

Statistical Learning Methods for Emerging Database Applications

Edward Chang
Associate Professor,
Electrical Engineering, UC Santa Barbara
CTO, VIMA Technologies

3/27/2003

DASFAA Tutorial, Kyoto

1

Useful Links

⌘ Related Publications

⊠ <http://www-db.stanford.edu/~echang/>

⌘ Software Free Trial

⊠ <http://www.imagebeagle.com>

⊠ Locate objectionable images on your hard drives

⊠ Before your boss finds it!!!

3/27/2003

DASFAA Tutorial, Kyoto

2

Outline

- ⌘ Statistical Learning
- ⌘ Emerging Applications Data Characteristics
- ⌘ Classical Models
- ⌘ Kernel Methods
 - ⊠ Linear Model View
 - ⊠ Nearest Neighbor View
 - ⊠ Geometric View
- ⌘ Dimension Reduction Methods

3/27/2003

DASFAA Tutorial, Kyoto

3

Statistical Learning

- ⌘ Program the computers to learn!
- ⌘ Computers improve performance with experience at some task
- ⌘ Example:
 - ⊠ Task: playing checkers
 - ⊠ Performance: % games it wins
 - ⊠ Experience: expert players

3/27/2003

DASFAA Tutorial, Kyoto

4

Statistical Learning

- ⌘ Task $\hat{Y} = f(U)$
 - ⊠ Represented by some model(s)
 - ⊠ Implies hypothesis
- ⌘ Performance
 - ⊠ Measured by error functions
- ⌘ Experience (L)
 - ⊠ Characterized by training data
- ⌘ Algorithm (Φ)

3/27/2003

DASFAA Tutorial, Kyoto

5

Supervised Learning

- ⌘ X: Data
 - ⊠ U: Unlabeled pool
 - ⊠ L: Labeled pool
- ⌘ G: Labels
 - ⊠ Regression
 - ⊠ Classification
- ⌘ Φ : Learning algorithm
- ⌘ $f = \Phi(L)$
- ⌘ $\hat{Y} = f(U)$

3/27/2003

DASFAA Tutorial, Kyoto

6

Learning Algorithms

- ⌘ Linear Model
- ⌘ K-NN
- ⌘ Neural Networks
- ⌘ Decision Trees
- ⌘ Kernel Methods
- ⌘ Etc.

3/27/2003

DASFAA Tutorial, Kyoto

7

Classical Model

- ⌘ N : Number of training instances
- ⌘ N^+ , N^-
- ⌘ D : Dimensionality
- ⌘ $N \gg D$ $N \rightarrow \infty$
 - ⊠ E.g., PAC learnability
- ⌘ $N^- \approx N^+$

3/27/2003

DASFAA Tutorial, Kyoto

8

Emerging DB Applications

- ⌘ $N < D$
- ⌘ $N^+ \ll N^-$
- ⌘ Examples
 - ⊠ Information Retrieval with relevance feedback
 - ⊠ Gene Profiling

3/27/2003

DASFAA Tutorial, Kyoto

9

Image Retrieval Demo

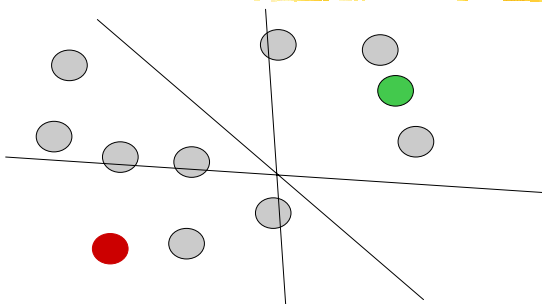
- ⌘ $N < D$
 - ⊠ $N < 50$
 - ⊠ $D = 150$
- ⌘ $N^+ \ll N^-$
- ⌘ ACM SIGMOD 01; ACM MM 01,02; IEEE CVPR 03

3/27/2003

DASFAA Tutorial, Kyoto

10

SVMactive

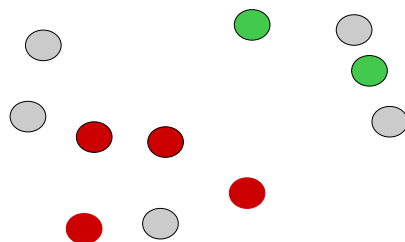


3/27/2003

DASFAA Tutorial, Kyoto

11

SVMactive

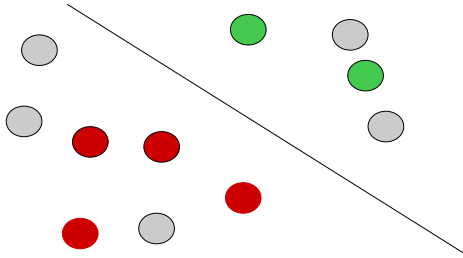


3/27/2003

DASFAA Tutorial, Kyoto

12

SVMactive

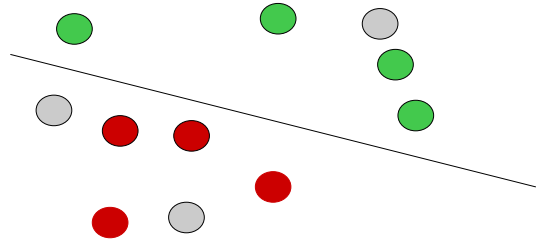


3/27/2003

DASFAA Tutorial, Kyoto

13

SVMactive

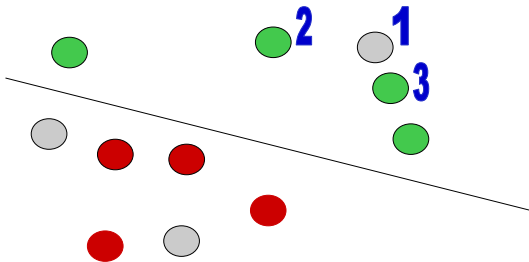


3/27/2003

DASFAA Tutorial, Kyoto

14

Ranking



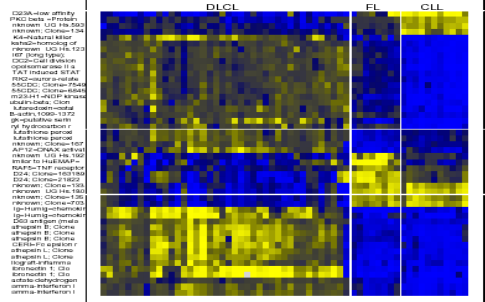
3/27/2003

DASFAA Tutorial, Kyoto

15

Gene Profiling Example

N = 59 cases, D = 4026 genes



3/27/2003

DASFAA Tutorial, Kyoto

16

Outline

- ⌘ Statistical Learning
- ⌘ Emerging Applications Data Characteristics
- ⌘ Classical Models (Classification)
- ⌘ Kernel Methods
 - ☒ Linear Model View
 - ☒ Nearest Neighbor View
 - ☒ Geometric View
- ⌘ Dimension Reduction Methods

