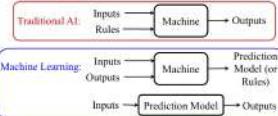


- Useful Resources: Datasets**
- UCI Repository: <http://www.ics.uci.edu/~mlearn/MLRepository.html>
 - Kaggle: <http://www.kaggle.com>
 - Distributed/parallel computing + machine learning → Distributed/parallel machine learning
 - Machine learning + hardware → AI chips

Machine Learning Definition

- Machine learning is a type of artificial intelligence that provides computers with the ability to learn from examples/experience without being explicitly programmed



$$DL \subseteq ML \subseteq AI$$

↓ of
methodology ↓ of
field

Course Objective

- To provide students with essential concepts and **principles** of machine learning algorithms
- To enable students to understand how to **revise** or **design** (beyond how to use) various machine learning algorithms to solve supervised learning and unsupervised learning problems

A Motivating Example

- Given a face image, to classify the face gender.



- Once upon a time, to develop an AI system to solve such a task, developers or domain experts need to provide rules and implement them in the system

If the face has long hair and does not have mustache, then this is a "female" face;
If the face has short hair and mustache, then this is a "male" face.

A Motivating Example (cont.)

- Limitations:
- Time consuming
- The defined rules may not be complete
- Not able to handle uncertainty

If the face has long hair and does not have mustache, then this is a "female" face;
If the face has short hair and mustache, then this is a "male" face.



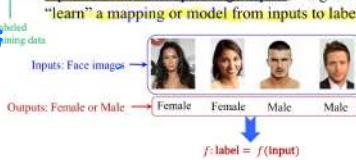
How to Represent an Example?



- Feature engineering (not machine learning focus)
- Representation learning (one of the crucial research topics in machine learning)
- Deep learning is the current most effective approach to representation learning

Supervised Learning

- The examples presented to computers are **pairs** of **inputs** and the corresponding **outputs**, the goal is to "learn" a mapping or model from inputs to labels



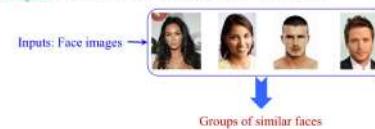
Outputs are discrete (i.e., categorical) values → classification
Labels are continuous values → regression

Example:

- Regression
- Classification
- Ensemble Learning

Unsupervised Learning

- The examples presented to computers are **a set of inputs without any outputs**, the goal is to "learn" an intrinsic structure of the examples, e.g., clusters of examples, density of the examples

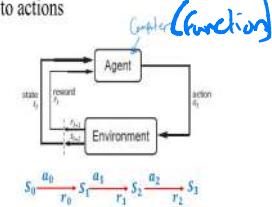


Example:

- Clustering
- Density estimation
- Dimensionality Reduction

Reinforcement Learning

- Learning by interacting with an environment to achieve a goal
- Objective: to learn an optimal policy mapping states to actions



Example:

- Deep Q-Network (DQN) (Cartoon)

Advanced paradigms:

- Semi-supervised learning
- Active learning
- Transfer learning

Various machine learning applications:
Face recognition, object recognition, text mining, activity recognition, stock price prediction, etc.

Various learning paradigms:
Supervised learning, unsupervised learning, reinforcement learning, other advanced learning.

Various types of methodologies:
Graphical models, deep learning, empirical risk minimization, entropy-based models, kernel methods, etc.

Various mathematical techniques:
Probability theory, linear algebra, calculus, optimization, information theory, functional analysis, etc.

