

UVA CS 6316/4501 – Fall 2016 Machine Learning

Lecture 9: Review of Regression

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

Where are we ? →

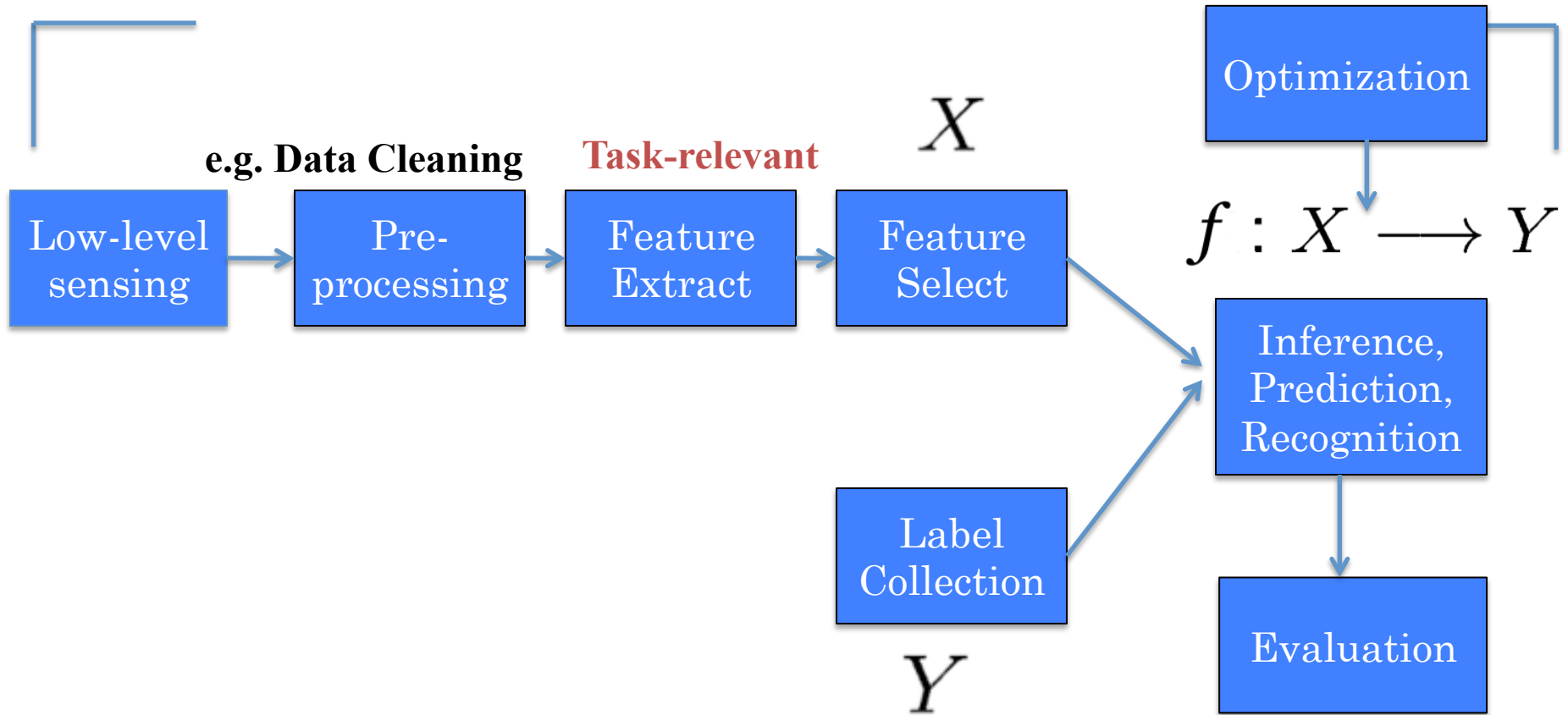
Five major sections of this course

- ❑ Regression (supervised)
- ❑ Classification (supervised)
- ❑ Unsupervised models
- ❑ Learning theory
- ❑ Graphical models

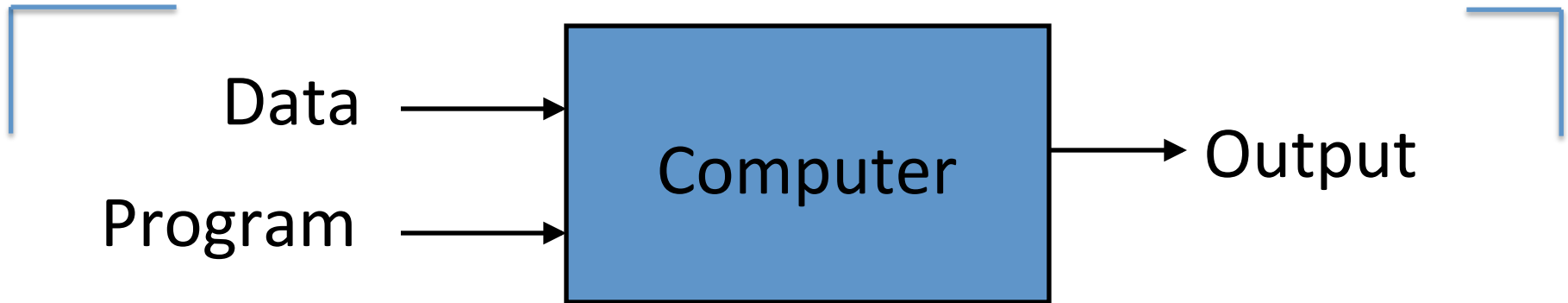
Today

- ❑ Review of basic pipeline
- ❑ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ❑ Feature Selection
- ❑ Model Selection

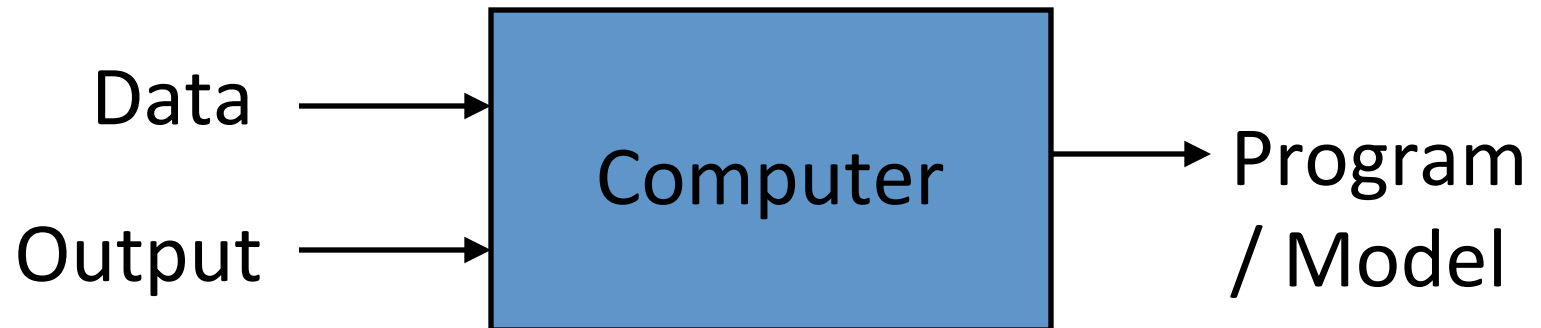
A Typical Machine Learning Pipeline



Traditional Programming



Machine Learning



e.g. SUPERVISED LEARNING

$$f : X \longrightarrow Y$$

- Find function to map **input** space X to **output** space Y

- **Generalisation**: learn function / hypothesis from **past data** in order to “explain”, “predict”, “model” or “control” **new** data examples

KEY

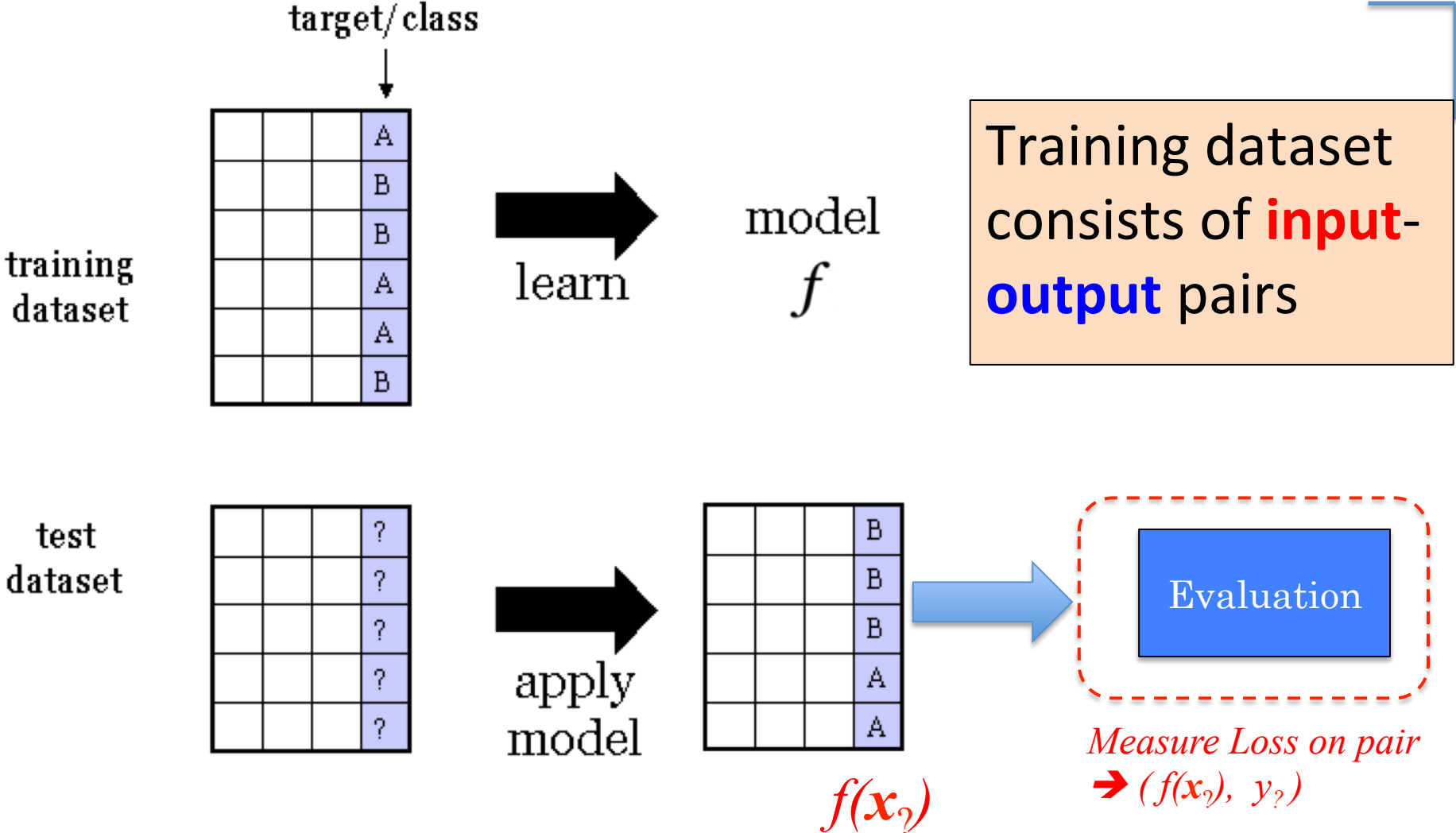
	X_1	X_2	X_3	Y
S_1				
S_2				
S_3				
S_4				
S_5				
S_6				