3/27/2003

DASFAA Tutorial, Kyoto

17

Linear Model

$$\mathbb{Y} = \beta_0 + \sum \beta_j X_j \quad (j = 1 \text{ to } p)$$

$$\mathbb{Y} = X^T \beta$$

$$\mathbb{RSS}(\beta) = (y - X\beta)^T (y - X\beta)$$

☒ RSS: Residual Sum of Square

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

3/27/2003

DASFAA Tutorial, Kyoto

18

Linear Model

Elements of Statistical Learning @Hastings, Tibshirani & Friedman 2001. Chapter 3

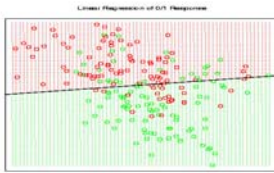


Figure 2.1: A classification example in two dimensions. The classes are coded as a binary variable—GREEN = 0, RED = 1—and then fit by linear regression. The line is the decision boundary defined by $x^T \beta = 0.5$. The red shaded region denotes that part of input space classified as RED, while the green region is classified as GREEN.

3/27/02

19

Maximum Likelihood

$$\mathbb{Y} = \beta_0 + \sum \beta_j X_j \quad (j = 1 \text{ to } p)$$

$$\mathbb{Y} = X^T \beta$$

$$\mathbb{Y} = X^T \beta + \epsilon$$

☒ ϵ (noise signals) are independent

☒ $\epsilon \rightarrow N(0, \sigma^2)$

☒ $P(y|\beta x)$ has a normal dist. with

☒ Mean at $y = \beta x$

☒ Variance σ^2

3/27/2003

DASFAA Tutorial, Kyoto

20

Linear Model

$$\mathbb{P}(y|\beta x) \rightarrow N(0, \sigma^2)$$

☒ Training

☒ Given $(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$

☒ Infer $P(\beta | x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n)$

☒ By Bayes rule, or

☒ Maximum Likelihood Estimate

3/27/2003

DASFAA Tutorial, Kyoto

21

Maximum Likelihood

☒ For what β is

☒ $P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n, \beta)$ maximized?

☒ $\prod P(y_i | \beta x_i)$ maximized?

☒ $\prod \exp(-\frac{1}{2}(y_i - \beta x_i / \sigma)^2)$ maximized?

☒ $\sum (-\frac{1}{2}(y_i - \beta x_i / \sigma)^2)$ maximized?

☒ $\sum (y_i - \beta x_i)^2$ minimized?

3/27/2003

DASFAA Tutorial, Kyoto

22

Least Square Linear Model

☒ Solution Method #1

$$\text{RSS}(\beta) = (y - X\beta)^T (y - X\beta)$$

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

☒ Solution Method #2 (for $D > N$)

☒ Gradient decent

☒ Perceptron

3/27/2003

DASFAA Tutorial, Kyoto

23

Other Linear Models

☒ LDA

☒ Find the projection direction which minimizes the overlap for two Gaussian distributions

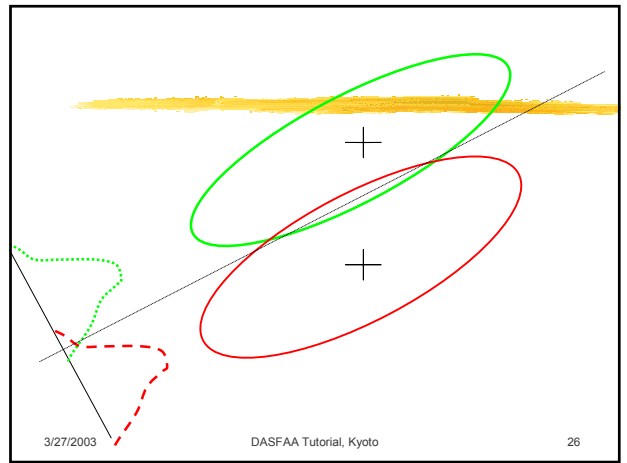
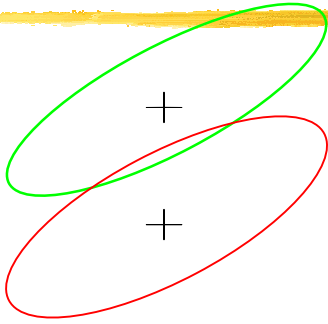
☒ Separating Hyperplane

