

Course ID 80250993 Dates: 9/16-12/23/2021 @ I-205
Tencent/VooV Meeting 521 4678 6257 (passwd 1205)



Chapter 4

Multi-category Classification and Nonlinear Classification

Xuegong Zhang
September 30, 2021



Xuegong Zhang

1



4.1

Multi-category Linear Classifiers

Xuegong Zhang

2

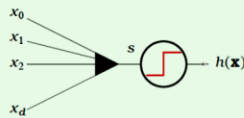
Recall: Three major types of linear machines



$$s = \sum_{i=0}^d w_i x_i$$

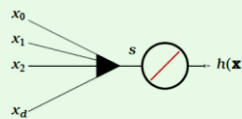
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



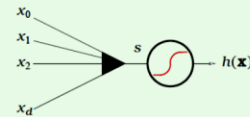
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

$$h(\mathbf{x}) = \theta(s)$$



Abu-Mostafa, Magdon-Ismael, Lin, *Learning from Data*, Lecture 9

→ Binary Classification

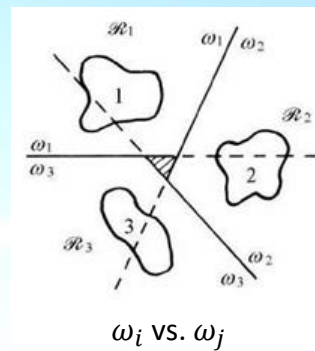
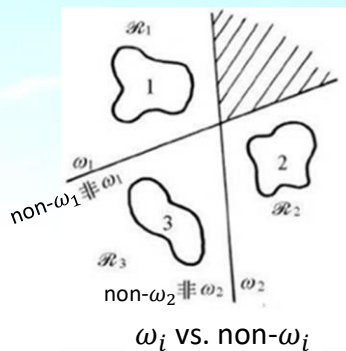
Xuegong Zhang

3

Multi-category classification



- Strategies for multi-class classification using binary methods
 - One-vs-rest, aka one-vs-all, one-over-all (OVA)
 - Pairwise classification



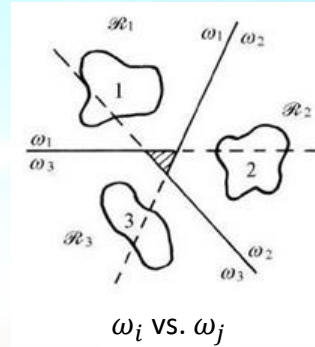
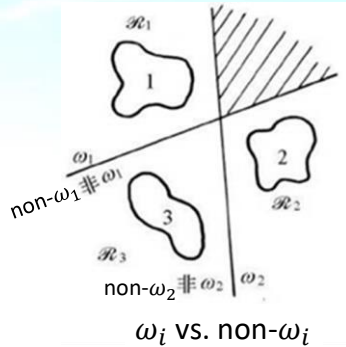
Xuegong Zhang

4



Challenges

- Need a “code-book” to index the binary classifiers
- Need comprehensive voting schemes to make decisions based on multiple binary classifiers
- There could be “undefined” voting results



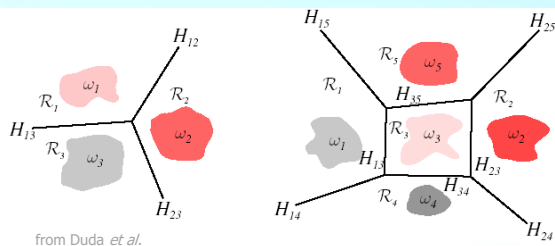
Xuegong Zhang

5

Multi-category Linear Discriminants



- Given C categories, define C discriminant functions $g_i(\mathbf{x}) = \mathbf{a}_i^T \mathbf{y}$
- Classify \mathbf{x} as a member of c_i if $g_i(\mathbf{x}) > g_j(\mathbf{x})$ for all $j \neq i$



Algorithm:

If $\mathbf{y}^k \in \omega_i$ but $\mathbf{a}_i(k)^T \mathbf{y}^k \leq \mathbf{a}_j(k)^T \mathbf{y}^k, j \neq i$

$$\text{then } \begin{cases} a_i(k+1) = a_i(k) + \rho_k y^k \\ a_j(k+1) = a_j(k) - \rho_k y^k \\ a_l(k+1) = a_l(k), l \neq i, j \end{cases}$$

——“Multi-class linear machine”

- Converge when data are linearly separable

Xuegong Zhang

6



Multiclass Fisher's Linear Discriminant



- a.k.a. Multi-class LDA
 - Within-class scatter matrix $S_w = \sum_k \sum_{x \in \omega_i} (x - m_k)(x - m_k)^T$
 - Between-class scatter matrix $S_b = \sum_k n_k (m_k - m)(m_k - m)^T$
- $$W: \max J(W) = \frac{|W^T S_b W|}{|W^T S_w W|} \quad \text{or} \quad \hat{w} = \operatorname{argmax}_w \frac{w^T S_b w}{w^T S_w w}$$
- Solution: the $\leq c - 1$ eigenvector solutions of $S_b w = \lambda S_w w$
 — reduction to $\leq c - 1$ dimensions, but not solving the classifiers

Xuegong Zhang

7



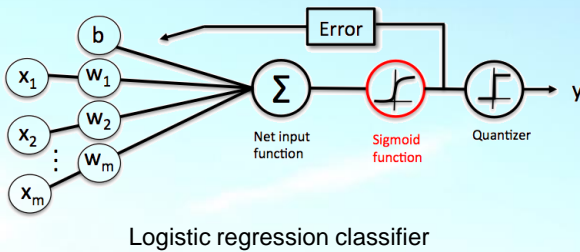
4.2 Multi-category Logistic Regression and Softmax

