

UVA CS 6316

– Fall 2015 Graduate: Machine Learning

Lecture 22: Review

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

11/16/15

1

Announcements: Rough Plan

- HW5
 - Out on Nov. 18th
 - Due on Dec. 7th
- Project Presentation
 - Due on Dec. 2nd midnight
 - Presentations @ Dec 3rd and Dec 4th
- Project Final Report
 - Due on Dec. 11th midnight

11/16/15

2

Announcements: Exam

- Open Note / Open Lectures
- No laptop / No Cell phone / No internet access / No electronic devices

- Covering contents till today
 - Practice with sample questions in HW4
 - HW4 due on Nov. 20th
 - Please review course slides carefully

Today

- History of AI & Machine Learning
- Review of ML methods covered so far

What are the goals of AI research?

Artifacts that THINK
like HUMANS

Artifacts that THINK
RATIONALLY

Artifacts that ACT
like HUMANS

Artifacts that ACT
RATIONALLY

11/16/15

5

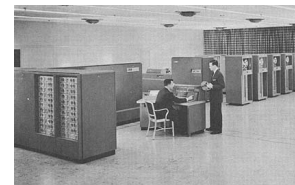
From: M.A. Papalaskar

A Bit of History

1940s

Advances in mathematical logic, information theory, concept of neural computation

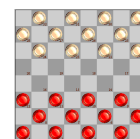
- ⌚ 1943: McCulloch & Pitts Neuron
- ⌚ 1948: Shannon: Information Theory
- ⌚ 1949: Hebbian Learning
 - ⌚ cells that fire together, wire together



1950s

Early computers. Dartmouth conference coins the phrase “artificial intelligence” and Lisp is proposed as the AI programming language

- ⌚ 1950: Turing Test
- ⌚ 1956: Dartmouth Conference
- ⌚ 1958: Friedberg: Learn Assembly Code
- ⌚ 1959: Samuel: Learning Checkers



11/16/15

6

From: M.A. Papalaskar

1960s

A.I. funding increased (mainly military). Famous quote: “Within a generation ... the problem of creating 'artificial intelligence' will substantially be solved.”

Early symbolic reasoning approaches.

- ⦿ Logic Theorist, GPS, Perceptrons
- ⦿ 1969: Minsky & Papert “Perceptrons”

1970s

A.I. “winter” – Funding dries up as people realize this is a hard problem!

Limited computing power and dead-end frameworks lead to failures.

- ⦿ eg: Machine Translation Failure

11/16/15

7
From: M.A. Papalaskar

1980s

Rule based “expert systems” used in medical / legal professions.

Bio-inspired algorithms (Neural networks, Genetic Algorithms).

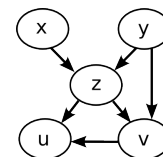
Again: A.I. promises the world – lots of commercial investment

Expert Systems (Mycin, Dendral, EMYCIN)

Knowledge Representation and reasoning:

Frames, Eurisko, Cyc, NMR, fuzzy logic

Speech Recognition (HEARSAY, HARPY, HWIM)

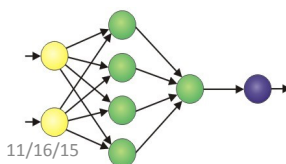


$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

Machine Learning:

- ⦿ 1982: Hopfield Nets, Decision Trees, GA & GP.

- ⦿ 1986: Backpropagation, Explanation-Based Learning



11/16/15

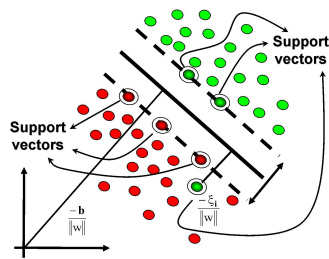
8
From: M.A. Papalaskar

1990s

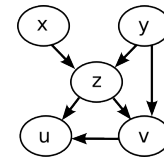
Some concrete successes begin to emerge. AI diverges into separate fields: Computer Vision, Automated Reasoning, Planning systems, Natural Language processing, **Machine Learning**...

...Machine Learning begins to overlap with statistics / probability theory.

- ☉ 1992: Koza & Genetic Programming
- ☉ 1995: Vapnik: Support Vector Machines



11/16/15



$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

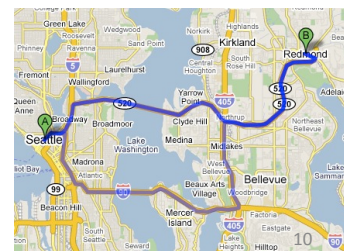
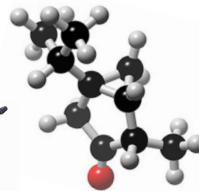
9

From: M.A. Papalaskar

2000s

First commercial-strength applications: Google, Amazon, computer games, route-finding, credit card fraud detection, spam filters, etc...

Tools adopted as standard by other fields e.g. biology



11/16/15

From: M.A. Papalaskar

2010s.... ??????

Deep Learning With massive amounts of computational power, machines can now recognize objects and translate speech in real time. Artificial intelligence is finally getting smart.	Temporary Social Media Messages that quickly self-destruct could enhance the privacy of online communications and make people freer to be spontaneous.	Prenatal DNA Sequencing Reading the DNA of fetuses will be the next frontier of the genomic revolution. But do you really want to know about the genetic problems or musical aptitude of your unborn child?	Additive Manufacturing Skeptical about 3-D printing? GE, the world's largest manufacturer, is on the verge of using the technology to make jet parts.	Baxter: The Blue-Collar Robot Rodney Brooks's newest creation is easy to interact with, but the complex innovations behind the robot show just how hard it is to get along with people.
Memory Implants A maverick neuroscientist believes he has deciphered the code by which the brain forms long-term memories. Next: testing a prosthetic implant for people suffering from long-term memory loss.	Smart Watches The designers of the Pebble watch realized that a mobile phone is more useful if you don't have to take it out of your pocket.	Ultra-Efficient Solar Power Doubling the efficiency of a solar cell would completely change the economics of renewable energy. Nanotechnology just might make it possible.	Big Data from Cheap Phones Collecting and analyzing information from simple cell phones can provide surprising insights into how people move about and behave – and even help us understand the spread of diseases.	Supergrids A new high-power circuit breaker could finally make highly efficient DC power grids practical.

How can we build more intelligent computer / machine ?

- Able to
 - **perceive the world**
 - **understand the world**
- This needs
 - Basic speech capabilities
 - Basic vision capabilities
 - Language understanding
 - User behavior / emotion understanding
 - **Able to think ??**

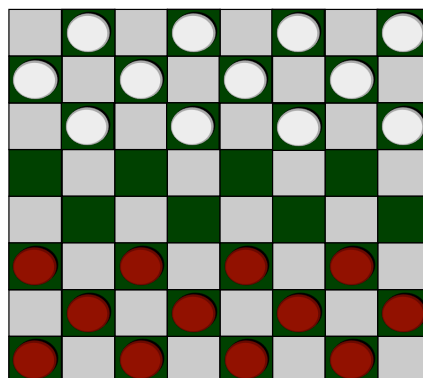
Plenty of Data

- **Text:** trillions of words of English + other languages
- **Visual:** billions of images and videos
- **Audio:** thousands of hours of speech per day
- **User activity:** queries, user page clicks, map requests, etc,
- **Knowledge graph:** billions of labeled relational triplets
-

Data-driven machine learning methods have made machines / computers much more intelligent

Samuel's definition of ML (1959)

- Arthur Samuel (1959). Machine Learning: Field of study that **gives computers the ability to learn without being explicitly programmed.**



Tom Mitchell (1998): Well-posed Learning Problem

*A computer program is said to learn from experience **E** with respect to some task **T** and some performance measure **P**, if its performance on **T**, as measured by **P**, improves with experience **E**.*

11/16/15

15

Defining the Learning Task

Improve on task, T, with respect to performance metric, P, based on experience, E.

T: Playing checkers
P: Percentage of games won against an arbitrary opponent
E: Playing practice games against itself

T: Recognizing hand-written words
P: Percentage of words correctly classified
E: Database of human-labeled images of handwritten words

T: Driving on four-lane highways using vision sensors
P: Average distance traveled before a human-judged error
E: A sequence of images and steering commands recorded while observing a human driver.

T: Determine which students like oranges or apples
P: Percentage of students' preferences guessed correctly
E: Student attribute data

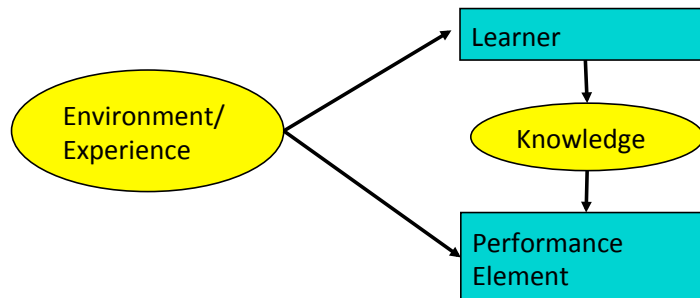
11/16/15

16

From: M.A. Papalaskar

Designing a Learning System

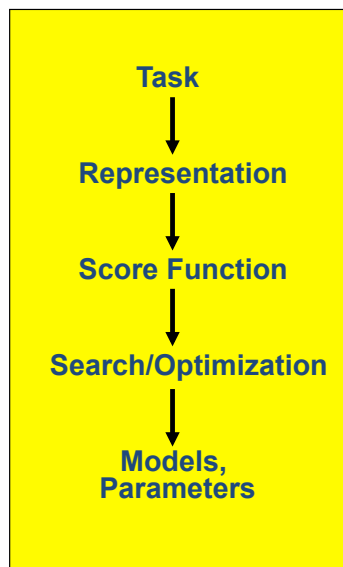
- Choose the **training experience**
- Choose exactly what is to be learned, i.e. the **target function**.
- Choose a **learning algorithm** to infer the target function from the experience.
- A learning algorithm will also determine a **performance measure**



17

11/16/15

Machine Learning in a Nutshell



ML grew out of
work in AI

*Optimize a
performance criterion
using example data or
past experience,*

*Aiming to generalize to
unseen data*

11/16/15

18

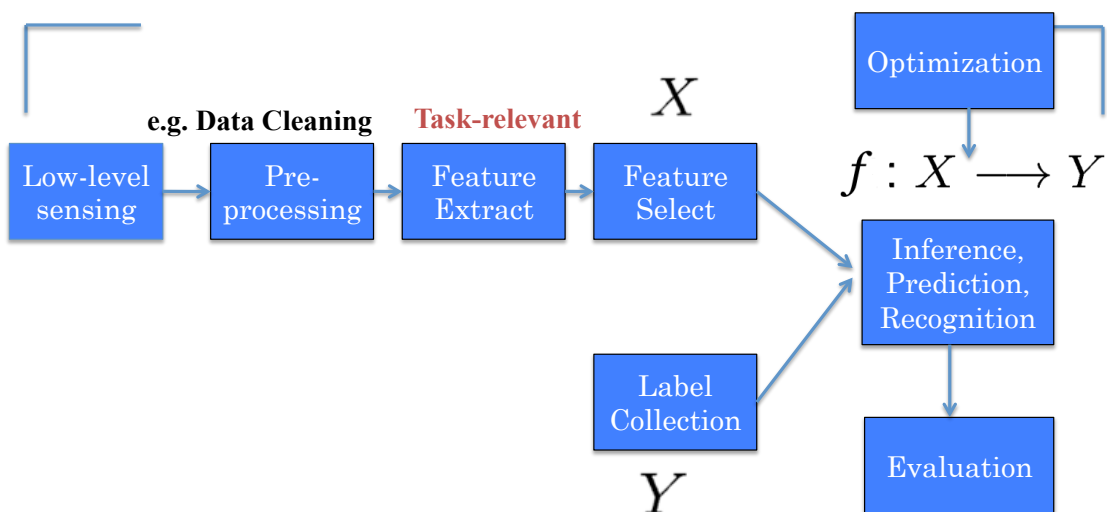
What we have covered for each

Task	
Representation	
Score Function	
Search/ Optimization	
Models, Parameters	

11/16/15

19

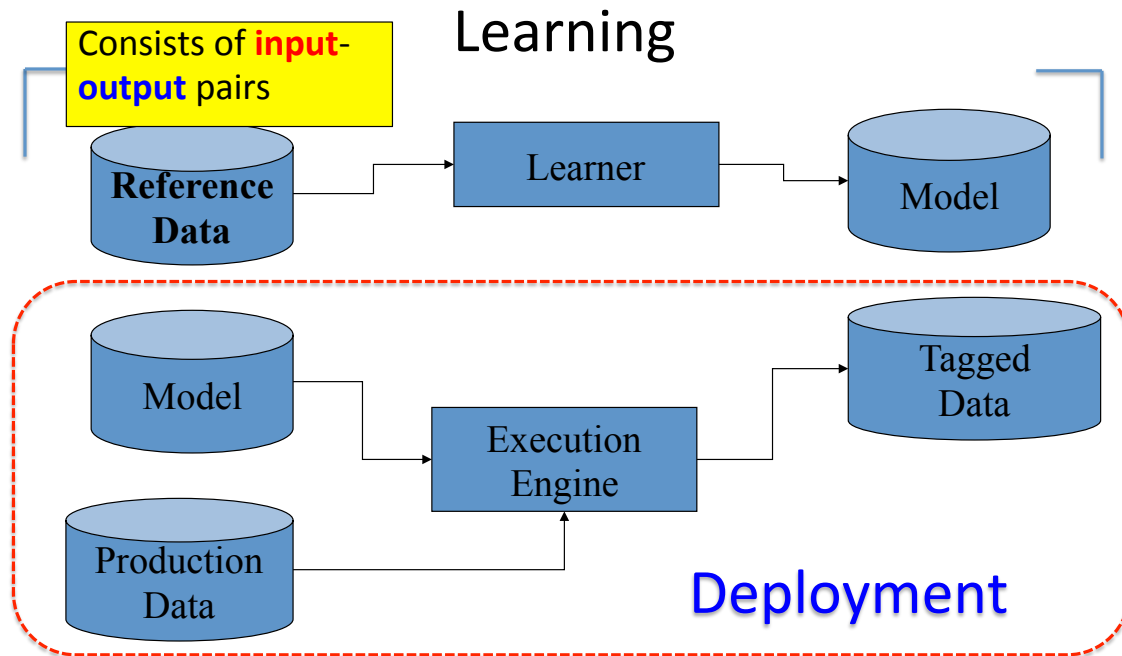
A Typical Machine Learning Pipeline



11/16/15

20

An **Operational** Model of Machine Learning



11/16/15

21

Today

- History of Machine Learning & AI
- Review of ML methods covered so far

11/16/15

22

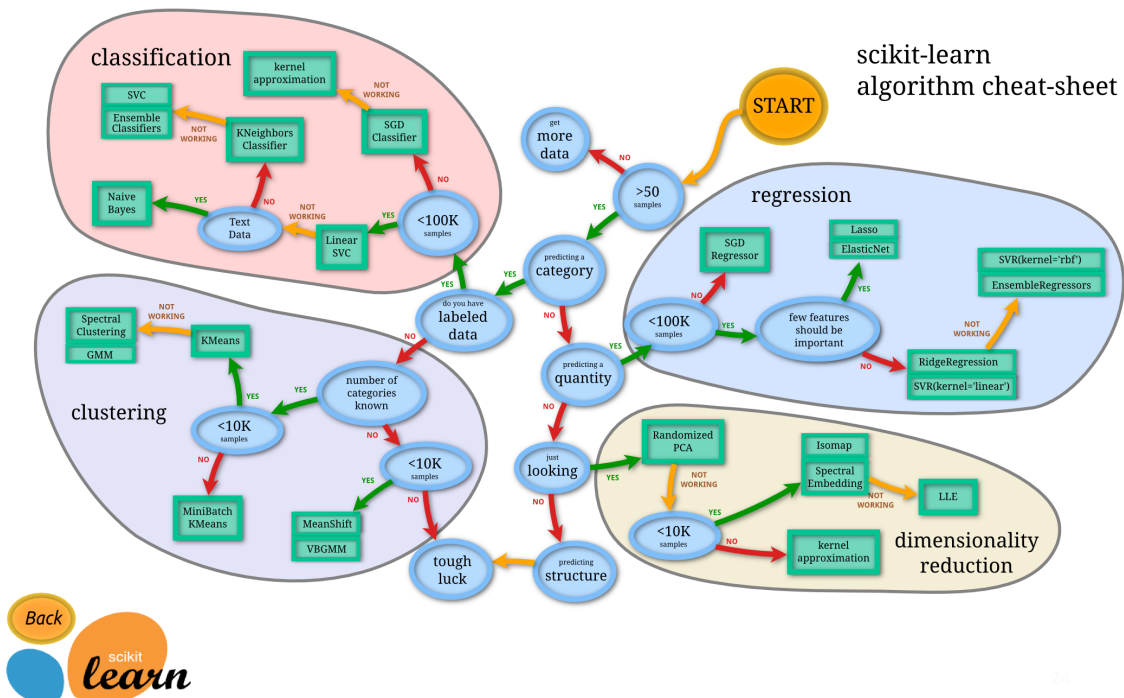
Where are we ? →

major sections of this course

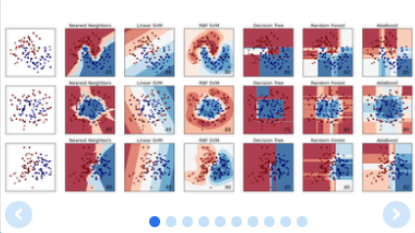
- ❑ Regression (supervised)
- ❑ Classification (supervised)
- ❑ Unsupervised models
- ❑ Learning theory
- ❑ Graphical models

http://scikit-learn.org/stable/tutorial/machine_learning_map/

Scikit-learn algorithm cheat-sheet



<http://scikit-learn.org/stable/>



scikit-learn

Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Classification

Identifying to which set of categories a new observation belong to.