Limitation of Supervised Learning

- Require sufficient labeled data to train a precise model (i.e., a model with good prediction performance)
- Sufficiency of labeled data is context-aware, depending on different kinds of applications and specific datasets
- When there is insufficient labeled data, can we still train a precise model?
- Advanced machine learning paradigms

Typical Learning Procedure

Two phases:

- Training phase
 - Given labeled training data $\{x_i, y_i\}$, for $i = 1, \dots, N$
 - Apply supervised learning algorithms on $\{x_i, y_i\}$ to learn a model f such that $f(x_i) = y_i$
- Testing or prediction phase
 - Given unseen test data x_i^* , for $i = 1, \dots, T$
 - Use the trained model f to make predictions $f(x_i^*)$'s

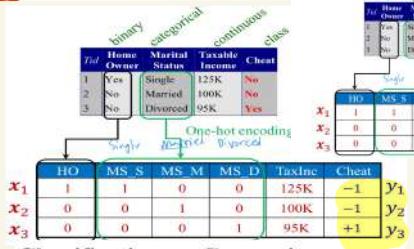
associate with curiosity, exploration
intrinsic motivation (improve itself)
desire of a learning system to acquire knowledge without any external reward or incentives
useful for continuous improve performance/adapt to new enviro
extrinsic motivation (
arise from external rewards/incentives.
designed to maximize a specific objective such as accuracy/performance on particular task.
useful in cases where ultimate goal is to optimize a specific outcome/achieve a particular result
dist: could lead to overfitting/suboptimal behavior if algorithm not properly designed/calibrated

Supervised Learning (cont.)

- In mathematics,
- Given: a set of (\mathbf{x}_i, y_i) for $i = 1, \dots, N$, where $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{im}]$ and y_i is a scalar,
 - Goal: to learn a mapping $f: \mathbf{x} \rightarrow \mathbf{y}$ by requiring $f(\mathbf{x}_i) = y_i$. Prediction model
 - The learned mapping f is expected to be able to make precise predictions on any unseen \mathbf{x}^* as $f(\mathbf{x}^*)$

Examples III

A vector of features



Classification v.s. Regression

- For classification, y is discrete
 - If y is binary, then binary classification
 - If y is nominal not binary, then multi-class classification
- For regression, y is continuous
- Can an example be assigned to more than one categories?

multi-label classification \Rightarrow can be for

multiple class

Common Evaluation Metrics

- Classification: Accuracy, error rate, F-score etc.

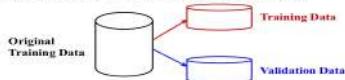
Actual	Prediction	
x1 +1	+1	✓
x2 -1	-1	✓
x3 +1	+1	✓
x4 -1	-1	✗
x5 +1	+1	✗
x6 -1	-1	✓
x7 +1	+1	✓
x8 -1	-1	✓

Accuracy = 5/8
Error rate = 3/8

- Regression: Mean Absolute Error (MAE), Root Mean Squared Error (RMSE)

Evaluation of Supervised Learning

- Recall: the learned predictive model should be able to make predictions on previously unseen data as accurately as possible
- Solution: the whole training data is divided into "training" and "test" sets, with "training" set used to learn the predictive model and "test" set used to validate the performance of the learned model



All Training data

	F1	F2	F3	F4	F5	F6		x_1
$y_1 +1$	1	1	0	0	1	0	...	x_1
$y_2 -1$	0	0	1	1	0	1	...	x_2
$y_N -1$	0	1	0	0	1	1	...	x_N
Used for training							Used for evaluation	
$y_1 +1$	1	0	0	1	0	0	x_1	$y_{l+1} = f(x_1)$
$y_2 -1$	0	1	1	0	1	0	x_2	$y_{l+1} = f(x_2)$
$y_l +1$	1	1	0	1	1	1	x_l	$y_{l+1} = f(x_l)$
Learn	$f: x \rightarrow y$							Evaluate for $i = l+1, \dots, N$
								Some evaluation metric
								$\hat{y}_{l+1} \dots \hat{y}_N$ vs. $y_{l+1} \dots y_N$

Example I



Female, Female, Male, Male

Image processing & Computer Vision tech

-1 -1 +1 +1

$y_1 = -1$ $x_1 = [x_{11}, x_{12}, \dots, x_{1m}]$

$y_2 = -1$ $x_2 = [x_{21}, x_{22}, \dots, x_{2m}]$

$y_3 = 1$ $x_3 = [x_{31}, x_{32}, \dots, x_{3m}]$

$y_4 = 1$ $x_4 = [x_{41}, x_{42}, \dots, x_{4m}]$

Example II

Compact: easy to operate; very good picture quality; looks sharp!