Xuegong Zhang

8



Binary logistic regression classifier



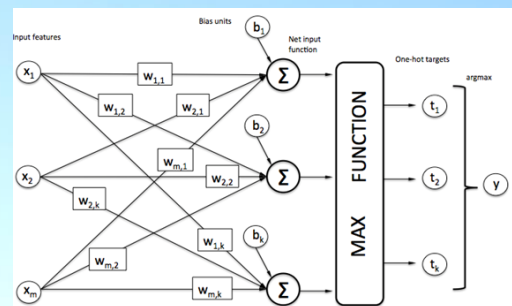
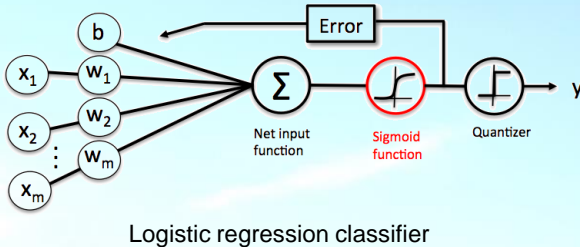
- Binary classification:
 - Weighted sum \rightarrow Logistic function \rightarrow Compare with threshold \rightarrow Classification

Xuegong Zhang

Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

9

Multiclass classification with logistic regression



- Binary classification:
 - Weighted sum \rightarrow Logistic function \rightarrow Compare with threshold \rightarrow Classification
- Multiclass classification:
 - Weighted sum \rightarrow Compare among peers (Max function) \rightarrow Classification

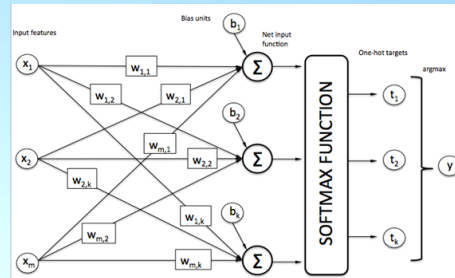
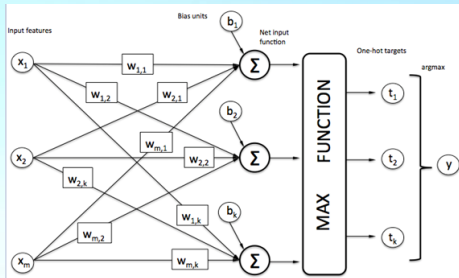
Xuegong Zhang

Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

10



Multiclass logistic regression and SoftMax



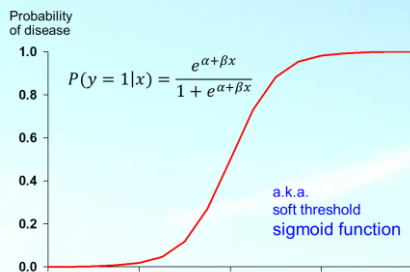
- Problem: no consideration on competition among class probabilities
- Solution: Max functions → SoftMax function

Xuegong Zhang

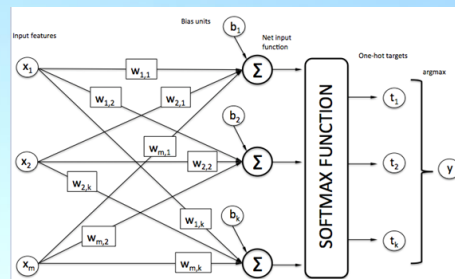
Modified from https://rasbt.github.io/mlxtend/user_guide/classifier/SoftmaxRegression/

11

SoftMax (Normalized exponential function)



$$P(y = 1|x) = \frac{e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}{1 + e^{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p}}$$



$$P(y = j|x) = \frac{e^{w_j \cdot x}}{\sum_{k=1}^K e^{w_k \cdot x}}, \quad j = 1, \dots, K$$

--- SoftMax function

Warning: many random names in ML 😊

Xuegong Zhang

12



4.3 Quadratic Discriminant

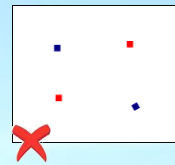
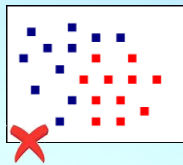
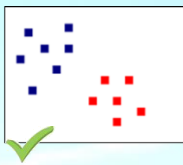
Xuegong Zhang

13

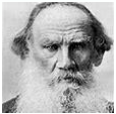
The need for nonlinear classifiers



- Linearly non-separable cases \rightarrow The need for nonlinear classifiers



- Linearly non-separable \neq nonlinearly separable
- “nonlinearities” are not the same



Leo Tolstoy:

“All happy families resemble one another, but each un-happy family is unhappy in its own way.”

Xuegong Zhang

14



Quadratic discriminant analysis (QDA)

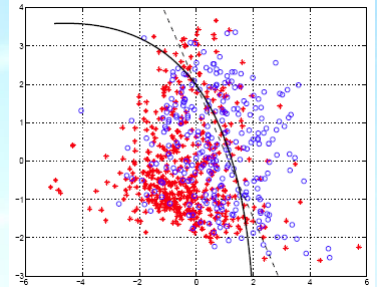
- Quadratic functions: might be the simplest nonlinear functions
- Quadratic discriminant:

$$g_i(\mathbf{x}) = k_i^2 - (\mathbf{x} - \hat{\mathbf{m}}_i)^T \hat{\Sigma}_i^{-1} (\mathbf{x} - \hat{\mathbf{m}}_i), \quad i = 1, \dots, C$$

$$\text{mean } \hat{\mathbf{m}}_i = \frac{1}{N_i} \sum_{j=1}^{N_i} \mathbf{x}_j$$

$$\text{covariance matrix } \hat{\Sigma}_i = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (\mathbf{x}_j - \mathbf{m}_i)(\mathbf{x}_j - \mathbf{m}_i)^T$$

- Decision: $\text{class}(\mathbf{x}) = \underset{i}{\operatorname{argmax}} g_i(\mathbf{x})$
- When to choose QDA (instead of LDA)?
 - Large training data set, roughly normal distribution
 - Covariance matrixes different between classes