Applications: Spam detection, Image recognition.

Algorithms: SVM, nearest neighbors, random forest, ...

Regression

Predicting a continuous value for a new example.

Applications: Drug response, Stock prices.

Algorithms: SVR, ridge regression, Lasso, ...

Clustering

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes

Algorithms: k-Means, spectral clustering, mean-shift, ...

Dimensionality reduction

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency

Algorithms: PCA, feature selection, non-negative matrix factorization.

Model selection

Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning

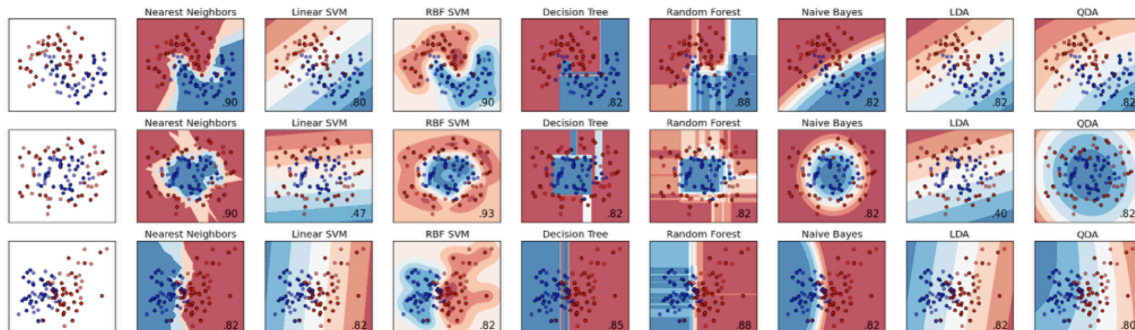
Modules: grid search, cross validation, metrics.

Preprocessing

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.

Modules: preprocessing, feature extraction.



- ✓ different assumptions on data
- ✓ different scalability profiles at training time
- ✓ different latencies at prediction time
- ✓ different model sizes (embedability in mobile devices)

What we have covered (I)

□ Supervised Regression models

- Linear regression (LR)
- LR with non-linear basis functions
- Locally weighted LR
- LR with Regularizations

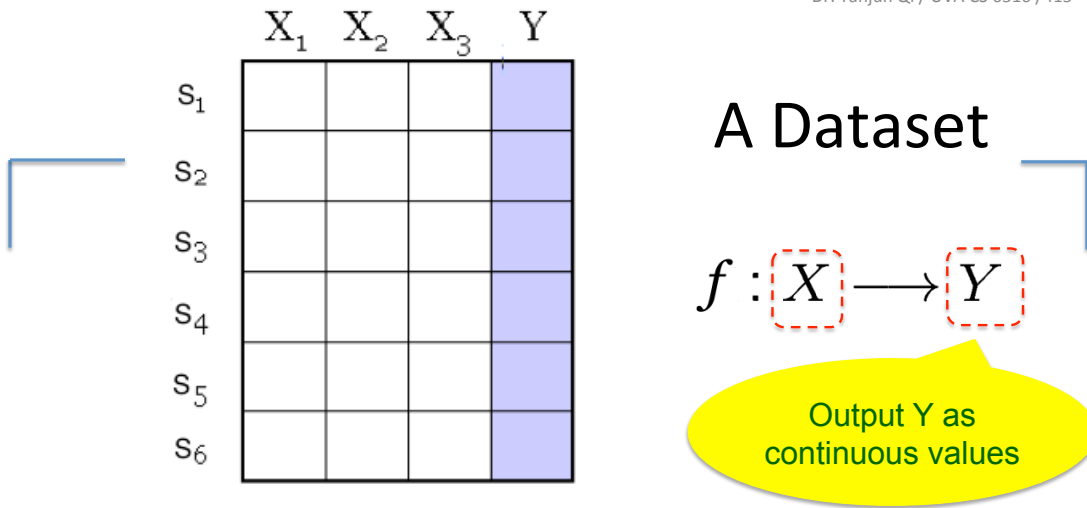
e.g. SUPERVISED LEARNING

$$f : X \longrightarrow Y$$

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

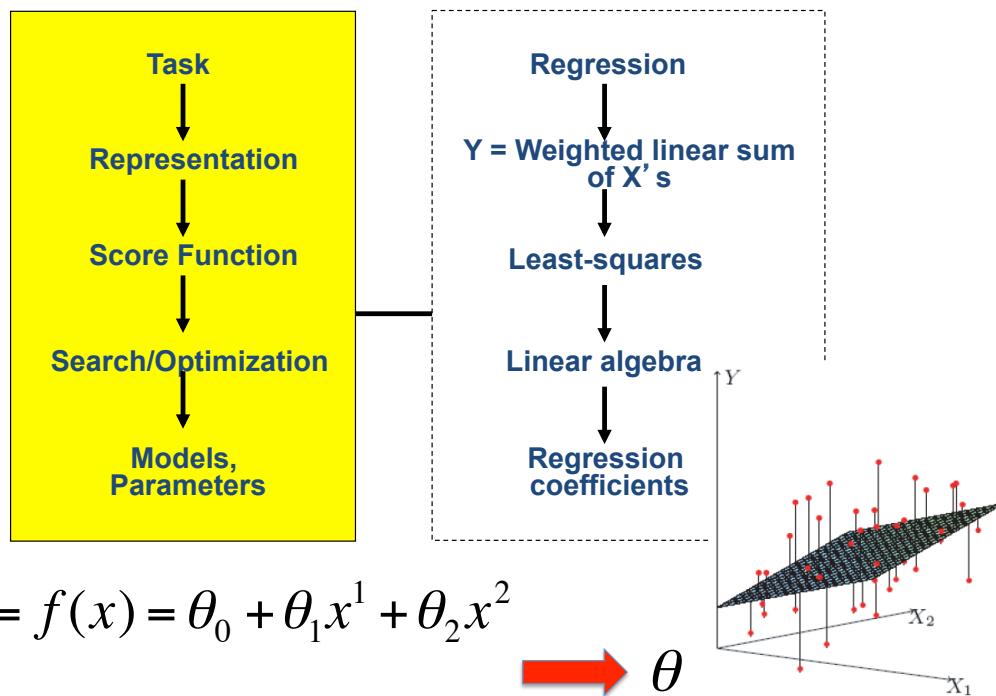
- **Generalisation**: learn function / hypothesis from **past data** in order to “explain”, “predict”, “model” or “control” **new** data examples

KEY



- **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]

(1) Multivariate Linear Regression



Method I: normal equations

- Write the cost function in matrix form:

$$\begin{aligned}
 J(\theta) &= \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2 \\
 &= \frac{1}{2} (X\theta - \bar{y})^T (X\theta - \bar{y}) \\
 &= \frac{1}{2} (\theta^T X^T X \theta - \theta^T X^T \bar{y} - \bar{y}^T X \theta + \bar{y}^T \bar{y})
 \end{aligned}
 \quad
 \mathbf{X} = \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ \vdots & \vdots & \vdots \\ - & \mathbf{x}_n^T & - \end{bmatrix}
 \quad
 \bar{\mathbf{y}} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

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

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

The normal equations

$$\Downarrow$$

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

11/16/15

31

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

11/16/15

32

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:

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

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

11/16/15

33

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?

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

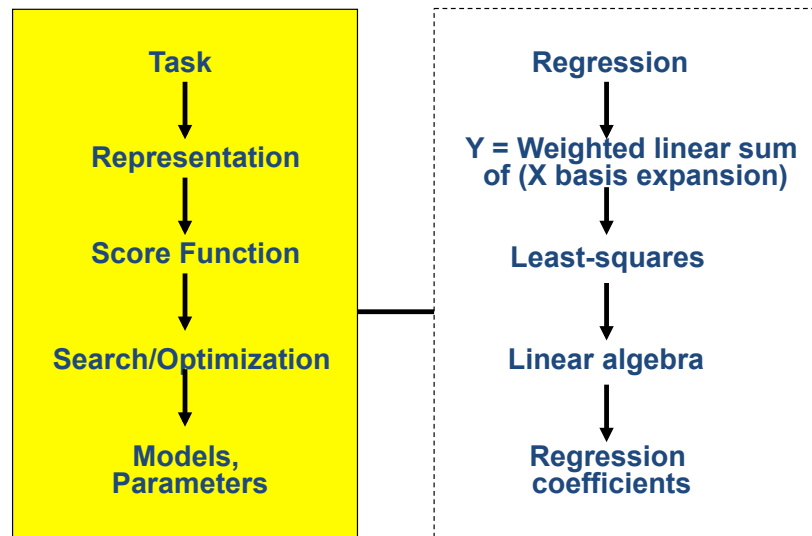
Newton's method
for Linear Regression

11/

9/9/14

25

(2) Multivariate Linear Regression with basis Expansion



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

11/16/15

35

(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:

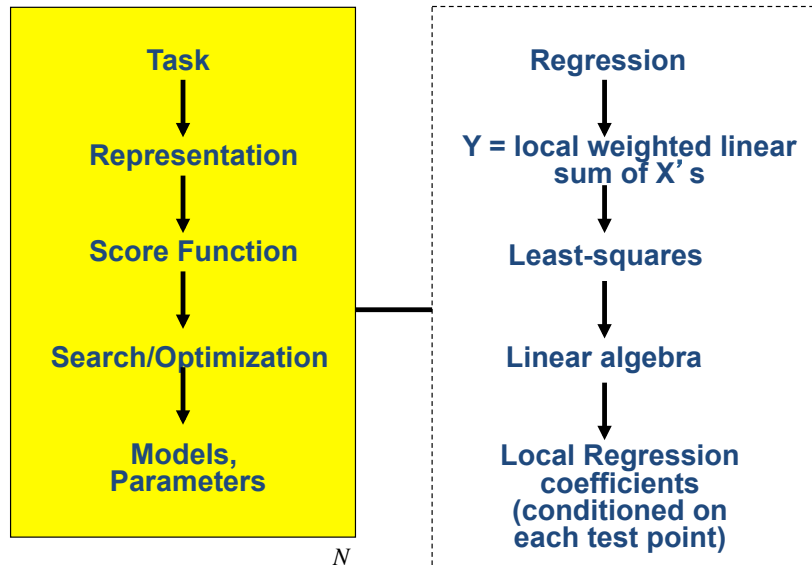
$$\varphi(x) := [1, x, x^2, x^3]$$

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

11/16/15

36

(3) Locally Weighted / Kernel Regression



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

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

11/16/15

37

(3) Locally weighted regression

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

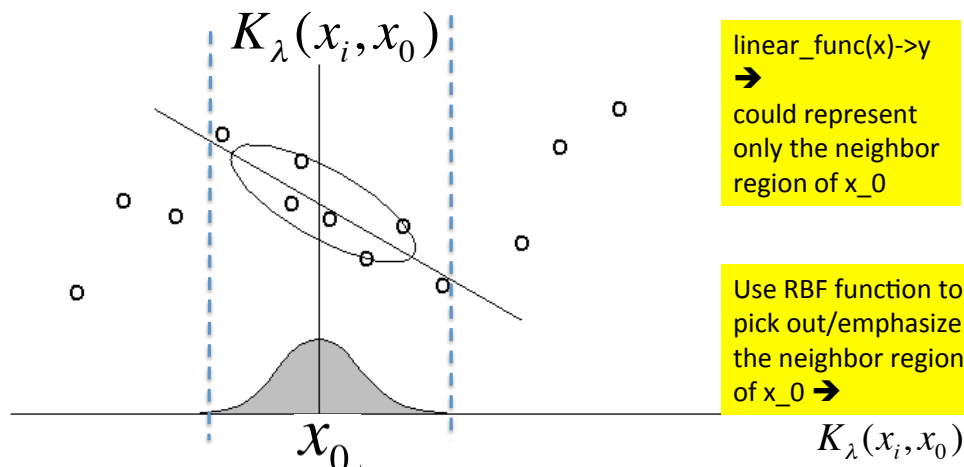
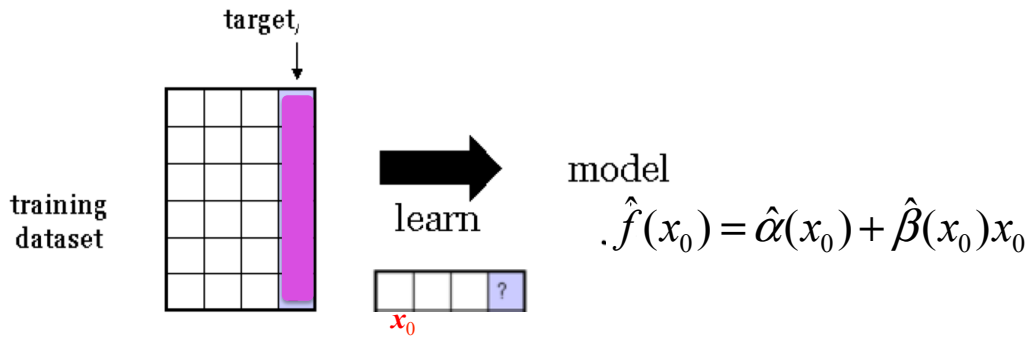


Figure 2: In locally weighted regression, points are weighted by proximity to the current x in question using a kernel. A regression is then computed using the weighted points.