3/27/2003

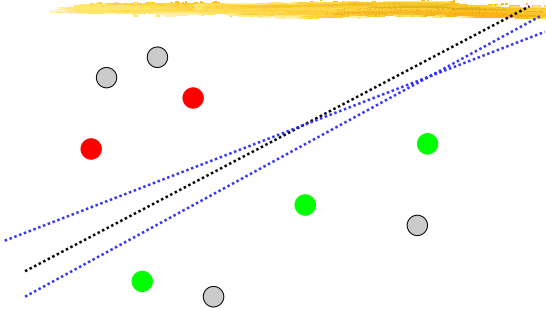
DASFAA Tutorial, Kyoto

24

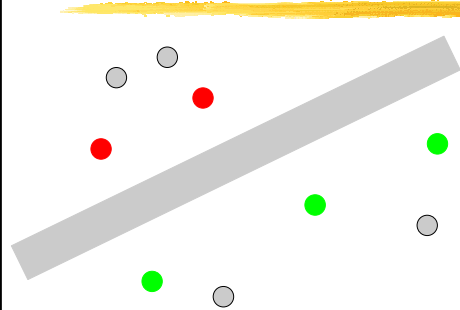
LDA



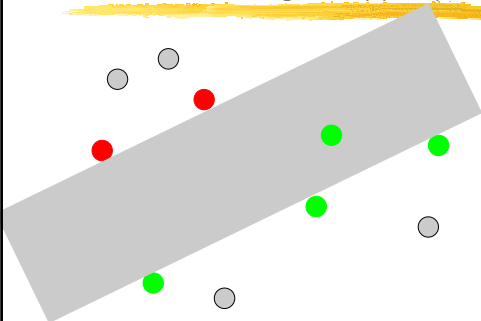
Separating Hyperplane



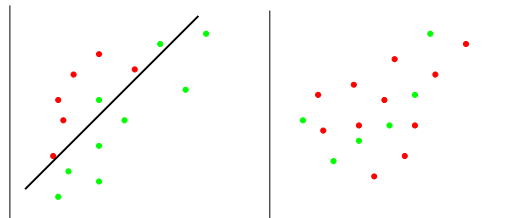
Separating Hyperplane



Maximum Margin Hyperplane



Linear Model Fits All Data?



How about Joining the Dots?

$$\mathbb{Y}(x) = 1/k \sum y_i,$$

$$\square x_i \in N_k(x)$$

$$\mathbb{K} = 1$$

3/27/2003

DASFAA Tutorial, Kyoto

31

Linear Models

$$\mathbb{N} \geq D$$

\square Least Square

\square LDA

$$\mathbb{D} > N$$

\square Perceptron

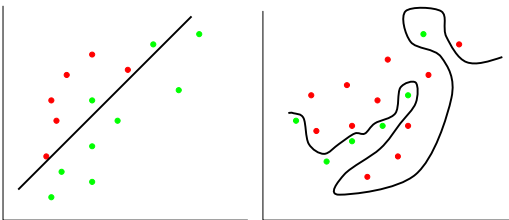
\square Maximum Hyperplane

3/27/2003

DASFAA Tutorial, Kyoto

32

Linear Model Fits All?



3/27/2003

DASFAA Tutorial, Kyoto

33

NN with k = 1

Elements of Statistical Learning @Harvard, MIT and Stanford 2001 Chapter 2

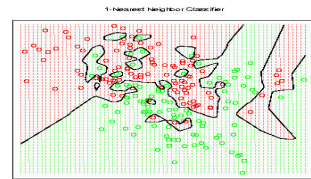


Figure 2.3: The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable (GREEN = 0, RED = 1), and then predicted by 1-nearest-neighbor classification. [C. Hastie, etc. 2001]

3/

Nearest Neighbor

\mathbb{F} Four Things Make a Memory Based Learner

\square A distance function

\square K: number of neighbors to consider?

\square A weighted function (optional)

\square How to fit with the local points?

3/27/2003

DASFAA Tutorial, Kyoto

35

Problems

\mathbb{F} Fitting Noise

\mathbb{F} Jagged Boundaries

3/27/2003

DASFAA Tutorial, Kyoto

36

Solutions

⌘ Fitting Noise

- ☒ Pick a Larger K ?

⌘ Jagged Boundaries

- ☒ Introducing Kernel as a weighting function

3/27/2003

DASFAA Tutorial, Kyoto

37

NN with $k = 15$

Elements of Statistical Learning @ Hastie, Tibshirani & Friedman 2001 Chapter 9

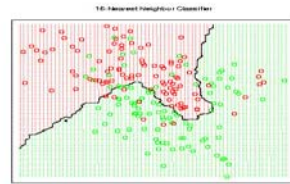


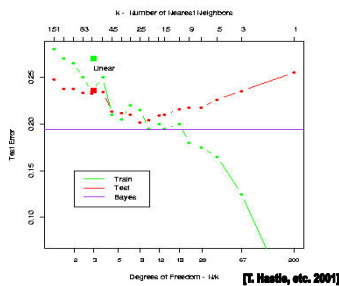
Figure 2.2: The same classification example in two dimensions as in Figure 2.1. The classes are coded as a binary variable ($\text{GREEN} = 0, \text{RED} = 1$) and then fit by 15-nearest-neighbor averaging as in (2.8). The predicted class is hence chosen by majority vote amongst the 15-nearest neighbors. [T. Hastie, etc. 2001]

3/27

3

NN

Elements of Statistical Learning @ Hastie, Tibshirani & Friedman 2001 Chapter 2



3

[T. Hastie, etc. 2001]

Solutions

⌘ Fitting Noise

- ☒ Pick a larger K ?

⌘ Jagged Boundaries

- ☒ Introducing Kernel as a weighting function

3/27/2003

DASFAA Tutorial, Kyoto

40

Nearest Neighbor -> Kernel Method

⌘ Four Things Make a Memory Based Learner

- ☒ A distance function
- ☒ K : number of neighbors to consider? All
- ☒ A weighted function: RBF kernels
- ☒ How to fit with the local points? Predict weights

3/27/2003

DASFAA Tutorial, Kyoto

41

Kernel Method

⌘ RBF Weighted Function

- ☒ Kernel width holds the key
- ☒ Use cross validation to find the "optimal" width

⌘ Fitting with the Local Points

- ☒ Where NN meets Linear Model

3/27/2003

DASFAA Tutorial, Kyoto

42

LM vs. NN

- ⌘ Linear Model
 - ⊠ $f(x)$ is approximated by a global linear function
 - ⊠ More stable, less flexible
- ⌘ Nearest Neighbor
 - ⊠ K-NN assumes $f(x)$ is well approximated by a locally constant function
 - ⊠ Less stable, more flexible
- ⌘ Between LM and NN
 - ⊠ The other models...

