Ch 2.2.3: Intro to classification Lecture 9 - CMSE 381

Prof. Mengsen Zhang

Michigan State University

:

Dept of Computational Mathematics, Science & Engineering

Mon, Sep 15, 2025

Announcements

CMSE381_F2025_Schedule : Schedule									
Lec #		Date	Topic	Reading	HW				
1	M	8/25	Intro / Python Review	1					
2	W	8/27	What is statistical learning	2.1					
3	F	8.29	Assessing Model Accuracy	2.2.1, 2.2.2					
	М	9/1	Labor Day - No Class						
4	W	9/3	Linear Regression	3.1					
5	F	9/5	More Linear Regression	3.1	HW #1 Due Sun 9/7				
6	М	9/8	Multi-linear Regression	3.2					
7	W	9/10	Probably More Linear Regression	3.3					
8	F	9/12	Last of the Linear Regression		HW #2 Due				
9	М	9/15	Intro to classification, Bayes classifier, KNN classifier	2.2.3	Sun 9/14				
10	W	9/17	Logistic Regression	4.1, 4.2, 4.3.1-3					
11	F	9/19	Multiple Logistic Regression / Multinomial Logistic Regression	4.3.4-5	HW #3 Due Sun 9/21				
	М	9/22	Project Day & Review						
	W	9/24	Midterm #1						
12	F	9/26	Leave one out CV	5.1.1, 5.1.2					
40	**	0/00	1. 4-1.4 007	F 4 0					

Last Time:

• Finished Linear Regression

Announcements:

- Homework #3 Due Sunday Sep 21
- Next Monday Review day
 - depends on what you ask!
 - Submit your questions Here
- Wed 9/24 Exam #1
 - ▶ Bring 8.5×11 sheet of paper
 - Handwritten both sides
 - Anything you want on it, but must be your work
 - ► You will turn it in

Covered in this lecture

- Ch 2.2.3
- Error rate (classification)
- Bayes Classifier
- K-NN classification

Section 1

Classification Overview

What is classification

Classification: When the response variable is qualitative

- Given feature vector X and qualitative response Y in the set S, the goal is to find a function (classifier) C(X) taking X as input and predicting its value for Y.
- We are more interested in estimating the probabilities that X belongs to each category

Some examples

- Predict whether a COVID19 vaccine will work on a patient given patient's age
- An online banking service wants to determine whether a transaction being performed is fraudulent on the basis of the user's IP address, past transactions, etc.

Section 2

Ch 2.2.3: Classification

Error rate

- Training data: $\{(x_1, y_1), \dots, (x_n, y_n)\}$ with y_i qualitative
- Estimate $\hat{y} = \hat{f}(x)$
- Indicator variable

Training error rate:

$$\frac{1}{n}\sum_{i=1}^n\mathrm{I}(y_i\neq\hat{y}_i)$$

Test error rate:

$$\operatorname{Ave}(\mathrm{I}(y_0 \neq \hat{y}_0))$$

Best ever classifier

We can't have nice things

Bayes Classifier:

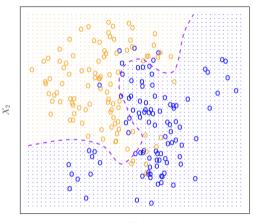
Give every observation the <u>highest</u> <u>probability</u> class given its predictor variables

$$\Pr(Y = j \mid X = x_0)$$

An example

- Survey students for amount of programming experience, and current GPA
- Try to predict if they will pass CMSE 381.
- If we have a survey of all students that could ever exist, we can determine the probability of failure given combo of those features.

Bayes decision boundary



 X_1

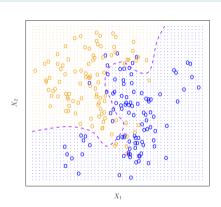
Bayes error rate

• Error at $X = x_0$

$$1 - \max_{j} \Pr(Y = j \mid X = x_0)$$

Overall Bayes error:

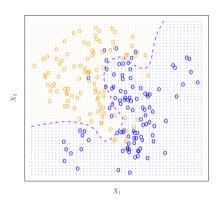
$$1 - E\left(\max_{j} \Pr(Y = j \mid X = x_0)\right)$$



12 / 20

Zhang (MSU-CMSE) Lec 8 Mon, Sep 15, 2025

The game

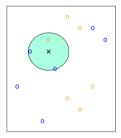


Test your understanding: PollEv

Section 3

K-Nearest Neighbors Classifier

K-Nearest Neighbors

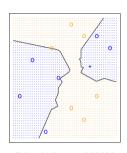


$$K = 3$$

- Fix K positive integer
- N(x) = the set of K closest neighbors to x
- Estimate conditional proability

$$\Pr(Y = j \mid X = x_0) = \frac{1}{K} \sum_{i \in N(x_0)} I(y_i = j)$$

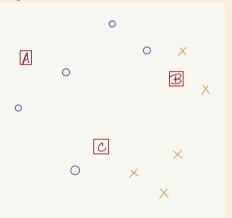
• Pick *j* with highest value



Black line: KNN decision boundary

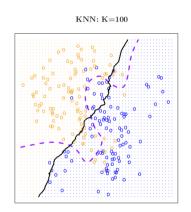
Example

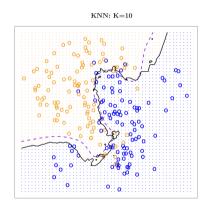
Here label is shown by O vs X. What are the knn predictions for points A, B and C for k = 1 or k = 3?

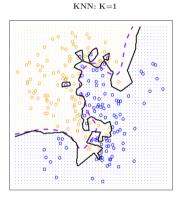


Point	k=1 Prediction	k = 3 Prediction	
Α			
В			
С			

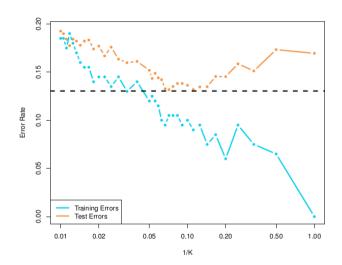
Tradeoff







More on tradeoff



Jupyter notebook

Next time

CMSE381_F2025_Schedule : Schedule

Lec #		Date	Topic	Reading	HW			
1	M	8/25	Intro / Python Review	1				
2	W	8/27	What is statistical learning	2.1				
3	F	8.29	Assessing Model Accuracy	2.2.1, 2.2.2				
	M	9/1	Labor Day - No Class					
4	W	9/3	Linear Regression	3.1				
5	F	9/5	More Linear Regression	3.1	HW #1 Due Sun 9/7			
6	M	9/8	Multi-linear Regression	3.2				
7	W	9/10	Probably More Linear Regression	3.3				
8	F	9/12	Last of the Linear Regression		HW #2 Due Sun 9/14			
9	М	9/15	Intro to classification, Bayes classifier, KNN classifier	2.2.3				
10	W	9/17	Logistic Regression	4.1, 4.2, 4.3.1-3				
11	F	9/19	Multiple Logistic Regression / Multinomial Logistic Regression	4.3.4-5	HW #3 Due Sun 9/21			
	M	9/22	Project Day & Review					
	W	9/24	Midterm #1					
12	F	9/26	Leave one out CV	5.1.1, 5.1.2				
40		0/00	1. 4-1-1 (017	E 4 0				