<https://online.stat.psu.edu/stat508/lesson/9/9.2/9.2.8>

Xuegong Zhang

15

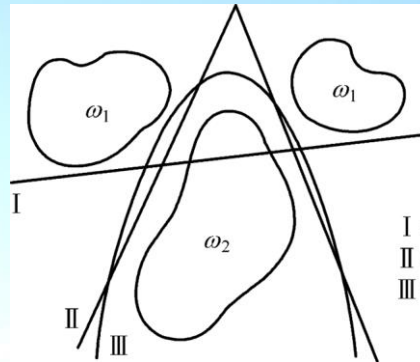
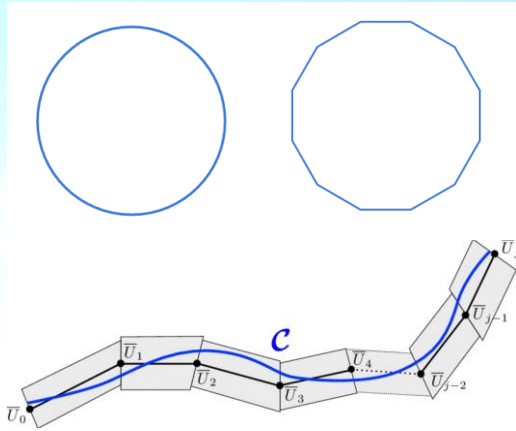


4.4 Piecewise Linear Classifiers

Xuegong Zhang

16

Piecewise linearity to approximate any nonlinearity

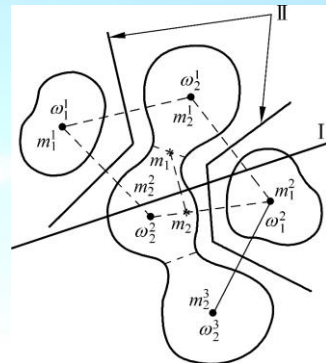
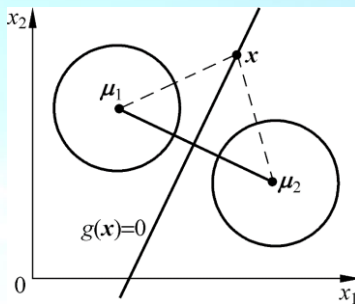


I : Linear classifier
II : Piecewise linear classifier
III : Quadratic classifier

Xuegong Zhang

17

Minimal distance classifier



- With given sub-class regions $R_i^l, l = 1, 2, \dots, L, i = 1, 2, \dots, c$

$$g_i(x) = \min_{l=1, \dots, L_i} \|x - m_i^l\|$$

- Decision:

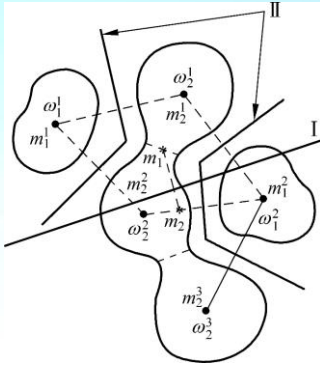
$$\text{class}(x) = \underset{i}{\operatorname{argmin}} g_i(x)$$

Xuegong Zhang

18



Piecewise linear discriminant



- With given sub-classes

$$\omega_i = \{\omega_i^1, \omega_i^2, \dots, \omega_i^{l_i}\}, i = 1, 2, \dots, c$$

- Linear discriminants

$$g_i^l(\mathbf{x}) = \mathbf{w}_i^l \cdot \mathbf{x} + \omega_{i0}^l, l = 1, \dots, l_i, i = 1, \dots, c$$

$$g_i(\mathbf{x}) = \max_{l=1, \dots, l_i} g_i^l(\mathbf{x}), i = 1, \dots, c$$

- Decision

$$\text{class}(\mathbf{x}) = \underset{i=1, \dots, c}{\operatorname{argmax}} g_i(\mathbf{x})$$

Xuegong Zhang

19

An algorithm for piecewise linear classification



① Initialization:

- Set: l_i subclasses in class ω_i , $i = 1, \dots, c$
- Weights of linear discriminant of subclass l_i at iteration t : $\alpha_i^l(t)$
- Initialize $\alpha_i^l(0)$, $i = 1, \dots, c$, $l = 1, \dots, l_i$

② For sample $\mathbf{y}_k \in \omega_j$, find $\alpha_j^m(t)^T \mathbf{y}_k = \max_{l=1, \dots, l_j} \{\alpha_j^l(t)^T \mathbf{y}_k\}$, check

- If $\alpha_j^m(t)^T \mathbf{y}_k > \alpha_i^l(t)^T \mathbf{y}_k, \forall i = 1, \dots, c, i \neq j, l = 1, \dots, l_i$, continue;
- Else if $\exists i \neq j, \alpha_j^m(t)^T \mathbf{y}_k \leq \alpha_i^n(t)^T \mathbf{y}_k$ in subclass n , find the subclass n with largest discriminant, do correction:

③ $\alpha_j^m(t+1) = \alpha_j^m(t) + \rho_k \mathbf{y}_k$

④ $\alpha_i^n(t+1) = \alpha_i^n(t) - \rho_k \mathbf{y}_k$

③ Repeat ② with the next sample, until convergence.