3/27/2003

DASFAA Tutorial, Kyoto

43

Decision Theories

- ⌘ Bias & Variance Tradeoff
- ⌘ Bayes Prediction
- ⌘ VC Dimensionality
- ⌘ PAC Learnability

3/27/2003

DASFAA Tutorial, Kyoto

44

Variance vs. Bias

- ⌘ $MSE(x_0) = E_T [f(x_0) - \hat{y}_0]^2$
 $= E_T [\hat{y}_0 - E_T(\hat{y}_0)]^2 + [E_T(\hat{y}_0) - f(x_0)]^2$
- ⌘ $Error = Var_T(\hat{y}_0) + Bias^2(\hat{y}_0)$

3/27/2003

DASFAA Tutorial, Kyoto

45

Outline

- ⌘ Statistical Learning
- ⌘ Emerging Applications Data Characteristics
- ⌘ Classical Models (Classification)
- ⌘ Kernel Methods
- ⌘ Dimension Reduction Methods

3/27/2003

DASFAA Tutorial, Kyoto

46

Where Are We and Where Am I Heading To ?

- ⌘ LM and NN
- ⌘ Kernel Method of Three Views
 - ⊠ LM view
 - ⊠ NN view
 - ⊠ Geometric view

3/27/2003

DASFAA Tutorial, Kyoto

47

Linear Model View

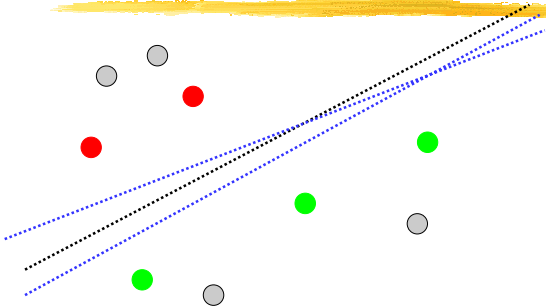
- ⌘ $Y = \beta_0 + \sum \beta X$
- ⌘ Separating Hyperplane
 - ⊠ $\text{Max}_{i \in \{1, \dots, l\}} C$
 - ⊠ Subject to $y_i f(x_i) \geq C$, or
 - ⊠ $y_i (\beta_0 + \beta x_i) \geq C$

3/27/2003

DASFAA Tutorial, Kyoto

48

Separating Hyperplane

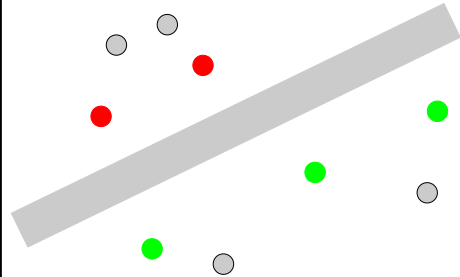


3/27/2003

DASFAA Tutorial, Kyoto

49

Separating Hyperplane

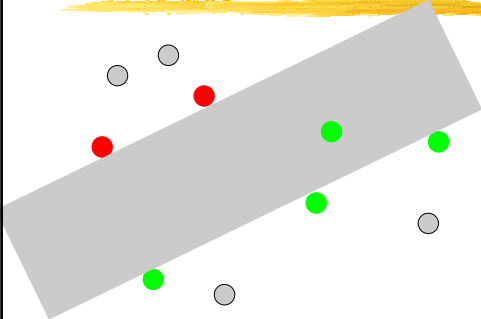


3/27/2003

DASFAA Tutorial, Kyoto

50

Maximum Margin Hyperplane



3/27/2003

DASFAA Tutorial, Kyoto

51

Classifier Margin

⌘ Margin

- ☑ Defined as width of the boundary before hitting a data object

⌘ Maximum Margin

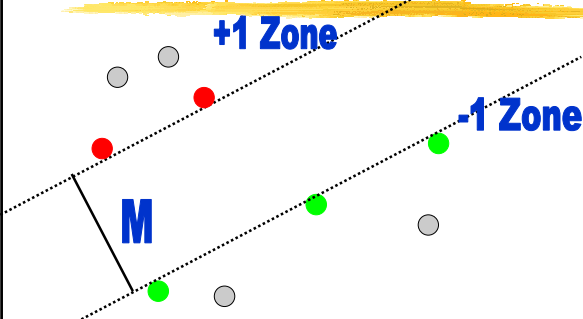
- ☑ Tends to minimize classification variance
- ☑ No formal theory for this yet

3/27/2003

DASFAA Tutorial, Kyoto

52

Separating Hyperplane



3/27/2003

DASFAA Tutorial, Kyoto

53

M's Mathematical Representation

⌘ Plus-plane

- ☑ $\{x: wx+b = +1\}$

⌘ Minus-plane

- ☑ $\{x: wx+b = -1\}$

⌘ $w \perp$ Plus-plane

- ☑ $w(u - v) = 0$, if u and v on plus-plane

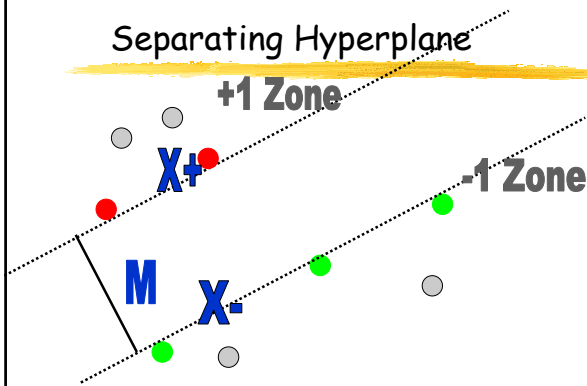
⌘ $w \perp$ Minus-plane

3/27/2003

DASFAA Tutorial, Kyoto

54

Separating Hyperplane



3/27/2003

DASFAA Tutorial, Kyoto

55

M

- ⌘ Let x^- be any point on minus-plane
- ⌘ Let x^+ be the closest plus-plane-point to x^-
- ⌘ $x^+ = x^- + \lambda w$, why
 - ⊠ The line $(x^+x^-) \perp$ minus-plane
- ⌘ $M = |x^+ - x^-|$

3/27/2003

DASFAA Tutorial, Kyoto

56

M

1. $w x^- + b = -1$
2. $w x^+ + b = 1$
3. $x^+ = x^- + \lambda w$
4. $M = |x^+ - x^-|$
5. $w(x^- + \lambda w) + b = 1$ (from 2 & 3)
6. $w x^- + b + \lambda w w = 1$
7. $\lambda w w = 2$

3/27/2003

DASFAA Tutorial, Kyoto

57

M

1. $\lambda w w = 2$
2. $\lambda = 2 / w w$
3. $M = |x^+ - x^-| = |\lambda w| = \lambda |w| = 2 / |w|$
4. Max M
 - ⌘ Gradient decent, simulated annealing, EM, Newton's method?