A Dataset

$$f : X \longrightarrow Y$$


- **Data/points/instances/examples/samples/records:** [rows]
- **Features/attributes/dimensions/independent variables/covariates/predictors/regressors:** [columns, except the last]
- **Target/outcome/response/label/dependent variable:** special column to be predicted [last column]

SUPERVISED LEARNING




Evaluation Metric

e.g. for linear regression models

training dataset 

$$\mathbf{X}_{train} = \begin{bmatrix} \text{--} & \mathbf{x}_1^T & \text{--} \\ \text{--} & \mathbf{x}_2^T & \text{--} \\ \vdots & \vdots & \vdots \\ \text{--} & \mathbf{x}_n^T & \text{--} \end{bmatrix} \quad \bar{\mathbf{y}}_{train} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

test dataset 

$$\mathbf{X}_{test} = \begin{bmatrix} \text{--} & \mathbf{x}_{n+1}^T & \text{--} \\ \text{--} & \mathbf{x}_{n+2}^T & \text{--} \\ \vdots & \vdots & \vdots \\ \text{--} & \mathbf{x}_{n+m}^T & \text{--} \end{bmatrix} \quad \bar{\mathbf{y}}_{test} = \begin{bmatrix} y_{n+1} \\ y_{n+2} \\ \vdots \\ y_{n+m} \end{bmatrix}$$

Evaluation Metric

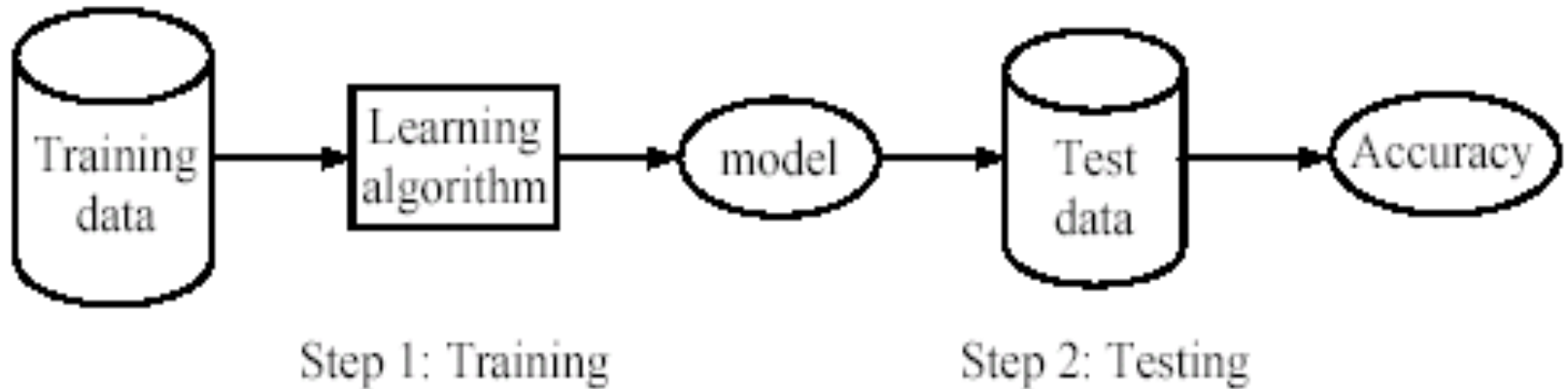
e.g. for linear regression models

- Testing MSE (mean squared error) to report:

$$J_{test} = \frac{1}{m} \sum_{i=n+1}^{n+m} (\mathbf{x}_i^T \boldsymbol{\theta}^* - y_i)^2$$

Evaluation Choice-I:

- ✓ **Training (Learning)**: Learn a model using the training data
- ✓ **Testing**: Test the model using **unseen test data** to assess the model accuracy



$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}},$$

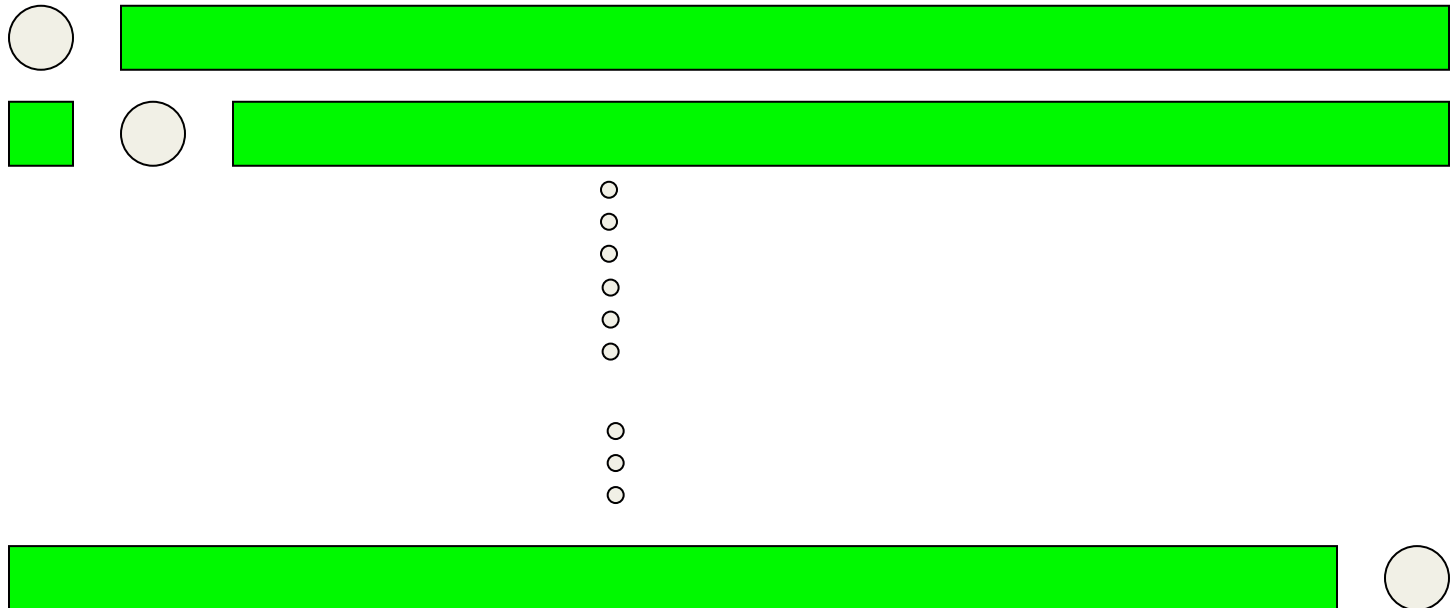
Evaluation Choice-II:

e.g. 10 fold Cross Validation

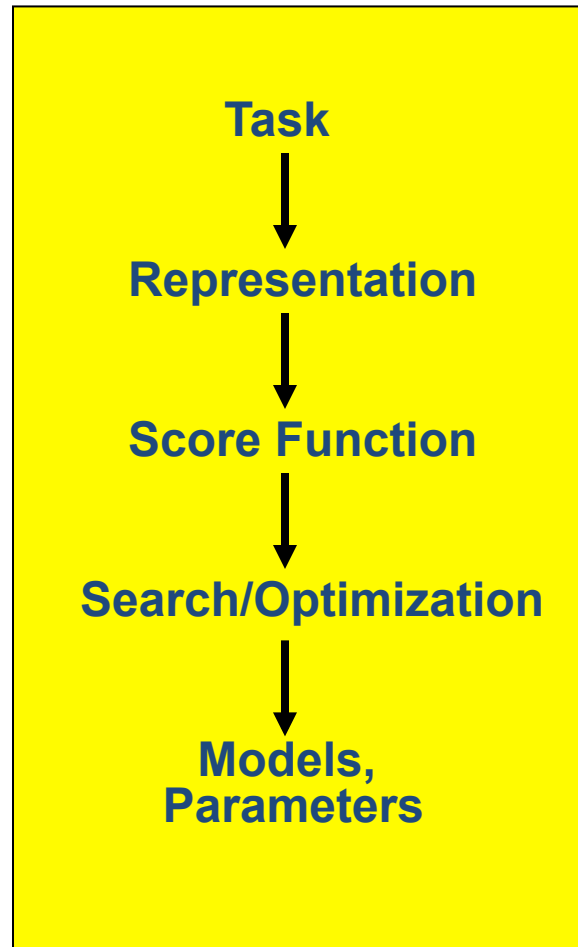
- Divide data into 10 equal pieces
- 9 pieces as training set, the rest 1 as test set
- Collect the scores from the diagonal

model	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10
1	train	train	train	train	train	train	train	train	train	test
2	train	train	train	train	train	train	train	train	test	train
3	train	train	train	train	train	train	train	test	train	train
4	train	train	train	train	train	train	test	train	train	train
5	train	train	train	train	train	test	train	train	train	train
6	train	train	train	train	test	train	train	train	train	train
7	train	train	train	test	train	train	train	train	train	train
8	train	train	test	train	train	train	train	train	train	train
9	train	test	train	train	train	train	train	train	train	train
10	test	train	train	train	train	train	train	train	train	train

e.g. Leave-one-out (n-fold cross validation)