11/16/15

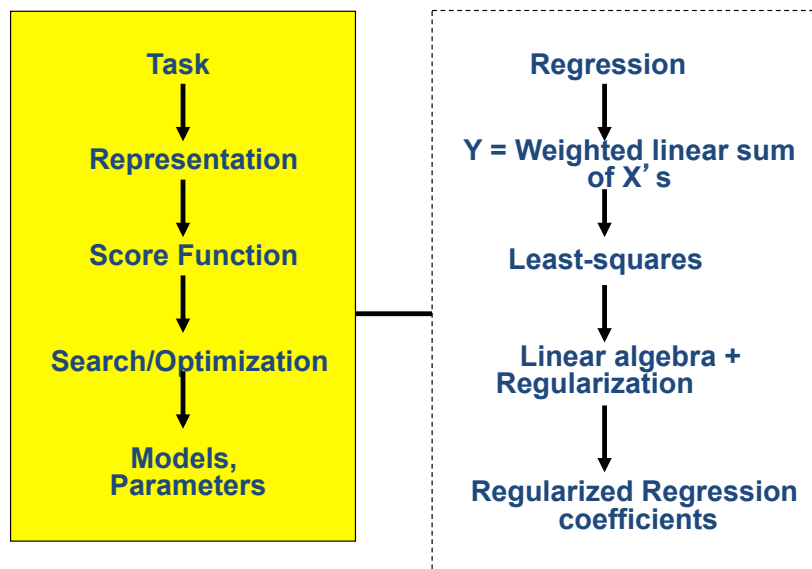
LEARNING of Locally weighted linear regression



→ 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

11/16/15

What we have covered (II)

- Supervised Classification models
 - Support Vector Machine
 - Bayes Classifier
 - Logistic Regression
 - K-nearest Neighbor
 - Random forest / Decision Tree
 - Neural Network (e.g. MLP)
 - *Feature selection

11/16/15

42

Three major sections for classification

- We can divide the large variety of classification approaches into **roughly three major types**

- Discriminative**
 - directly estimate a decision rule/boundary
 - e.g., **logistic regression**, support vector machine, decisionTree
- Generative:**
 - build a generative statistical model
 - e.g., **naïve bayes classifier**, Bayesian networks
- Instance based classifiers**
 - Use observation directly (no models)
 - e.g. **K nearest neighbors**

11/16/15

43

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

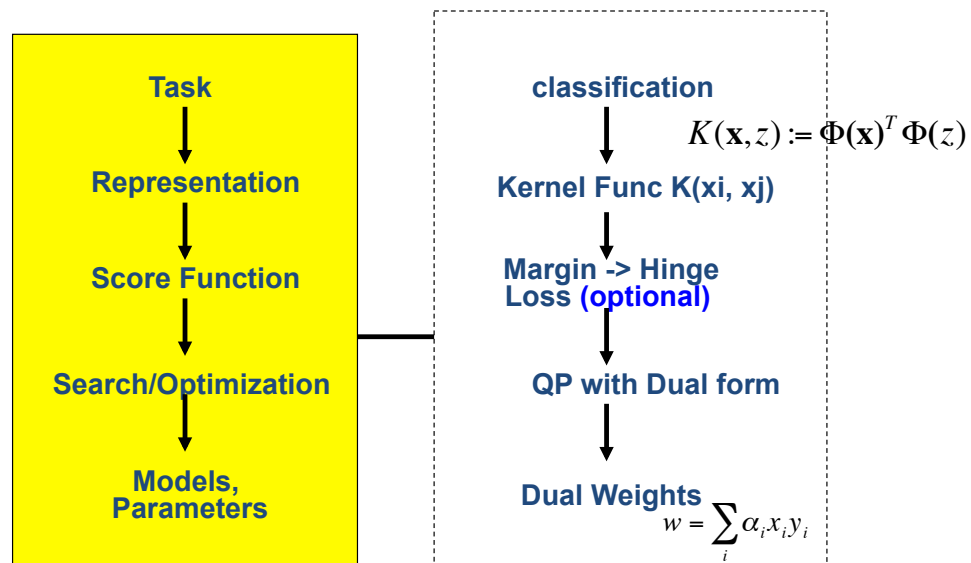
Output as Discrete
Class Label
 C_1, C_2, \dots, C_L

- 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]

11/16/15

44

(1) Support Vector Machine



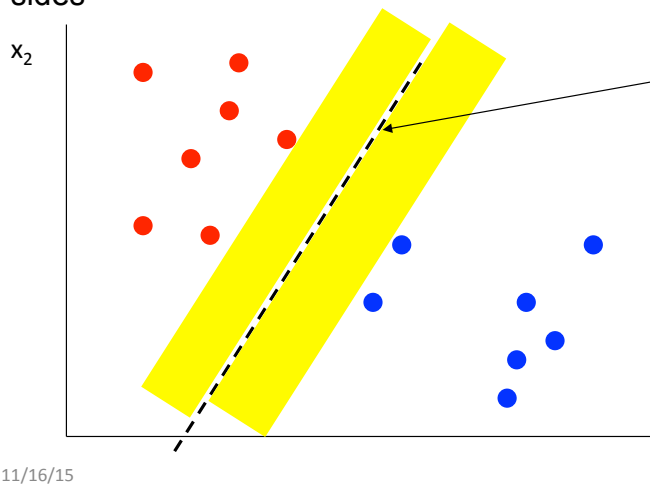
$$\underset{w, b}{\operatorname{argmin}} \sum_{i=1}^p w_i^2 + C \sum_{i=1}^n \varepsilon_i$$

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

11/16/15

(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

11/16/15

x₁

46

X_1	X_2	X_3	Y

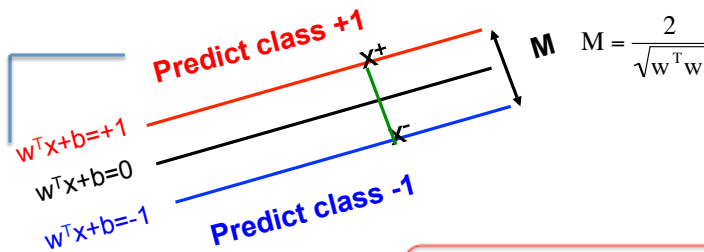
A Dataset for binary classification

$$f : X \rightarrow Y$$

Output as Binary Class Label:
1 or -1

- **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]

When linearly Separable Case: Optimization Step i.e. learning optimal parameter for SVM



1. Correctly classifies all points
2. 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

$$w^T x + b \leq -1$$

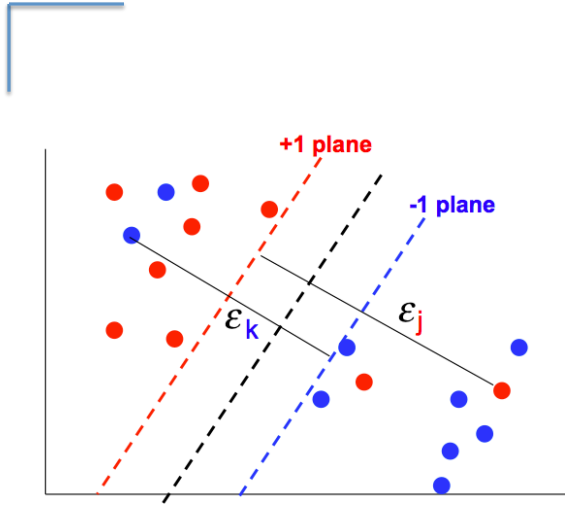
A total of n constraints if we have n input samples

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

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

SVM as a QP problem

Final optimization for non linearly separable case



The new optimization problem is:

$$\min_w \frac{w^T w}{2} + \sum_{i=1}^n C \epsilon_i$$

hyperpara

subject to the following inequality constraints:

For all x_i in class + 1

$$w^T x_i + b \geq 1 - \epsilon_i$$

For all x_i in class - 1

$$w^T x_i + b \leq -1 + \epsilon_i$$

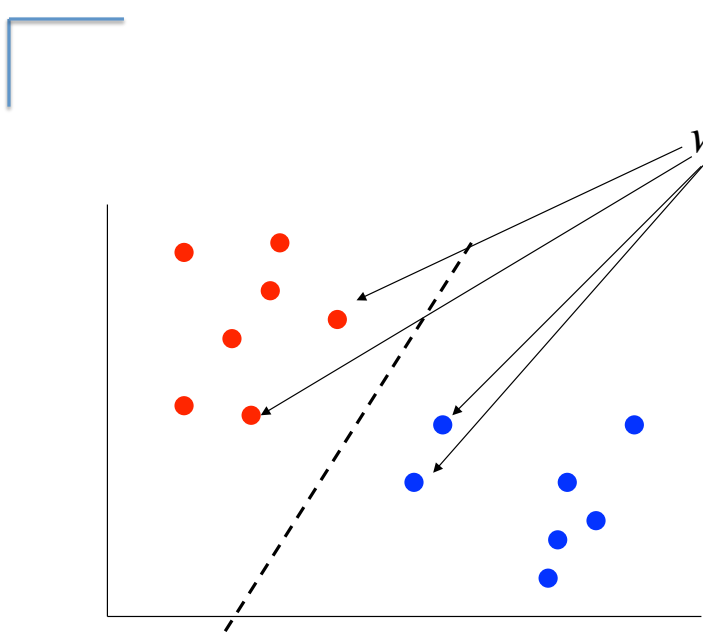
For all i

$$\epsilon_i \geq 0$$

} A total of n constraints

} Another n constraints

Dual SVM - interpretation



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

training points

$\alpha_i \geq 0$

For α_i that are 0, no influence

Dual SVM for linearly separable case

Our dual target function: $\max_{\alpha} \sum_{i=1}^{n_{train}} \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$

Handwritten notes: n_{train} above the sum, $\sum_{i=1}^{n_{train}}$ and $\sum_{j=1}^{n_{train}}$ with arrows pointing to the indices.

Training

$$\sum_i \alpha_i y_i = 0$$

$$\alpha_i \geq 0 \quad \forall i$$

Dot product for all training samples

Dot product with training samples

To evaluate a new sample \mathbf{x}_{ts} we need to compute:

Testing

$$\mathbf{w}^T \mathbf{x}_{ts} + b = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_{ts} + b$$

Handwritten notes: \mathbf{x}_{ts} circled, $\mathbf{x}_i^T \mathbf{x}_{ts}$ circled, α_i circled.

~~$\alpha_i \geq 0$~~
 $\alpha_i > 0$

$$\hat{y}_{ts} = \text{sign} \left(\sum_{i \in SV} \alpha_i y_i (\mathbf{x}_i^T \mathbf{x}_{ts}) + b \right)$$

Dual formulation for linearly non separable case

Dual target function:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C > \alpha_i \geq 0, \forall i$$

Hyperparameter C should be tuned through k-folds CV

The only difference is that the α are now bounded

To evaluate a new sample \mathbf{x}_j we need to compute:

$$\mathbf{w}^T \mathbf{x}_j + b = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x}_j + b$$

This is very similar to the optimization problem in the linear separable case, except that there is an upper bound C on α_i now

Once again, efficient algorithm exist to find α_i

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

However, we can obtain dramatic savings by noting that

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

$O(m)$

$K(x, z)$
We only need m operations!

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

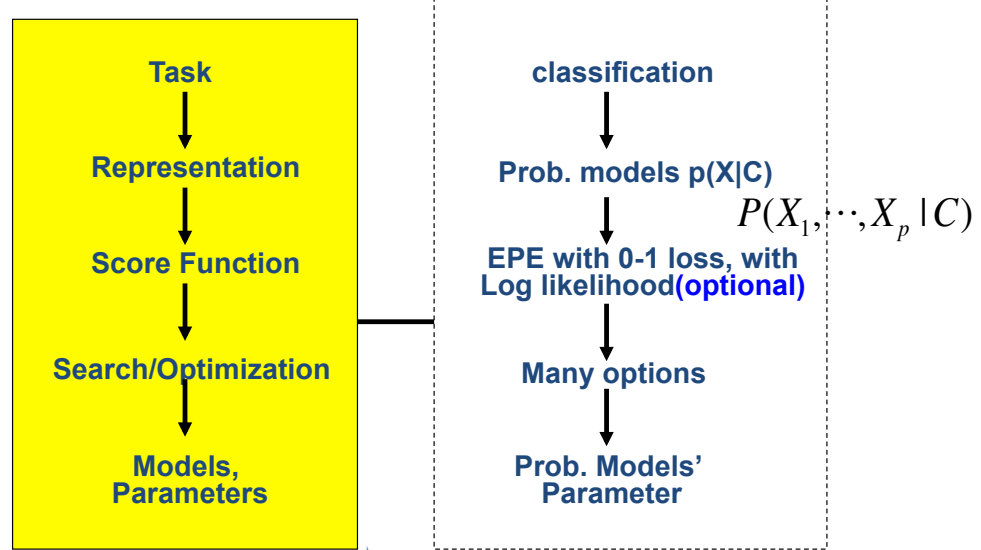
$$K(x, z) = (x^T z + 1)^2$$

11/16/15

53

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

(2) Bayes Classifier



Bernoulli Naive

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

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

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

11/16/15

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

- **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]

11/16/15

55

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

11/16/15

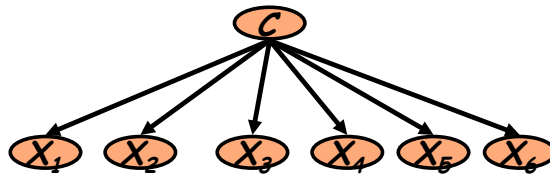
56

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!**

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



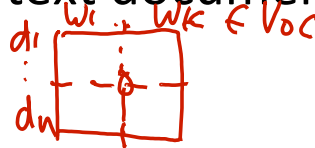
11/16/15

57

Adapt from Prof. Ke Chen NB slides

Probabilistic Models of text documents

$\Pr(D | C = c)$? $D = (w_1, w_2, \dots, w_k)$



Two
Previous
models

$\Pr(W_1 = \text{true}, W_2 = \text{false}, \dots, W_k = \text{true} | C = c)$

Multivariate Bernoulli Distribution

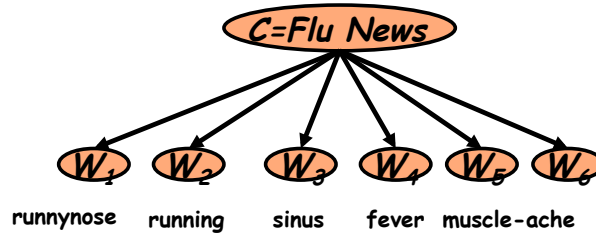
$\Pr(W_1 = n_1, W_2 = n_2, \dots, W_k = n_k | C = c)$

Multinomial Distribution

11/16/15

58

(2.1) : Multivariate Bernoulli for text



- **Conditional Independence Assumption:** Features (word presence) are *independent* of each other given the class variable:

this is naïve

$$\Pr(W_1 = true, W_2 = false, \dots, W_k = true | C = c) = P(W_1 = true | C) \cdot P(W_2 = false | C) \cdot \dots \cdot P(W_k = true | C)$$

Ber

- Multivariate Bernoulli model is appropriate for **binary feature variables**

(2.2) Multinomial Naïve Bayes as Stochastic Language Models

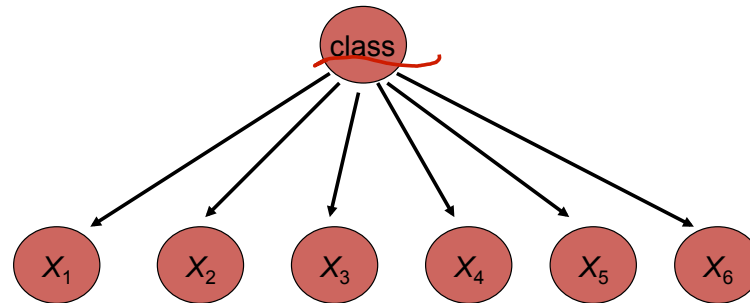
Model C1	
0.2	the
0.01	boy
0.0001	said
0.0001	likes
0.0001	black
0.0005	dog
0.01	garden

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

the	boy	likes	black	dog
0.2	0.01	0.0001	0.0001	0.0005
0.2	0.0001	0.02	0.1	0.01

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

Multinomial Naïve Bayes = a class conditional unigram language model



- Think of X_i as the word on the i^{th} position in the document string
- Effectively, the probability of each class is done as a class-specific unigram language model

11/16/15

Adapt From Manning' textCat tutorial ⁶¹

(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

11/16/15

62

Naïve Gaussian means ?

