# L4: Logistic Regression

Shan Wang

Lingnan College, Sun Yat-sen University

# Last lecture

- Linear regression method
  - Model: $y = \boldsymbol{\theta}'\boldsymbol{x}$
  - Strategy
    - Least squared error: $\min \frac{1}{N}\sum_{i=1}^{N}(y_i - \boldsymbol{\theta}'\boldsymbol{x}_i)^2$
    - Maximum likelihood: $\max \sum_{i=1}^{N} \log P(y_i|\boldsymbol{x}_i, \theta)$
  - Algorithm
    - Normal equation: $\widehat{\boldsymbol{\theta}} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{y}$
    - Gradient descent method: $\boldsymbol{\theta}_{\text{new}} = \boldsymbol{\theta}_{\text{old}} - \eta \frac{\partial R(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$
- Regularization
  - Lasso: $\lambda\|\boldsymbol{\theta}\|_1$
  - Ridge: $\lambda\|\boldsymbol{\theta}\|_2^2$
    - Relationship between Ridge and MAP
- Application: quality of wine

# Regression vs Classification

- **Linear regression** would predict a **continuous** outcome
  - $y = \boldsymbol{\theta}'\boldsymbol{x}$

- To predict the quality of care
  - The dependent variable is modelled as a binary variable
  - 1 if low-quality care, 0 if high-quality care

- This is a **categorical variable**
  - Typically a small number of possible outcomes, 2 (low-quality care and high-quality care) in this case

- How can we extend the idea of linear regression to situations where the outcome variable is categorical?
  - Only want to predict 1 or 0
  - Could round outcome to 1 or 0
  - But we can do better with **logistic regression**

# Course outline

- Supervised learning
  - Linear regression
  - **Logistic regression**
  - SVM and kernel
  - Tree models

- Deep learning
  - Neural networks
  - Convolutional NN
  - Recurrent NN

- Unsupervised learning
  - Clustering
  - PCA
  - EM

- Reinforcement learning
  - MDP
  - ADP
  - Deep Q-Network

# This Lecture

- Classification problem

- Logistic regression method
  - Model
  - Strategy
  - Algorithm

- Prediction & evaluation

- Multi-class logistic regression

- Application: diabetes care

Reference: VE 445, Shuai LI (SJTU); OPIM 326 Daniel Zheng (SMU)
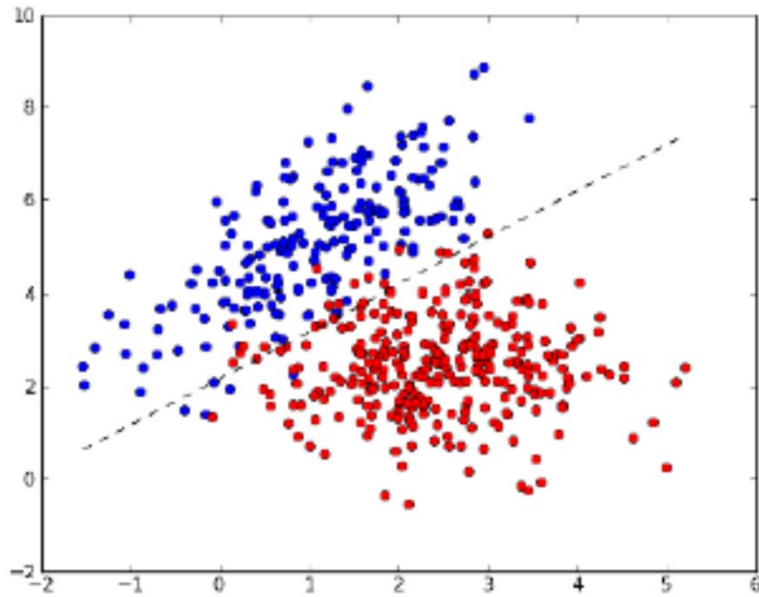
# Classification problem

# What is classification?

- Compared to regression problem, which predicts the labels from many numerical features