Xuegong Zhang

20



1-minute break



Xuegang Zhang

21



4.5 Nearest Neighbor Classifiers

Xuegang Zhang

22

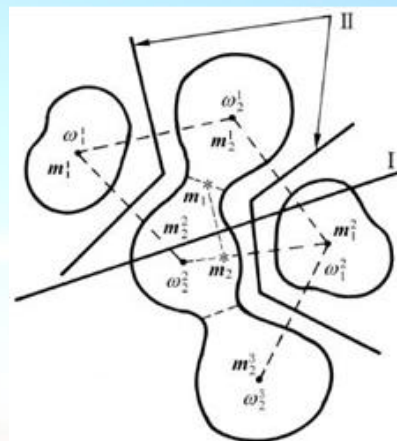
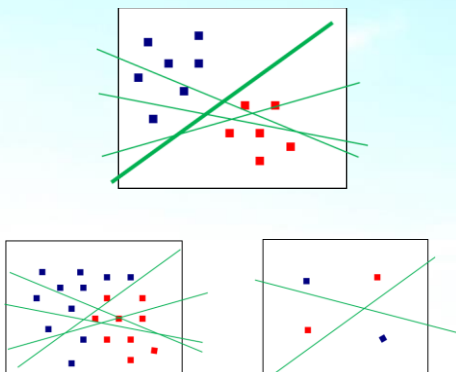
4.5.1 Nearest Neighbor Methods

Xuegong Zhang

23

Piece-wise linear classifier

- One way for nonlinearity: Piece-wise linear classifier
→ Classification of multiple sub-classes



Xuegong Zhang

24



Extreme case: → Nearest-Neighbor (NN) method

- Sample set: $S_N = \{(x_1, \theta_1), (x_2, \theta_2), \dots, (x_N, \theta_N)\}$
 - x_i : samples, labels: $\theta_i = \{1, 2, \dots, c\}$
- Distance measure $\delta(x_i, x_j)$, e.g. $\delta(x_i, x_j) = \|x_i - x_j\|$
- Nearest Neighbor Decision (1-NN):

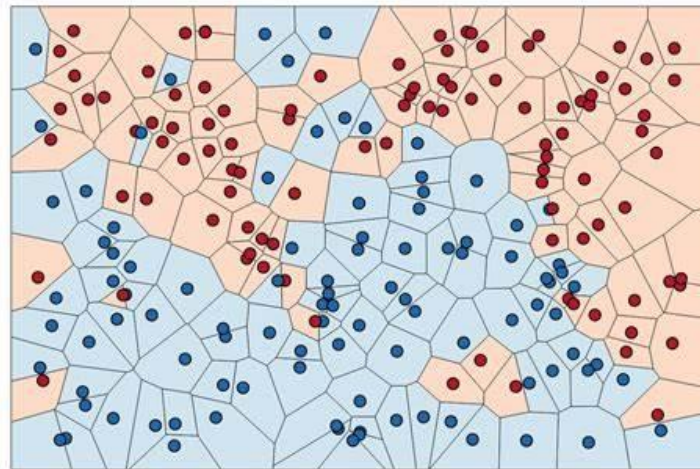
If $\delta(x, x') = \min_{j=1, \dots, N} \delta(x, x_j)$ and $\omega(x') = \theta'$, then $\hat{\omega}_1(x) = \theta'$
- Or:

Discriminant function of ω_i : $g_i(x) = \min_{k=1, \dots, N_i} \|x - x_i^k\|$

Decision: If $g_j(x) = \min_{i=1, \dots, c} g_i(x)$, then $x \in \omega_j$

Xuegong Zhang

25



Xuegong Zhang

26



Some popular distance measurements

- Minkovski Metric (of order s): $\delta(x_k, x_l) = [\sum_{i=1}^d |x_{ki} - x_{li}|^s]^{\frac{1}{s}}$
- Euclidean Distance: $\delta_E(x_k, x_l) = [(x_k - x_l)^T(x_k - x_l)]^{\frac{1}{2}}$
- City-Block Distance: $\delta(x_k, x_l) = \sum_{i=1}^d |x_{ki} - x_{li}|$
- Chobyshev Distance: $\delta(x_k, x_l) = \max_i |x_{ki} - x_{li}|$
- Squared Distance: $\delta(x_k, x_l) = (x_k - x_l)^T Q (x_k - x_l)$
- Nonlinear distances, e.g., $\delta(x_k, x_l) = \begin{cases} H & \text{if } \delta_E(x_k, x_l) \geq T \\ 0 & \text{if } \delta_E(x_k, x_l) < T \end{cases}$

Xuegong Zhang

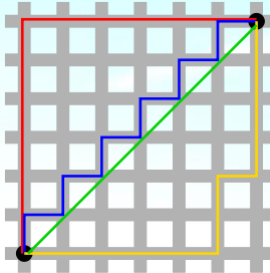
27

Manhattan Distance vs. Euclidean Distance

- Euclidean (L2) Distance: $\delta_E(x_k, x_l) = [(x_k - x_l)^T(x_k - x_l)]^{\frac{1}{2}}$
- City-Block (L1) Distance: $\delta(x_k, x_l) = \sum_{i=1}^d |x_{ki} - x_{li}|$

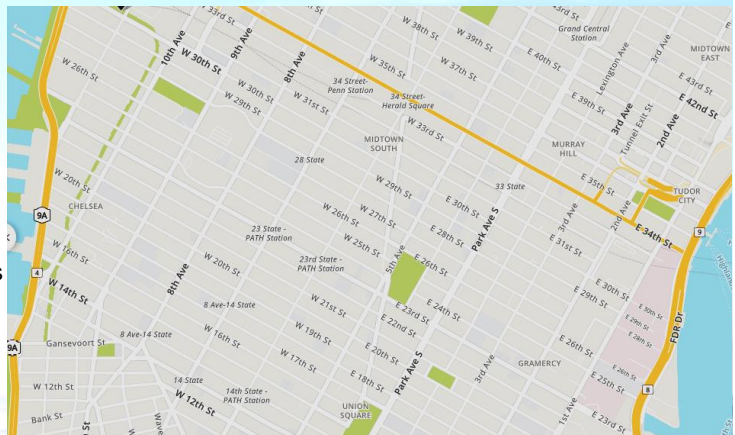


a.k.a. Manhattan distance, taxicab distance



Taxicab geometry versus Euclidean distance: In taxicab geometry, the red, yellow, and blue paths all have the same shortest path length of 12. In Euclidean geometry, the green line has length $6\sqrt{2} \approx 8.49$, and is the unique shortest path.