$O(kp^2 + kp)$
 \sum_k

Not Naive

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

Naive

p-dim random vector

$$P(X_1, X_2, \dots, X_p | C = c_j) = P(X_1 | C) P(X_2 | C) \dots P(X_p | C)$$

$$= \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp \left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2} \right)$$

$O(kp + kp)$

Diagonal Matrix

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

Each class' covariance matrix is diagonal

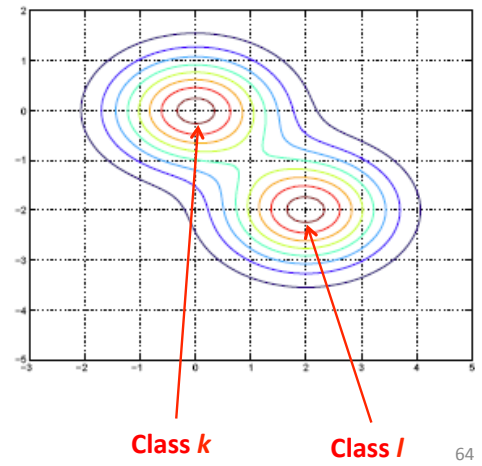
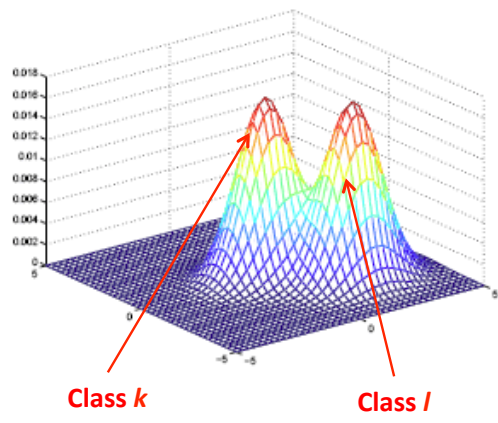
(2.4) LDA (Linear Discriminant Analysis)

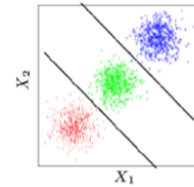
$O(p^2 + kp)$

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





Optimal Classification

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

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

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

- Note

Linear Discriminant Function for LDA

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

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

$$\log \frac{P(C = k | X)}{P(C = l | X)} = \log \frac{P(X | C = k)}{P(X | C = l)} + \log \frac{P(C = k)}{P(C = l)}$$

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

$$+ x^T \Sigma^{-1} (\mu_k - \mu_l)$$

Equals to zero

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

(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:

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

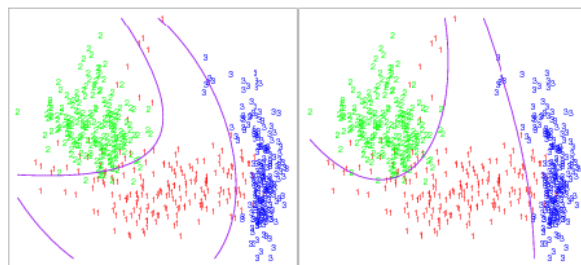
- ▶ Decision boundaries are quadratic equations in x .
- ▶ QDA fits the data better than LDA, but has more parameters to estimate.

11/16/15

67

(2.6) LDA on Expanded Basis

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



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 $x_1, x_2, x_1 x_2, x_1^2, x_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.

11/16/15

68

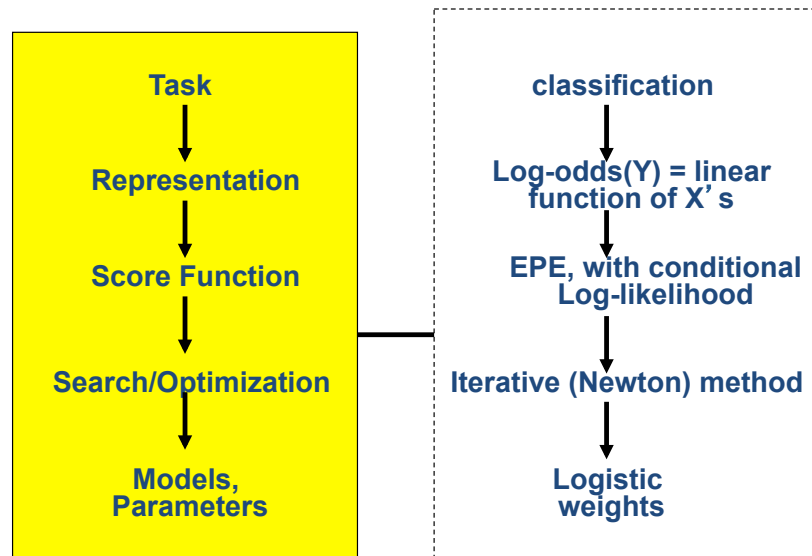
(2.7) Regularized Discriminant Analysis

- ▶ 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.

(3) Logistic Regression



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

Logistic Regression—when?

Logistic regression models are appropriate for target variable coded as 0/1.

We only observe “0” and “1” for the target variable—but we think of the target variable conceptually as a probability that “1” will occur.

Binary
 $p(y=1|x)$
logistic

This means we use Bernoulli distribution to model the target variable with its Bernoulli parameter $p(y=1|x)$ predefined.

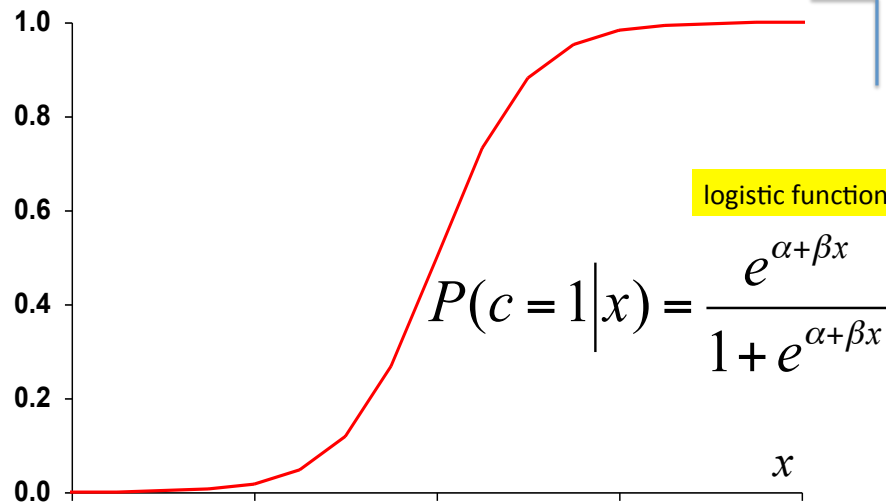
The main interest → predicting the probability that an event occurs (i.e., the probability that $p(y=1|x)$).

Discriminative

Logistic regression models for binary target variable coded 0/1.

e.g.
Probability of disease

$P(C=1|X)$



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

Logit function

$$\ln \left[\frac{P(c = 1|x)}{P(c = 0|x)} \right] = \ln \left[\frac{P(c = 1|x)}{1 - P(c = 1|x)} \right] = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

Decision Boundary → equals to zero

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.

How?

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

$p(y_i | x_i)$

x_i are $(p+1)$ -dimensional input vector with leading entry 1
 β is a $(p+1)$ -dimensional vector

11/16/15 We want to maximize the log-likelihood in order to estimate β

73

Discriminative vs. Generative

Generative approach

- Model the joint distribution $p(X, C)$ using $p(X | C = c_k)$ and $p(C = c_k)$

Class prior

Discriminative approach

- Model the conditional distribution $p(c | X)$ directly

e.g.,

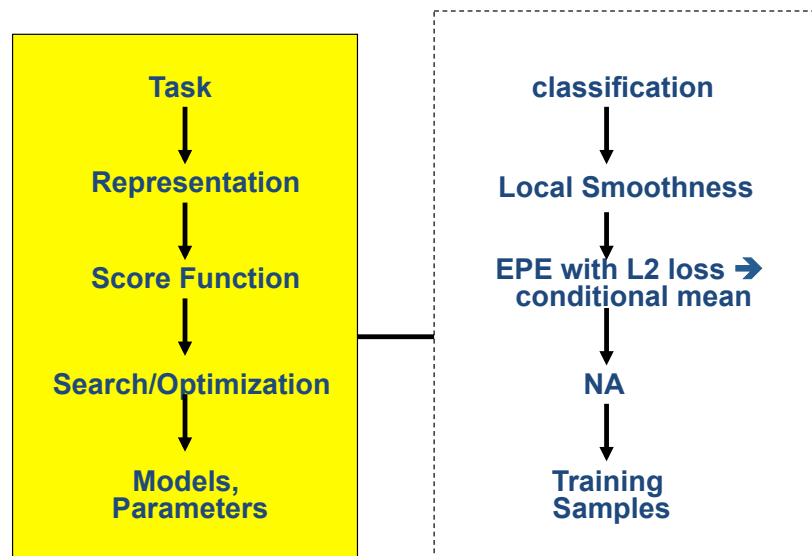
$$\frac{1}{1 + e^{-(\beta_0 + \beta_1 * X)}}$$

Discriminative vs. Generative

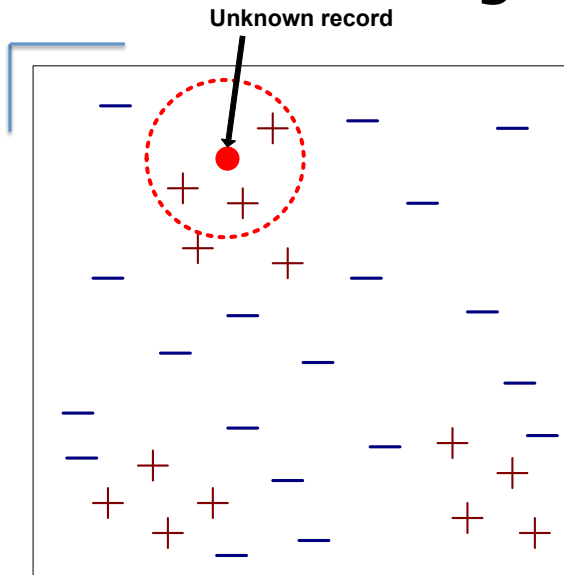
- Empirically, **generative** classifiers approach their asymptotic error faster than discriminative ones
 - Good for small training set
 - Handle missing data well (EM)
- Empirically, **discriminative** classifiers have lower asymptotic error than generative ones
 - Good for larger training set

Dr. Yanjun Qi / UVA CS 6316 / f15

(4) K-Nearest Neighbor



Nearest neighbor classifiers



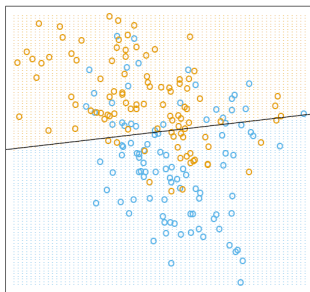
To classify unknown sample:

1. Compute distance to other training records
2. Identify k nearest neighbors
3. Use class labels of nearest neighbors to determine the class label of unknown record (e.g., by taking majority vote)

11/16/15

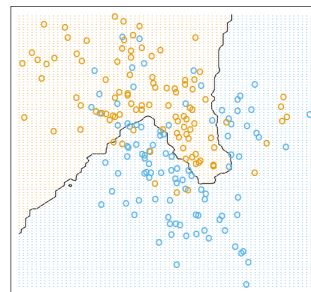
77

Decision boundaries in global vs. local models

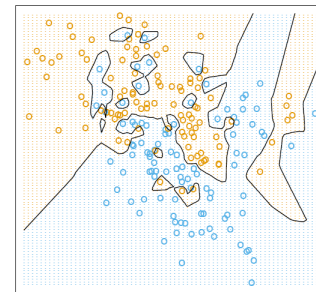


linear regression

- global
- stable
- can be inaccurate



15-nearest neighbor



1-nearest neighbor

- local
- accurate
- unstable

What ultimately matters: **GENERALIZATION**

11/16/15

78

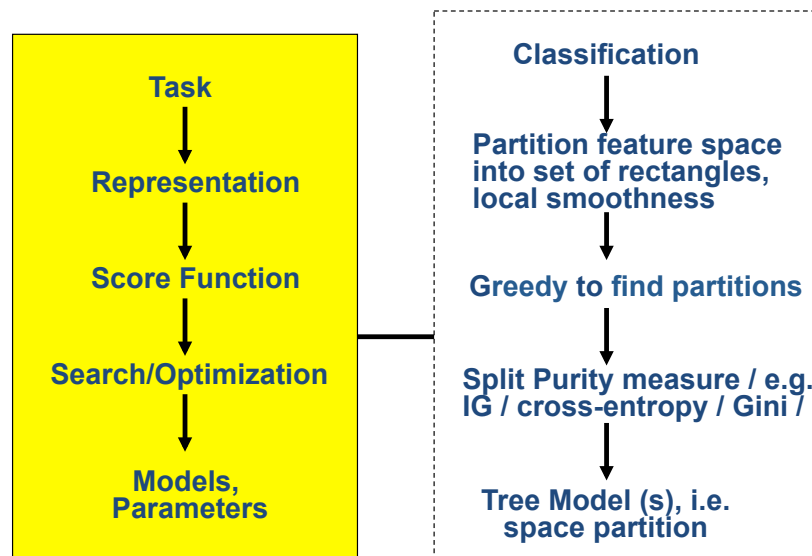
Nearest neighbor classification

- 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.

11/16/15

79

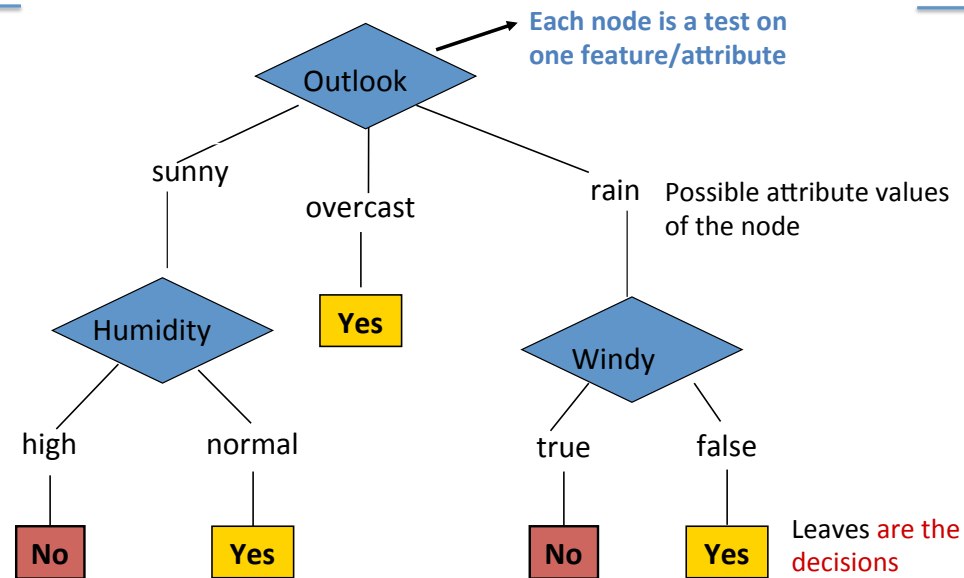
(5) Decision Tree / Random Forest



11/16/15

80

Anatomy of a decision tree



11/16/15

81

Decision trees

- Decision trees represent a disjunction of conjunctions of constraints on the attribute values of instances.