- Many applications
  - Spam Detection: Predicting if an email is Spam or not based on word frequencies
  - Credit Card Fraud: Predicting if a given credit card transaction is fraud or not based on their previous usage
  - Health: Predicting if a given mass of tissue is benign or malignant
  - Marketing: Predicting if a given user will buy an insurance product or not
  - …

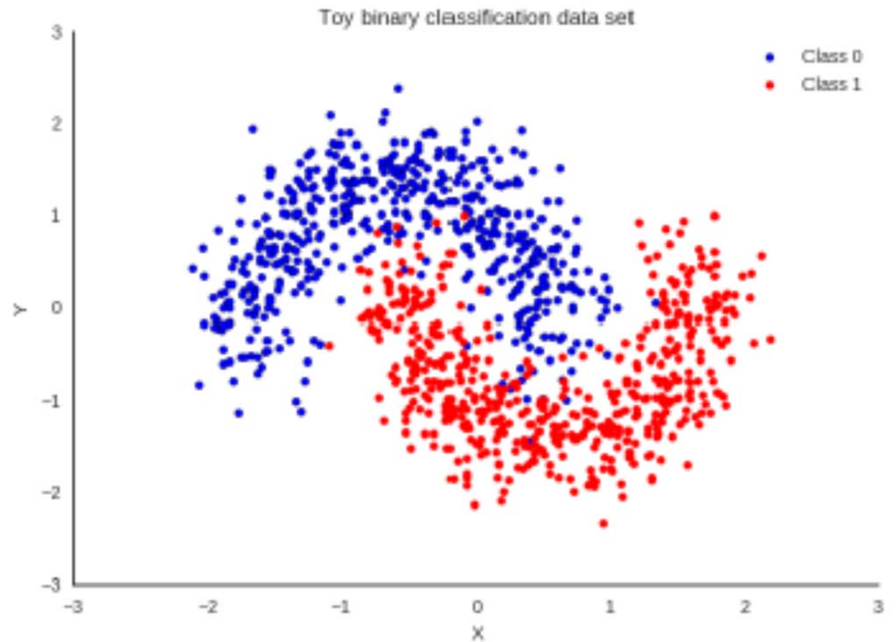# Classification problem definition

- Given:
  - A description of an instance $x \in X$
  - A fixed set of categories: $C = \{c_1, c_2, \ldots, c_K\}$

- Determine:
  - The category of $x$: $f(x) \in C$ where $f(x)$ is a categorization function whose domain is $X$ and whose range is $C$
  - If the category set binary, i.e. $C = \{0, 1\}$ ({false, true}, {negative, positive}) then it is called binary classification

# Binary classification



Linearly separable

Nonlinearly separable

# Linear discriminative model

- Modeling the dependence of unobserved variables on observed ones

- a.k.a. conditional models

- Deterministic: $y = f_\theta(x)$

- Probabilistic: $p_\theta(y|x)$

- For binary classification
    - $p_\theta(y = 1|x)$
    - $p_\theta(y = 0|x) = 1 - p_\theta(y = 1|x)$

# Logistic regression model

# What is logistic regression

- Logistic regression is a binary classification model
  - $p_\theta(y = 1|x) = \sigma(\theta'x) = \frac{e^{\theta'x}}{1+e^{\theta'x}}$
  - $p_\theta(y = 0|x) = \frac{1}{1+e^{\theta'x}}$
  - $\sigma(x)$ is the logistic function or the sigmoid function

- Interpretation
  - Odds: $\frac{p}{1-p}$
  - $p = \frac{e^{\theta'x}}{1+e^{\theta'x}} \leftrightarrow \log(\frac{p}{1-p}) = \theta'x$
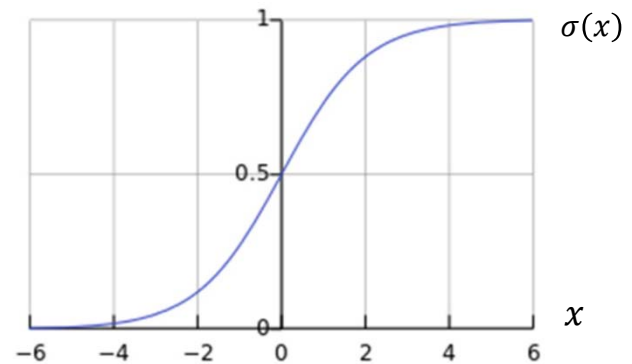    - Logistic regression is a linear regression for the log odds

*Intercept term is omitted in this lecture, you may imagine the first term in $x$ is always 1.

