

UVA CS 4501 - 001 / 6501 – 007

Introduction to Machine Learning and Data Mining

Lecture 17: Review / Bias-Variance Tradeoff

Yanjun Qi / Jane, , PhD

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

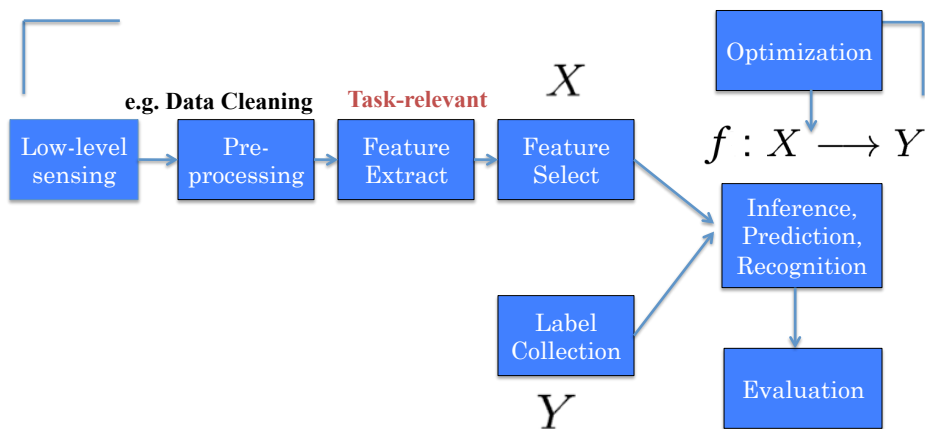
Today

- ❑ Review of basic pipeline
- ❑ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ❑ Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
- ❑ Model Selection

10/30/14

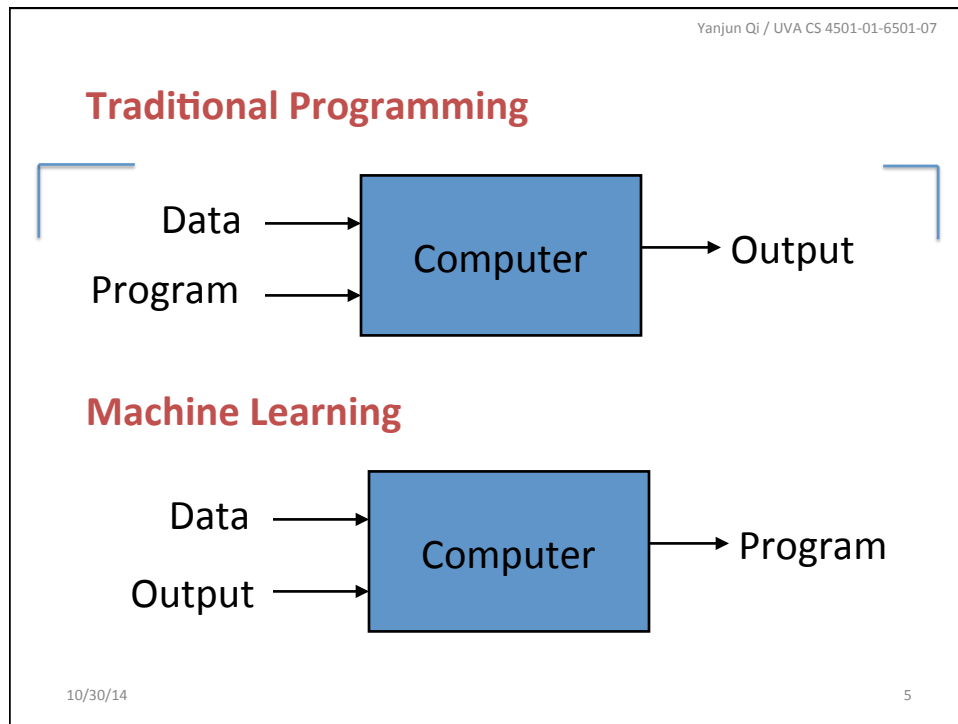
3

A Typical Machine Learning Pipeline



10/30/14

4



Yanjun Qi / UVA CS 4501-01-6501-07

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

10/30/14 6

Yanjun Qi / UVA CS 4501-01-6501-07

	X ₁	X ₂	X ₃	Y
S ₁				
S ₂				
S ₃				
S ₄				
S ₅				
S ₆				

A Dataset

$$f : X \rightarrow Y$$

- **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]

10/30/147

Yanjun Qi / UVA CS 4501-01-6501-07

SUPERVISED LEARNING

target/class

↓

		A
		B
		B
		A
		A
		B

learn

→

model

f

Training dataset consists of **input-output** pairs

test dataset

		?
		?
		?
		?
		?

apply model

→

		B
		B
		B
		A
		A

$f(x_?)$

→

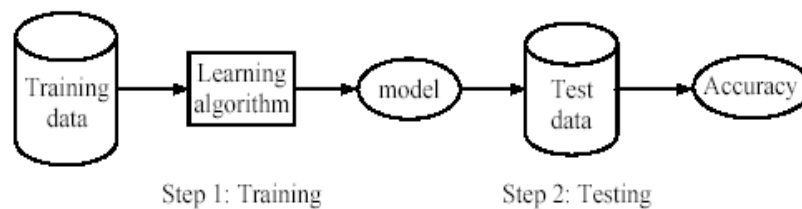
Evaluation

Measure Loss on pair $\rightarrow (f(x_?), y_?)$

10/30/148

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



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

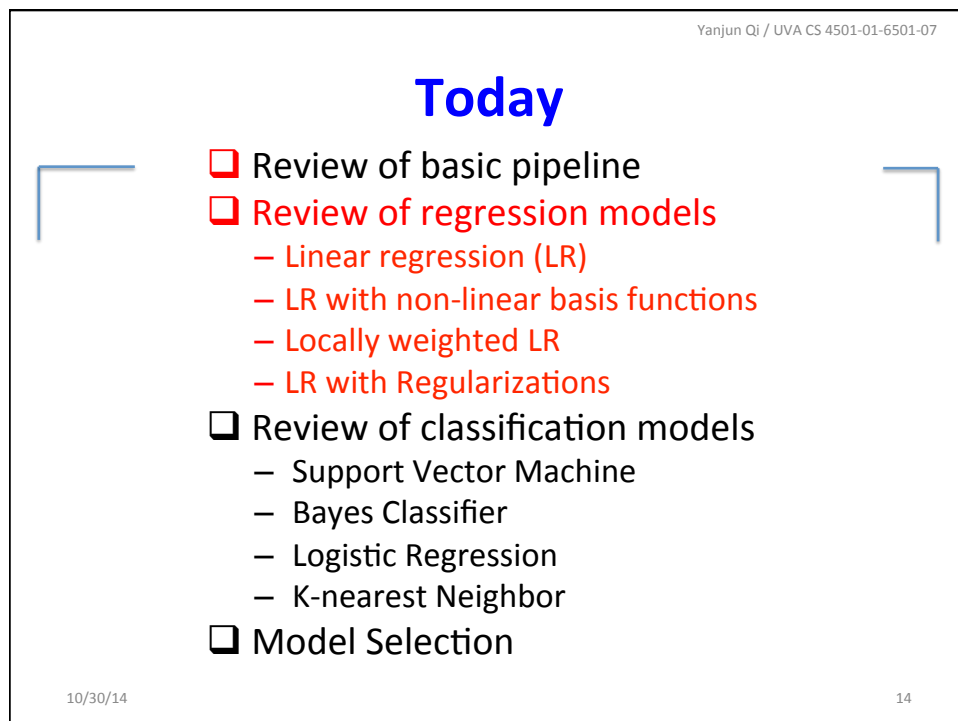
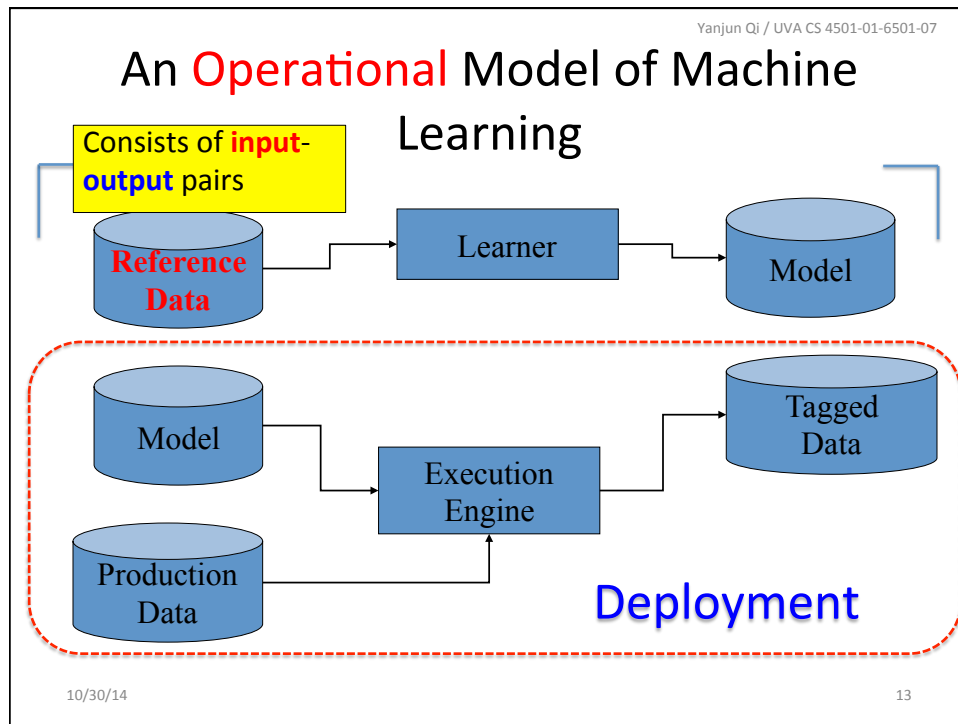
10/30/14