- (Outlook ==overcast)
- OR
- ((Outlook==rain) and (Windy==false))
- OR
- ((Outlook==sunny) and (Humidity=normal))
- => yes play tennis

11/16/15

82

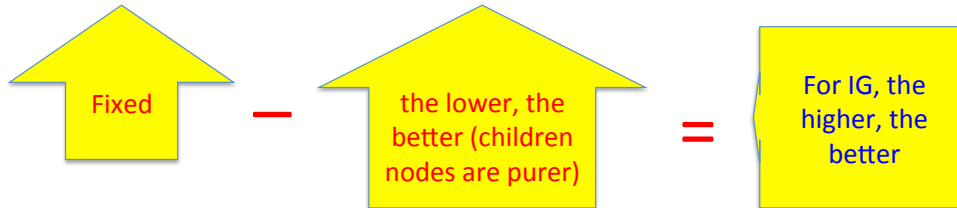
Information gain

- $IG(X_i, Y) = H(Y) - H(Y|X_i)$

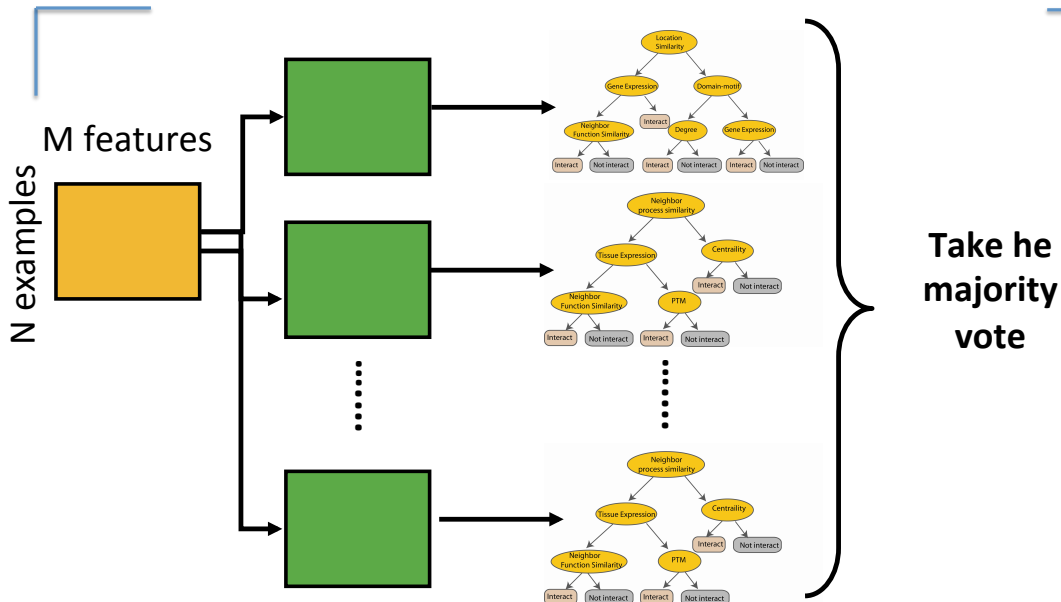
Reduction in uncertainty by knowing a feature X_i

Information gain:

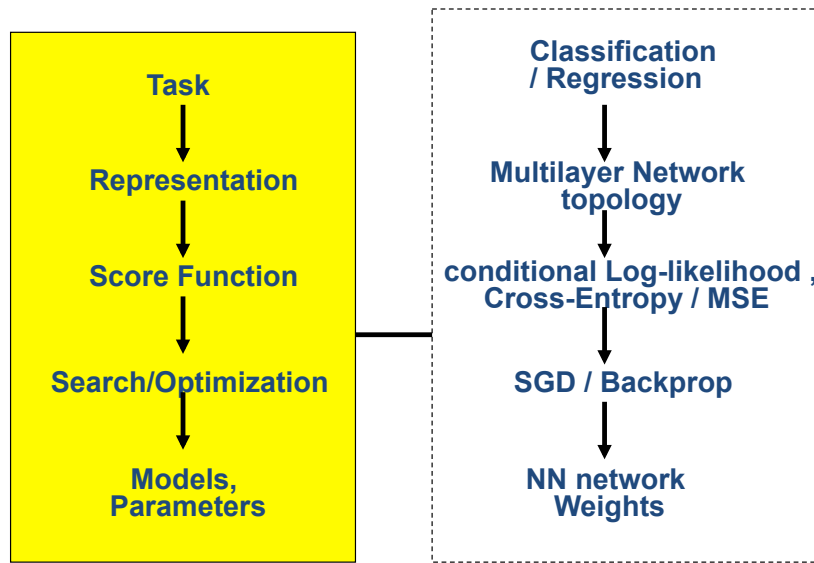
= (information before split) – (information after split)
 = entropy(parent) – [average entropy(children)]



Random Forest Classifier



(6) Neural Network

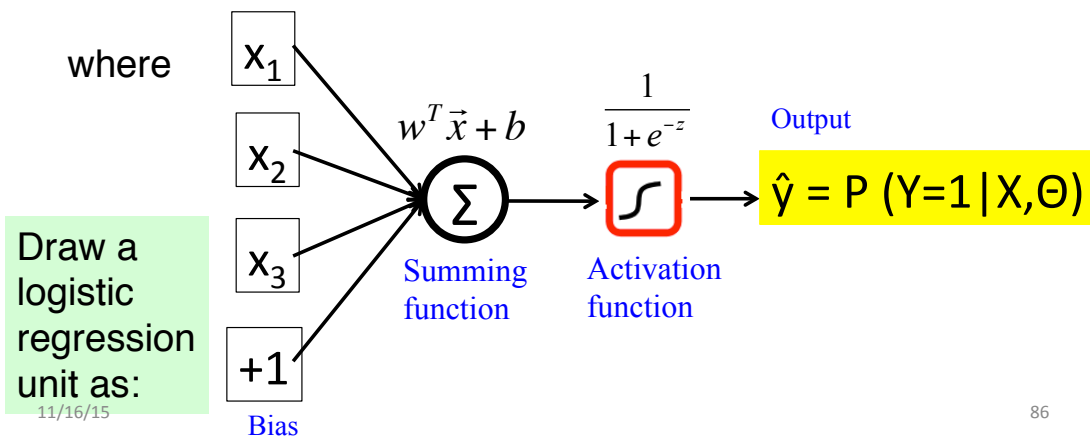


Logistic regression

Logistic regression could be illustrated as a module

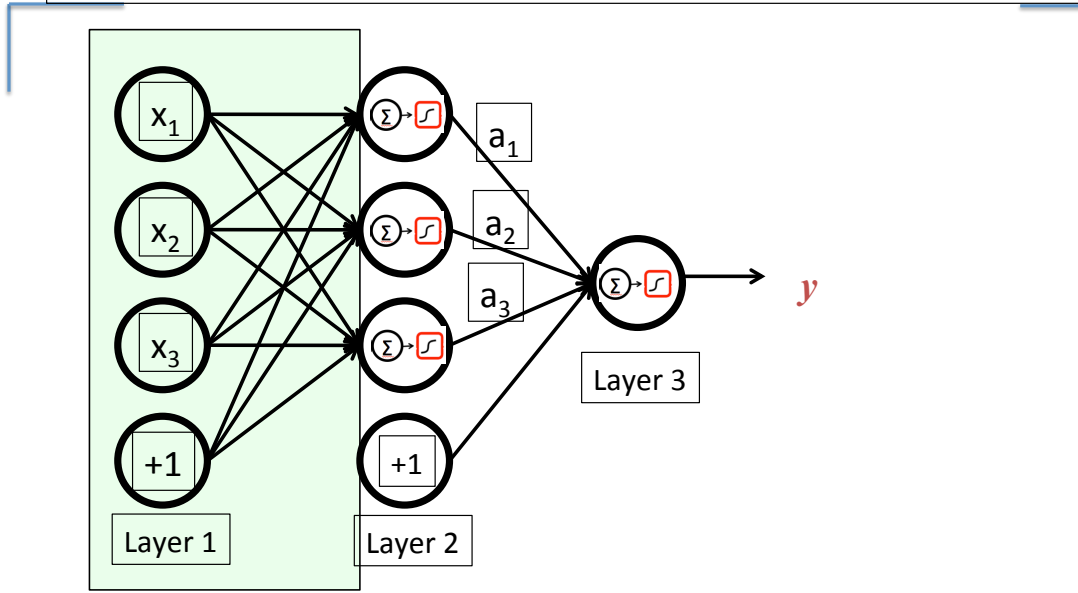
On input x , it outputs \hat{y} :

where



Multi-Layer Perceptron (MLP)

String a lot of logistic units together. Example: 3 layer network:

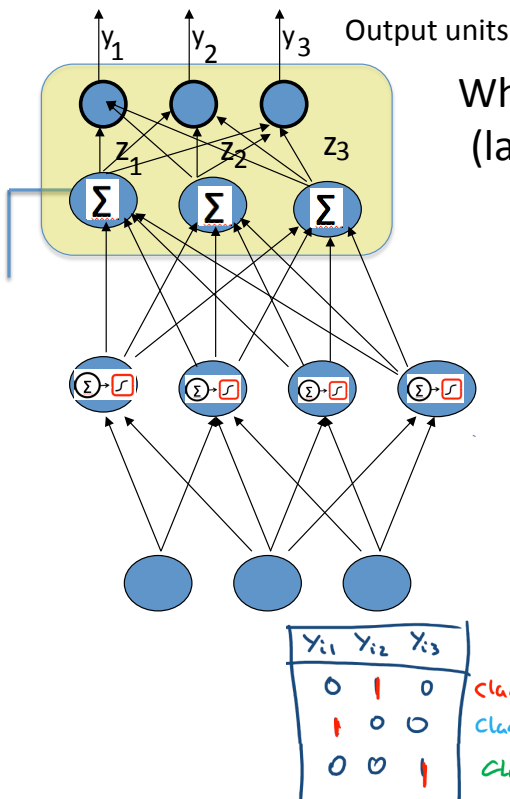


11/16/15 input

hidden

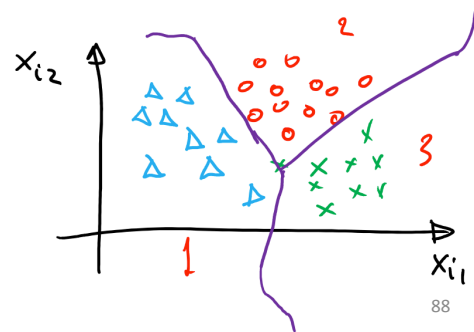
output

87



When for multi-class classification (last output layer: softmax layer)

When multi-class output, last layer is softmax output layer → a multinomial logistic regression unit

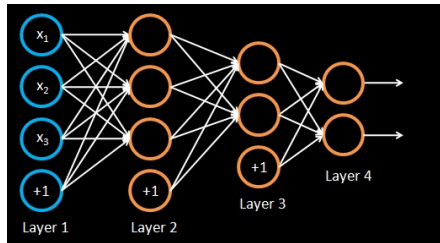
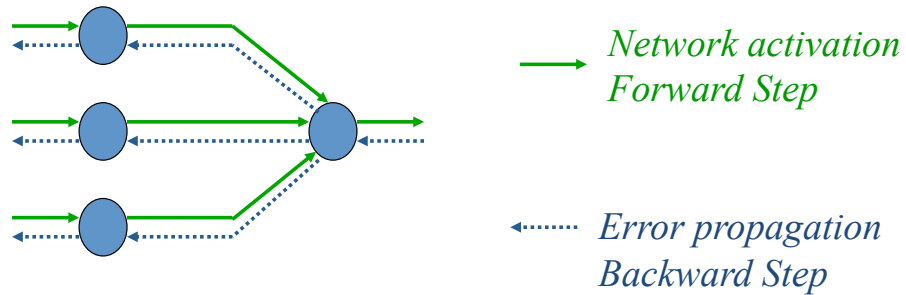


11/16/15

88

Backpropagation

- Back-propagation training algorithm



11/16/15

89

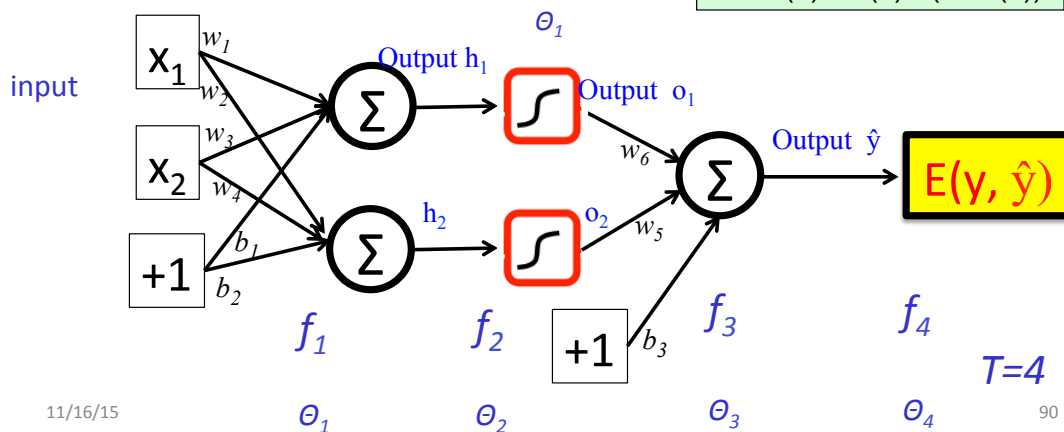
to train this layered network. The stacked layers in our network can be written in a more general form of multi-level functions:

$$l_{\mathbf{x}} = \mathbf{f}_T(\mathbf{f}_{T-1}(\dots(\mathbf{f}_1(\mathbf{x}))\dots)),$$

where $l_{\mathbf{x}}$ denotes the loss on a single example \mathbf{x}

For instance → for regression

for sigmoid unit o ,
its derivative is,
 $o'(h) = o(h) * (1 - o(h))$



11/16/15

90

$\mathbf{f}_i, i \in [1, T]$, the derivative for updating its parameter set θ_i is using the delta rule:

$$\frac{\partial l}{\partial \theta_i} = \frac{\partial f_T}{\partial f_i} \times \frac{\partial f_i}{\partial \theta_i},$$

and the first factor on the right can be recursively calculated:

$$\frac{\partial f_T}{\partial f_i} = \frac{\partial f_T}{\partial f_{i+1}} \times \frac{\partial f_{i+1}}{\partial f_i}.$$

Note that \mathbf{f} and θ are usually vectors

so $\frac{\partial f_T}{\partial f_{i+1}}$ and $\frac{\partial f_i}{\partial \theta_i}$ are Jacobian matrices, and “ \times ” is matrix multiplication.

e.g.