# Logistic function properties

- Properties for the logistic function

  - $\sigma(x) = \dfrac{e^x}{1+e^x}$

    - Bounded in $(0,1)$
    - $\sigma(x) \to 1$ when $x \to \infty$
    - $\sigma(x) \to 0$ when $x \to -\infty$

  - $\sigma'(x) = \dfrac{\partial \frac{e^x}{1+e^x}}{\partial x}$

    $= \dfrac{e^x(1+e^x) - e^x e^x}{(1+e^x)^2} = \dfrac{e^x}{1+e^x} \times \dfrac{1}{1+e^x}$

    $= \sigma(x)(1-\sigma(x))$

# Logistic regression strategy

- Goal: we want to choose the right $\boldsymbol{\theta}$, to make the best prediction

- Which loss function to use?

# Entropy

- Entropy
  - A measure of the uncertainty
  - Suppose the are $K$ classes, class $k$ with $p_k$
    - Entropy$= -\sum_{k=1}^{K} p_k \log p_k$
  - Uniform distribution has maximum entropy

- Cross entropy
  - To calculate the difference between two probability distributions
    - E.g., The true distribution and predicted distribution
  - Cross entropy$= -\sum_{k=1}^{K} p_k \log q_k$
    - $p_k$: true label distribution
    - $q_k$: predicted label distribution

# Cross entropy loss for logistic regression

- Loss function for data point $(\boldsymbol{x}, y)$ with prediction model $p_\theta(\cdot \,|x)$ is

$$L(y, \boldsymbol{x}, p_\theta) =$$

$$= -1_{y=1} \log p_\theta(1|\boldsymbol{x}) - 1_{y=0} \log p_\theta(0|\boldsymbol{x})$$

$$= -y \log p_\theta(1|\boldsymbol{x}) - (1 - y) \log(1 - p_\theta(1|\boldsymbol{x}))$$

- Where $p_\theta(1|\boldsymbol{x}) = \sigma(\boldsymbol{\theta}'\boldsymbol{x})$

- Minimize the empirical error $\widehat{R}(\boldsymbol{\theta})$

$$\widehat{\boldsymbol{\theta}} = \mathrm{argmin} -\frac{1}{N} \sum_{i=1}^{N} [y_i \log p_\theta(1|\boldsymbol{x}_i) + (1 - y_i) \log(1 - p_\theta(1|\boldsymbol{x}_i))]$$

# Logistic regression algorithm

# Solve the $\widehat{\boldsymbol{\theta}}$

- $\hat{R}(\boldsymbol{\theta})$ is convex in $\boldsymbol{\theta}$

- Gradient descent:

- $$\frac{\partial L(y,x,p_\theta)}{\partial \boldsymbol{\theta}} = -y \frac{1}{\sigma(\boldsymbol{\theta}'x)} \sigma(\boldsymbol{\theta}'x)\big(1 - \sigma(\boldsymbol{\theta}'x)\big)x$$

$$-(1-y)\frac{-1}{1-\sigma(\boldsymbol{\theta}'x)}\sigma(\boldsymbol{\theta}'x)\big(1 - \sigma(\boldsymbol{\theta}'x)\big)$$

$$= -\big(y - \sigma(\boldsymbol{\theta}'x)\big)x \qquad \sigma'(x) = \sigma(x)(1-\sigma(x))$$

- $$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \boldsymbol{\eta}\big(y - \sigma(\boldsymbol{\theta}'x)\big)x$$