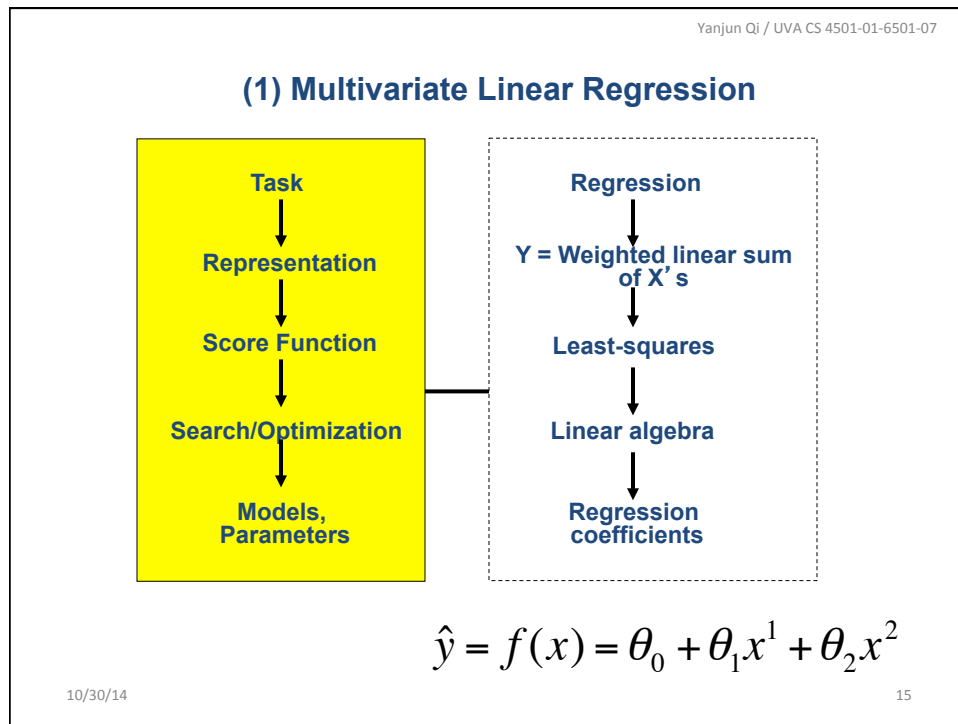
Evaluation Choice-II: e.g. 10 fold Cross Validation

- 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

10/30/14





Yanjun Qi / UVA CS 4501-01-6501-07

(1) Linear Regression (LR)

$f: X \rightarrow Y$

→ e.g. Linear Regression Models

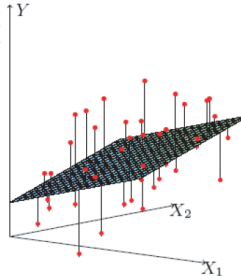
$$\hat{y} = f(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2$$

→ To minimize the cost function:

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

→

θ



10/30/14

Yanjun Qi / UVA CS 4501-01-6501-07

Our goal:

• We can represent the whole Training set:

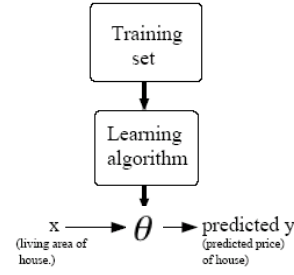
$$\mathbf{X} = \begin{bmatrix} \text{--} & \mathbf{x}_1^T & \text{--} \\ \text{--} & \mathbf{x}_2^T & \text{--} \\ \vdots & \vdots & \vdots \\ \text{--} & \mathbf{x}_n^T & \text{--} \end{bmatrix} = \begin{bmatrix} x_1^0 & x_1^1 & \dots & x_1^{p-1} \\ x_2^0 & x_2^1 & \dots & x_2^{p-1} \\ \vdots & \vdots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

• Predicted output for each training sample:

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

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

10/30/14



Yanjun Qi / UVA CS 4501-01-6501-07

Method I: normal equations

• Write the cost function in matrix form:

$$J(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \boldsymbol{\theta} - y_i)^2$$

$$= \frac{1}{2} (\mathbf{X} \boldsymbol{\theta} - \bar{\mathbf{y}})^T (\mathbf{X} \boldsymbol{\theta} - \bar{\mathbf{y}})$$

$$= \frac{1}{2} (\boldsymbol{\theta}^T \mathbf{X}^T \mathbf{X} \boldsymbol{\theta} - \boldsymbol{\theta}^T \mathbf{X}^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T \mathbf{X} \boldsymbol{\theta} + \bar{\mathbf{y}}^T \bar{\mathbf{y}})$$

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

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

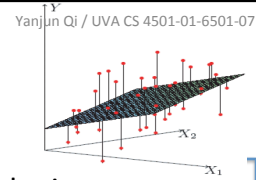
\Rightarrow $\mathbf{X}^T \mathbf{X} \boldsymbol{\theta} = \mathbf{X}^T \bar{\mathbf{y}}$

The normal equations

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

10/30/14

Probabilistic Interpretation of Linear Regression



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

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

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:

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

- By independence (among samples) assumption:

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

10/30/14

19

Probabilistic Interpretation of Linear Regression (cont.)

YanJun Qi / UVA CS 4501-01-6501-07

- Hence the log-likelihood is:

$$l(\theta) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

- Do you recognize the last term?

Yes it is:

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

- Thus under independence assumption, residual means square is equivalent to MLE of ϑ !

10/30/14

20

Method II: LR with batch Steepest descent / Gradient descent

$$\theta_t = \theta_{t-1} - \alpha \nabla J(\theta_{t-1}) \quad \text{For the t-th epoch}$$

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

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

– 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:

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

- This is known as the Least-Mean-Square update rule, or the Widrow-Hoff learning rule
- This is actually a "**stochastic**", "**coordinate**" descent algorithm
- This can be used as a **on-line** algorithm

Yanjun Qi / UVA CS 4501-01-6501-07

Method IV: Newton's method for optimization

- The most basic **second-order** optimization algorithm
- Updating parameter with

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

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

WHY ???
Normal Eq?

Newton's method for Linear Regression

10/ 9/9/14
25

Yanjun Qi / UVA CS 4501-01-6501-07

(2) Multivariate Linear Regression with basis Expansion

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

Models,
Parameters

Regression

↓

Y = Weighted linear sum
of (X basis expansion)

↓

Least-squares

↓

Linear algebra

↓

Regression
coefficients

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

10/30/14
24

(2) LR with polynomial basis functions

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

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

where the $\phi_j(x)$ are fixed basis functions (also define $\phi_0(x) = 1$)

- E.g.: polynomial regression:

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

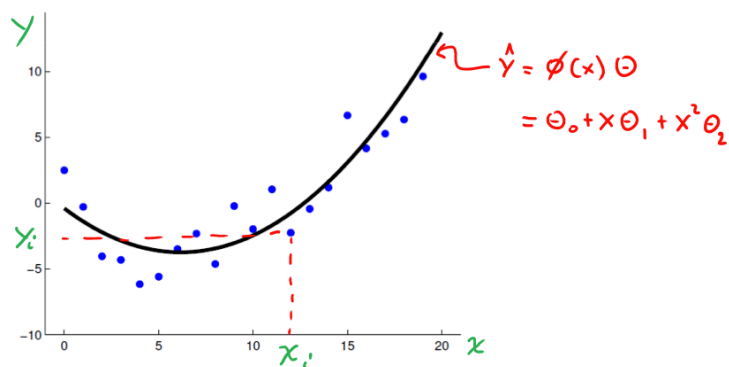
$$\theta^* = (\varphi^T \varphi)^{-1} \varphi^T \bar{y}$$