It is also quite blurry in very dark settings. I will never buy HP again

Bag-of-words representation

F1	F2	F3	F4	F5	F6	...	x_1
1	1	0	0	1	0	...	x_1

F1	F2	F3	F4	F5	F6	...	x_2
0	0	1	1	0	1	...	x_2

F1	F2	F3	F4	F5	F6	...	x_N
0	1	0	0	1	1	...	x_N

Dictionary

F1	F2	F3	F4	F5	F6	...	x_1
compact	easy	quite	blurry	good	never_buy	...	x_1

Inductive Learning

- How to learn a predictive model $f(\cdot)$
- Bayesian Classifiers \rightarrow probability
- Decision Trees \rightarrow predict behavior by asking questions
- Artificial Neural Networks
- Support Vector Machines \rightarrow geometry
- Regularized Regression Model

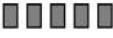
Typical Learning Procedure

Two phases:

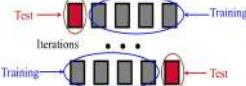
- Training phase
 - Given labeled training data (\mathbf{x}_i, y_i) , for $i = 1, \dots, N$
 - Apply supervised learning algorithms on $\{\mathbf{x}_i, y_i\}$ to learn a model f such that $f(\mathbf{x}_i) = y_i$
- Testing or prediction phase
 - Given unseen test data \mathbf{x}_i^* , for $i = 1, \dots, T$
 - Use the trained model f to make predictions $f(\mathbf{x}_i^*)$'s

Cross Validation

- k-fold cross-validation: partition data into k subsets of the same size

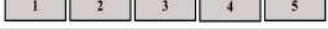


- Hold aside one group for testing and use the rest to build model



5-fold Cross-Validation

Partition into 5 subsets



For a specific setting of hyper-parameters H_3 , e.g., 10 convolutional layers followed by 2 fully connected layers, etc.

- $f(x; \theta_1, H_1)$
- $f(x; \theta_2, H_1)$
- $f(x; \theta_3, H_1)$
- $f(x; \theta_4, H_2)$
- $f(x; \theta_5, H_2)$
- $f(x; \theta_6, H_2)$
- $f(x; \theta_7, H_3)$
- $f(x; \theta_8, H_3)$
- $f(x; \theta_9, H_3)$
- $f(x; \theta_{10}, H_3)$

Evaluation score s_1 for the hyper-parameters setting H_1

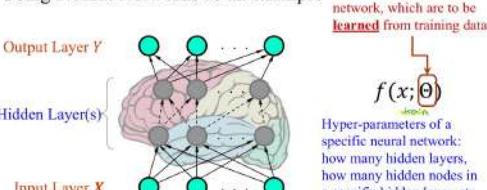
Based on evaluation scores of other hyper-parameters settings, H_2, H_3, \dots , to choose the best one

Once the best hyper-parameters setting H_* is determined, all training data is used to train the neural network $f(x; \theta_*, H_*)$ to deploy it for use

Notes on Cross-Validation

- Using Neural Networks as an example

Parameters of a neural network, which are to be learned from training data



Hyper-parameters of a specific neural network: how many hidden layers, how many hidden nodes in a specific hidden layer, etc.

Note: hyper-parameters need to be manually determined, not learned from training data!

Inductive Learning

- How to learn a predictive model $f(\mathbf{x})$:
 - Bayesian Classifiers
 - Decision Trees
 - Artificial Neural Networks
 - Support Vector Machines
 - Regularized Regression Model

Uncertainty in Prediction

Recall: in supervised learning, given a set of (\mathbf{x}_i, y_i) for $i = 1, \dots, N$, the goal is to learn a mapping $f: \mathbf{x} \rightarrow y$ by requiring $f(\mathbf{x}) = y$.