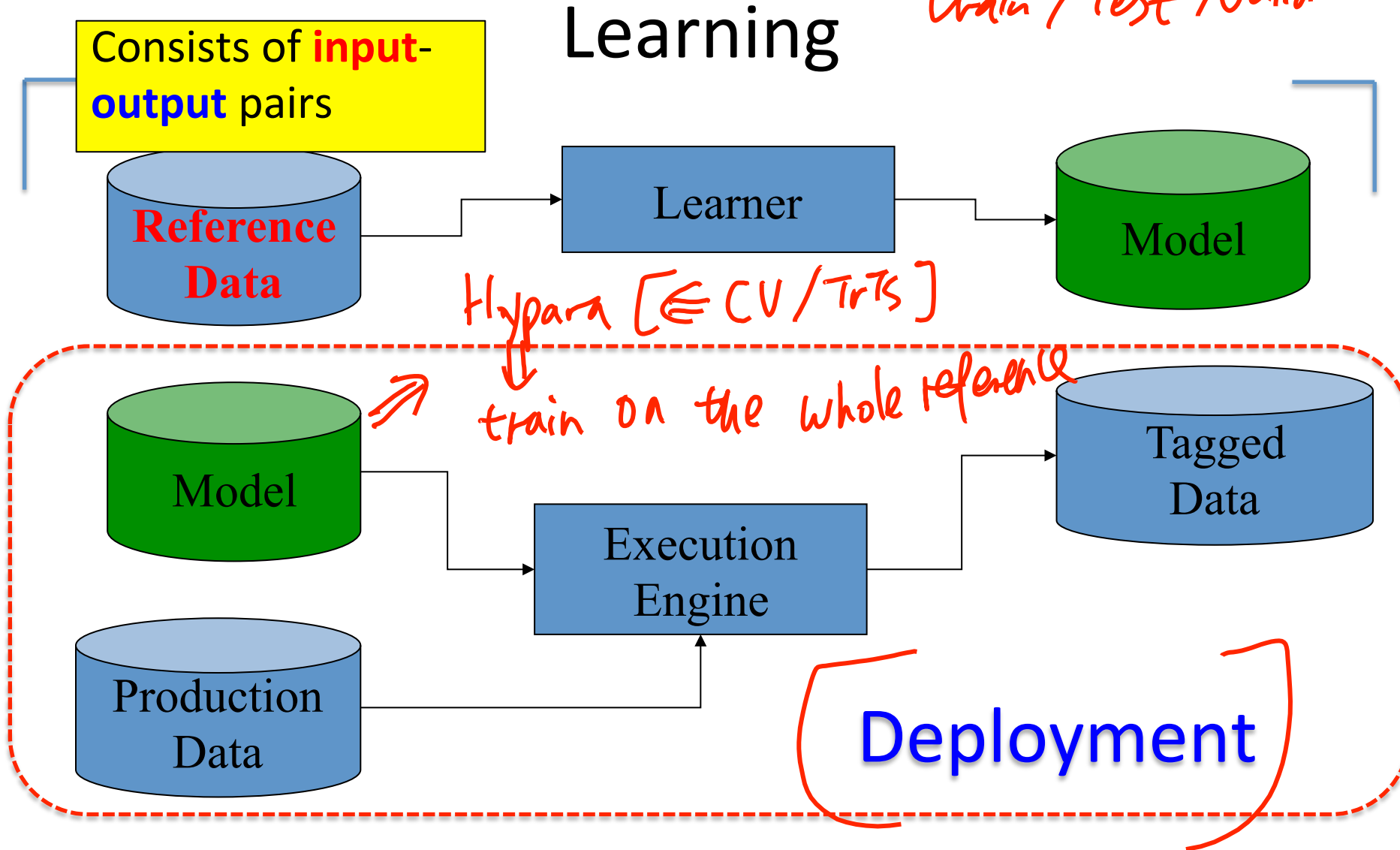


Machine Learning in a Nutshell



An **Operational** Model of Machine

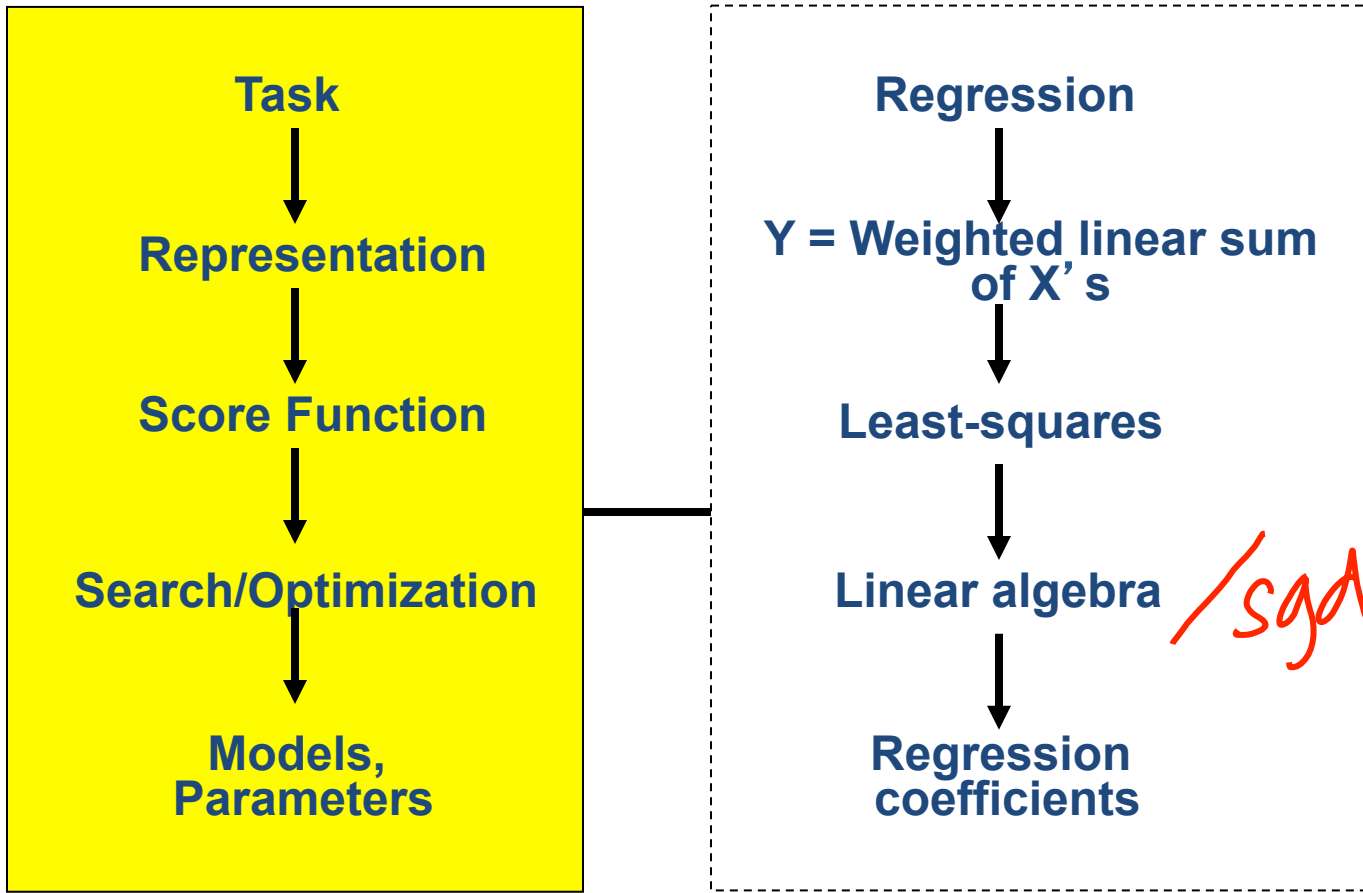
train / test / valid



Today

- ❑ Review of basic pipeline
- ❑ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ❑ Feature Selection
- ❑ Model Selection

(1) Multivariate Linear Regression



$$\hat{y} = f(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2$$

(1) Linear Regression (LR)

$$f: X \longrightarrow Y$$

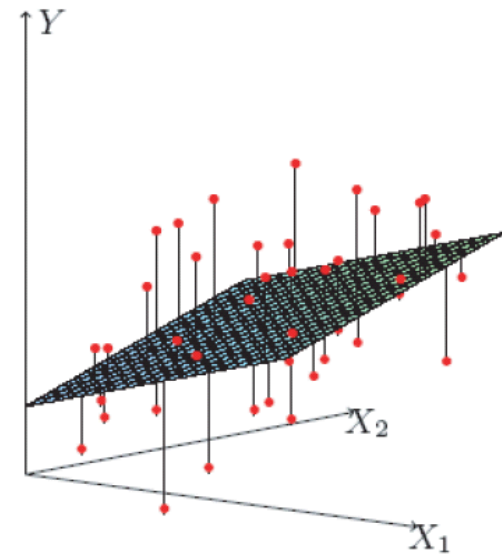
→ e.g. Linear Regression Models

$$\hat{y} = f(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2$$

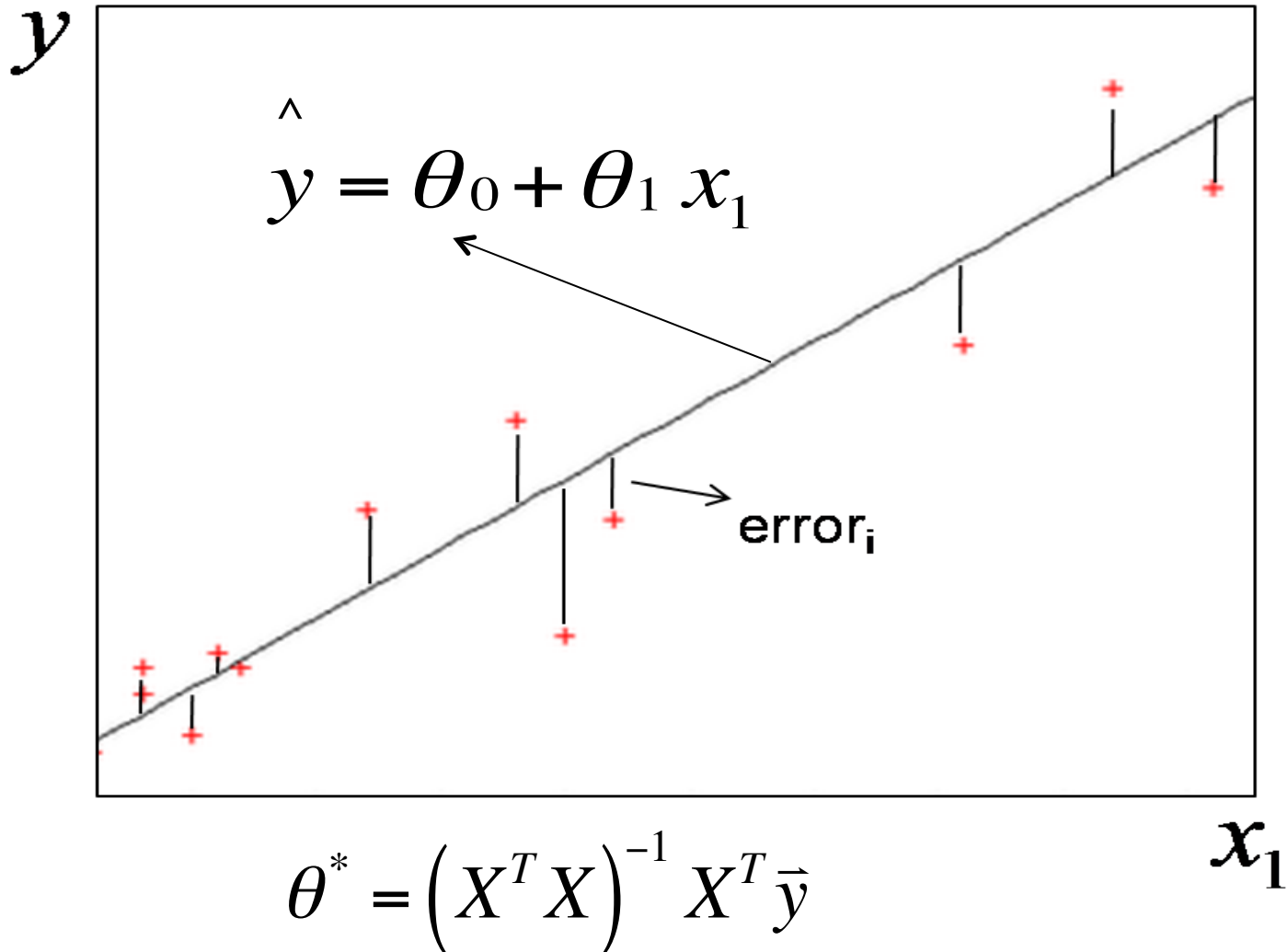
→ To minimize the “least square” cost function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i(\bar{x}_i) - y_i)^2$$

→ θ



Linear regression (1D example)