Xuegong Zhang

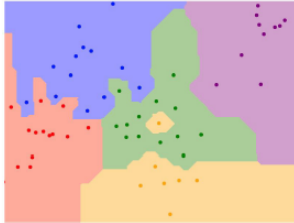




K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance

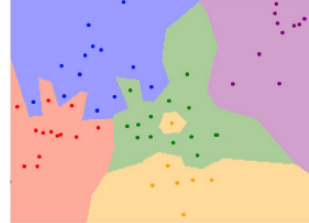
$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$



K = 1

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$



K = 1

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 33

April 5, 2018

Xuegong Zhang

29

k-Nearest Neighbors (k-NN) method



- A generalization of 1-NN

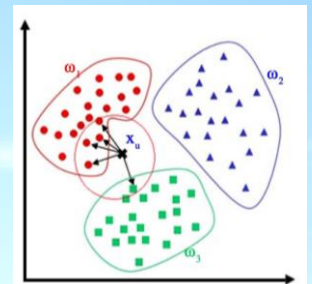
$$S_N = \{(x_1, \theta_1), (x_2, \theta_2), \dots, (x_N, \theta_N)\}, \quad \theta_i = \{1, 2, \dots, c\}$$

- Discriminant for ω_i :

$$g_i(x) = k_i$$

where $k_i, i = 1, \dots, c$ is the number of samples belonging to ω_i among the k nearest neighbors of x

- Decision: If $g_j(x) = \max_{i=1, \dots, c} g_i(x)$, then $x \in \omega_j$



Xuegong Zhang

30

4.5.2

Asymptotic Errors of Nearest Neighbor Methods

Xuegong Zhang

31

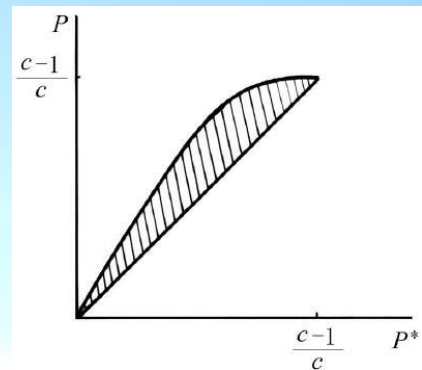
Conclusion:

- Error rate P_1 of NN method for i.i.d. samples, comparing with the best possible error rate

$$P^* \leq P_1 \leq P^* \left(2 - \frac{c}{c-1} P^* \right)$$

where P^* is the Bayesian error (smallest error given the distributions),

and P_1 is the **asymptotic error** of NN



Xuegong Zhang

32

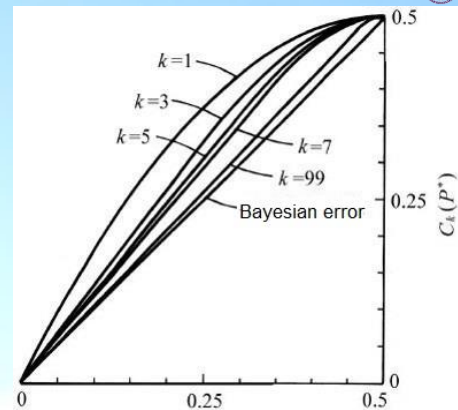
- Asymptotic error of k-NN

$$P^* \leq P_k \leq P^* \left(2 - \frac{c}{c-1} P^* \right)$$

Or:

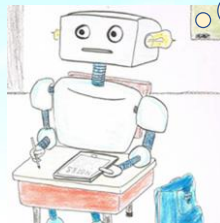
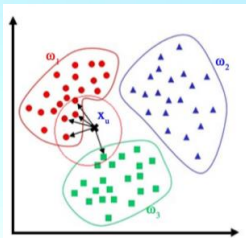
$$P^* \leq P_k \leq 2P^*$$

- When $N \rightarrow \infty$, $k \uparrow$, $P_k \downarrow$
given that $k \ll N$



Xuegong Zhang

33



So simple, so easy, and so good!



Who told me to learn? No need! I just need to memorize training samples.

Xuegong Zhang

34

单选题 1分

设置

4-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



提交

35



4.5.3

Improved Nearest Neighbor Methods

Xufeng Zhang

36



- So simple, so easy, and so good!

—— Wait! No free lunch!



- Potential challenges:
 - Huge memory and computation
 - “Voting democracy” can be risky, esp. with noisy samples
 - Rejection when voting is too close
 - Weight voting, ...
 - May not work well when N is too small

Xuegong Zhang

37



• Computational Issues

- Need to store all training samples → Memory Cost
- Need to compare with all samples → Computation Cost

Q: What is the time complexity of finding the nearest neighbor in N training samples?



• Fast k-NN algorithms

→ Basic idea: Reduce the online time by offline preparation beforehand

Xuegong Zhang

38

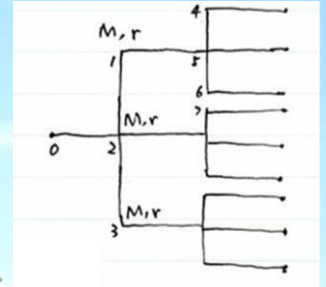


- Basic idea of fast k-NN algorithms

→ Reduce the online time by offline preparation beforehand

- The Branch-Bound (BB) Algorithm

- Organize samples as subsets in a tree structure
- Use few parameters to represent a subset (node)
- Compare with nodes instead of all samples
- Only compare with individual samples at the end nodes



Xuegong Zhang

39

Branch-Bound (BB) Algorithm (for 1-NN)



- Symbols:

\mathcal{X}_p : Sample subset of node p

N_p : Number of samples in \mathcal{X}_p

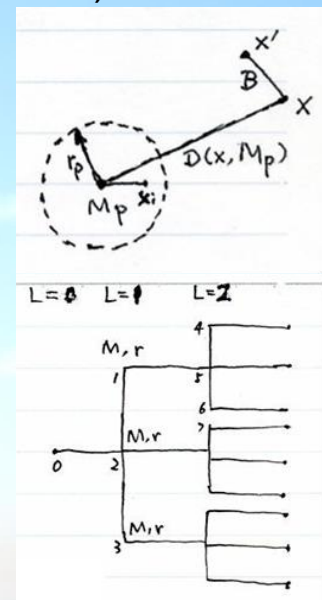
M_p : Sample mean of \mathcal{X}_p

$r_p = \max_{x_i \in \mathcal{X}_p} D(x_i, M_p)$: the farthest distance in \mathcal{X}_p from its center

B : Record on the current nearest distance

- Two procedures:

1. Divide samples into subsets (e.g., using [clustering algorithms](#)), compute and store the M_p, r_p and all $D(x_i, M_p)$ of \mathcal{X}_p
2. Using the BB algorithm to search through the tree



Xuegong Zhang

40



- Rule 1:

For new sample x and node \mathcal{X}_p :

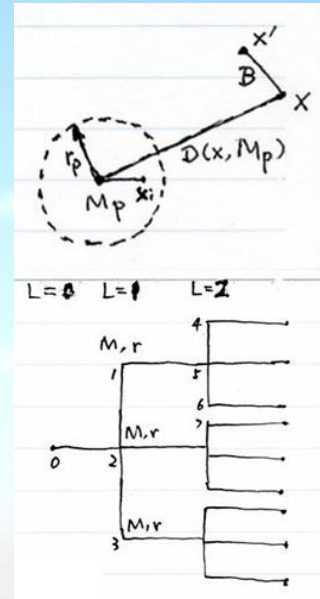
If $D(x, M_p) > B + r_p$, the nearest neighbor of x cannot be in \mathcal{X}_p .

- Rule 2:

For new sample x and training sample $x_i \in \mathcal{X}_p$:

If $D(x, M_p) > B + D(x_i, M_p)$,

x_i is NOT the nearest neighbor of x .



Xuegong Zhang

41

- BB Algorithm for the nearest neighbor:

1° (Initialize)

Set $B = \infty, L = 1, p = 0$.

2° (Expand)

Put children of current node in the Active Table of the current level, and compute and store their $D(x, M_p)$.

3° (Check)

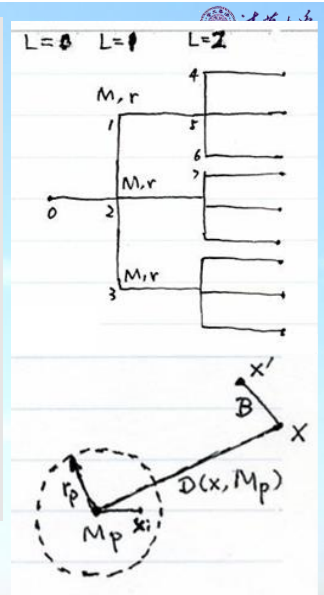
For each node in the Active Table, if $D(x, M_p) > B + r_p$, remove it from the Table (according to Rule 1).

4° (Trace back)

If no node in the Active Table, go back to the upper level and set $L = L - 1$.

If $L == 0$, stop; otherwise go to step 3°.

If there are nodes in the Active Table, continue to 5°.



Xuegong Zhang

42

5° (Choose nearest node)

Choose the nearest node with smallest $D(x, M_p)$ in the Active Table, say, p' , set it as the current node.

If current level L is the end level, go to 6°; otherwise, set $L = L + 1$, go to 2°.

6° (Check)

For every x_i in current node p' ,

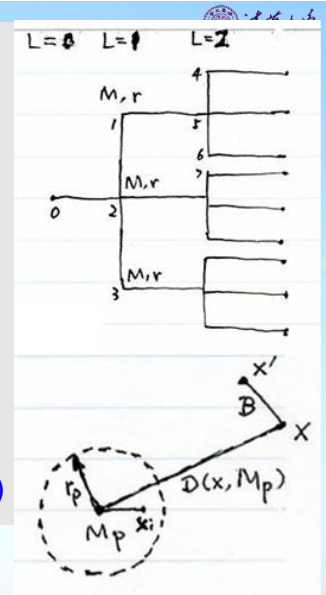
If $D(x, M_p) > B + D(x_i, M_p)$ then not the nearest (Rule 2)

Otherwise, compute $D(x, x_i)$. If $D(x, x_i) < B$, set $NN = i, B = D(x, x_i)$

After all x_i in p' being checked, go to 3°.

- When algorithm ends, output the x_{NN} and $D(x, x_{NN}) = B$ as the nearest neighbor and the nearest distance.

(for k-NN, we only need to modify step 6°)



Xuegong Zhang

43



Examples of Improved Nearest Neighbor Methods

Xuegong Zhang

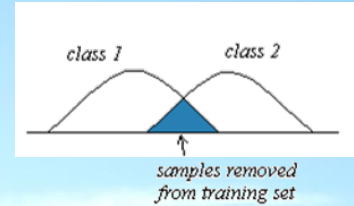
44

Editing Nearest Neighbor Method



- Consideration

- Samples in the overlapping region of two classes can confuse the algorithm
- Removing them should help to improve



Xuegong Zhang

45

Editing Nearest Neighbor Method



- How to know samples in the confusing zone?
- Idea:
 - Pre-classification to detect samples in confusing zone
 - Edit: remove samples in confusing zone
 - Can be done for multiple rounds if we have sufficient samples
- Asymptotic error of 1-NN using 1-NN pre-classification editing