10/30/14

25

e.g. polynomial regression

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



10/30/14

26
Dr. Nando de Freitas's tutorial slide

(2) LR with radial-basis functions

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

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

where the $\phi_j(x)$ are fixed basis functions (also define $\phi_0(x) = 1$)

- E.g.: LR with RBF regression:

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

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

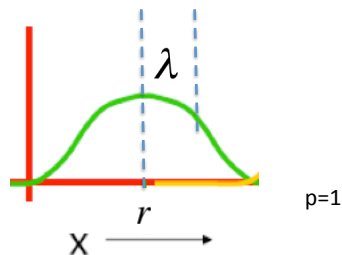
10/30/14

27

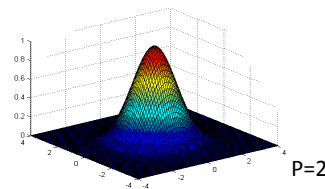
RBF = radial-basis function: a function which depends only on the radial distance from a centre point

Gaussian RBF →
$$K_{\lambda}(\underline{x}, r) = \exp\left(-\frac{(\underline{x}-r)^2}{2\lambda^2}\right)$$

as distance from the centre r increases, the output of the RBF decreases



10/30/14



28

Yanjun Qi / UVA CS 4501-01-6501-07

(2) Linear regression with RBF basis functions (predefined centres)

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

$$\hat{y} = \theta_0 + e^{-\|x-1\|^2} \theta_1 + e^{-\|x-2\|^2} \theta_2 + e^{-\|x-4\|^2} \theta_3$$

The green curve is a weighted sum of the red curves

Dr. Nando de Freitas's tutorial slide

Yanjun Qi / UVA CS 4501-01-6501-07

$$k_{\lambda=1}(x,1) = e^{-\|x-1\|^2}$$

$$k_{\lambda=1}(x,2) = e^{-\|x-2\|^2}$$

$$k_{\lambda=1}(x,4) = e^{-\|x-4\|^2}$$

$$\varphi(x) := [1, e^{-\|x-1\|^2}, e^{-\|x-2\|^2}, e^{-\|x-4\|^2}]$$

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

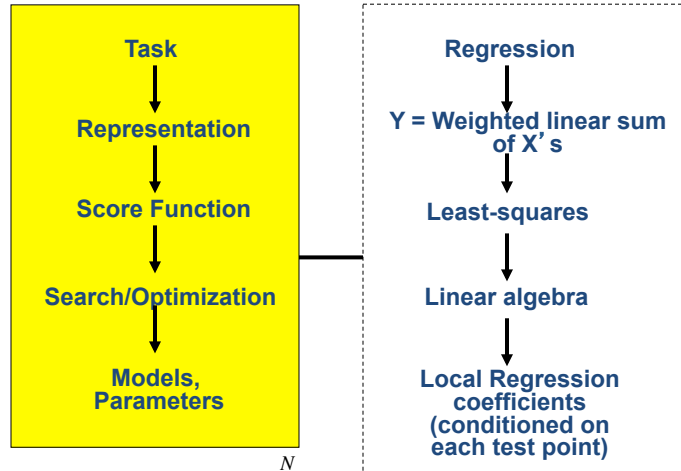
$$\hat{y} = \varphi(x)\theta$$

$$= \theta_0 + \theta_1 \exp(-(x-1)^2) + \theta_2 \exp(-(x-2)^2) + \theta_3 \exp(-(x-4)^2)$$

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

30

(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, ...

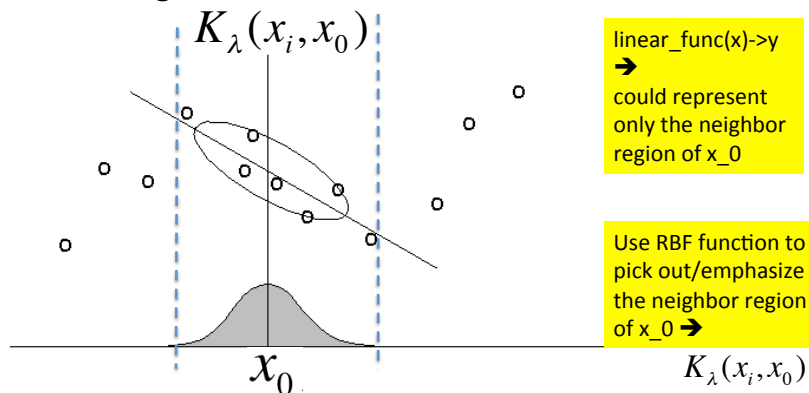
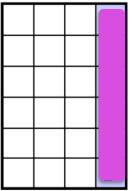


Figure 2: In locally weighted regression, points are weighted by proximity to the current x in question using a kernel. A regression is then computed using the weighted points.

Yanjun Qi / UVA CS 4501-01-6501-07


LEARNING of Locally weighted linear regression

target




training dataset

learn



model

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


x_0

→ Separate weighted least squares
at each target point x_0

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

10/30/1433

Yanjun Qi / UVA CS 4501-01-6501-07

(4) Regularized multivariate linear regression

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

Models,
Parameters

Regression

↓

Y = Weighted linear sum
of X' s

↓

Least-squares

↓

Linear algebra +
Regularization

↓

Regression
coefficients

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

10/30/1434

(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:
$$\min J(\beta) = \sum (Y - \hat{Y})^2$$
 - LASSO estimation:
$$\min J(\beta) = \sum_{i=1}^n (Y - \hat{Y})^2 + \lambda \sum_{j=1}^p |\beta_j|$$
 - Ridge regression estimation:
$$\min J(\beta) = \sum_{i=1}^n (Y - \hat{Y})^2 + \lambda \sum_{j=1}^p \beta_j^2$$
- Error on data + Regularization

10/30/14

35/54

(4) 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

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

to minimize $J(\beta)$, take derivative and set to zero

10/30/14

36

Today

- ☐ Review of basic pipeline
- ☐ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ☐ Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
- ☐ Model Selection / Bias Variance Tradeoff

Where are we ? →

Three major sections for classification

- We can divide the large variety of classification approaches into roughly three major types
- ➔ 1. Discriminative
 - directly estimate a decision rule/boundary
 - e.g., logistic regression, support vector machine, decisionTree
 - ➔ 2. Generative:
 - build a generative statistical model
 - e.g., naïve bayes classifier, Bayesian networks
 - ➔ 3. Instance based classifiers
 - Use observation directly (no models)
 - e.g. K nearest neighbors

YanJun Qi / UVA CS 4501-01-6501-07

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

Output as Discrete Class Label
 C_1, C_2, \dots, C_L

- **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]

10/30/14
39

YanJun Qi / UVA CS 4501-01-6501-07

(1) Support Vector Machine

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

Models,
Parameters

classification

↓

$K(\mathbf{x}, \mathbf{z}) := \Phi(\mathbf{x})^T \Phi(\mathbf{z})$

Kernel Func $K(x_i, x_j)$

↓

Margin + Hinge Loss (optional)

↓

QP with Dual form

↓

Dual Weights
 $w = \sum_i \alpha_i x_i y_i$

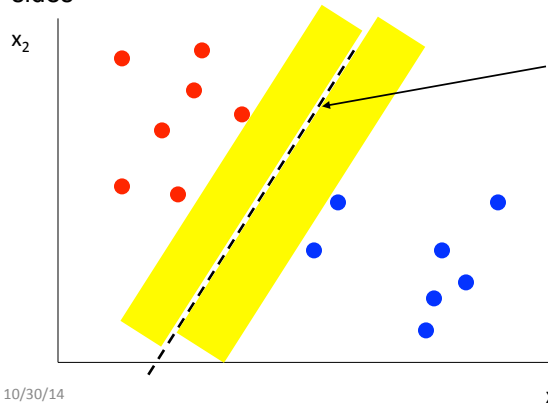
$$\operatorname{argmin}_{\mathbf{w}, b} \sum_{i=1}^p w_i^2 + C \sum_{i=1}^n \epsilon_i$$

subject to $\forall \mathbf{x}_i \in D_{train} : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1 - \epsilon_i$

10/30/14

(1) SVM as Max margin classifiers

- Instead of fitting all points, focus on boundary points
- Learn a boundary that leads to the largest margin from points on both sides



Why?