3/27/2003

DASFAA Tutorial, Kyoto

58

Max M

- ⌘ Max $M = 2 / |w|$
- ⌘ Min $|w| / 2$
- ⌘ Min $|w|^2 / 2$
 - ⊠ subject to $y_i(x_i w + b) \geq 1$
 - ⊠ $i = 1, \dots, N$
- ⌘ Quadratic criterion with linear inequality constraints

3/27/2003

DASFAA Tutorial, Kyoto

59

Max M

- ⌘ Min $|w|^2 / 2$
 - ⊠ subject to $y_i(x_i w + b) \geq 1$
 - ⊠ $i = 1, \dots, N$
- ⌘ $L_p = \min_{w,b} |w|^2 / 2 + \sum_{i=1..N} \alpha_i [y_i(x_i w + b) - 1]$
- ⌘ $w = \sum_{i=1..N} \alpha_i y_i x_i$
- ⌘ $0 = \sum_{i=1..N} \alpha_i y_i$

3/27/2003

DASFAA Tutorial, Kyoto

60

Wolfe Dual

$$\text{⌘ } Ld = \sum_{i=1..N} \alpha_i - 1/2 \sum_{i,j=1..N} \alpha_i \alpha_j y_i y_j x_i \cdot x_j$$

⌘ Subject to

$$\text{⊠ } \alpha_i \geq 0$$

$$\text{⊠ } \alpha_i [y_i(x_i \cdot w + b) - 1] = 0$$

⊠ KKT conditions

$$\text{⊠ } \alpha_i > 0, y_i(x_i \cdot w + b) = 1 \text{ (Support Vectors)}$$

$$\text{⊠ } \alpha_i = 0, y_i(x_i \cdot w + b) > 1$$

Class Prediction

$$\text{⌘ } y_q = w \cdot x_q + b$$

$$\text{⌘ } w = \sum_{i=1..N} \alpha_i y_i x_i$$

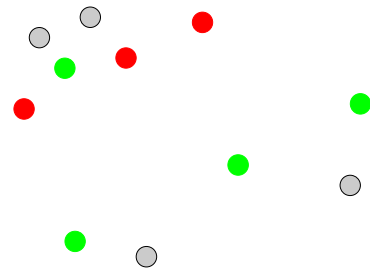
$$\text{⌘ } y_q = \text{sign}(\sum_{i=1..N} \alpha_i y_i (x_i \cdot x_q) + b)$$

Non-seperatable Classes

⌘ Soft Margin Hyperplane

⌘ Basis Expansion

Non-separating Case



Soft Margin SVMs

⌘ Min $|w|^2/2$

⊠ subject to $y_i(x_i \cdot w + b) \geq 1$

⊠ $i = 1, \dots, N$

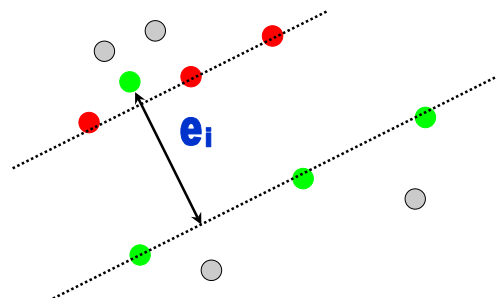
⌘ Min $|w|^2/2 + C \sum \epsilon_i$

⊠ $x_i \cdot w + b \geq 1 - \epsilon_i$ if $y_i = 1$

⊠ $x_i \cdot w + b \leq -1 + \epsilon_i$ if $y_i = -1$

⊠ $\epsilon_i \geq 0$

Non-separating Case



Wolfe Dual

$$\text{⌘ } Ld = \sum_{i=1..N} \alpha_i - 1/2 \sum_{i,j=1..N} \alpha_i \alpha_j \gamma_i \gamma_j x_i x_j$$

⌘ Subject to

$$\text{⊠ } C \geq \alpha_i \geq 0$$

$$\text{⊠ } \sum \alpha_i \gamma_i = 0$$

⊠ KKT conditions

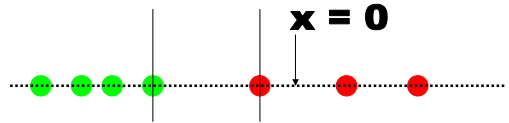
$$\text{⌘ } \gamma_i = \text{sign}(\sum_{i=1..N} \alpha_i \gamma_i (x_i - x_q) + b)$$

3/27/2003

DASFAA Tutorial, Kyoto

67

Basis Function

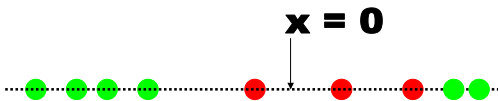


3/27/2003

DASFAA Tutorial, Kyoto

68

Harder 1D Example



3/27/2003

DASFAA Tutorial, Kyoto

69

Basis Function

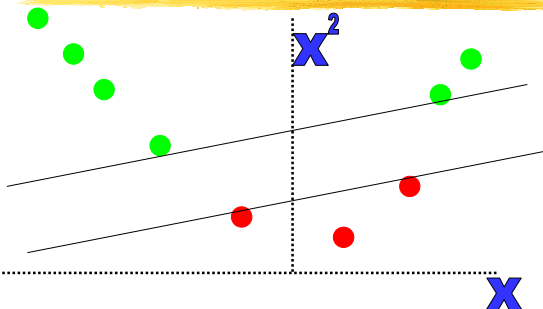
$$\text{⌘ } \Phi(X) = (x, x^2)$$

3/27/2003

DASFAA Tutorial, Kyoto

70

Harder 1D Example



3/27/2003

DASFAA Tutorial, Kyoto

71

Some Basis Functions

$$\text{⌘ } \Phi(X) = \sum \gamma_m h_m(X)$$

$$\text{⊠ } h_m(X) \mathbb{R}^p \rightarrow \mathbb{R}$$

⌘ Common Functions

⊠ Polynomial

⊠ Radial basis functions

⊠ Sigmoid functions

3/27/2003

DASFAA Tutorial, Kyoto

72

Wolfe Dual

$$\text{⌘ } L_d = \sum_{i=1..N} \alpha_i - 1/2 \sum_{i,j=1..N} \alpha_i \alpha_j y_i y_j \Phi(x_i) \Phi(x_j)$$

⌘ Subject to

$$\text{⊠ } C \geq \alpha_i \geq 0$$

$$\text{⊠ } \sum \alpha_i y_i = 0$$

⊠ KKT conditions

$$\text{⌘ } y_q = \text{sign} \left(\sum_{i=1..N} \alpha_i y_i (\Phi(x_i) \cdot \Phi(x_q)) + b \right)$$

$$\text{⌘ } K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$$

⊠ Kernel function!