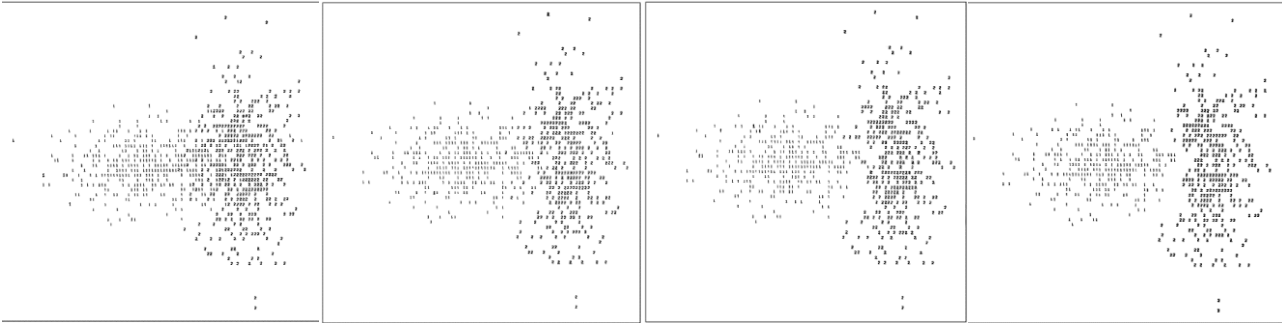
$$P_1^E(e) = \frac{P_1(e)}{2[1 - P_1(e)]}$$

$$\text{ref. } P^* \leq P_1 \leq 2P^*$$

Xuegong Zhang

46

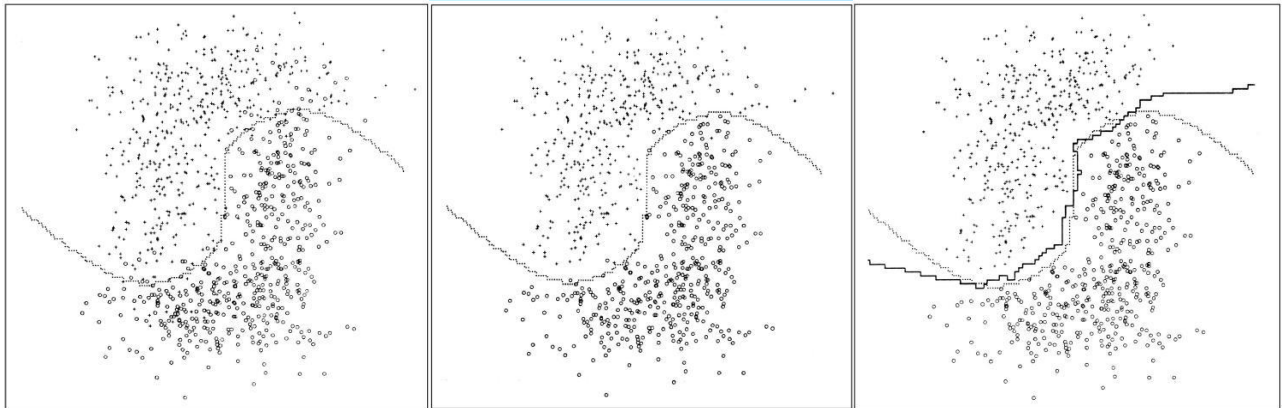
- Examples (multi-edit)



Xuegong Zhang

47

- Examples (multi-edit)



Xuegong Zhang

48



Condensed Nearest Neighbor (CNN) Method

- Find a few representative samples to significantly reduce storage
 - Separate \mathcal{X}^N as \mathcal{X}_S and \mathcal{X}_G
 - Start with only 1 sample in \mathcal{X}_S , all others in \mathcal{X}_G
 - Consider each sample in \mathcal{X}_G , if correctly classified with \mathcal{X}_S then stay, otherwise move to \mathcal{X}_S
 - ...
 - Use only samples in \mathcal{X}_S as the final set

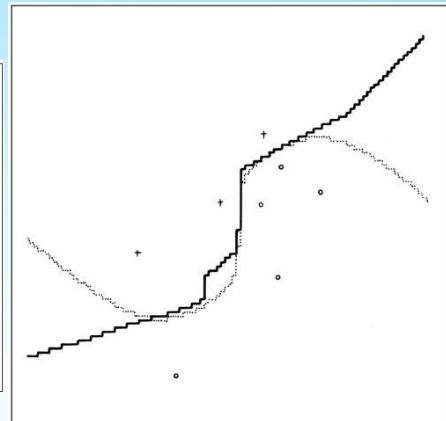
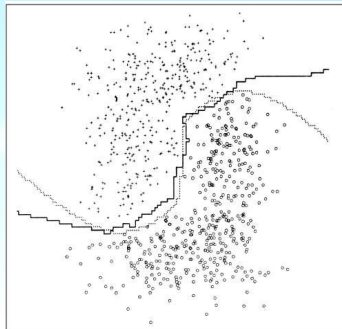
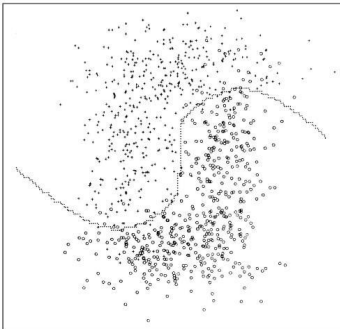
Xuegong Zhang

49



Condensed Nearest Neighbor (CNN) Method

- Can be combined with editing NN



Xuegong Zhang

50

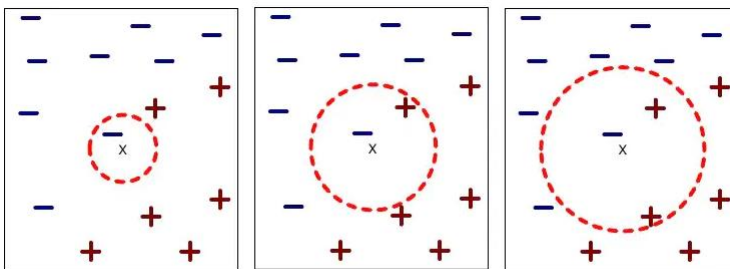
4.5.4

Discussions on Nearest Neighbor Methods

Xuegong Zhang

51

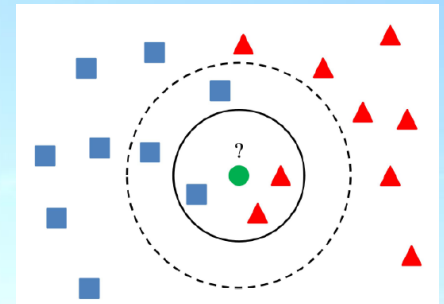
k-Nearest Neighbors (k-NN)