$$\frac{\partial f_4}{\partial f_3} = \frac{\partial (y - \hat{y})^2}{\partial f_3} = \frac{\partial (y - \hat{y})}{\partial f_3} = \frac{\partial (y - \hat{y})}{\partial f_3}$$

output error

11/16/15

Dr. Qi's CIKM 2012 paper/talk ⁹¹

for $j = 1$ to MaxIter **do**

if converge **then**

 break

end if

$\mathbf{x}, y \leftarrow$ random sampled data point and label

 calculate loss $l(\mathbf{x}; y)$

 cumulative $\leftarrow 1$

for $i = T$ to 1 **do**

$$\frac{\partial l}{\partial \theta_i} \leftarrow \text{cumulative} * \frac{\partial f_i}{\partial \theta_i}$$

$$\theta_i \leftarrow \theta_i - \lambda \frac{\partial l}{\partial \theta_i}$$

$$\text{cumulative} \leftarrow \text{cumulative} * \frac{\partial f_{i+1}}{\partial f_i}$$

end for

11/16/15 **end for**

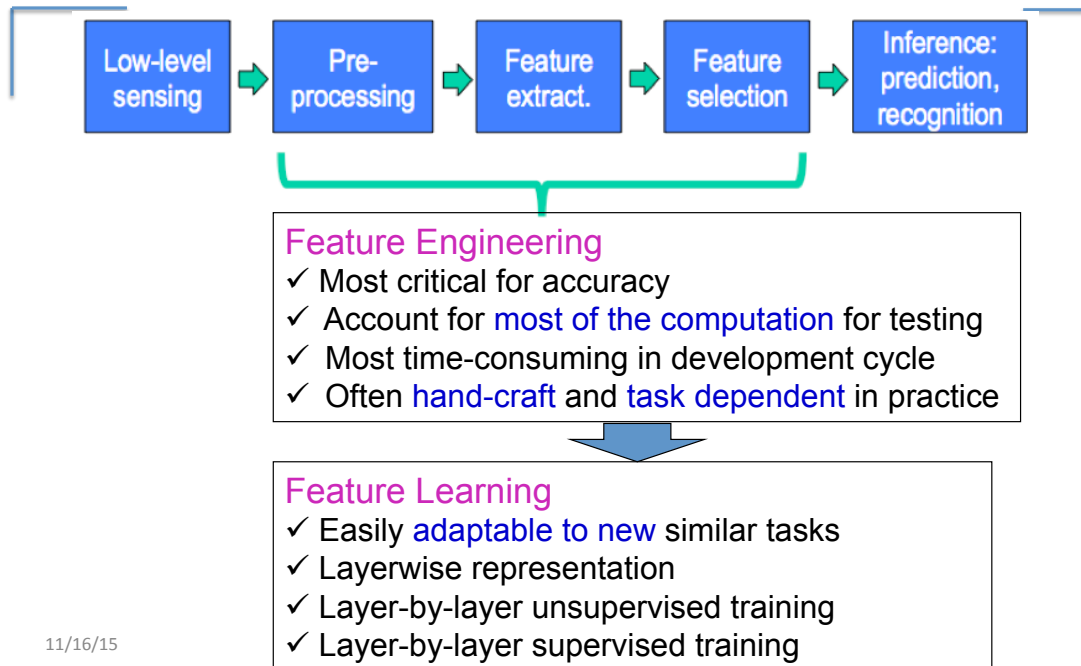
Dr. Yanjun Qi / UVA CS 6316 / f15

*Error propagation
Backward Step*



Dr. Qi's CIKM 2012 paper/talk ⁹²

Deep Learning Way: Learning features / Representation from data



11/16/15

DESIGN ISSUES for Deep NN

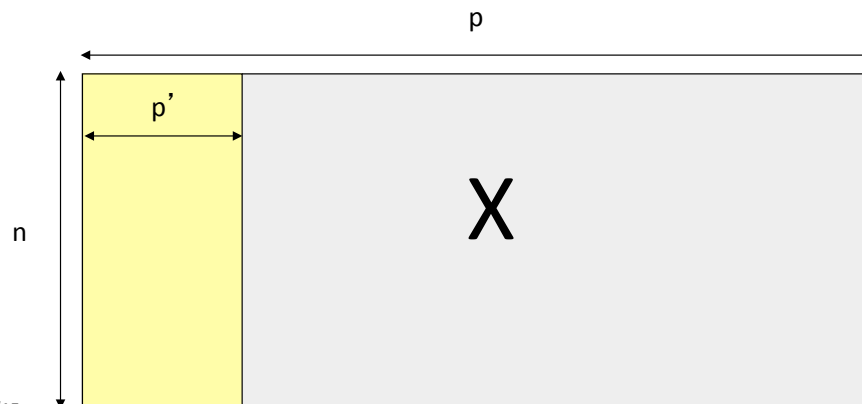
- Data representation
- Network Topology
- Network Parameters
- Training
 - Scaling up with **graphics processors**
 - Scaling up with **Asynchronous SGD**
- Validation

11/16/15

94

(7) Feature Selection

- **Thousands to millions of low level features:** select the most relevant one to build **better, faster, and easier to understand** learning machines.

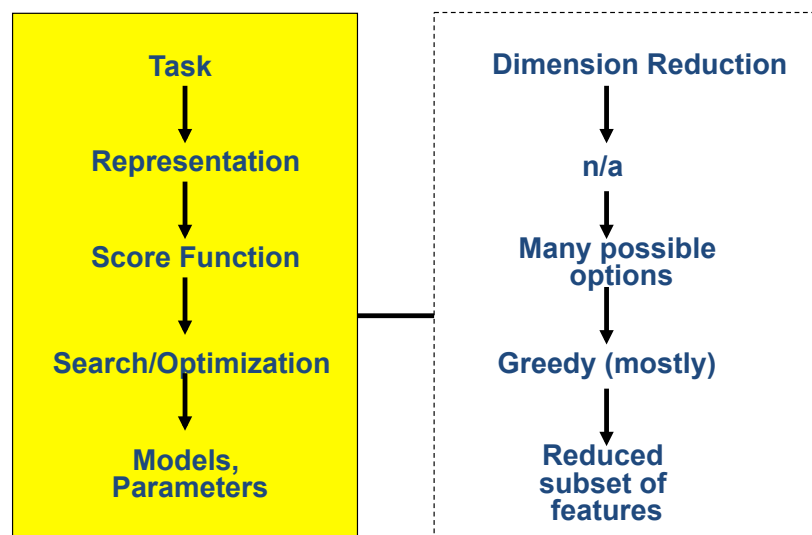


11/16/15

From Dr. **Isabelle Guyon**

95

(7) Feature Selection (not covered)



11/16/15

96

Feature Selection

- **Filtering approach:**
ranks features or feature subsets independently of the predictor (classifier).
 - ...using univariate methods: consider **one** variable at a time
 - ...using multivariate methods: consider **more than one** variables at a time

- **Wrapper approach:**
uses a classifier to assess (many) features or feature subsets.

- **Embedding approach:**
uses a classifier to build a (single) model with a subset of features that are internally selected.

11/16/15

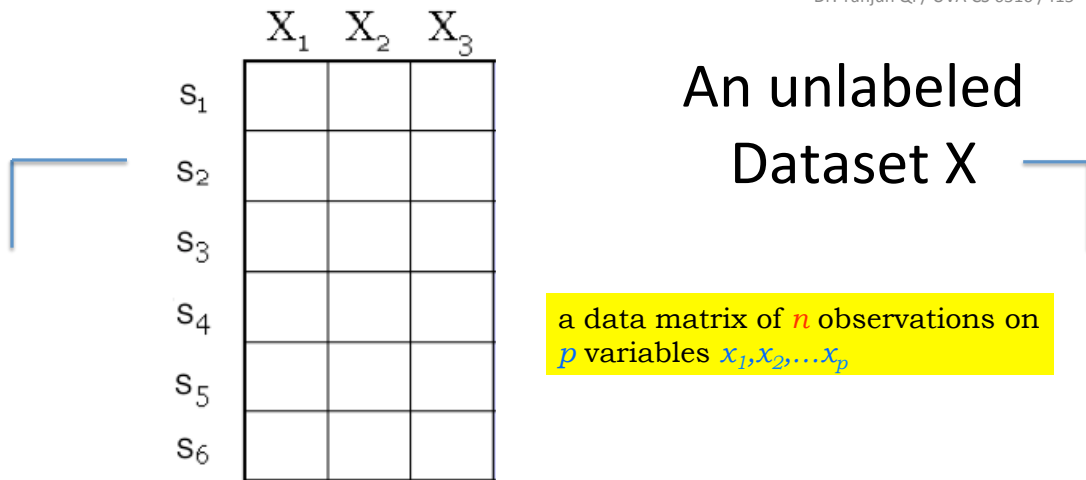
97/54

What we have covered (III)

- **Unsupervised models**
 - Dimension Reduction (PCA)
 - Hierarchical clustering
 - K-means clustering
 - GMM/EM clustering

11/16/15

98



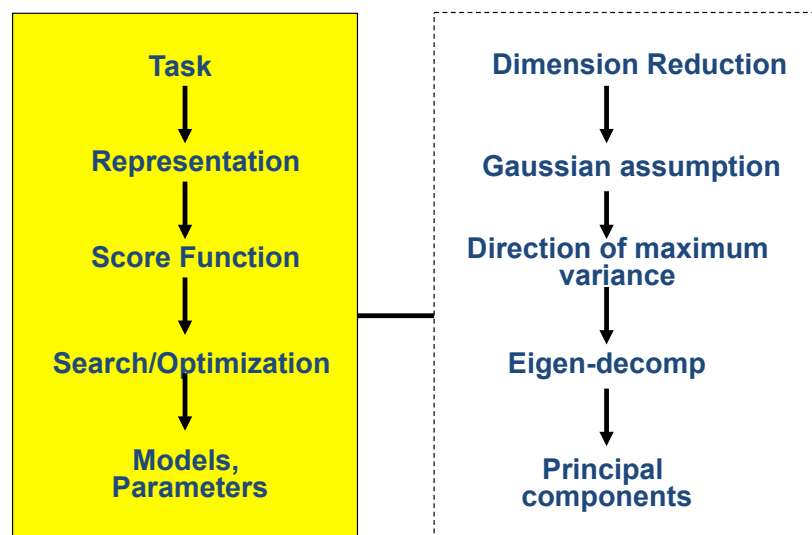
Unsupervised learning = learning from raw (unlabeled, unannotated, etc) data, as opposed to supervised data where a label of examples is given

- **Data**/points/instances/examples/samples/records: [rows]
- **Features**/attributes/dimensions/independent variables/covariates/predictors/regressors: [columns]

11/16/15

99

(0) Principal Component Analysis

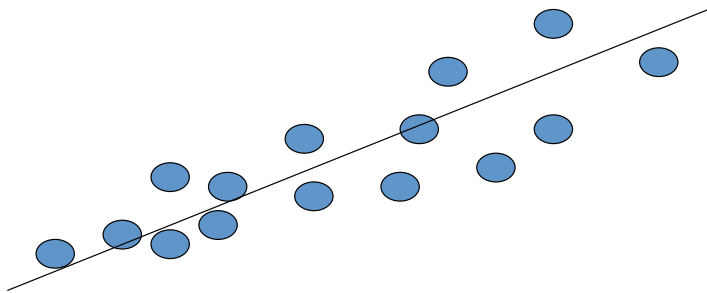


11/16/15

100

Algebraic Interpretation – 1D

- Given n points in a p dimensional space, for large p , how does one project on to a 1 dimensional space?



- Choose a line that fits the data so **the points are spread out well along the line**

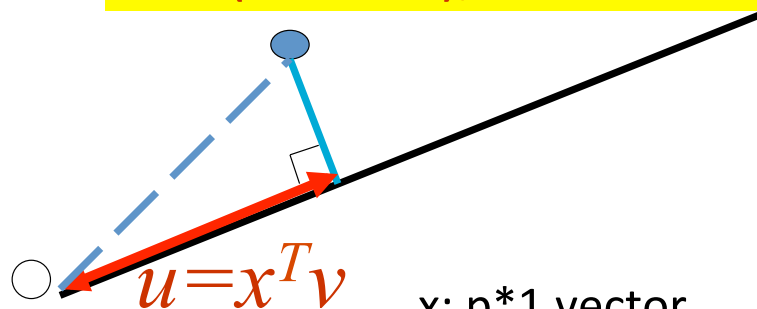
11/16/15

From Dr. S. Narasimhan¹⁰¹

Algebraic Interpretation – 1D

- Minimizing sum of squares of distances to the line is the same as **maximizing the sum of squares of the projections on that line**, thanks to Pythagoras.

$$\max(v^T X^T X v), \text{ subject to } v^T v = 1$$



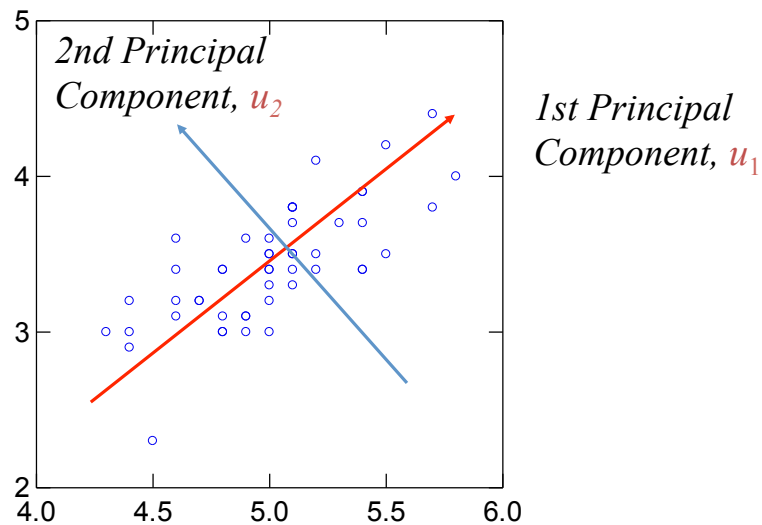
x : $p \times 1$ vector
 v : $p \times 1$ vector

assuming data
is centered

11/16/15

102

PCA Eigenvectors → Principal Components



11/16/15

103

PCA & Gaussian Distributions.

- PCA is similar to learning a Gaussian distribution for the data.
- Dimension reduction occurs **by ignoring the directions in which the covariance is small.**

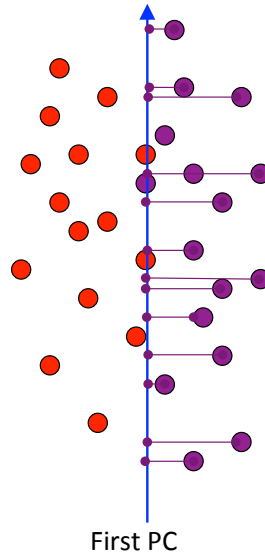
11/16/15

104

PCA Limitations

- The direction of maximum variance is not always good for classification

not ideal for discrimination



11/16/15

From Prof. Derek Hoiem¹⁰⁵

What we have covered (III)

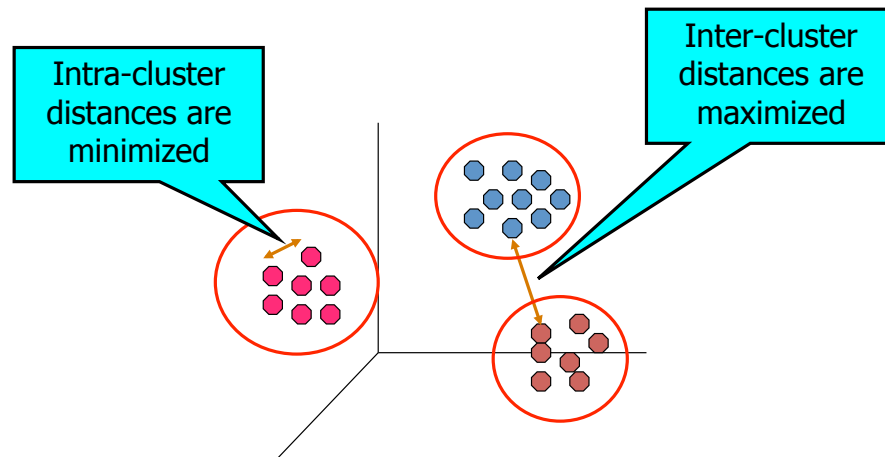
- Unsupervised models
 - Dimension Reduction (PCA)
 - Hierarchical clustering
 - K-means clustering
 - GMM/EM clustering

11/16/15

106

What is clustering?