3/27/2003

DASFAA Tutorial, Kyoto

73

Quadratic Basis Functions

$$\text{⌘ } \Phi(X) = \{1, x_i, x_i x_j\}, ij = 1..p$$

⊠ (p+1)(p+2) terms

⊠ P² terms

⊠ O(P²) computational cost

⌘ It is equivalent to $(x_i x_j + 1)^2$

⊠ O(p) computational cost

⌘ Total Cost

⊠ O(N²p)

3/27/2003

DASFAA Tutorial, Kyoto

74

Dot Product Saves the Day

$$\text{⌘ } O(N^2 p)$$

⌘ Quadratic

$$\text{⊠ } O(N^2 p^2)$$

⌘ Cubic

$$\text{⊠ } O(N^2 p^3)$$

⌘ Quartic

$$\text{⊠ } O(N^2 p^4)$$

3/27/2003

DASFAA Tutorial, Kyoto

75

Quiz

⌘ What is a polynomial kernel degree d function's signature?

$$\text{⌘ } (x_i x_j + 1)^d$$

3/27/2003

DASFAA Tutorial, Kyoto

76

Nearest Neighbor View

⌘ Z , a set of zero mean jointly Gaussian random variables,

⊠ Each Z_i corresponds to one example X_i

$$\text{⊠ } \text{Cov}(z_i, z_j) = K(x_i, x_j)$$

⌘ y_i , the label of z_i , +1 or -1

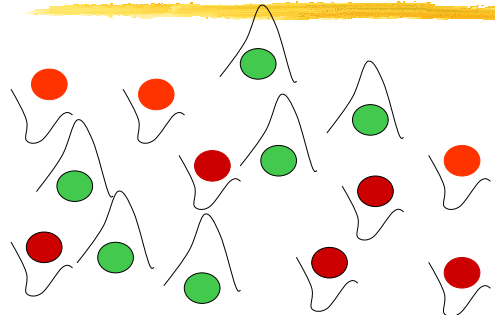
$$\text{⊠ } P(y_i | z_i) = \sigma(y_i, z_i)$$

3/27/2003

DASFAA Tutorial, Kyoto

77

Training Data



3/27/2003

DASFAA Tutorial, Kyoto

78

General Kernel Classifier [Jaakkola, etc. 99]

⌘ MAP Classification for x_t

- ⊠ $y_t = \text{sign}(\sum \alpha_i y_i K(x_t, x_i))$
- ⊠ $K(x_i, x_j) = \text{Cov}(z_i, z_j)$ (some similarity function)

⌘ Supervised Training: Compute α_i

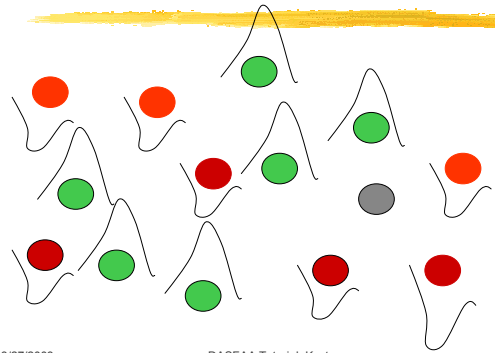
- ⊠ Given X and y , and
- ⊠ An error function such as
$$J(\alpha) = -\frac{1}{2} \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j) + \sum F(\alpha_i)$$

3/27/2003

DASFAA Tutorial, Kyoto

79

Leave One Out



3/27/2003

DASFAA Tutorial, Kyoto

80

SVMs

$$\text{⌘ } y_t = \text{sign}(\sum \alpha_i y_i K(x_t, x_i))$$

⌘ (y_i, x_i) training data, α_i nonnegative, and kernel K positive definite

⌘ α_i is obtained by maximizing

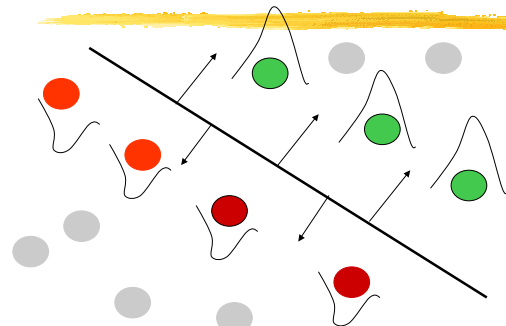
- ⊠ $J(\alpha) = -\frac{1}{2} \sum \alpha_i \alpha_j y_i y_j K(x_i, x_j) + \sum F(\alpha_i)$
- ⊠ $F(\alpha_i) = \alpha_i$
- ⊠ $\alpha_i \geq 0, \sum y_i \alpha_i = 0$

3/27/2003

DASFAA Tutorial, Kyoto

81

SVMs



3/27/2003

DASFAA Tutorial, Kyoto

82

Important Insight

$$\text{⌘ } K(x_i, x_j) = \text{Cov}(z_i, z_j)$$

⌘ To design of a kernel is to design a similarity function that produces a **positive definite** covariance matrix on the training instances

3/27/2003

DASFAA Tutorial, Kyoto

83

Basis Function Selection

⌘ Three General Approaches

- ⊠ Restriction methods
 - ⊠ Limit the class of functions
- ⊠ Selection methods
 - ⊠ Scan the dictionary adaptively (Boosting)
- ⊠ Regularization methods
 - ⊠ Use the entire dictionary but restrict coefficients (Ridge Regression)

3/27/2003

DASFAA Tutorial, Kyoto

84

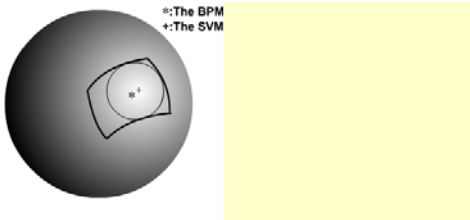
Overfitting?