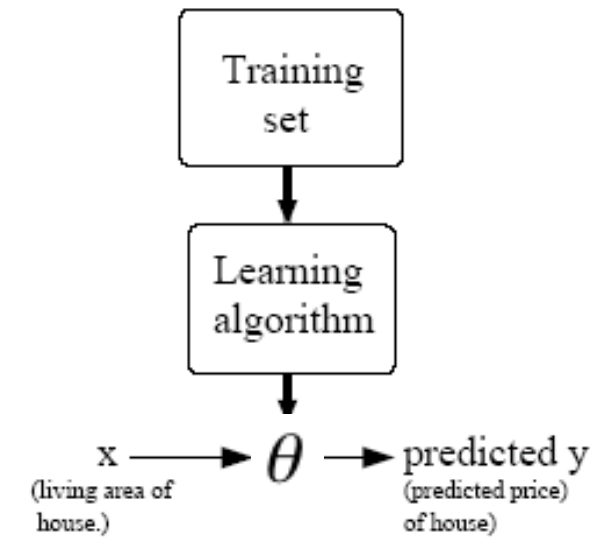


- We can represent the whole Training set:

$$\mathbf{X} = \begin{bmatrix} \text{--} & \mathbf{x}_1^T & \text{--} \\ \text{--} & \mathbf{x}_2^T & \text{--} \\ \vdots & \vdots & \vdots \\ \text{--} & \mathbf{x}_n^T & \text{--} \end{bmatrix} = \begin{bmatrix} x_1^0 & x_1^1 & \dots & x_1^{p-1} \\ x_2^0 & x_2^1 & \dots & x_2^{p-1} \\ \vdots & \vdots & \vdots & \vdots \\ x_n^0 & x_n^1 & \dots & x_n^{p-1} \end{bmatrix}$$

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

- Predicted output for each training sample:



$$\begin{bmatrix} f(\mathbf{x}_1^T) \\ f(\mathbf{x}_2^T) \\ \vdots \\ f(\mathbf{x}_n^T) \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \theta \\ \mathbf{x}_2^T \theta \\ \vdots \\ \mathbf{x}_n^T \theta \end{bmatrix} = \mathbf{X}\theta$$

Our goal:

Method I: normal equations

- Write the cost function in matrix form:

$$\begin{aligned}
 J(\theta) &= \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2 \\
 &= \frac{1}{2} (X\theta - \bar{\mathbf{y}})^T (X\theta - \bar{\mathbf{y}}) \\
 &= \frac{1}{2} (\theta^T X^T X \theta - \theta^T X^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T X \theta + \bar{\mathbf{y}}^T \bar{\mathbf{y}})
 \end{aligned}$$

$$\mathbf{X} = \begin{bmatrix} \text{--} & \mathbf{x}_1^T & \text{--} \\ \text{--} & \mathbf{x}_2^T & \text{--} \\ \vdots & \vdots & \vdots \\ \text{--} & \mathbf{x}_n^T & \text{--} \end{bmatrix} \quad \bar{\mathbf{y}} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

To minimize $J(\theta)$, take derivative and set to zero:

$$\Rightarrow X^T X \theta = X^T \bar{\mathbf{y}}$$

The normal equations

$$\Downarrow$$

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

Method II: LR with batch Steepest descent / Gradient descent

$$\theta_t = \theta_{t-1} - \alpha \nabla J(\theta_{t-1})$$

For the t-th epoch

$$\nabla_{\theta} J = \left[\frac{\partial}{\partial \theta_1} J, \dots, \frac{\partial}{\partial \theta_k} J \right]^T = - \sum_{i=1}^n (y_i - \bar{\mathbf{x}}_i^T \theta) \mathbf{x}_i$$

$$\theta_j^{t+1} = \theta_j^t + \alpha \sum_{i=1}^n (y_i - \bar{\mathbf{x}}_i^T \theta^t) x_i^j$$

– This is as a **batch** gradient descent algorithm

Method III: LR with Stochastic GD →

- From the batch steepest descent rule:

$$\theta_j^{t+1} = \theta_j^t + \alpha \sum_{i=1}^n (y_i - \bar{\mathbf{x}}_i^T \theta^t) x_i^j$$

- For a single training point, we have:

$$\longrightarrow \theta^{t+1} = \theta^t + \alpha (y_i - \bar{\mathbf{x}}_i^T \theta^t) \bar{\mathbf{x}}_i$$

- a "**stochastic**", "**coordinate**" descent algorithm
- This can be used as an **on-line** algorithm

Method IV: Newton's method for optimization

- The most basic **second-order** optimization algorithm

$$\theta_{k+1} = \theta_k - \mathbf{H}_K^{-1} \mathbf{g}_k$$

- Updating parameter with

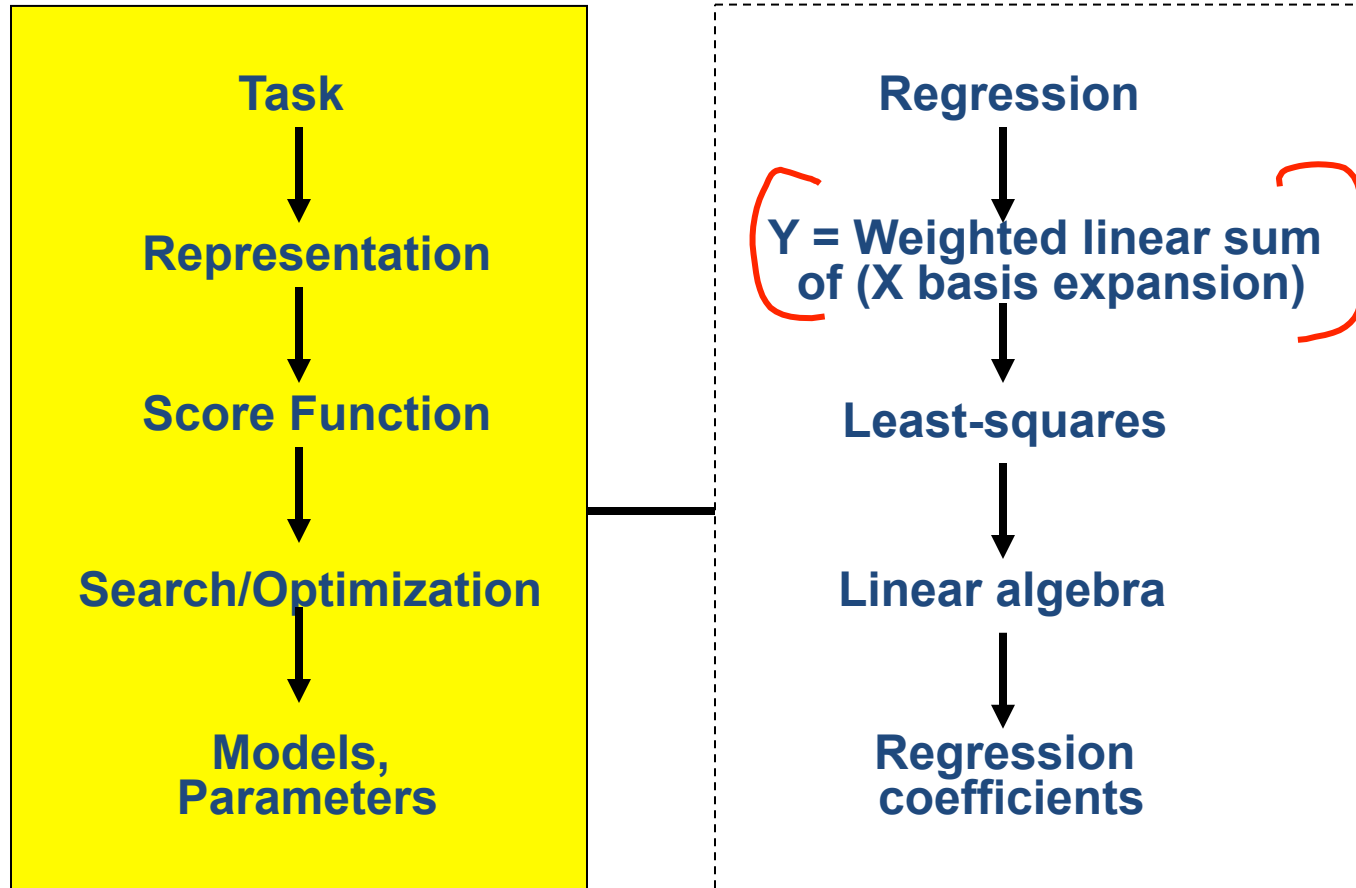
$$\begin{aligned} \Rightarrow \theta^{t+1} &= \theta^t - \mathbf{H}^{-1} \nabla f(\theta) \\ &= \theta^t - (\mathbf{X}^T \mathbf{X})^{-1} [\mathbf{X}^T \mathbf{X} \theta^t - \mathbf{X}^T \vec{y}] \end{aligned}$$

WHY ???
Normal Eq?

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

Newton's method
for Linear Regression

(2) Multivariate Linear Regression with basis Expansion



$$\hat{y} = \theta_0 + \sum_{j=1}^m \theta_j \varphi_j(x) = \varphi(x)\theta$$

(2) LR with polynomial basis functions

- LR does not mean we can only deal with linear relationships

$$y = \theta_0 + \sum_{j=1}^m \theta_j \varphi_j(x) = \varphi(x)\theta$$

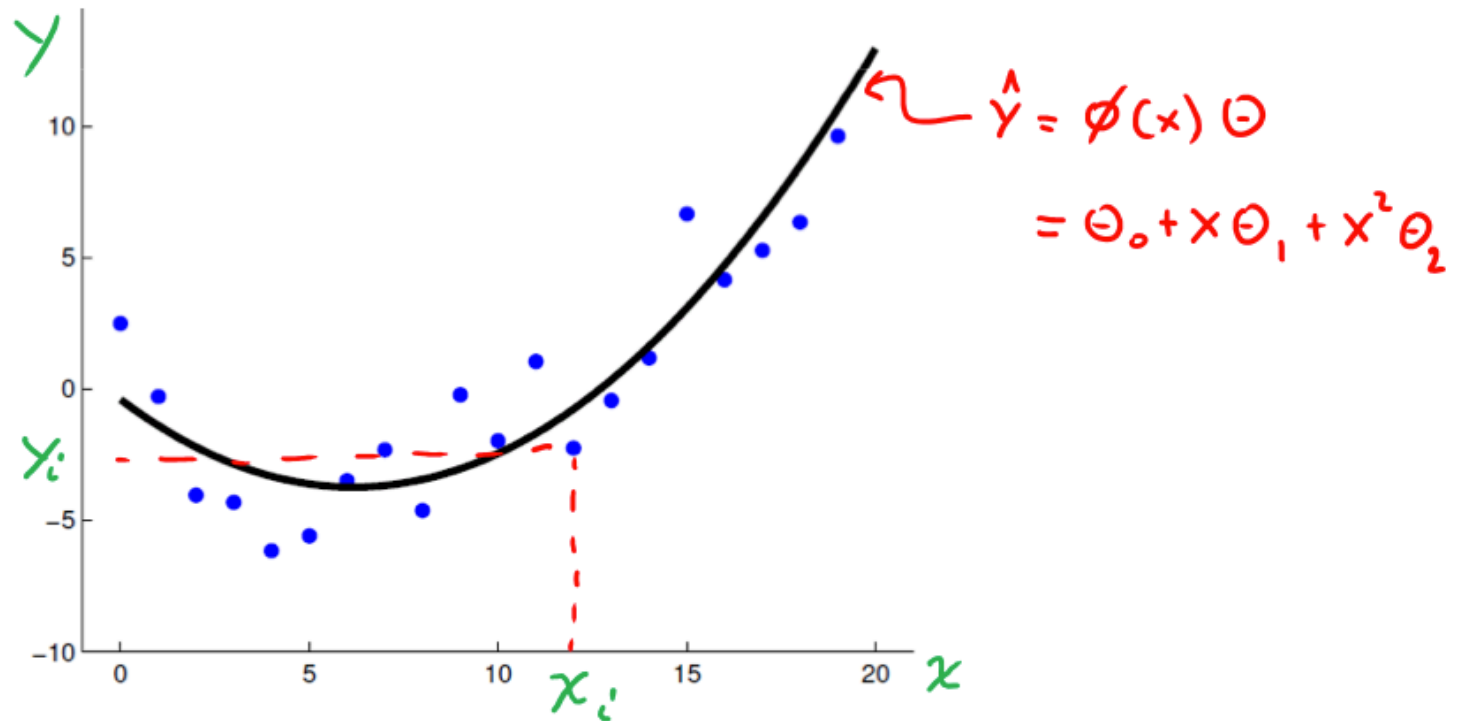
- E.g.: polynomial regression:

$$\varphi(x) := [1, x, x^2, x^3]$$

$$\theta^* = (\varphi^T \varphi)^{-1} \varphi^T \bar{y}$$

e.g. polynomial regression

For example, $\phi(x) = [1, x, x^2]$



LR with radial-basis functions

- LR does not mean we can only deal with linear relationships

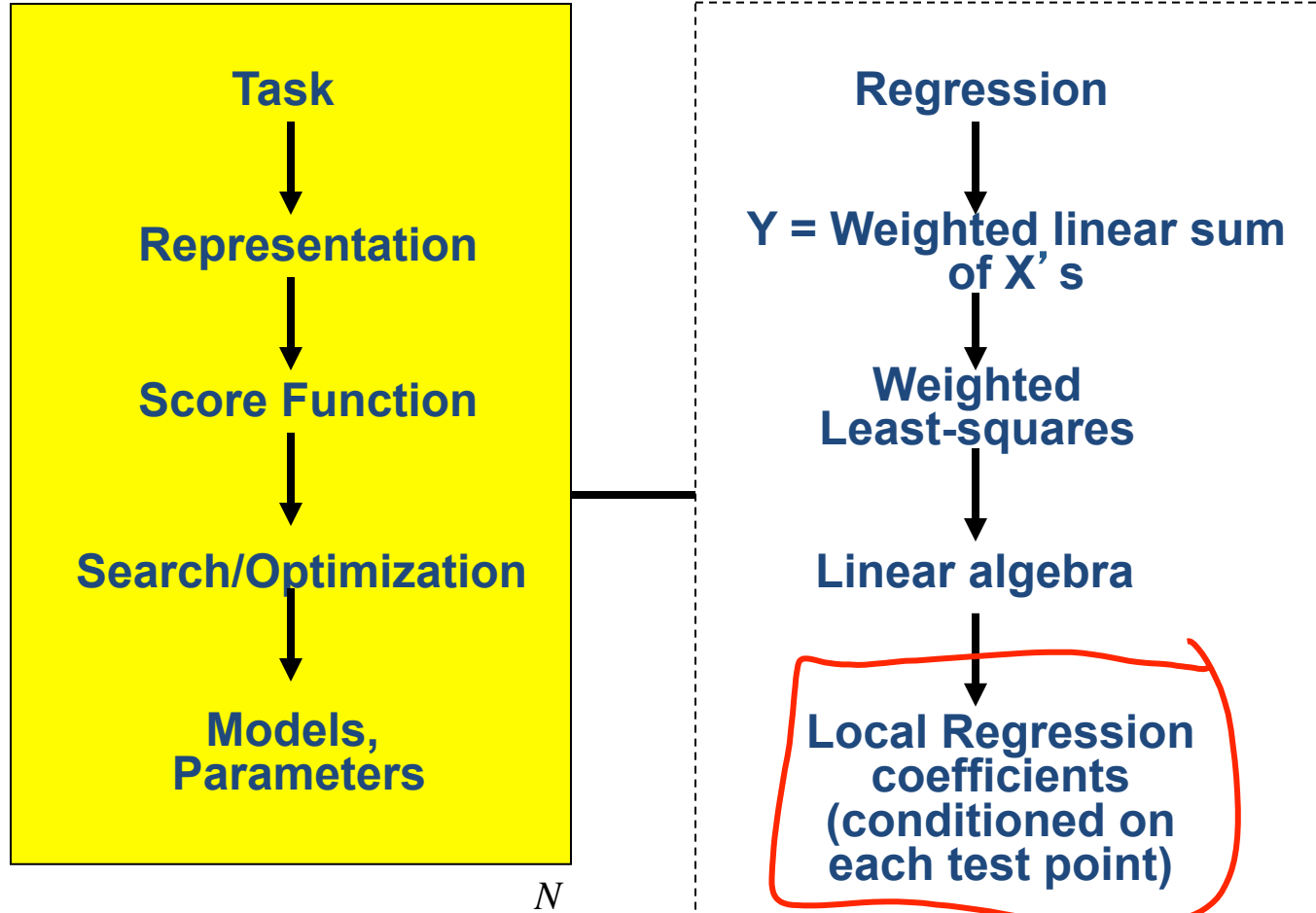
$$\hat{y} = \theta_0 + \sum_{j=1}^m \theta_j \varphi_j(x) = \varphi(x)\theta$$

- E.g.: LR with RBF regression: $K_\lambda(\underline{x}, r) = \exp\left(-\frac{(\underline{x} - r)^2}{2\lambda^2}\right)$

$$\varphi(x) := [1, K_{\lambda=1}(x, 1), K_{\lambda=1}(x, 2), K_{\lambda=1}(x, 4)]$$

$$\theta^* = (\varphi^T \varphi)^{-1} \varphi^T \vec{y}$$

(3) Locally Weighted / Kernel Regression

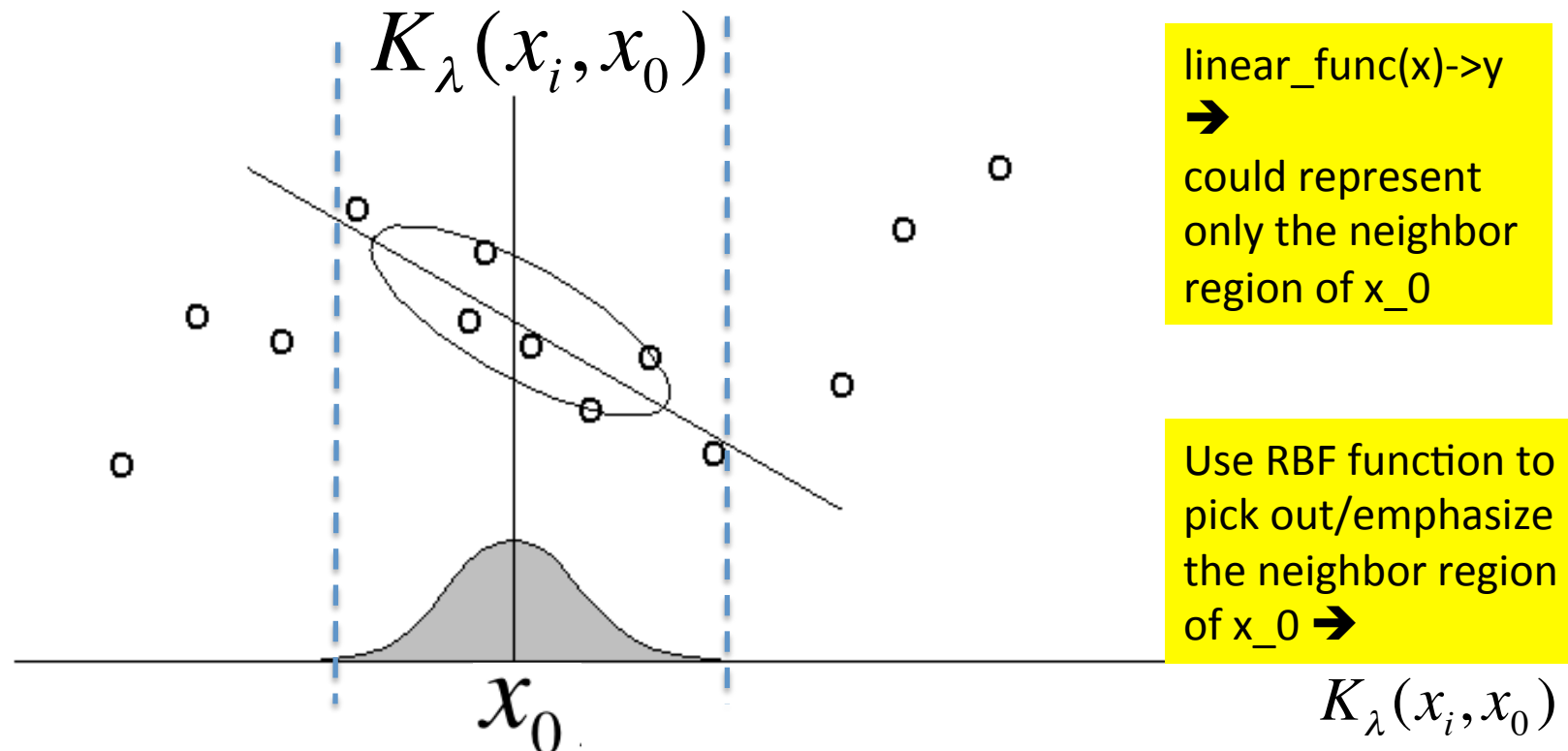


$$\min_{\alpha(x_0), \beta(x_0)} \sum_{i=1}^N K_{\lambda}(x_i, x_0) [y_i - \alpha(x_0) - \beta(x_0)x_i]^2$$

$$\hat{f}(x_0) = \hat{\alpha}(x_0) + \hat{\beta}(x_0)x_0$$

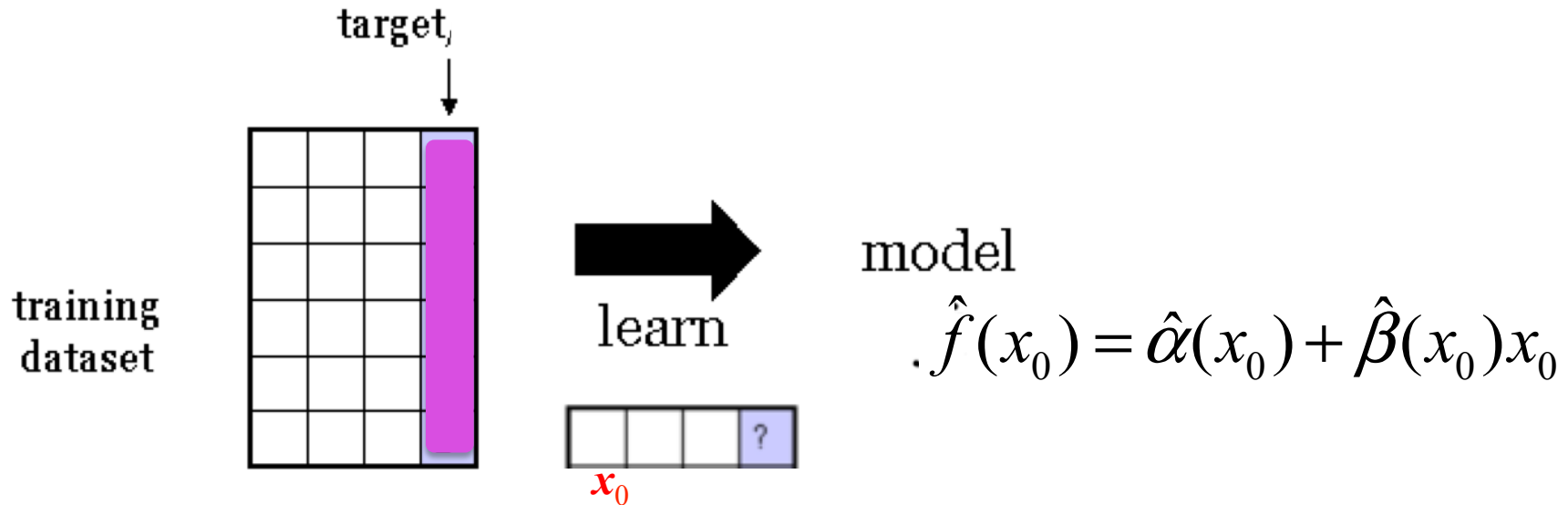
(3) Locally weighted regression

- *aka* locally weighted regression, locally linear regression, LOESS, ...