- Find groups (clusters) of data points such that data points in a group will be similar (or related) to one another and different from (or unrelated to) the data points in other groups



11/16/15

107

Issues for clustering

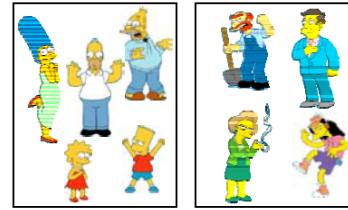
- What is a natural grouping among these objects?
 - Definition of "groupness"
- What makes objects "related"?
 - Definition of "similarity/distance"
- Representation** for objects
 - Vector space? Normalization?
- How many** clusters?
 - Fixed a priori?
 - Completely data driven?
 - Avoid "trivial" clusters - too large or small
- Clustering **Algorithms**
 - Partitional algorithms
 - Hierarchical algorithms
- Formal** foundation and convergence

11/16/15

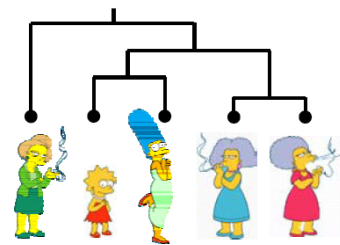
108

Clustering Algorithms

- **Partitional algorithms**
 - Usually start with a random (partial) partitioning
 - Refine it iteratively
 - K means clustering
 - Mixture-Model based clustering



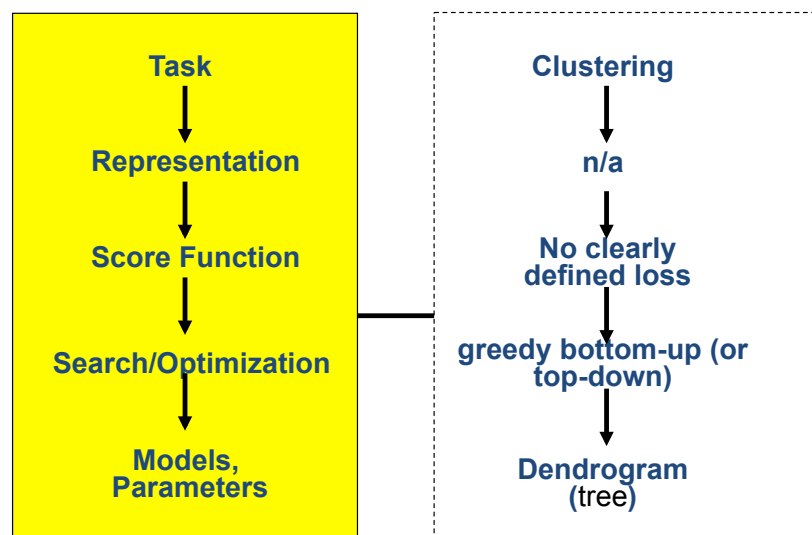
- **Hierarchical algorithms**
 - Bottom-up, agglomerative
 - Top-down, divisive



11/16/15

109

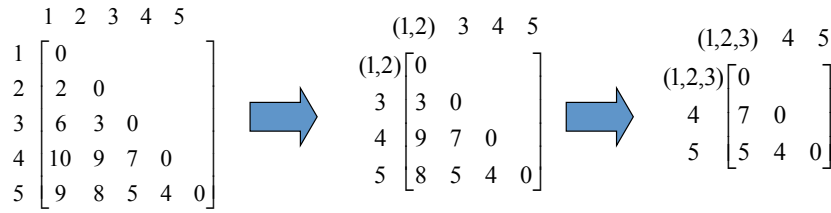
(1) Hierarchical Clustering



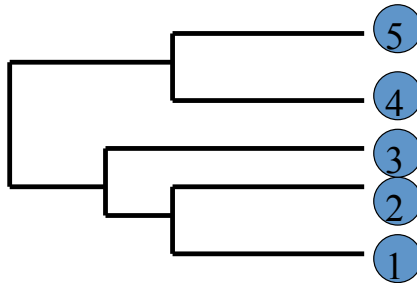
11/16/15

110

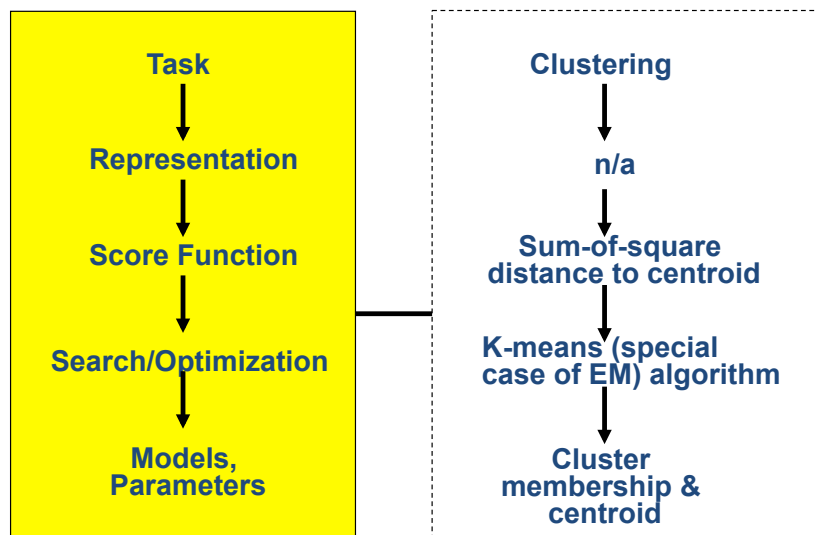
Example: single link



$$d_{(1,2,3),(4,5)} = \min\{d_{(1,2,3),4}, d_{(1,2,3),5}\} = 5$$

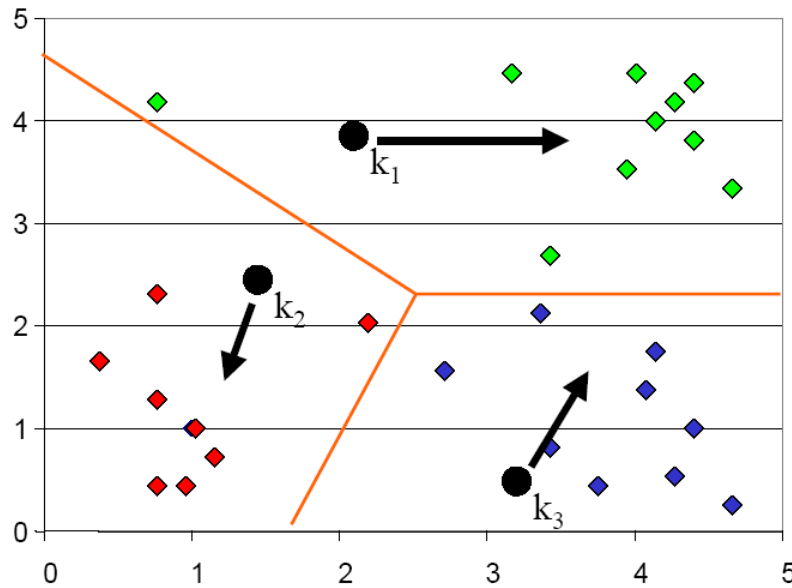


(2) K-means Clustering

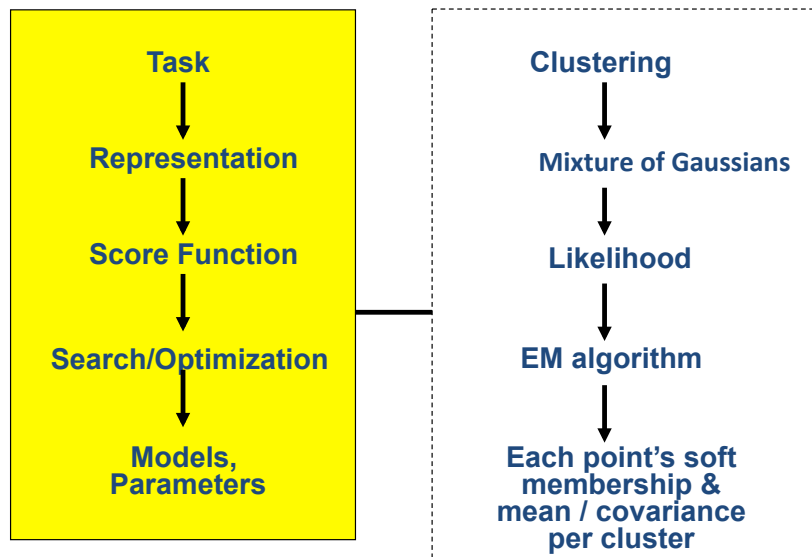


K-means Clustering: Step 2

- Determine the membership of each data points



(3) GMM Clustering



$$\sum_i \log p(x = x_i) = \sum_i \log \left[\sum_{\mu_j} p(\mu = \mu_j) \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{\|x_i - \mu_j\|_2}{2\sigma^2}\right) \right]$$

Expectation-Maximization

Dr. Yanjun Qi / UVA CS 6316 / f15

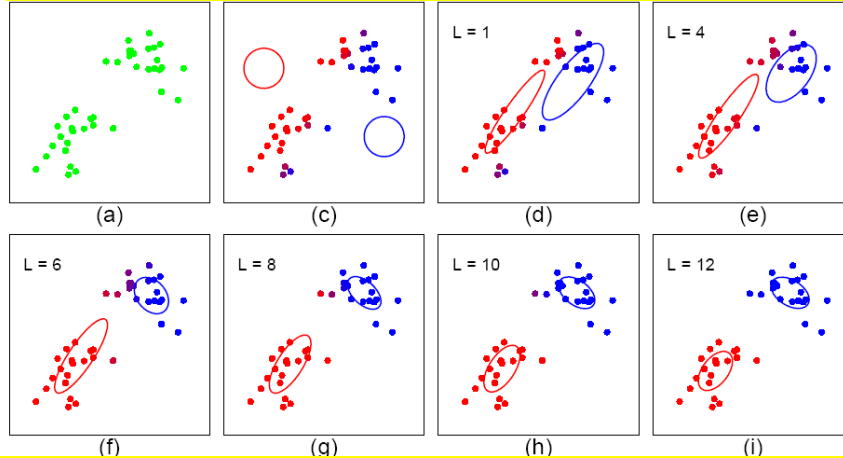
for training GMM

- Start:

- "Guess" the centroid m_k and covariance S_k of each of the K clusters

- Loop

each cluster, revising both the mean (centroid position) and covariance (shape)



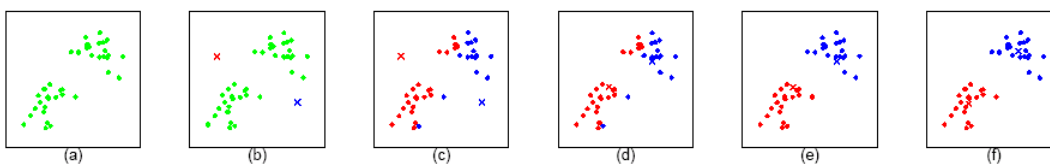
11/16/15

For each point, revising its proportions belonging to each of the K clusters

© Eric Xing @ CMU, 2006-2008

Compare: K-means

- The EM algorithm for mixtures of Gaussians is like a "soft version" of the K-means algorithm.
- In the K-means "E-step" we do **hard** assignment:
- In the **K-means** "M-step" we update the means as the **weighted sum** of the data, but now the weights are 0 or 1:



116

What we have covered (IV)

□ Learning theory / Model selection

- K-folds cross validation
- Expected prediction error
- Bias and variance tradeoff

11/16/15

117

(1) Evaluation Choice: 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













11/16/15

118

CV-based Model Selection

We're trying to decide which algorithm to use.

- We train each model and make a table...

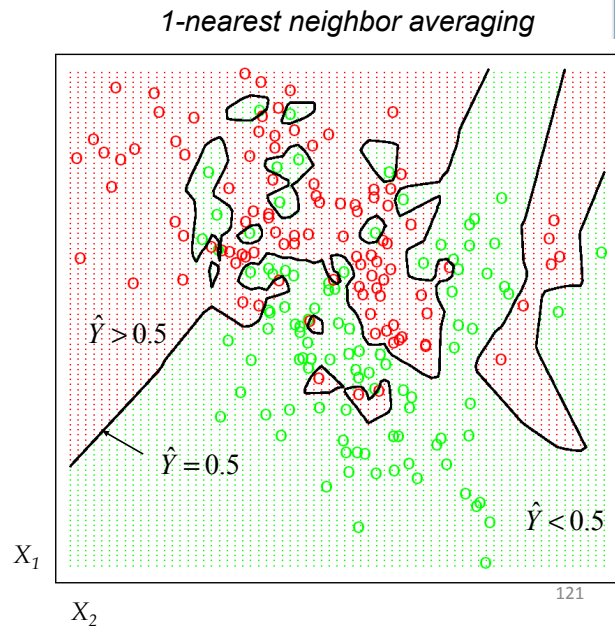
i	f_i	TRAINERR	10-FOLD-CV-ERR	Choice
1	f_1			
2	f_2			
3	f_3			✓
4	f_4			
5	f_5			
6	f_6			

Which kind of cross-validation ?

	Downside	Upside
Test-set	Variance: unreliable estimate of future performance	Cheap
Leave-one-out	Expensive. Has some weird behavior	Doesn't waste data
10-fold	Wastes 10% of the data. 10 times more expensive than test set	Only wastes 10%. Only 10 times more expensive instead of R times.
3-fold	Wastier than 10-fold. Expensivier than test set	Slightly better than test-set
R-fold	Identical to Leave-one-out	

(2) 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)



11/16/15

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):

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

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

Consider
population
distribution

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

11/16/15

122

Expected prediction error (EPE)

Consider joint 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 :

Conditional mean

$$\hat{f}(x) = E(Y | X = x)$$

e.g. KNN



NN methods are the direct implementation (approximation)

- For 0-1 loss: $L(k, \ell) = 1 - d_{kl}$

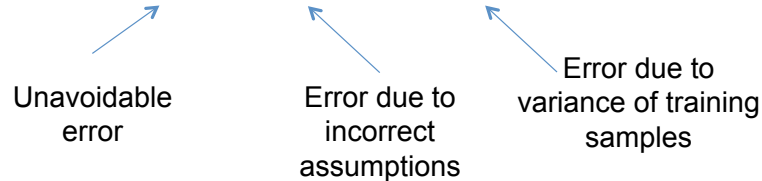


Bayes Classifier

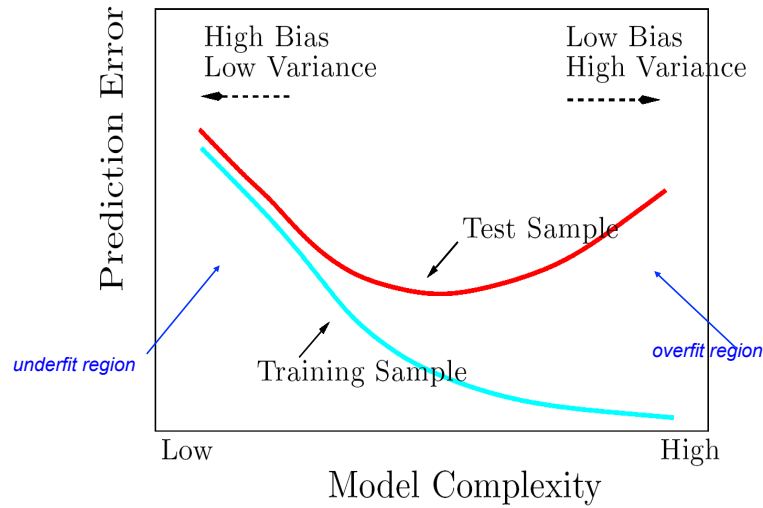
$$\hat{f}(X) = C_k \text{ if } \Pr(C_k | X = x) = \max_{g \in C} \Pr(g | X = x)$$