- ⌘ Probably Not
- ⌘ Because
 - ☑ N free parameters (not D)
 - ☑ Maximizing margin

Geometrical View

- ⌘ $S = w \cdot X + b$
- ⌘ $|w| = 1, b = 0$
- ⌘ $V = \{w \mid \exists i f(x_i) > 0; i = 1..n, |w| = 1\}$
- ⌘ SVM is the center of the largest sphere contained in V

SVMs



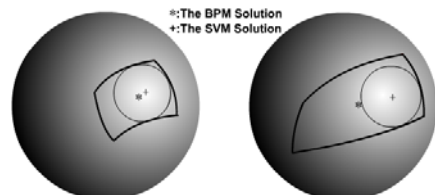
BPMs

- ⌘ Bayes Objective Function
 - ☑ $\hat{S}_t = \text{Bayes}_Z(X_t)$
 - $= \text{argmin}_{S_i \in S} E_{H|Z=x} [l(H(x), S_i)]$
- ⌘ BPMs [Herbrich, etc. 2001]
 - ☑ $A_{bp} = \text{argmin}_{h \in H} E_x [E_{H|Z=x} [l(H(x), h(x))]]$

BPMs

- ⌘ Linear Classifier
- ⌘ Input X Posses Spherical Gaussian Density
- ⌘ BP is the Center of Mass of the Version Space

BPMs vs. SVMs



BPMs

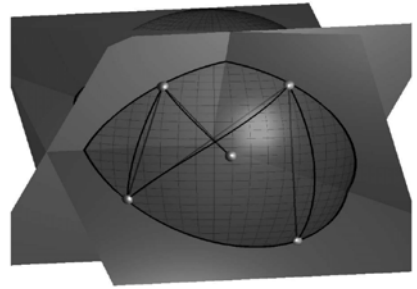
- ⌘ Use SVMs to find a good h in H
- ⌘ Find the BP
 - ⊠ Billiard Algorithm [Herbrich, etc. 2001]
 - ⊠ Perceptron Algorithm [Herbrich, etc. 2001]

3/27/2003

DASFAA Tutorial, Kyoto

91

Billiard Ball Algorithm (R. Herbrich)



3/27/2003

DASFAA Tutorial, Kyoto

92

Outline

- ⌘ Statistical Learning
- ⌘ Emerging Applications Data Characteristics
- ⌘ Classical Models (Classification)
- ⌘ Kernel Methods
- ⌘ Dimension Reduction Methods

3/27/2003

DASFAA Tutorial, Kyoto

93

Dimensionality Curse

- ⌘ D: Data Dimension
- ⌘ When D increases
 - ⊠ Nearest neighbors are not local
 - ⊠ All points are equally distanced

3/27/2003

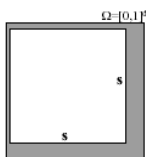
DASFAA Tutorial, Kyoto

94

Sparse High-D Space

[C. Aggarwal, etc. ICDT 2001]

⌘ Hyper-cube Range Queries



$$P^d[s] = s^d$$

3/27/2003

DASFAA Tutorial, Kyoto

95

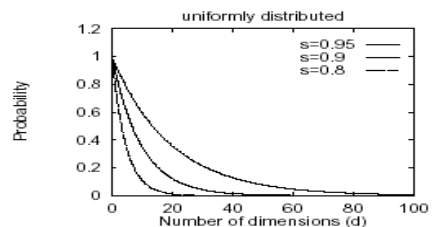


Figure 2: The probability function $P^d[s]$.

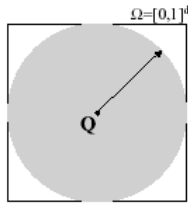
3/27/2003

DASFAA Tutorial, Kyoto

96

Sparse High-D Space

⌘ Spherical Range Queries

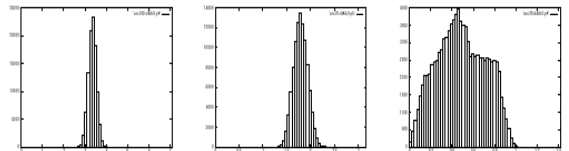


$$P[R \in sp^d(Q, 0.5)] = \frac{\sqrt{\pi^d} \cdot (0.5)^d}{\Gamma(\frac{d}{2} + 1)}$$

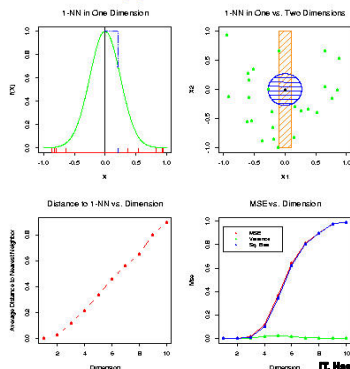
d	$P[R \in sp^d(Q, 0.5)]$	$N(d)$
2	0.785	1.273
4	0.308	3.242
10	0.002	401.5
20	$2.461 \cdot 10^{-8}$	$40^6 631^6 627$
40	$3.278 \cdot 10^{-21}$	$3.050 \cdot 10^{60}$
100	$1.868 \cdot 10^{-70}$	$5.353 \cdot 10^{69}$

Table 2: Probability that a point lies within the largest range query inside Ω , and the expected database size.

Dimensionality Curse



(a) 50 Dimensions (b) 10 Dimensions (c) 2 Dimensions
Figure 3: Distance Distribution of Uniform Data



So?

⌘ Is nearest neighbor *estimate* cursed in high-D spaces?

- Yes!
- When D is large and N is relatively small, the estimate is off!!

Are We Doomed?