10/5/16 **Figure 2:** In locally weighted regression, points are weighted by proximity to the current x in question using a kernel. A regression is then computed using the weighted points.

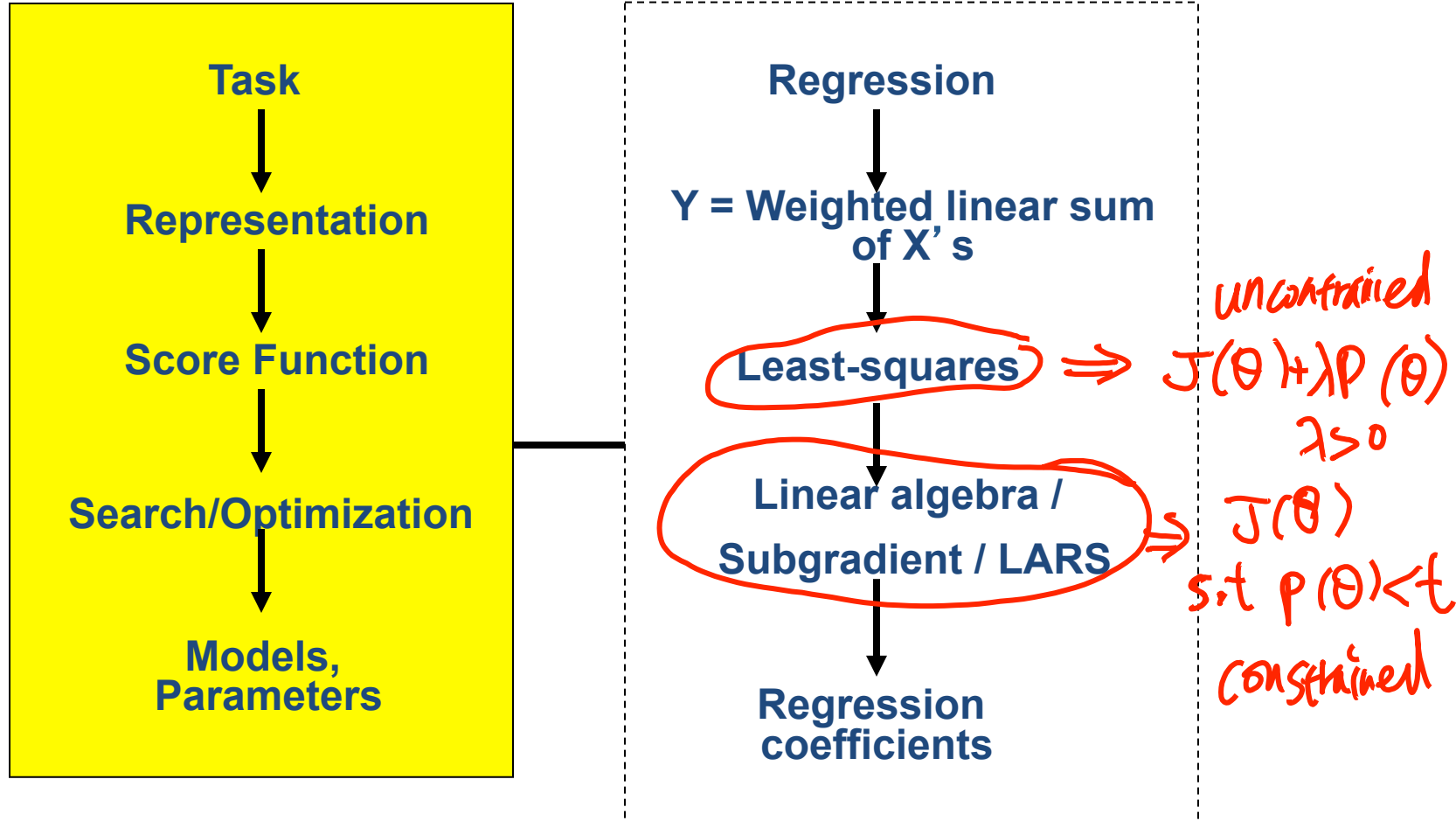
LEARNING of Locally weighted linear regression



→ Separate weighted least squares
at each target point x_0

$$\min_{\alpha(x_0), \beta(x_0)} \sum_{i=1}^N K_{\lambda}(x_i, x_0) [y_i - \alpha(x_0) - \beta(x_0)x_i]^2$$

(4) Regularized multivariate linear regression



$$\min J(\beta) = \sum_{i=1}^n \left(Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

(4) LR with Regularizations / Regularized multivariate linear regression

- Basic model

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_p x_p$$

- LR estimation:

$$\min J(\beta) = \sum \left(Y - \hat{Y} \right)^2$$

- LASSO estimation:

$$\min J(\beta) = \sum_{i=1}^n \left(Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- Ridge regression estimation:

$$\min J(\beta) = \sum_{i=1}^n \left(Y - \hat{Y} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

Error on data

+

Regularization

33/54

LR with Regularizations / Ridge Estimator

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \cdots + \hat{\beta}_p x_p$$

$$\beta^* = (X^T X + \lambda I)^{-1} X^T \bar{y}$$

- The ridge estimator is solution from RSS (regularized sum of square errors)

$$\hat{\beta}^{ridge} = \arg \min J(\beta) = \arg \min (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

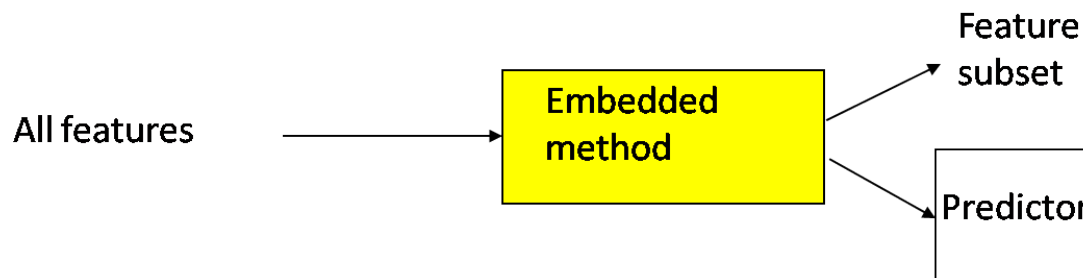
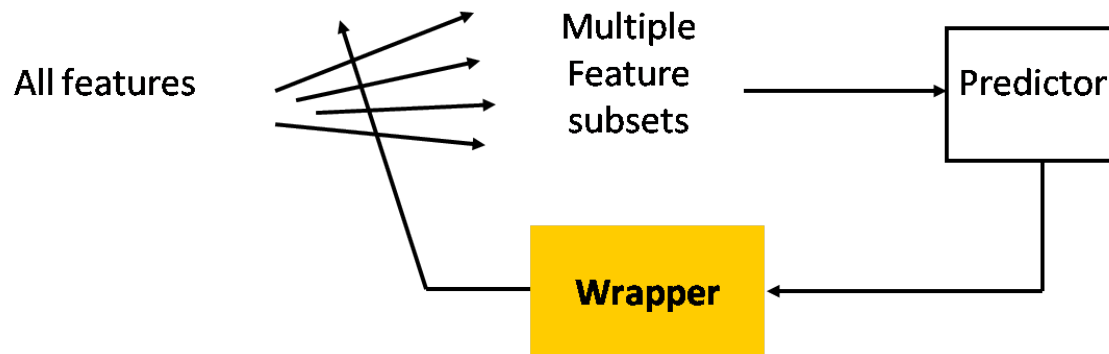
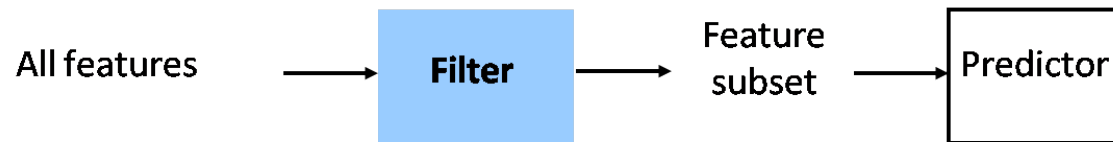
to minimize, take derivative and set to zero

Today

- ❑ Review of basic pipeline
- ❑ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ❑ **Feature Selection**
- ❑ Review of Model Selection

Summary: filters vs. wrappers vs. embedding

- **Main goal:** rank subsets of useful features



$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X} \mathbf{y}$$

$$\Rightarrow \begin{matrix} \mathbf{X}^T \mathbf{X} \\ p \times n & n \times p \end{matrix} : O(np^2)$$

$$\Rightarrow \underbrace{(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}}_{p \times p} : O(p^3)$$

$$\Rightarrow \mathbf{X} \mathbf{y} : O(np)$$

Computationally,

choose to
make $p \downarrow$
if we can
★

image at:

www.butterflyeffe.com/.../OccamsRazor.html
Remove frame

The principle of Occam's razor



Occam's razor:

states that the explanation of any phenomenon should make as few assumptions as possible, eliminating those that make no difference to any observable predictions of the theory

Occam's razor: law of parsimony

Basically it says that explanations must not include elements that have nothing to do with the phenomenon under analysis.