- Intuitive, 'makes sense'
- Some theoretical support
- Works well in practice

A Dataset for binary classification

$$f : X \rightarrow Y$$

Output as Binary Class Label:
1 or -1

X ₁	X ₂	X ₃	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]

Yanjun Qi / UVA CS 4501-01-6501-07

When linearly Separable Case

- The decision boundary should be as far away from the data of both classes as possible

- Correctly classifies all points
- Maximizes the margin (or equivalently minimizes $w^T w$)

W is a p-dim vector; b is a scalar

$w^T x + b = 1$

$w^T x + b = 0$

$w^T x + b = -1$

Class -1

Class 1

10/30/14 43

Yanjun Qi / UVA CS 4501-01-6501-07

Optimization Step

i.e. learning optimal parameter for SVM

- Correctly classifies all points
- Maximizes the margin (or equivalently minimizes $w^T w$)

Min $(w^T w)/2$
 subject to the following constraints:

For all x in class + 1
 $w^T x + b \geq 1$

For all x in class - 1
 $w^T x + b \leq -1$

A total of n constraints if we have n input samples

$$\operatorname{argmin}_{w,b} \sum_{i=1}^p w_i^2$$

subject to $\forall x_i \in D_{train} : y_i (x_i \cdot w + b) \geq 1$

SVM as a QP problem

10/30/14 44

Yanjun Qi / UVA CS 4501-01-6501-07

Dual formulation

Two optimization problems: For the separable and non separable cases

<p>Min $(w^T w)/2$</p> <p>For all x in class + 1</p> <p>$w^T x + b \geq 1$</p> <p>For all x in class - 1</p> <p>$w^T x + b \leq -1$</p>	<p>$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \epsilon_i$ → Hinge loss (optional)</p> <p>For all x_i in class + 1</p> <p>$w^T x + b \geq 1 - \epsilon_i$</p> <p>For all x_i in class - 1</p> <p>$w^T x + b \leq -1 + \epsilon_i$</p> <p>For all i</p> <p>$\epsilon_i \geq 0$</p>
----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

- Instead of solving these QPs directly we will solve a dual formulation of the SVM optimization problem
- The main reason for switching to this type of representation is that it would allow us to use a neat trick that will make our lives easier (and the run time faster)

10/30/1445

Yanjun Qi / UVA CS 4501-01-6501-07

A Geometrical Interpretation

$$w = \sum_i \alpha_i x_i y_i$$

For those α_i that are 0, no influence

10/30/1446

The kernel trick


How many operations do we need for the dot product?

$$\Phi(x)^T \Phi(z) = \sum_i 2x_i z_i + \sum_i x_i^2 z_i^2 + \sum_i \sum_{j=i+1} 2x_i x_j z_i z_j + 1$$

m
m
m(m-1)/2
 $\approx m^2$

$K(\mathbf{x}, \mathbf{z}) := \Phi(\mathbf{x})^T \Phi(\mathbf{z})$

However, we can obtain dramatic savings by noting that

$$\begin{aligned} \Phi(x)^T \Phi(z) &= (x^T z + 1)^2 = (x \cdot z + 1)^2 = (x \cdot z)^2 + 2(x \cdot z) + 1 \\ &= \left(\sum_i x_i z_i \right)^2 + \sum_i 2x_i z_i + 1 \\ &= \sum_i 2x_i z_i + \sum_i x_i^2 z_i^2 + \sum_i \sum_{j=i+1} 2x_i x_j z_i z_j + 1 \end{aligned}$$


We only need m operations!

So, if we define the **kernel function** as follows, there is no need to carry out $\phi(\cdot)$ explicitly

$$K(\mathbf{x}, \mathbf{z}) = (x^T z + 1)^2 \quad 47$$

10/30/14

Today

- Review of basic pipeline
- Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
- Model Selection

10/30/14

48

Yanjun Qi / UVA CS 4501-01-6501-07

$\operatorname{argmax}_k P(C = k | X) = \operatorname{argmax}_k P(X, C) = \operatorname{argmax}_k P(X|C)P(C)$
(2) Bayes Classifier

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

Models,
Parameters

classification

↓

Prob. models $p(X|C)$

↓

$P(X_1, \dots, X_p | C)$

EPE, with Log
likelihood(optional)

↓

Many options

↓

Prob. Models'
Parameter

Bernoulli Naive

$p(W_i = \text{true} | c_k) = p_{i,k}$

10/30/14

Gaussian Naive

$$\hat{P}(X_j | C = c_k) = \frac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left(-\frac{(X_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right)$$

Multinomial

$$P(W_1 = n_1, \dots, W_v = n_v | c_k) = \frac{N!}{n_{1k}! n_{2k}! \dots n_{vk}!} \theta_{1k}^{n_{1k}} \theta_{2k}^{n_{2k}} \dots \theta_{vk}^{n_{vk}}$$

Yanjun Qi / UVA CS 4501-01-6501-07

	X_1	X_2	X_3	C

A Dataset for classification

$f : X \rightarrow C$

Output as Discrete Class Label
 C_1, C_2, \dots, C_L

Generative

$$\operatorname{argmax}_C P(C | X) = \operatorname{argmax}_C P(X, C) = \operatorname{argmax}_C P(X|C)P(C)$$

- 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]

10/30/14
50

(2) Bayes classifier

- Treat each attribute and class label as random variables.
- Given a sample \mathbf{x} with attributes (x_1, x_2, \dots, x_p) :
 - Goal is to predict class C .
 - Specifically, we want to find the value of C_i that maximizes $p(C_i | x_1, x_2, \dots, x_p)$.

- Bayes classification

$$P(C | \mathbf{X}) \propto P(\mathbf{X} | C)P(C) = P(X_1, \dots, X_p | C)P(C)$$

Difficulty: learning the joint probability $P(X_1, \dots, X_p | C)$

10/30/14

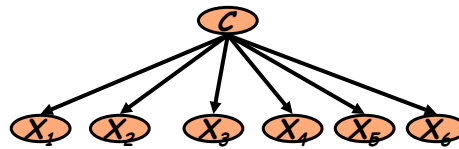
51

(2.1) Naïve Bayes Classifier

Difficulty: learning the joint probability $P(X_1, \dots, X_p | C)$

- Naïve Bayes classification
 - Assumption that **all input attributes are conditionally independent!**

$$\begin{aligned} P(X_1, X_2, \dots, X_p | C) &= P(X_1 | X_2, \dots, X_p, C)P(X_2, \dots, X_p | C) \\ &= P(X_1 | C)P(X_2, \dots, X_p | C) \\ &= P(X_1 | C)P(X_2 | C) \cdots P(X_p | C) \end{aligned}$$



10/30/14

52

Adapt from Prof. Ke Chen NB slides

(2.2) Multinomial Naïve Bayes as Stochastic Language Models

Model C1	Model C2	the	boy	likes	black	dog
0.2 the	0.2 the	0.2	0.01	0.0001	0.0001	0.0005
0.01 boy	0.0001 boy	0.2	0.0001	0.02	0.1	0.01
0.0001 said	0.03 said					
0.0001 likes	0.02 likes					
0.0001 black	0.1 black					
0.0005 dog	0.01 dog					
0.01 garden	0.0001 garden					

$P(s|C2) P(C2) > P(s|C1) P(C1)$

10/30/14

53

(2.3) Gaussian Naïve Bayes Classifier

- Continuous-valued Input Attributes
 - Conditional probability modeled with the normal distribution

$$\hat{P}(X_j | C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

μ_{ji} : mean (average) of attribute values X_j of examples for which $C = c_i$

σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

- **Learning Phase:** for $\mathbf{X} = (X_1, \dots, X_p)$, $C = c_1, \dots, c_L$
Output: $p \times L$ normal distributions and $P(C = c_i) \quad i = 1, \dots, L$
- **Test Phase:** for $\mathbf{X}' = (X'_1, \dots, X'_p)$
 - Calculate conditional probabilities with all the normal distributions
 - Apply the MAP rule to make a decision

10/30/14

54

Naïve Gaussian means ?

Not Naïve

$$P(X_1, X_2, \dots, X_p | C) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

Naïve

$$P(X_1, X_2, \dots, X_p | C = c_k) = P(X_1 | c_k) P(X_2 | c_k) \cdots P(X_p | c_k) = \prod_j \frac{1}{\sqrt{2\pi}\sigma_{j,k}} \exp \left(-\frac{(X_j - \mu_{j,k})^2}{2\sigma_{j,k}^2} \right)$$

Diagonal Matrix

$$\boldsymbol{\Sigma}_{c_k} = \boldsymbol{\Lambda}_{c_k}$$

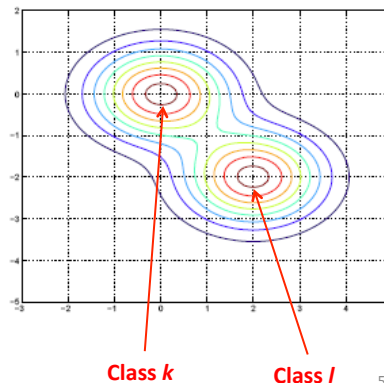
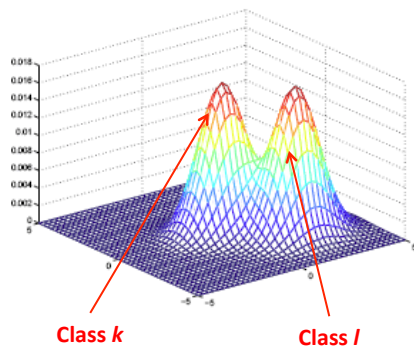
Each class' covariance matrix is diagonal