(3) Bias-Variance Trade-off

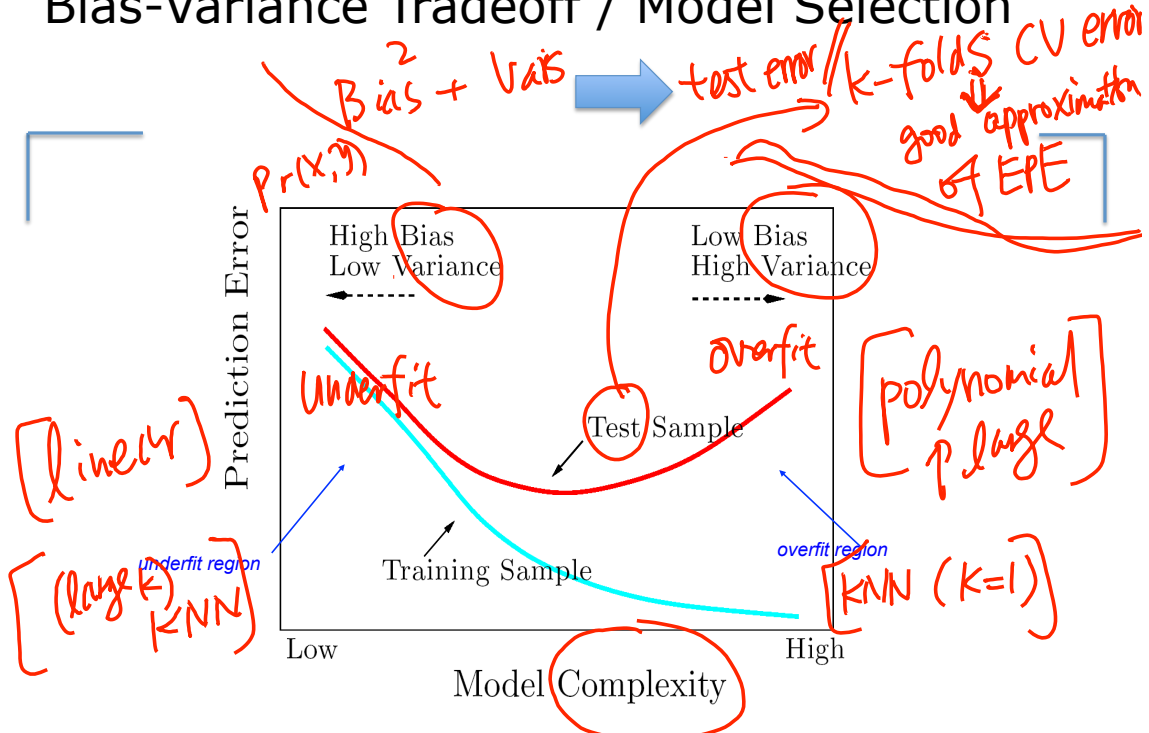
$$E(\text{MSE}) = \text{noise}^2 + \text{bias}^2 + \text{variance}$$



Bias-Variance Tradeoff / Model Selection

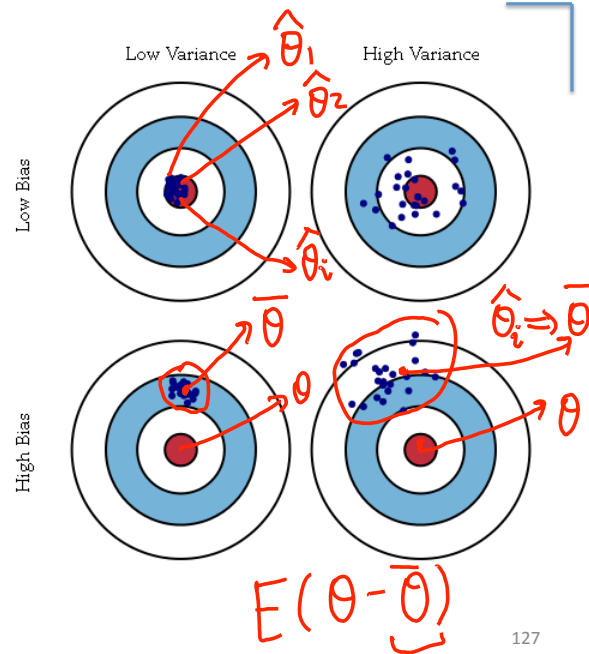


Bias-Variance Tradeoff / Model Selection



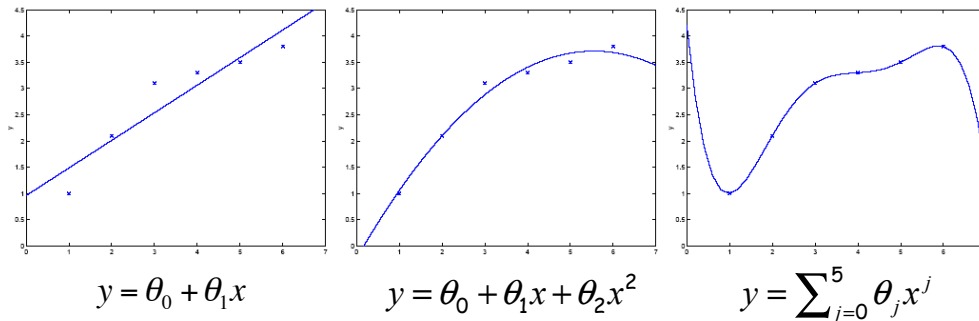
Model “bias” & Model “variance”

- 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



Which function f to choose?

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



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

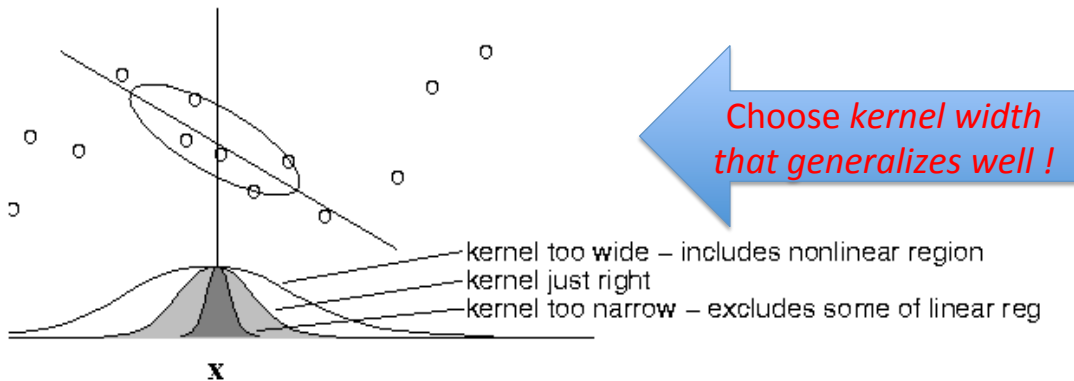
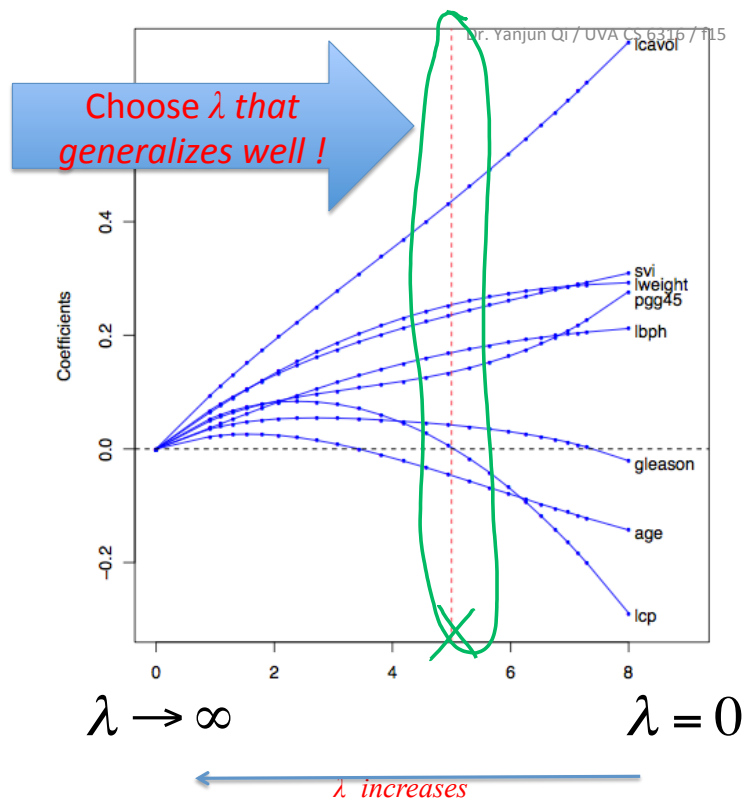


Figure 3: The estimator variance is minimized when the kernel includes as many training points as can be accommodated by the model. Here the linear LOESS model is shown. Too large a kernel includes points that degrade the fit; too small a kernel neglects points that increase confidence in the fit.

Which regularization parameter to choose ?

Ridge Regression

when varying λ ,
how θ_j varies.



Regularization path of a Lasso Estimator

Choose λ that generalizes well!

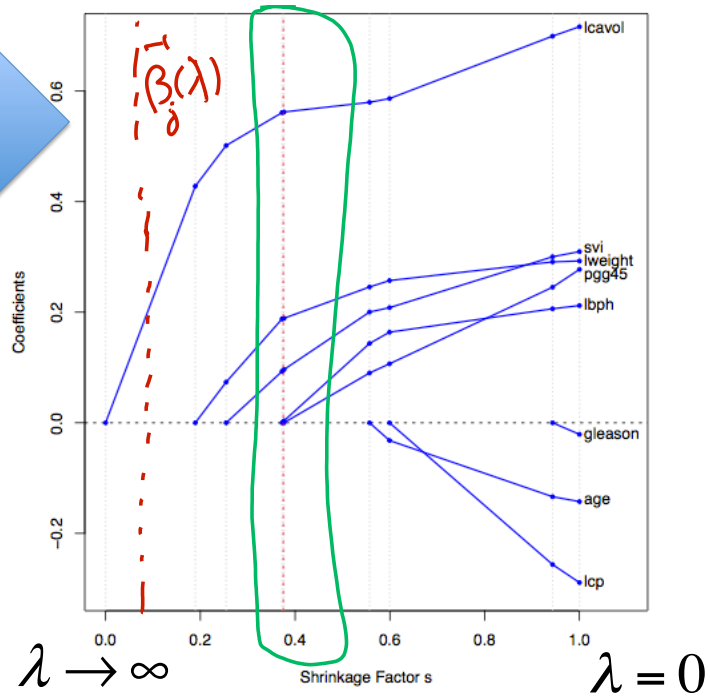


FIGURE 3.10. Profiles of lasso coefficients, as the tuning parameter t is varied. Coefficients are plotted versus $s = t / \sum_1^p |\hat{\beta}_j|$. A vertical line is drawn at $s = 0.36$, the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for 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.

11/16/15

need to make assumptions that are able to generalize

Dr. Yanjun Qi / UVA CS 6316 / f15

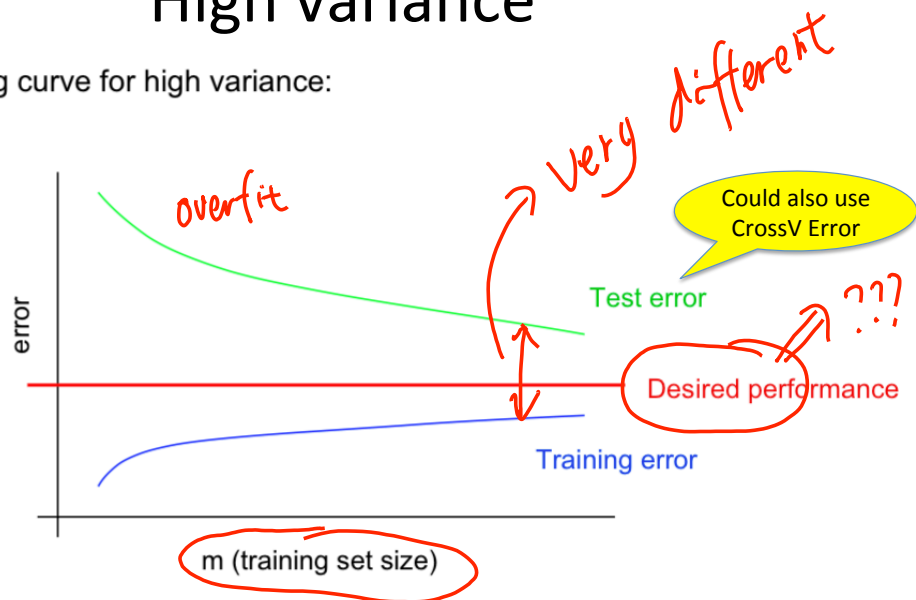
- Components of generalization error
 - **Bias:** how much the average model over all training sets differ from the true model?
 - Error due to inaccurate assumptions/simplifications made by the model
 - **Variance:** how much models estimated from different training sets differ from each other
- **Underfitting:** model is too “simple” to represent all the relevant class characteristics
 - High bias and low variance
 - High training error and high test error
- **Overfitting:** model is too “complex” and fits irrelevant characteristics (noise) in the data
 - Low bias and high variance
 - Low training error and high test error

11/16/15

132
Slide credit: L. Lazebnik

High variance

Typical learning curve for high variance:



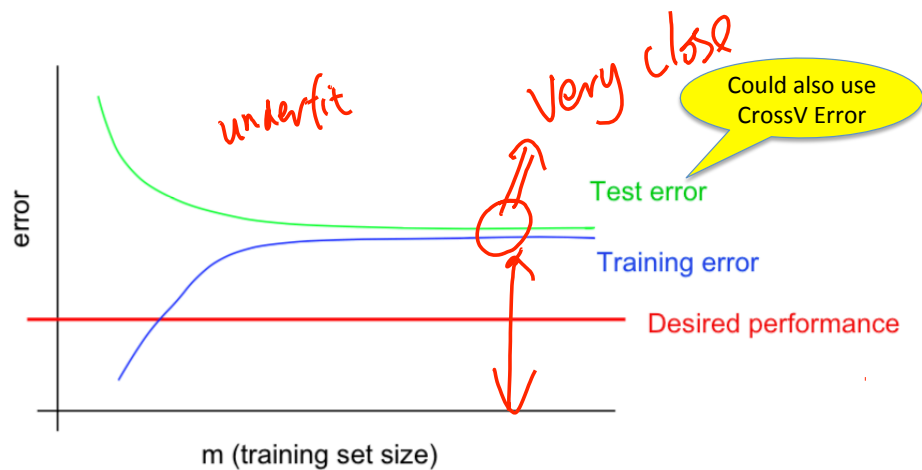
- Test error still decreasing as m increases. Suggests larger training set will help.
- Large gap between training and test error.
- **Low training error and high test error**

11/16/15

33
Slide credit: A. Ng

High bias

Typical learning curve for high bias:



- Even training error is unacceptably high.
- Small gap between training and test error.
- **High training error and high test error**

11/16/15

134
Slide credit: A. Ng

What we have covered for each component

Task	Regression, classification, clustering, dimen-reduction
Representation	Linear func, nonlinear function (e.g. polynomial expansion), local linear, logistic function (e.g. $p(c x)$), tree, multi-layer, prob-density family (e.g. Bernoulli, multinomial, Gaussian, mixture of Gaussians), local func smoothness, kernel matrix
Score Function	MSE, Hinge (margin), log-likelihood, EPE (e.g. L2 loss for KNN, 0-1 loss for Bayes classifier), cross-entropy, cluster points distance to centers, variance,
Search/ Optimization	Normal equation, gradient descent, stochastic GD, Newton, Linear programming, Quadratic programming (quadratic objective with linear constraints), greedy, EM, asyn-SGD, eigenDecomp
Models, Parameters	Regularization (e.g. L1, L2)

References

- Hastie, Trevor, et al. The elements of statistical learning. Vol. 2. No. 1. New York: Springer, 2009.
- Prof. M.A. Papalaskar's slides
- Prof. Andrew Ng's slides