...not as often stated:

“The simplest explanation is the correct one”

As it is not about simplicity or complexity

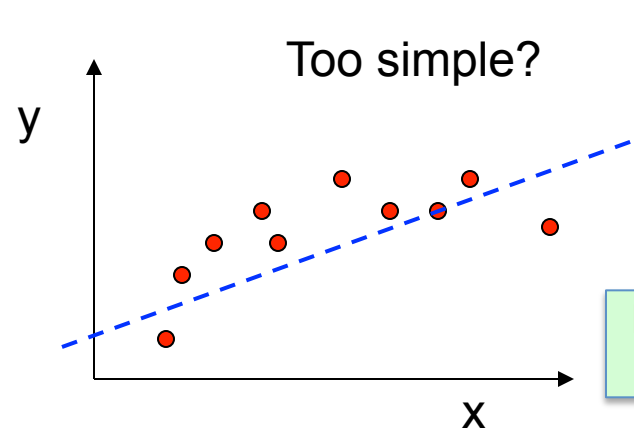
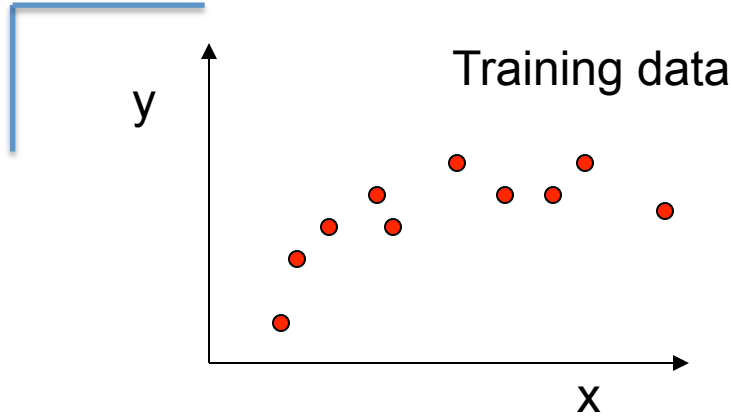
William of Occam 1288 - 1348

parsimony: extreme unwillingness to spend money or use resources.

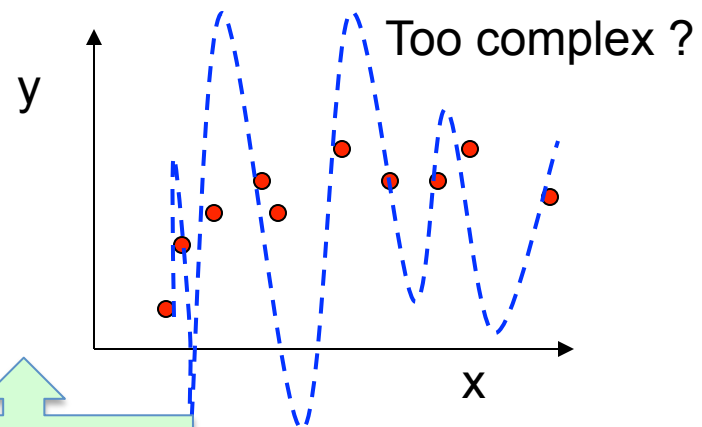
Today

- ❑ Review of basic pipeline
- ❑ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ❑ Feature Selection
- ❑ Review of Model Selection

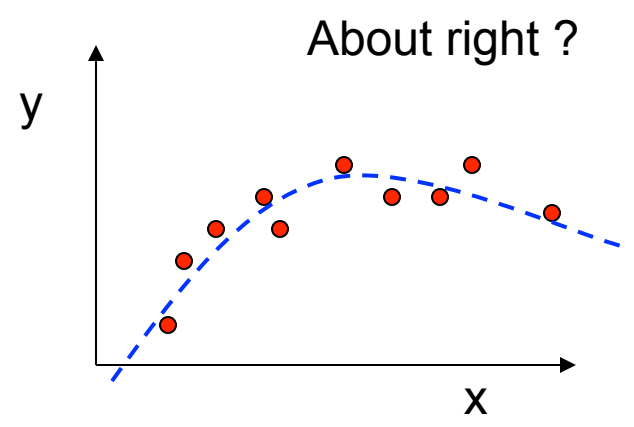
Regression: Complexity versus Goodness of Fit



Low Variance / High Bias



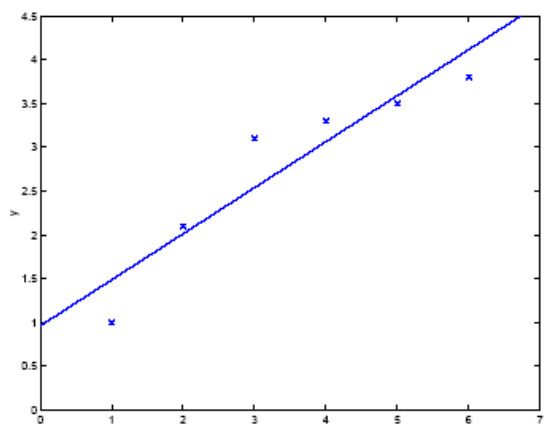
Low Bias / High Variance



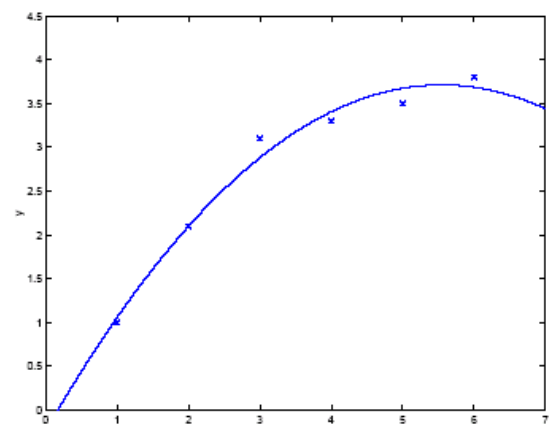
What ultimately matters: **GENERALIZATION**

Which function *f to choose?*

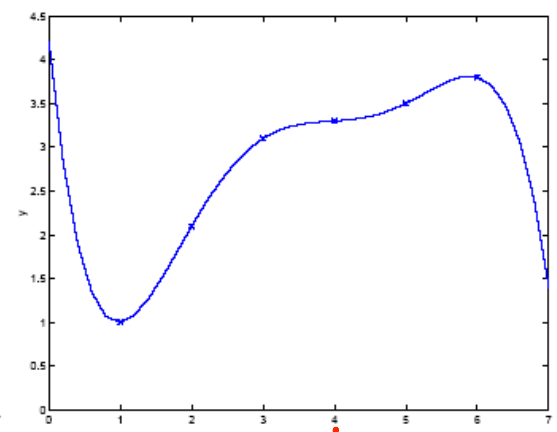
Many possible choices , e.g. LR with polynomial basis functions



$$y = \theta_0 + \theta_1 x$$



$$y = \theta_0 + \theta_1 x + \theta_2 x^2$$



$$y = \sum_{j=0}^{st} \theta_j x^j$$

Generalisation: learn function / hypothesis from **past data** in order to “explain”, “predict”, “model” or “control” **new data** examples

Choose f that generalizes well !

Which kernel width to choose ?

e.g. locally weighted LR

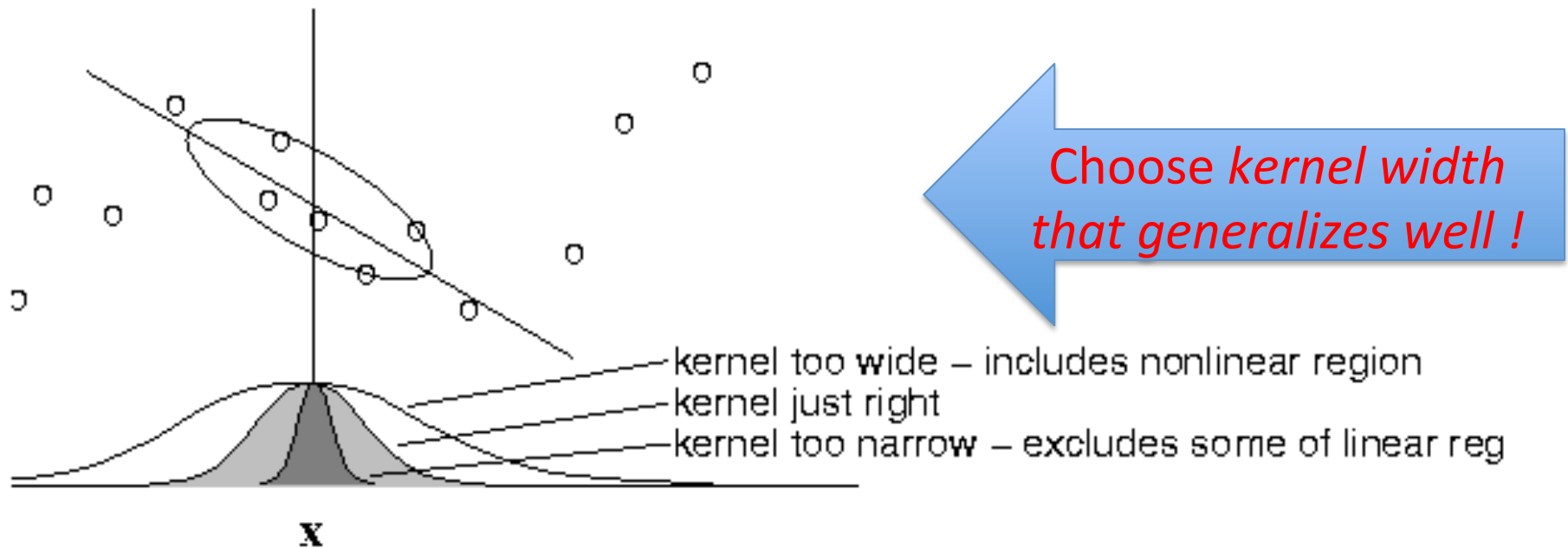


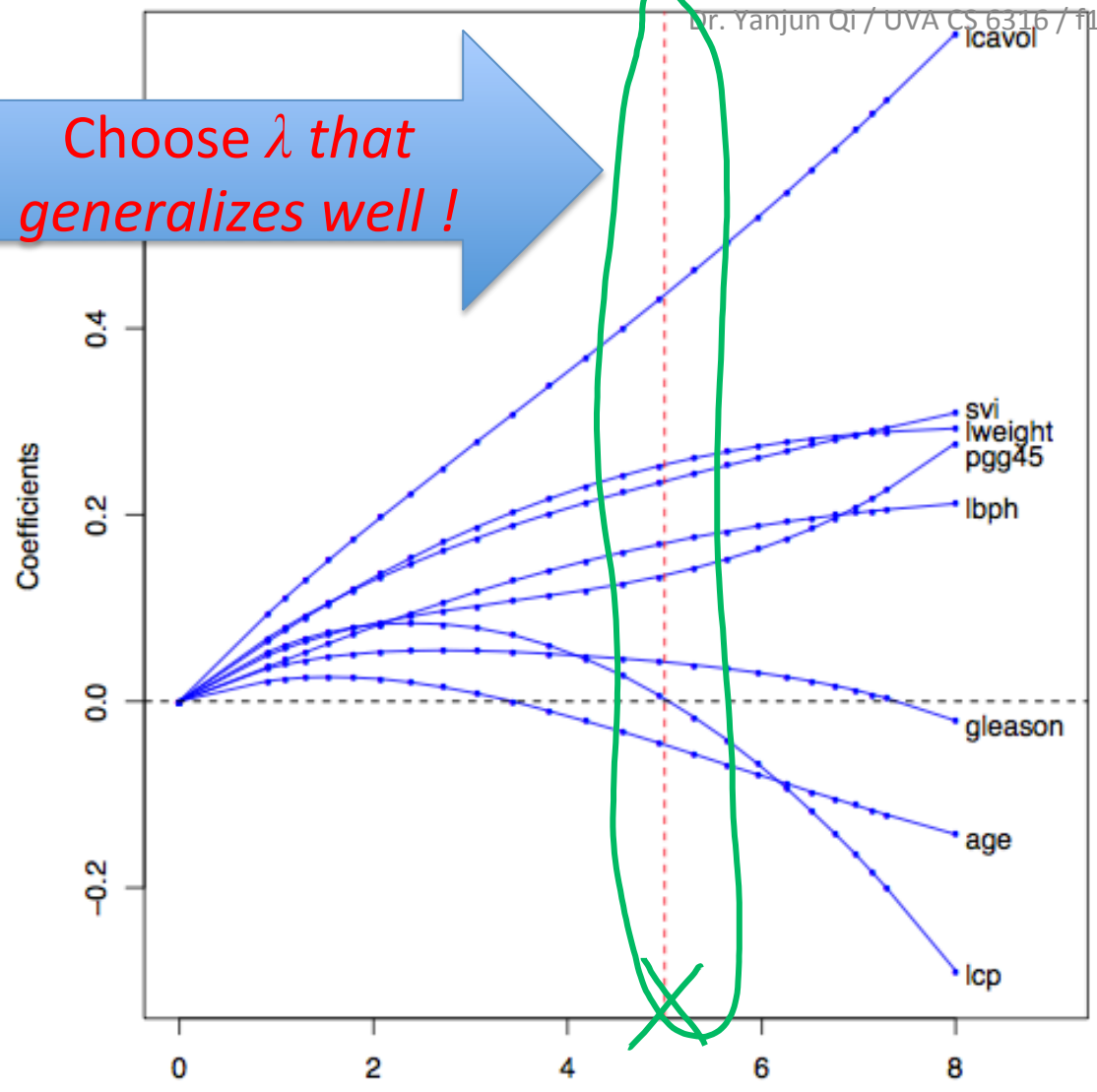
Figure 3: The estimator variance is minimized when the kernel includes as many training points as can be accommodated by the model. Here the linear LOESS model is shown. Too large a kernel includes points that degrade the fit, too small a kernel neglects points that increase confidence in the fit.