(2.4) LDA (Linear Discriminant Analysis)

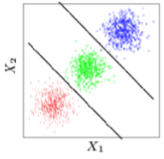
Linear Discriminant Analysis : $\boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}, \forall k$

Each class' covariance matrix is the same

The Gaussian Distribution are shifted versions of each other



Yanjun Qi / UVA CS 4501-01-6501-07



Optimal Classification

$$\operatorname{argmax}_k P(C_k | X) = \operatorname{argmax}_k P(X, C) = \operatorname{argmax}_k P(X | C) P(C)$$

$$= \operatorname{argmax}_k \left[-\log \left(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \right) - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

$$= \operatorname{argmax}_k \left[-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

- Note

Linear Discriminant Function for LDA

$$-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k - \frac{1}{2} x^T \Sigma^{-1} x$$

10/30/14 57

Yanjun Qi / UVA CS 4501-01-6501-07

Define Linear Discriminant Function

$$\delta(x) = -\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log \pi_k$$

→ The **Decision Boundary Between class k and l**, $\{x : \delta_k(x) = \delta_l(x)\}$, is a linear line/plane

$$\log \frac{P(C_k | X)}{P(C_l | X)} = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{P(C_k)}{P(C_l)}$$

$$= \log \frac{\pi_k}{\pi_l} - \frac{1}{2} (\mu_k + \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) + x^T \Sigma^{-1} (\mu_k - \mu_l) \quad (4.9)$$

Equals to zero

Boundary points X : when $P(c_k | X) = P(c_l | X)$, the left linear equation = 0, a linear line / plane

10/30/14 58

(2.5) QDA (Quadratic Discriminant Analysis)

- ▶ Estimate the covariance matrix Σ_k separately for each class k , $k = 1, 2, \dots, K$.

- ▶ Quadratic discriminant function:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k .$$

- ▶ Classification rule:

$$\hat{G}(x) = \arg \max_k \delta_k(x) .$$

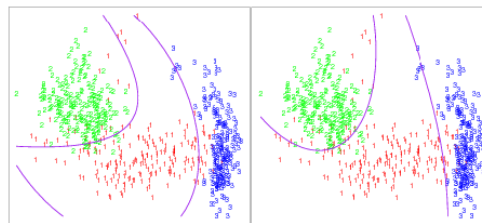
- ▶ Decision boundaries are quadratic equations in x .
- ▶ QDA fits the data better than LDA, but has more parameters to estimate.

10/30/14

59

(2.6) LDA on Expanded Basis

- ▶ Expand input space to include X_1X_2 , X_1^2 , and X_2^2 .
- ▶ Input is five dimensional: $X = (X_1, X_2, X_1X_2, X_1^2, X_2^2)$.



LDA with
quadratic basis
Versus
QDA

Figure 4.6: Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space $x_1, x_2, x_{12}, x_1^2, x_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

10/30/14

60

(2.7) Regularized Discriminant Analysis

- ▶ A compromise between LDA and QDA.
- ▶ Shrink the separate covariances of QDA toward a common covariance as in LDA.
- ▶ Regularized covariance matrices:

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma}.$$

- ▶ The quadratic discriminant function $\delta_k(x)$ is defined using the shrunken covariance matrices $\hat{\Sigma}_k(\alpha)$.
- ▶ The parameter α controls the complexity of the model.

Today

- Review of basic pipeline
- Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
- Model Selection

Yanjun Qi / UVA CS 4501-01-6501-07

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

X

 \rightarrow

C

Output as Discrete
Class Label
 C_1, C_2, \dots, C_L

Discriminative $\rightarrow P(C | X) \quad C = c_1, \dots, c_L$

- **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]

10/30/1463

Yanjun Qi / UVA CS 4501-01-6501-07

(3) Logistic Regression

Task

↓

Representation

↓

Score Function

↓

Search/Optimization

↓

**Models,
Parameters**

classification

↓

**Log-odds(Y) = linear
function of X' s**

↓

**EPE, with conditional
Log-likelihood**

↓

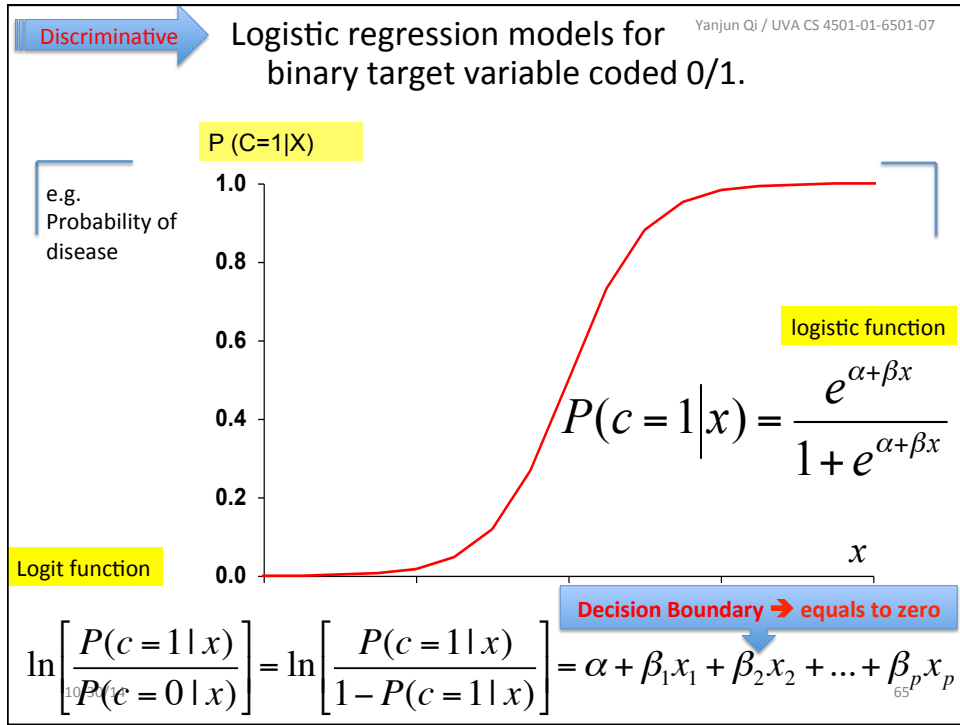
Iterative (Newton) method

↓

**Logistic
weights**

$$P(c = 1 | x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

10/30/1464



Yanjun Qi / UVA CS 4501-01-6501-07

MLE for Logistic Regression Training

Let's fit the logistic regression model for $K=2$, i.e., number of classes is 2

Training set: $(x_i, y_i), i=1, \dots, N$

For Bernoulli distribution
 $p(y | x)^y (1 - p)^{1-y}$

(conditional) Log-likelihood:

$$\begin{aligned}
 l(\beta) &= \sum_{i=1}^N \{\log \Pr(Y = y_i | X = x_i)\} \\
 &= \sum_{i=1}^N y_i \log(\Pr(Y = 1 | X = x_i)) + (1 - y_i) \log(\Pr(Y = 0 | X = x_i)) \\
 &= \sum_{i=1}^N (y_i \log \frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)} + (1 - y_i) \log \frac{1}{1 + \exp(\beta^T x_i)}) \\
 &= \sum_{i=1}^N (y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i)))
 \end{aligned}$$

x_i are $(p+1)$ -dimensional input vector with leading entry 1
 β is a $(p+1)$ -dimensional vector
 $y_i = 1$ if $C_i = 1$; $y_i = 0$ if $C_i = 0$

10/30/14 We want to maximize the log-likelihood in order to estimate β 66

Today

- ☐ Review of basic pipeline
- ☐ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ☐ Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - **K-nearest Neighbor**
- ☐ Model Selection