# Prediction & evaluation

# Label prediction

- Logistic regression provides the probability
  - $p_\theta(y = 1|x) = \sigma(\theta'x) = \frac{e^{\theta'x}}{1+e^{\theta'x}}$
  - $p_\theta(y = 0|x) = \frac{1}{1+e^{\theta'x}}$
- The final label of an instance is decided by setting a threshold $h$
  - If $p_\theta(y = 1|x) > h$, predict 1
  - Otherwise, predict 0

# How to choose threshold?

- In general, 0.5 is a good choose

- If you have precision-recall trade-off
  - Higher threshold
    - More False Negative
    - Less False Positive
    - Higher precision
    - Lower recall
  - Lower threshold
    - Less False Negative
    - More False Positive
    - Lower precision
    - Higher recall

|       | Prediction | |
|-------|------------|--------------|
|       | **1**        | **0**          |
| **1** | True Positive | False Negative |
| **0** | False Positive | True Negative |

Label

- **Precision**: the ratio of true class 1 cases in those with prediction 1

$$\text{Prec} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

|       | Prediction | |
|-------|------------|--------------|
|       | **1**        | **0**          |
| **1** | True Positive | False Negative |
| **0** | False Positive | True Negative |

Label

- **Recall**: the ratio of cases with prediction 1 in all true class 1 cases

$$\text{Rec} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

# Model evaluation

- Accuracy
  - How many are correctly classified

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN}$$

- Recall
  - How many positive ones are correctly classified

$$Precision = \frac{TP}{TP + FP}$$

- Precision
  - How many are correctly classified among those are labelled as positive

$$Recall = \frac{TP}{TP + FN}$$

- F1 measure

$$F1 - measure = \frac{2TP}{2TP + FP + FN}$$

# ROC curve

- ROC curve:
  - Receiver Operating Characteristic (ROC) Curve
  - X axis: false positive rate (1-specificity)=FP/(TN+FP)
  - Y axis: true positive rate (recall/sensitivity)=TP/(TP+FN)
  - A performance measurement for classification problem at various thresholds settings



Bottom left corner, threshold=1, every thing is predicted to be negative, so TPR=0 and FPR=0

Decrease the threshold, more true positives, more false positives

Top right corner, threshold=0, every thing is predicted to be positive, so TPR=1 and FPR=1

# AUC

- Area Under ROC Curve (AUC)
  - The higher, the better
  - Tells how much the model is capable of distinguishing between classes
  - Perfect classifier get AUC=1 and random classifier get AUC = 0.5

# PR curve and AUPR

- PR curve:
  - The precision recall curve
  - X axis: recall=TP/(TP+FN)
  - Y axis: precision=TP/(TP+FP)

- AUPR
  - Area Under PR curve
  - Can handle imbalanced datas
  - another plot to measure the performance of binary classifier
  - Usually, the classifiers gets lower AUPR value than AUC value

**Precision-recall curves – examples**

# Quick bonus question

Do you have any other idea to handle imbalanced data?

i.e., one class is frequent, the other is infrequent

# Multi-class logistic regression

# Multi-class logistic regression

# Model

- Class $C = \{c_1, c_2, \dots c_K\}$
- The probability of class $k$

$$p_\theta(y = c_k | \boldsymbol{x}) = \frac{e^{\boldsymbol{\theta}'_k \boldsymbol{x}}}{\sum_{k=1}^{K} e^{\boldsymbol{\theta}'_k \boldsymbol{x}}} \text{ for } k = 1, \dots, K$$

(softmax function)

- Parameters
  - $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K$
  - Can be normalized to be K-1 groups of parameters
    - $p_\theta(y = c_k | \boldsymbol{x}) = \frac{e^{\boldsymbol{\theta}'_k \boldsymbol{x}}}{1 + \sum_{k=1}^{K-1} e^{\boldsymbol{\theta}'_k \boldsymbol{x}}} \text{ for } k = 1, \dots, K-1$
    - $p_\theta(y = c_K | \boldsymbol{x}) = \frac{1}{1 + \sum_{k=1}^{K-1} e^{\boldsymbol{\theta}'_k \boldsymbol{x}}}$