an example
(ESL Fig3.8),

Ridge Regression

when varying
 λ , how θ_j
varies.

Choose λ that
generalizes well!



$\lambda \rightarrow \infty$

$\lambda = 0$

λ increases

Choose λ that generalizes well!

Regularization path of a Lasso Estimator

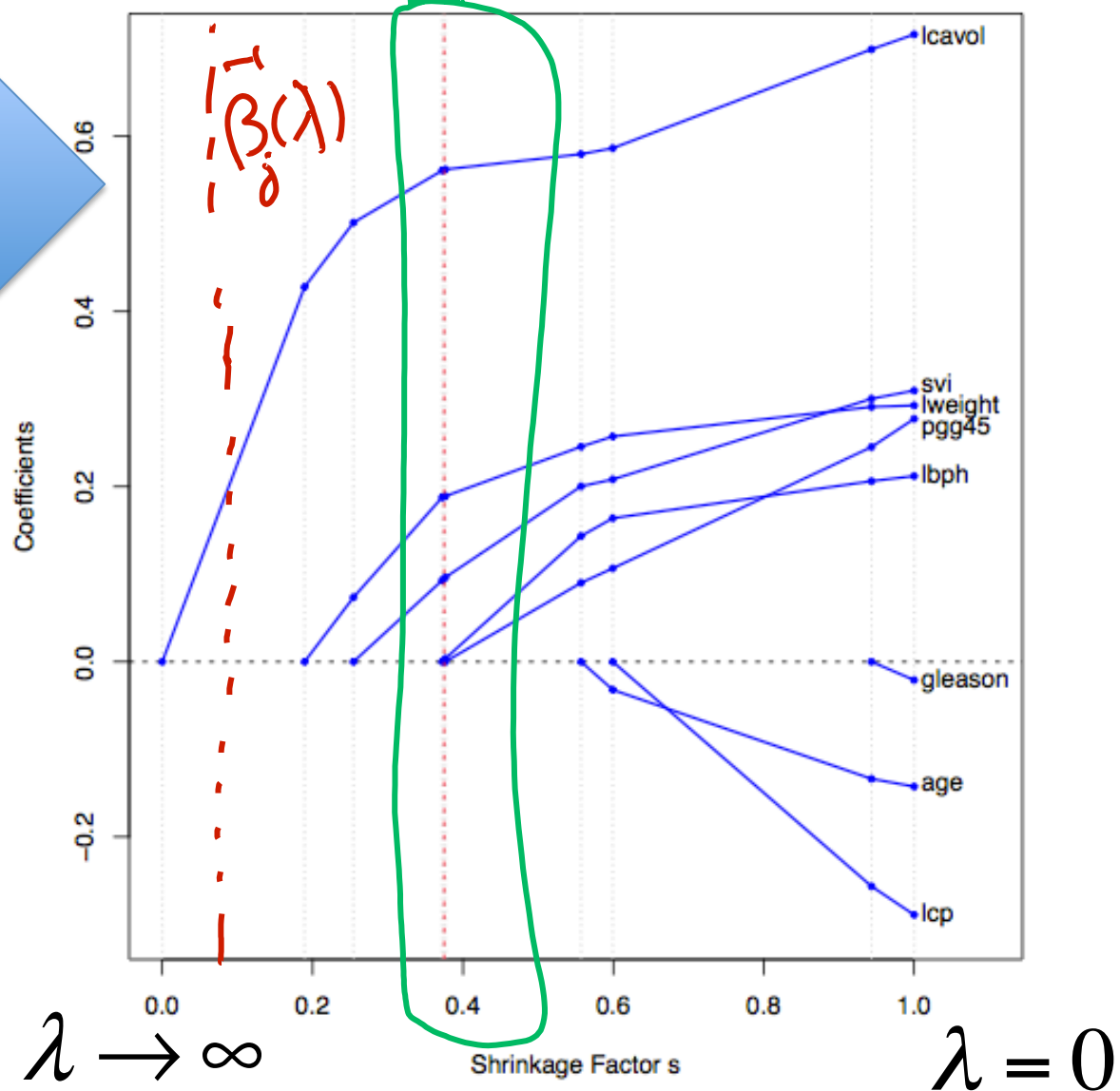
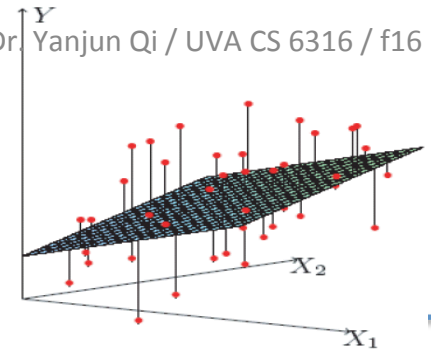


FIGURE 3.10. Profiles of lasso coefficients, as the tuning parameter t is varied. Coefficients are plotted versus $s = t / \sum_1^p |\hat{\beta}_j|$. A vertical line is drawn at $s = 0.36$, the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for ridge do not. The profiles are piece-wise linear, and so are computed only at the points displayed; see Section 3.4.4 for details.

Probabilistic Interpretation of Linear Regression (LATER)



- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

error data on each point

where ε is an error term of unmodeled effects or random noise

- Now assume that ε follows a Gaussian $N(0, \sigma)$, then we have:

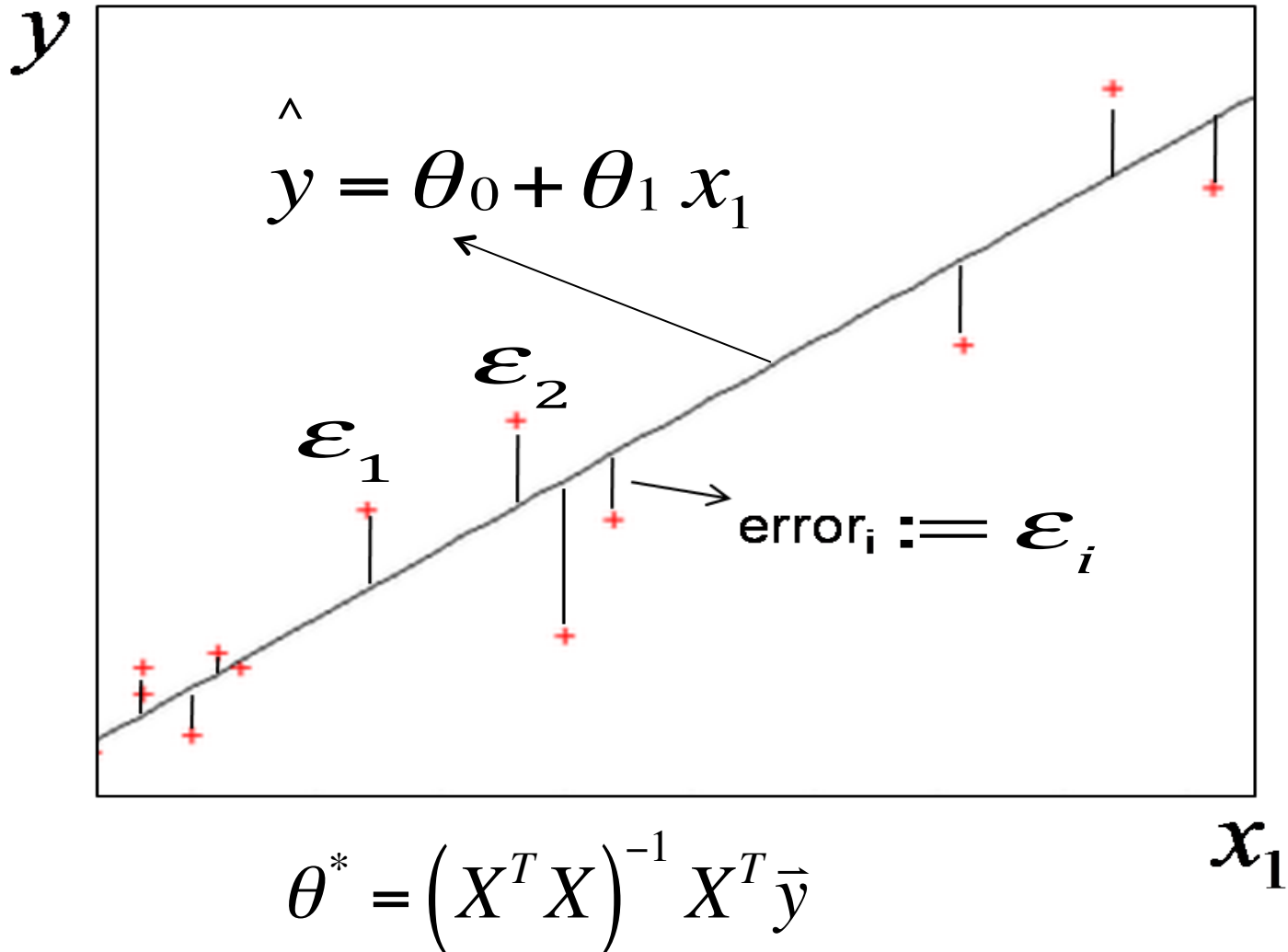
$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

Many more variations of LR from this perspective, e.g. binomial / poisson (LATER)

- By iid (among samples) assumption:

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

Linear regression (1D example)



References

- ❑ Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- ❑ Prof. Alexander Gray's slides