10/30/14

67

X_1	X_2	X_3	C

A Dataset for
classification

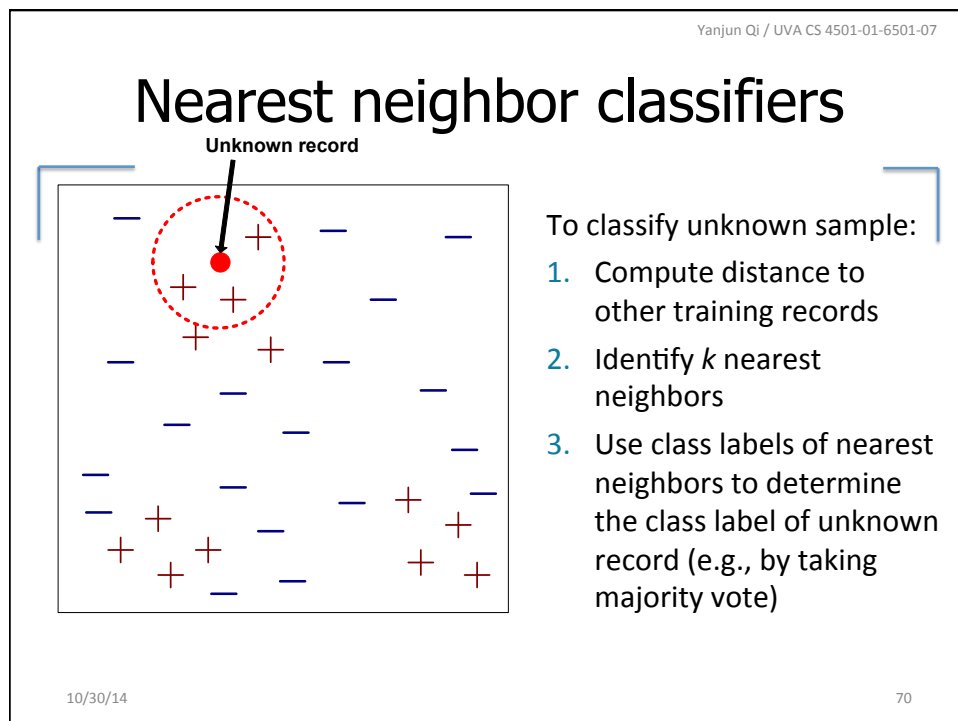
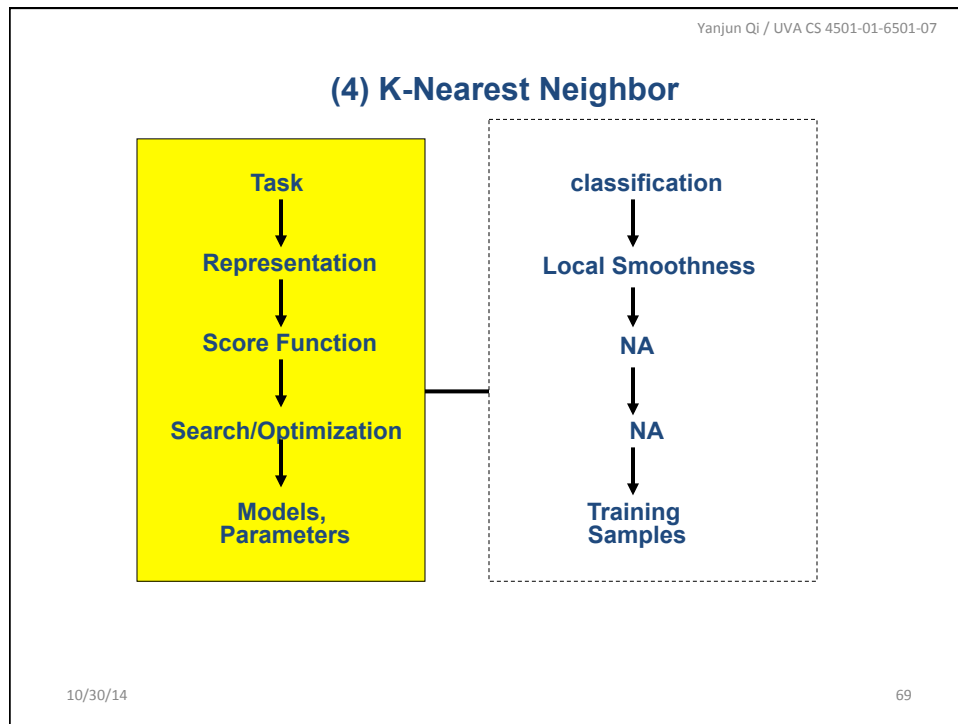
$$f : X \rightarrow C$$

Output as Discrete
Class Label
 C_1, C_2, \dots, C_L

- **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]

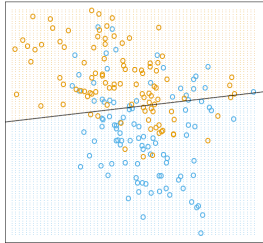
10/30/14

68



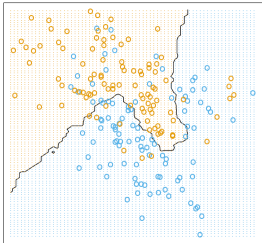
Yanjun Qi / UVA CS 4501-01-6501-07

Decision boundaries in global vs. local models

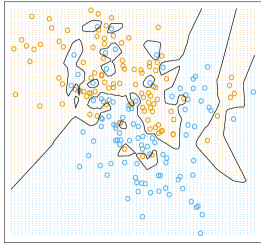


linear regression

- global
- stable
- can be inaccurate



15-nearest neighbor



1-nearest neighbor

- local
- accurate
- unstable

What ultimately matters: **GENERALIZATION**

10/30/14 71

Yanjun Qi / UVA CS 4501-01-6501-07

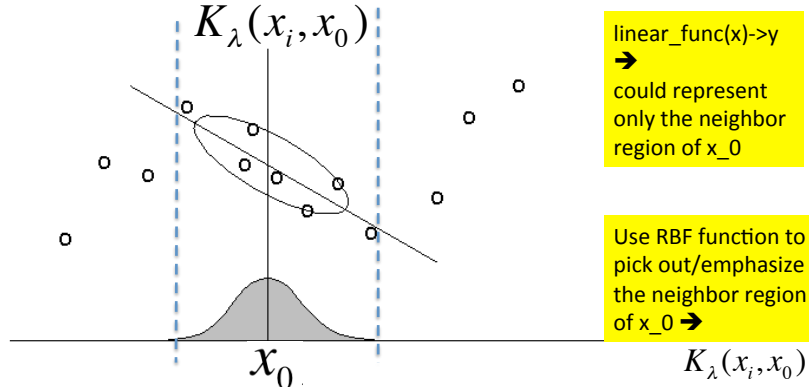
Nearest neighbor classification

- k -Nearest neighbor classifier is a **lazy** learner
 - Does not build model explicitly.
 - Unlike **eager** learners such as decision tree induction and rule-based systems.
 - Classifying unknown samples is relatively expensive.
- k -Nearest neighbor classifier is a **local** model, vs. **global** model of linear classifiers.

10/30/14 72

Vs. Locally weighted regression

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



10/30/14 **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.

Vs. Locally weighted regression

- Separate weighted least squares **at each target point x_0 :**

x_0 ?

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

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

$$K_{\tau}(\mathbf{x}_i, \mathbf{x}_0) = \exp\left(-\frac{(\mathbf{x}_i - \mathbf{x}_0)^2}{2\tau^2}\right)$$

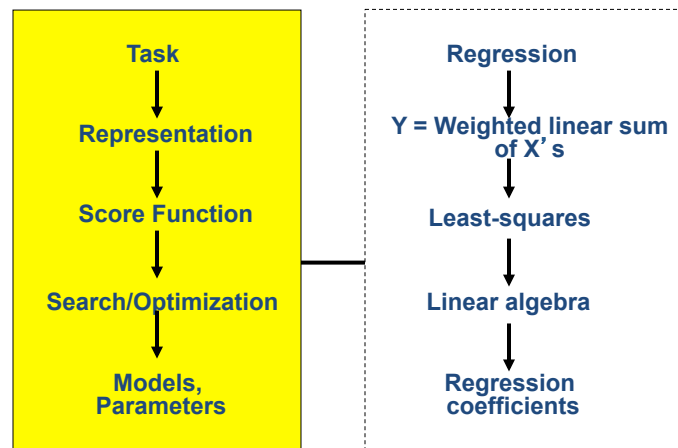
Today

- ❑ Review of basic pipeline
- ❑ Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- ❑ Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
- ❑ Model Selection / Bias Variance Tradeoff