# Strategy

- Cross entropy

$$L(y, \boldsymbol{x}, p_\theta) = -\sum_{k=1}^{K} 1_{y=c_k} \log p_\theta(c_k|\boldsymbol{x})$$

  - Empirical error

$$\hat{R}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, ..., \boldsymbol{\theta}_K) = -\sum_{i=1}^{N} \log p_\theta(y_i|\boldsymbol{x})$$

- Likelihood

$$\mathcal{L}(y, \boldsymbol{x}, p_\theta) = \prod_{i=1}^{N} p_\theta(y_i|\boldsymbol{x})$$

  - Take log:

$$\log \mathcal{L}(y, \boldsymbol{x}, p_\theta) = \sum_{i=1}^{N} \log p_\theta(y_i|\boldsymbol{x})$$

Minimize empirical error with cross entropy loss
⇕
Maximize the log likelihood

# Algorithm

- Gradient

$$\frac{\partial}{\partial \boldsymbol{\theta}_k} \log p_\theta(y = c_k | \boldsymbol{x}) = \frac{\partial}{\partial \boldsymbol{\theta}_k} \log \frac{e^{\boldsymbol{\theta}'_k \boldsymbol{x}}}{\sum_{k=1}^{K} e^{\boldsymbol{\theta}'_k \boldsymbol{x}}}$$

$$= \frac{\partial}{\partial \boldsymbol{\theta}_k} \log e^{\boldsymbol{\theta}'_k \boldsymbol{x}} - \frac{\partial}{\partial \boldsymbol{\theta}_k} \log \sum_{k=1}^{K} e^{\boldsymbol{\theta}'_k \boldsymbol{x}}$$

$$= \boldsymbol{x} - \frac{e^{\boldsymbol{\theta}'_k \boldsymbol{x}} \boldsymbol{x}}{\sum_{k=1}^{K} e^{\boldsymbol{\theta}'_k \boldsymbol{x}}} = \boldsymbol{x}(1 - p_\theta(y = c_k | \boldsymbol{x}))$$

# Application

Diabetes Care

# Service Quality vs. Efficiency

- Key trade-off in many service systems
  - Healthcare

- Improve operational efficiency without sacrificing quality

- How to assess quality?
  - Quality of decisions vs. outcomes (risk)
  - Critical decisions are often made by people with expert knowledge like physicians

- How to incorporate quality into operations management models?
  - Need quantitative and objective assessment of quality

# Healthcare Quality Assessment

- Importance: Critical in improving care quality and efficiency of healthcare operations
  - Timely intervention to revert poor quality care
  - Capacity shortage in healthcare systems
- Challenge: No single set of guidelines for defining quality of healthcare
- How?
  - Health professionals are experts in quality of care assessment

# Experts Assessment

- Healthcare quality assessment through expert opinions
  - Expert physicians can evaluate quality by examining a patient's records
  - This process is time consuming and inefficient
  - Experts are limited by memory and time
  - They cannot assess quality for millions of patients
- Similar practice in other industries
  - Accreditation in education, manufacturing, etc.

# Replicating Expert Assessment

- Can we develop analytical tools that replicate expert assessment?

- Learn from expert human judgment
  - Develop a model, interpret results, and adjust the model

- Make predictions/evaluations on a large scale

# Claims Data

**Medical Claims**

Diagnosis, Procedures, Doctor/Hospital, Cost

**Pharmacy Claims**

Drug, Quantity, Doctor, Medication Cost

- Electronically available
- Standardized

- Not 100% accurate
- Under-reporting is common
- Claims for hospital visits can be vague

# Building a model

- We use a **logistic regression**
  - Predicts an outcome variable, or *dependent/response variable*
  - Using a set of *independent/explanatory variables*

- Dependent variable: Poor care or not
  - is equal to 1 if the patient had poor care, and equal to 0 if the patient had good care

- Independent variables:
  - Number of Office Visits (OfficeVisits)
  - Number of Narcotics Prescribed (Narcotics)
  - Etc.