(a) 1-nearest neighbor

(b) 2-nearest neighbor

(c) 3-nearest neighbor

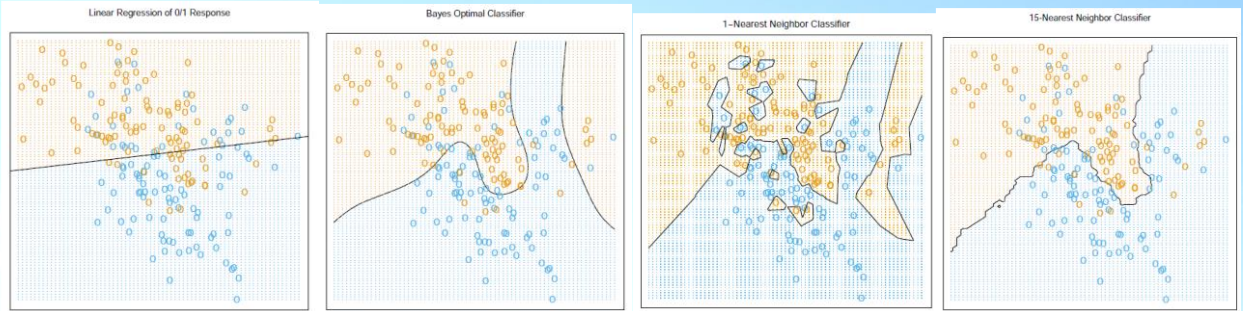


Xuegong Zhang

52



About the choice of k in k -NN



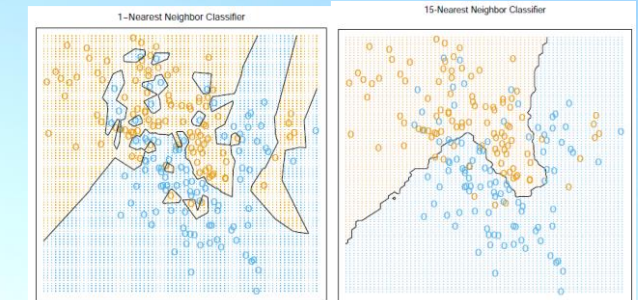
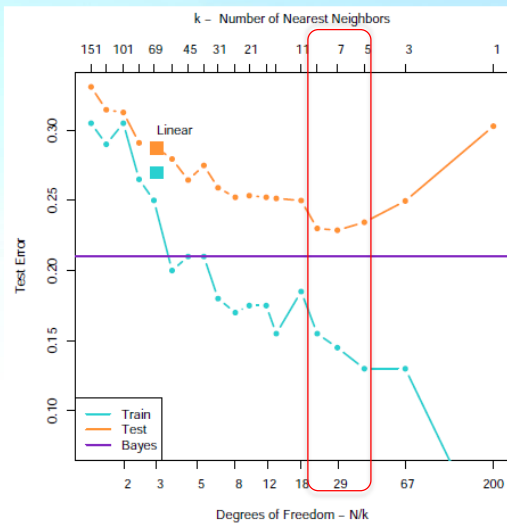
Ref. Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning*, Springer

Xuegong Zhang

53



About the choice of k in k -NN



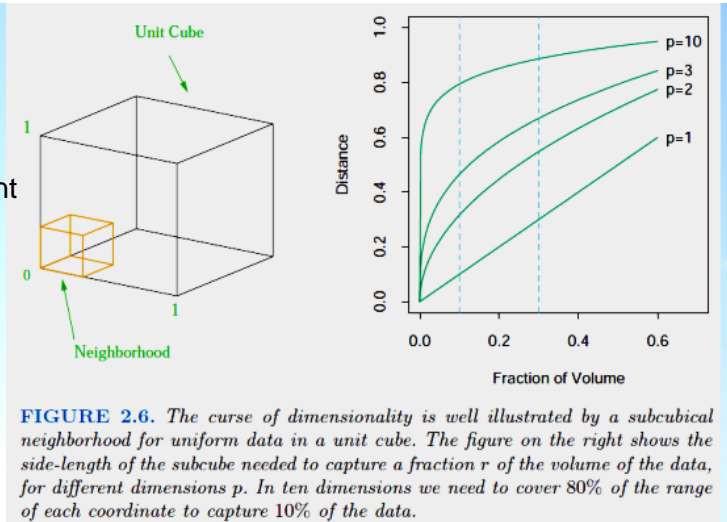
Ref. Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning*, Springer

Xuegong Zhang

54

• The curse of dimensionality

- In high dimension space, your nearest neighbor can be far away
- You may need tremendous amount of samples to make NN work



Ref. Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning*, Springer

Xuegong Zhang

55

Discussion



- In what scenario would you prefer using k-NN?
 - Dimensionality, Sample size
 - Expectation on the complexity of classification
 - Interpretability
- How do you choose k in a practical application?
 - Rule of thumb
 - Exploratory experiments
- What distance/preprocessing would you use in a practical application?
 - Characteristics of features
 - Relative importance of features



Xuegong Zhang

56

Homework



- Computer exercises (Ex2)
 - Find a package of KNN
 - Describe its algorithm
 - Write your own code of MLP
 - Experiment on the medical dataset



- Deadline:
 - Oct. 13 (Sunday), 23:00

Xuegong Zhang

57

单选题 1分

设置

5-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



Xuegong Zhang

提交

58