10/30/14

75

(1) Multivariate Linear Regression

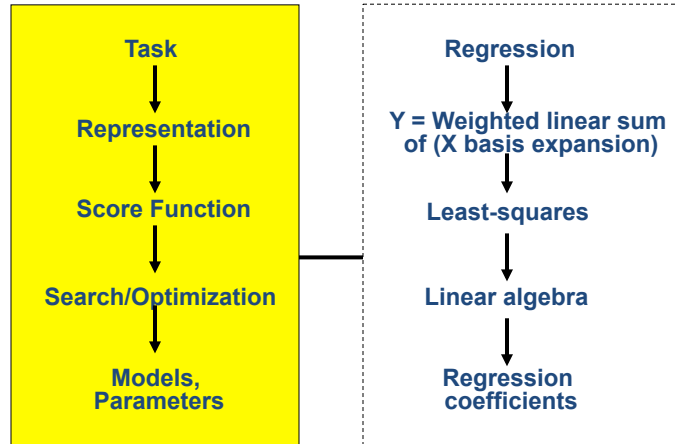


$$\hat{y} = f(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2$$

10/30/14

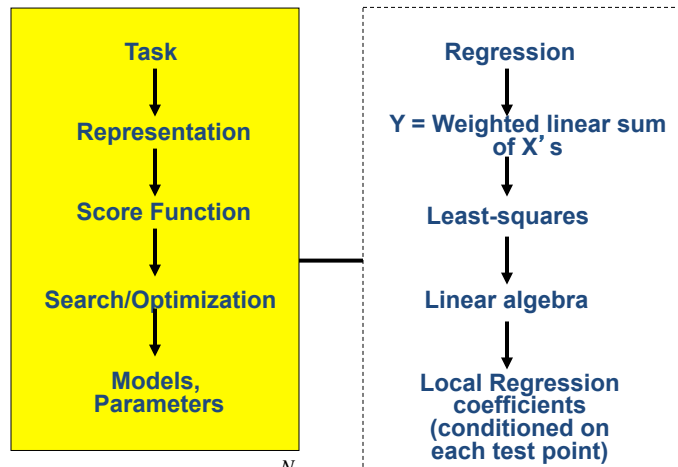
76

(2) Multivariate Linear Regression with basis Expansion



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

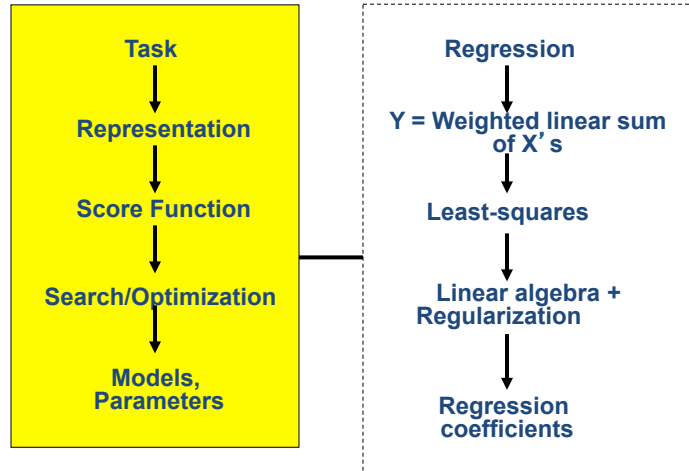
(3) Locally Weighted / Kernel Regression



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

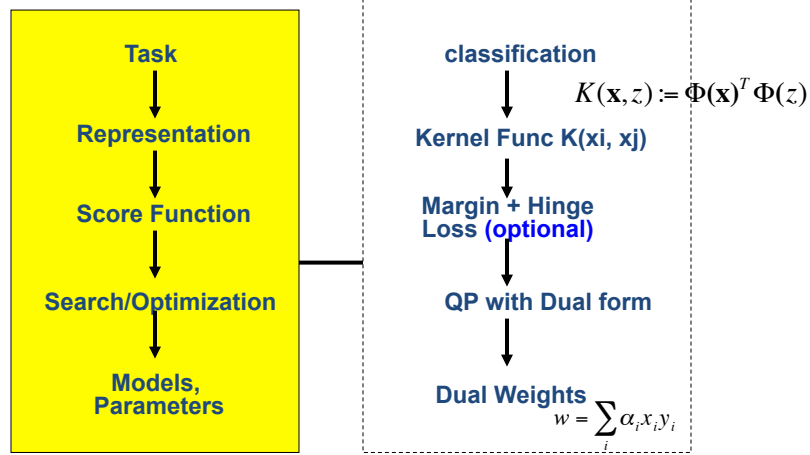
$$\hat{f}(x_0) = \hat{\alpha}(x_0) + \hat{\beta}(x_0)x_0$$

(4) Regularized multivariate linear regression



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

(1) Support Vector Machine



$$\operatorname{argmin}_{w,b} \sum_{i=1}^p w_i^2 + C \sum_{i=1}^n \epsilon_i$$

subject to $\forall x_i \in D_{train} : y_i (\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1 - \epsilon_i$

YanJun Qi / UVA CS 4501-01-6501-07

$$\operatorname{argmax}_k P(C = k | X) = \operatorname{argmax}_k P(X, C) = \operatorname{argmax}_k P(X|C)P(C)$$

(2) Bayes Classifier

Task
↓
Representation
↓
Score Function
↓
Search/Optimization
↓
Models,
Parameters

classification
↓
Prob. models $p(X|C)$
↓
EPE, with Log
likelihood(optional)
↓
Many options
↓
Prob. Models'
Parameter

Bernoulli Naive $p(W_i = \text{true} | c_k) = p_{i,k}$ Gaussian Naive

10/30/14 Multinomial

$$\hat{P}(X_j | C = c_k) = \frac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left(-\frac{(X_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right)$$

$$P(W_1 = n_1, \dots, W_v = n_v | c_k) = \frac{N!}{n_{1k}! n_{2k}! \dots n_{vk}!} \theta_{1k}^{n_{1k}} \theta_{2k}^{n_{2k}} \dots \theta_{vk}^{n_{vk}}$$

YanJun Qi / UVA CS 4501-01-6501-07

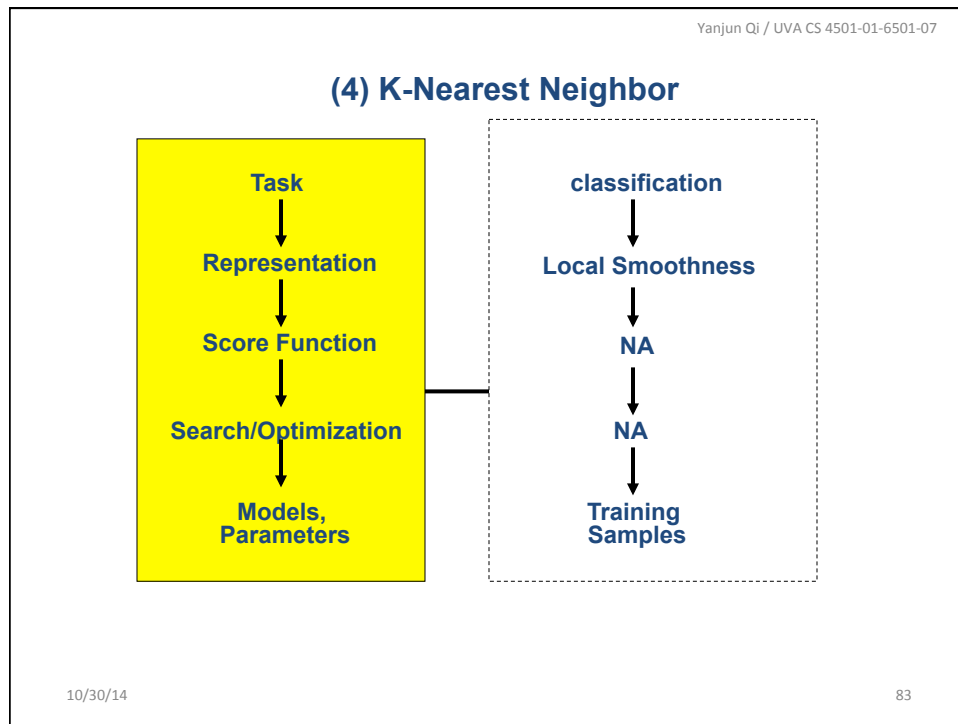
(3) Logistic Regression

Task
↓
Representation
↓
Score Function
↓
Search/Optimization
↓
Models,
Parameters

classification
↓
Log-odds(Y) = linear
function of X' s
↓
EPE, with conditional
Log-likelihood
↓
Iterative (Newton) method
↓
Logistic
weights

$$P(c = 1 | x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

10/30/14
82



Yanjun Qi / UVA CS 4501-01-6501-07

Today

- Review of basic pipeline
- Review of regression models
 - Linear regression (LR)
 - LR with non-linear basis functions
 - Locally weighted LR
 - LR with Regularizations
- Review of classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
- Model Selection / Bias Variance Tradeoff**

10/30/14 84

Yanjun Qi / UVA CS 4501-01-6501-07

e.g. Training Error from KNN, Lesson Learned

- When $k = 1$,
- No misclassifications (on training): **Overtraining**
- Minimizing training error is not always good (e.g., 1-NN)

1-nearest neighbor averaging

10/30/14

Yanjun Qi / UVA CS 4501-01-6501-07

Statistical Decision Theory

- Random input vector: X
- Random output variable: Y
- Joint distribution: $\Pr(X, Y)$
- Loss function $L(Y, f(X))$
- Expected prediction error (EPE):
- $$\text{EPE}(f) = \mathbb{E}(L(Y, f(X))) = \int L(y, f(x)) \Pr(dx, dy)$$

e.g. = $\int (y - f(x))^2 \Pr(dx, dy)$

Consider population distribution

10/30/14

e.g. Squared error loss (also called L2 loss)

86

Expected prediction error (EPE)