# Model for Healthcare Quality

- Plot of the independent variables
  - Number of Office Visits (OfficeVisits)
  - Number of Narcotics Prescribed (Narcotics)

- Red are poor care

- Green are good care

- Are these variables predictive of good care or poor care?
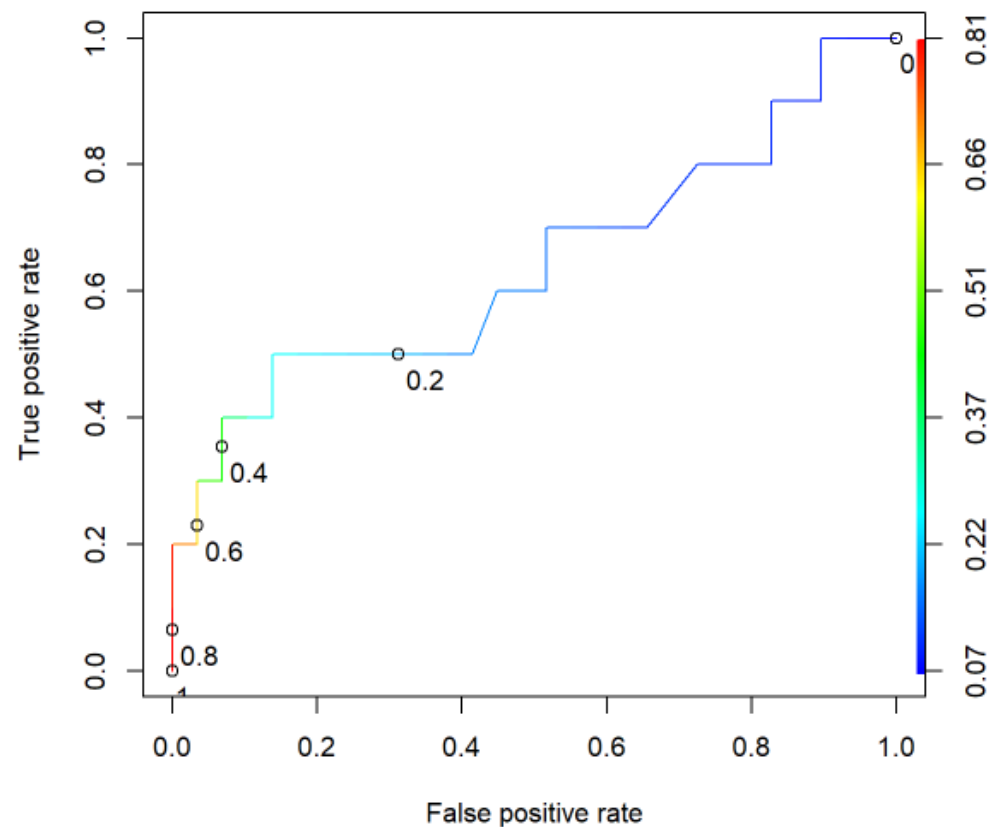
# Let's get our hands dirty!

# ROC and AUC

# Making Predictions

- Out-of-sample ROC
- Out-of-sample AUC = 0.64

# Lecture 4 wrap-up

✓Classification problem

✓Logistic regression method

  ✓Model

  ✓Strategy

  ✓Algorithm

✓Prediction & evaluation

✓Multi-class logistic regression

✓Application: diabetes care

# Next lecture

- Supervised learning
  - Linear regression
  - Logistic regression
  - SVM and kernel
  - Tree models

- Deep learning
  - Neural networks
  - Convolutional NN
  - Recurrent NN

- Unsupervised learning
  - Clustering
  - PCA
  - EM

- Reinforcement learning
  - MDP
  - ADP
  - Deep Q-Network

# Questions?

Shan Wang (王杉)

https://wangshan731.github.io/