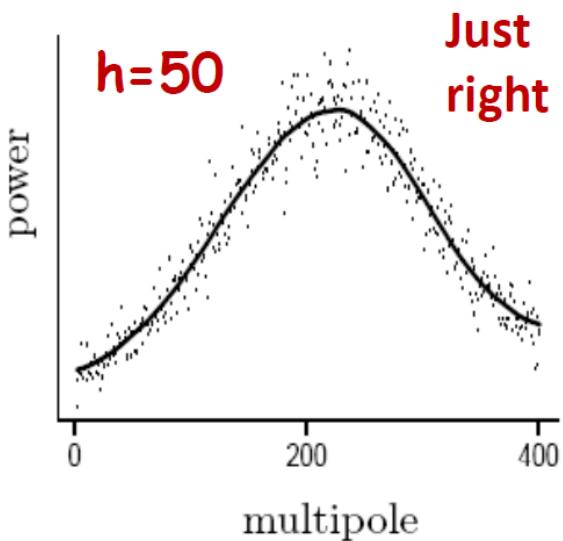


Image Source:  
Larry's book – All  
of Nonparametric  
Statistics



Choice of kernel is  
not that important

# Kernel Regression as Weighted Least Squares



$$\min_f \sum_{i=1}^n w_i(f(X_i) - Y_i)^2$$

$$w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

Weighted Least Squares

Kernel regression corresponds to locally constant estimator obtained from (locally) weighted least squares

i.e. set  $f(X_i) = \beta$  (a constant)



# Kernel Regression as Weighted Least Squares

set  $f(X_i) = \beta$  (a constant)

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

↓  
constant

$$w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

$$\frac{\partial J(\beta)}{\partial \beta} = 2 \sum_{i=1}^n w_i (\beta - Y_i) = 0$$

Notice that  $\sum_{i=1}^n w_i = 1$

$$\Rightarrow \hat{f}_n(X) = \hat{\beta} = \sum_{i=1}^n w_i Y_i$$





To minimize  $J$ , let's find its derivatives with respect to  $\theta$ .

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (\vec{X}\theta - \vec{y})^T (\vec{X}\theta - \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} \left( (\vec{X}\theta)^T \vec{X}\theta - (\vec{X}\theta)^T \vec{y} - \vec{y}^T (\vec{X}\theta) + \vec{y}^T \vec{y} \right)$$

$$a^T b = b^T a = \frac{1}{2} \nabla_{\theta} \left( \theta^T (\vec{X}^T \vec{X}) \theta - \vec{y}^T (\vec{X}\theta) - \vec{y}^T (\vec{X}\theta) \right)$$

$$\nabla_x b^T x = b = \frac{1}{2} \nabla_{\theta} \left( \theta^T (\vec{X}^T \vec{X}) \theta - 2(\vec{X}^T \vec{y})^T \theta \right)$$

$$\nabla_x x^T A x = 2Ax = \frac{1}{2} (2\vec{X}^T \vec{X}\theta - 2\vec{X}^T \vec{y})$$

$$= \vec{X}^T \vec{X}\theta - \vec{X}^T \vec{y}$$



$$\theta = (X^T X)^{-1} X^T y$$

$(X^T X)^{-1}$  is inverse of matrix  $X^T X$ .

Octave: **pinv (X' \*X) \*X' \*y**



## Normal equation

$$\theta = (X^T X)^{-1} X^T y$$

- What if  $X^T X$  is non-invertible? (singular/  
degenerate)
- Octave: **pinv (X' \*X) \*X' \*y**

# Local Linear/Polynomial Regression



$$\min_f \sum_{i=1}^n w_i(f(X_i) - Y_i)^2$$

$$w_i(X) = \frac{K\left(\frac{X-X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X-X_i}{h}\right)}$$

## Weighted Least Squares

Local Polynomial regression corresponds to locally polynomial estimator obtained from (locally) weighted least squares

i.e. set  $f(X_i) = \beta_0 + \beta_1(X_i - X) + \frac{\beta_2}{2!}(X_i - X)^2 + \dots + \frac{\beta_p}{p!}(X_i - X)^p$   
**(local polynomial of degree p around X)**