Consider sample population distribution

$$EPE(f) = E(L(Y, f(X))) = \int L(y, f(x)) \Pr(dx, dy)$$

- For L2 loss: e.g. $\int (y - f(x))^2 \Pr(dx, dy)$

under L2 loss, best estimator for EPE (Theoretically) is :

e.g. KNN → **Conditional mean** $f(x) = E(Y | X = x)$
 NN methods are the direct implementation (approximation)

- For 0-1 loss: $L(k, \ell) = 1 - \delta_{kl}$
 $\hat{G}(X) = C_k$ if $\Pr(C_k | X = x) = \max_{g \in C} \Pr(g | X = x)$

→ **Bayes Classifier**

Decomposition of EPE

$$Y = f(X) + \epsilon, \epsilon \sim (0, \sigma^2)$$

- When additive error model:

- Notations

- Output random variable: Y
- Prediction function: f
- Prediction estimator: \hat{f}

$$\begin{aligned} EPE(x_0) &= E[(Y - \hat{f})^2 | X = x_0] \\ &= E[((Y - f) + (f - \hat{f}))^2 | X = x_0] \\ &= E[\underbrace{(Y - f)^2}_{\epsilon} | X = x_0] + \underbrace{E[(f - \hat{f})^2 | X = x_0]}_{MSE} \\ &= \sigma^2 + Var(\hat{f}) + Bias^2(\hat{f}) \end{aligned}$$

MSE component of \hat{f} in estimating f

σ^2 ↑ Irreducible / Bayes error

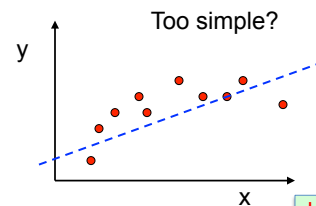
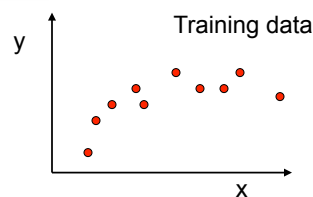
BIAS AND VARIANCE TRADE-OFF

- θ : true value (normally unknown)
- $\hat{\theta}$: estimator
- $\bar{\theta} = E[\hat{\theta}]$ (mean, i.e. expectation of the estimator)

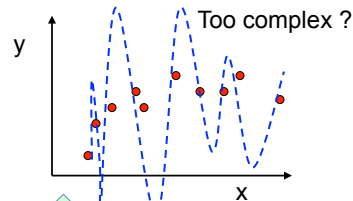
- Bias $E[(\bar{\theta} - \theta)^2]$
 - measures **accuracy** or **quality** of the estimator
 - low bias implies on average we will accurately estimate true parameter or func from training data
- Variance $E[(\hat{\theta} - \bar{\theta})^2]$
 - Measures **precision** or **specificity** of the estimator
 - Low variance implies the estimator does not **change** much as **the training set varies**

10/30/14

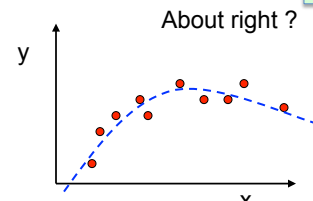
Regression: Complexity versus Goodness of Fit



Low Variance /
High Bias

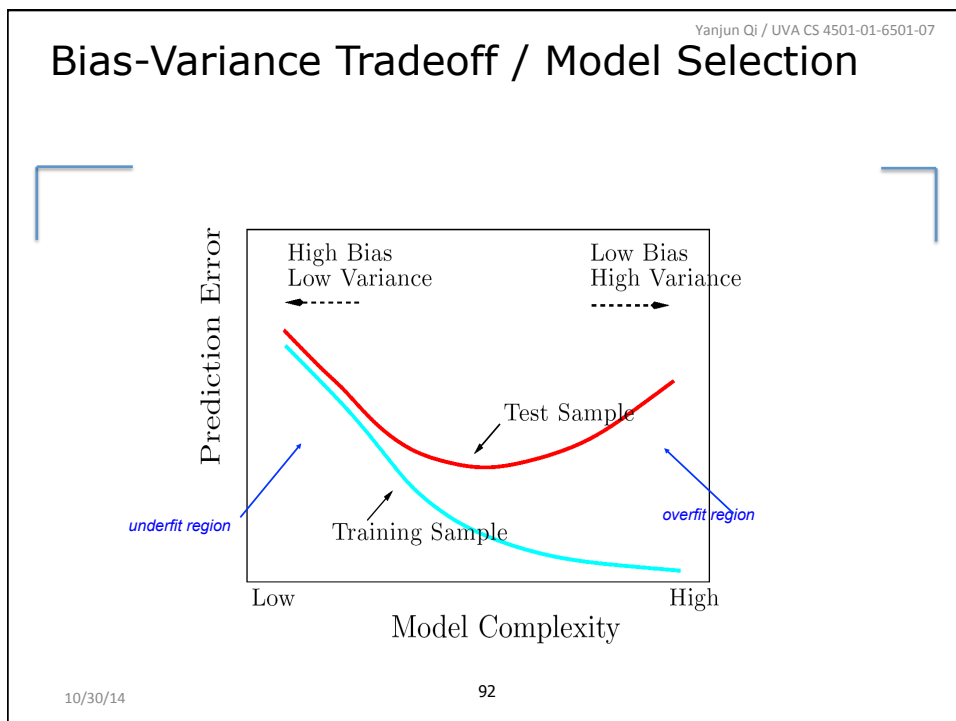
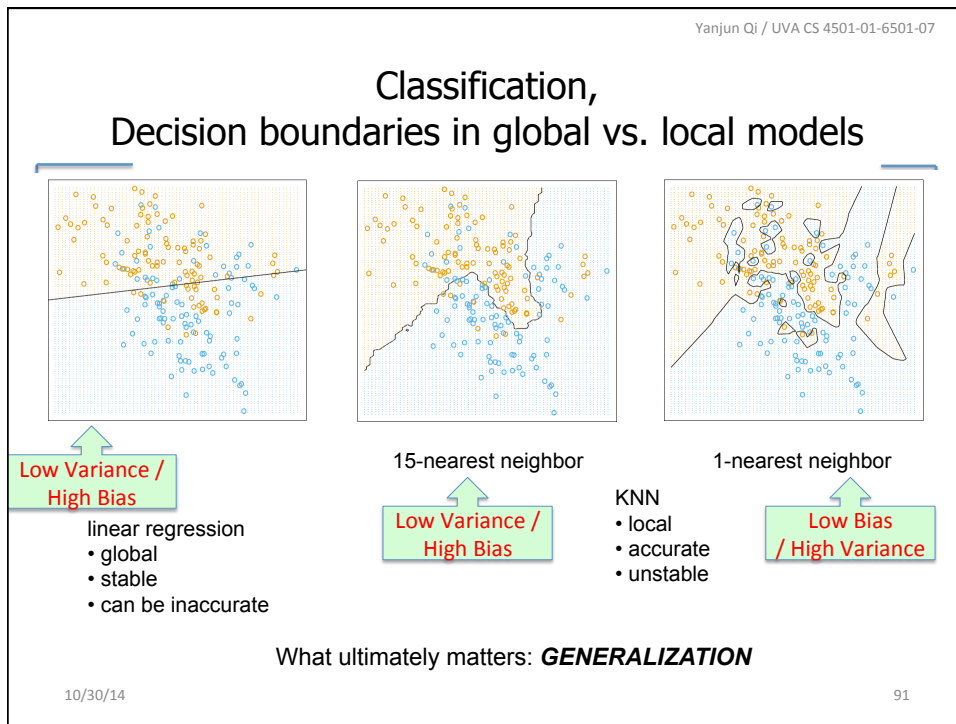


Low Bias
/ High Variance



What ultimately matters: **GENERALIZATION**

90

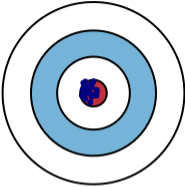
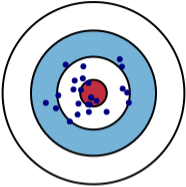


Yanjun Qi / UVA CS 4501-01-6501-07

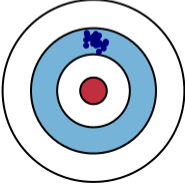
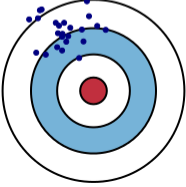
Model “bias” & Model “variance”

- Middle RED:
 - TRUE function
- Error due to bias:
 - How far off in general from the middle red
- Error due to variance:
 - How wildly the blue points spread

Low Bias

High Bias

10/30/14 93

Yanjun Qi / UVA CS 4501-01-6501-07

References

- Prof. Tan, Steinbach, Kumar’s “Introduction to Data Mining” slide
- Prof. Andrew Moore’s slides
- Prof. Eric Xing’s slides
- Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.

10/30/14 94

Midterm

- Open Note / Open Book
- No laptop / No Cell phone / No internet access
- Easier than sample questions in HW4