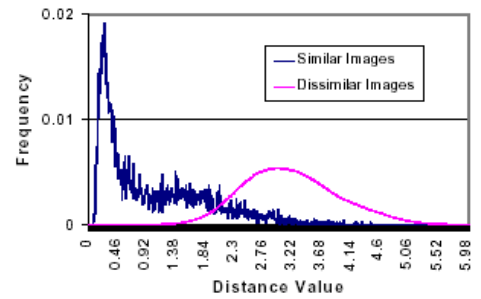
- ⌘ How does the curse affect classification?
- ⌘ Similar objects tend to cluster together
- ⌘ Classification makes binary prediction

3/27/2003

DASFAA Tutorial, Kyoto

103

Distribution of Distances



(a) $m = 144$

3/

34

Some Solutions to High-D

- ⌘ Restricted Estimators
 - ⊠ Specifying the nature of local neighborhood
- ⌘ Adaptive Feature Reduction
 - ⊠ PCA, LDA
- ⌘ Dynamic Partial Function

3/27/2003

DASFAA Tutorial, Kyoto

105

Three Major Paradigms

- ⌘ Preserve data description in a lower dimensional space
 - ⊠ PCA
- ⌘ Maximize discriminability in a lower dimensional space
 - ⊠ LDA
- ⌘ Activate only similar channels
 - ⊠ DPF

3/27/2003

DASFAA Tutorial, Kyoto

106

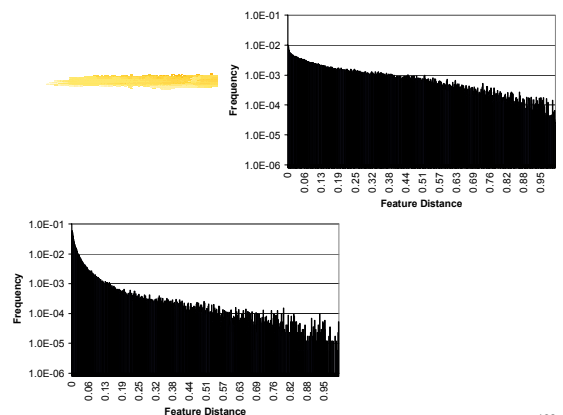
Minkowski Distance

- ⌘ Objects P and Q
- ⌘ $D = (\sum_M (p_i - q_i)^n)^{1/n}$
- ⌘ Similar images are similar in all M features

3/27/2003

DASFAA Tutorial, Kyoto

107

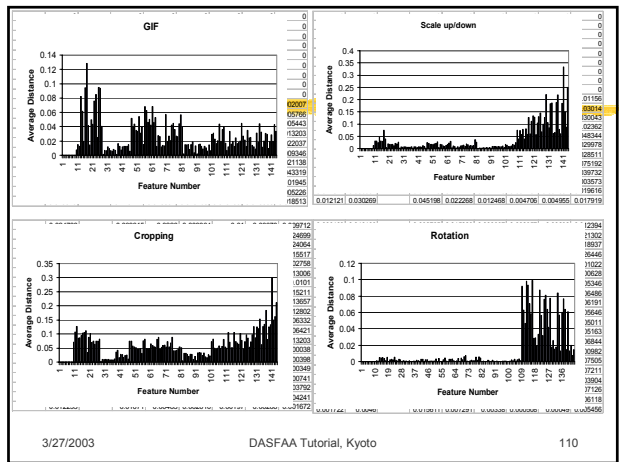


108

Weighted Minkowski Distance

$$D = (\sum_M w_i (p_i - q_i)^n)^{1/n}$$

Similar images are similar in the same subset of the M features



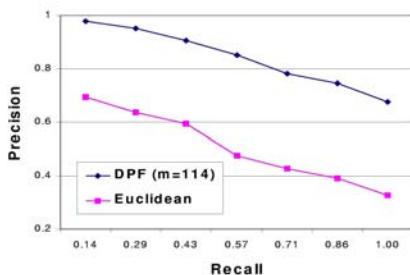
Similarity Theories

- Objects are similar in all respects (Richardson 1928)
- Objects are similar in some respects (Tversky 1977)
- Similarity is a process of determining respects, rather than using predefined respects (Goldstone 94)

DPF

- Which Place is Similar to Kyoto?
- Partial
- Dynamic
- Dynamic Partial Function

Precision/Recall



Summary

- Statistical Learning
- Emerging Applications Data Characteristics
- Classical Models (Classification)
 - Kernel Methods
 - Linear Model View
 - Nearest Neighbor View
 - Geometric View
- Dimension Reduction Methods

Emerging DB Applications

⌘ $N < D$

⌘ $N^+ \ll N^-$

⌘ Examples

- ☒ Information Retrieval with relevance feedback
- ☒ Gene Profiling
- ☒ Bioinformatics

Useful Links

⌘ Related Publications

☒ <http://www-db.stanford.edu/~echang/>

⌘ Software Free Trial

☒ <http://www.imagebeagle.com>

☒ Locate objectionable images on your hard drives

☒ Before your boss finds it!!!

References

1. The Elements of Statistical Learning, T. Hastie, R. Tibshirani, and J. Friedman, Springer, N.Y., 2001
2. Machine Learning, T. Mitchell, 1997
3. High-dimensional Data Analysis, D. Donoho, American Math. Society Lecture, 2000
4. Support Vector Machine Active Learning for Image Retrieval, S. Tong and E. Chang, ACM MM, 2001
5. Dynamic Partial Function, B. Li and E. Chang, ACM Journal, 2003
6. Pattern Discovery in Sequences under a Markov Assumption, D. Chudova and P. Smyth, ACM KDD 2002
7. Bayes Point Machines, R. Herbrich, T. Graepel and C. Campbell, Journal of Machine Learning Research, 2001
8. The Nature of Statistical Learning Theory, V. Vapnik, Springer, N.Y., 1995
9. Probabilistic Kernel Regression Models, T. Jaakkola and D. Haussler, Conference of AI and Statistics, 1999
10. Support Vector Machines, Lecture Notes, A. Moore, CMU
11. On the Surprising Behavior of Distance Metrics in High-dimensional Space, C. Aggarwal, A. Hinneburg, and D. Keim, ICDD 2001