In many applications, the mapping or relationship f between the input features and the output labels is non-deterministic (uncertain).

For example, suppose you are asked to predict a result of the final of a football cup between Team 1 and Team 2; what team will win?

Notes on Bayesian Classifiers

Why not computing $P(y = c|\mathbf{x})$ for each class directly from the training data, but using Bayes rule to estimate $P(x|y) = c$ and $P(y = c)$ instead?

To estimate $P(y = c|\mathbf{x})$, one needs to consider each combination of values of \mathbf{x} .

E.g., suppose \mathbf{x} is m -dimensional and each feature has binary values, then the total number of possible value combinations of \mathbf{x} is 2^m — require a huge size of training dataset and time consuming, not practical!

The form of $P(\mathbf{x}|y = c)$ can be decomposed based on some assumptions and probability properties — need no to consider all possible combinations.

Bayesian Classifiers

- From a probability point of view, the mapping $f: \mathbf{x} \rightarrow y$ can be modeled as a conditional probability $P(y|\mathbf{x})$.
- Bayesian classifiers aim to learn the mapping $f: \mathbf{x} \rightarrow y$ for supervised learning in the form of conditional probability $P(y|\mathbf{x})$, such that for any input \mathbf{x}^* , one can use $P(y = c|\mathbf{x}^*)$ to predict the probability of \mathbf{x}^* belonging to class c , where $c \in \{0, \dots, C - 1\}$.
- How to estimate $P(y = c|\mathbf{x}^*)$ for different classes?
- How to make use of $P(y = c|\mathbf{x}^*)$ to make a prediction?
- We first review some important probability concepts

Naïve Bayes Classifiers (L3)

- How to estimate $P(\mathbf{x}|y = c)$ from the training data?
- Assume that the features are conditionally independent given the class label:

$$\begin{aligned} P(\mathbf{x}|y = c) &= \prod_{i=1}^d P(x_i|y = c) \\ \text{where } \mathbf{x} &= [x_1, x_2, \dots, x_d] \\ P(x_1, x_2, \dots, x_d|y = c) &= \prod_{i=1}^d P(x_i|y = c) \end{aligned}$$

Bayesian Belief Networks (L4)

- A more general approach to modeling the independence and conditional independence among \mathbf{x} and y , s.t. the computation of $P(\mathbf{x}, y) = P(\mathbf{x}|y)P(y)$ is tractable.
- Use a graphical representation of the probabilistic relationships among features (\mathbf{x}) and output class (y).

Bayes Rule or Bayes Theorem

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

This is induced from product rule:

$$P(A, B) = P(A|B) \times P(B) = P(B|A) \times P(A)$$

Generalized to the case when \mathbf{A} and \mathbf{B} are a set of variables:

$$\begin{aligned} P(A_1 \dots A_k | B_1 \dots B_p) &= \frac{P(B_1 \dots B_p | A_1 \dots A_k)}{P(B_1 \dots B_p)} \\ &= \frac{P(B_1 \dots B_p | A_1 \dots A_k)P(A_1 \dots A_k)}{P(B_1 \dots B_p)} \end{aligned}$$

Define Variables

- Let Y be the random variable that represents the result of the match (0, 1, 2)
 - $Y = 0$: Manchester United wins the match
 - $Y = 1$: Manchester City wins the match
 - $Y = 2$: Draw
- Let X be the random variable that represents the team hosting the match (0 or 1)
 - $X = 0$: Manchester United hosts the match
 - $X = 1$: Manchester City hosts the match
- To estimate $P(Y = 0|X = 0)$, $P(Y = 1|X = 0)$, and $P(Y = 2|X = 0)$

$$\begin{aligned} P(Y = 1|X = 0) &\stackrel{\text{Bayes rule}}{=} \frac{P(X = 0|Y = 1) \times P(Y = 1)}{P(X = 0)} \\ &\quad \boxed{\substack{\text{Prior: } P(X = 0) \\ \text{Sum rule: } P(X) = \sum_i P(X, Y_i)}} \\ &\quad \boxed{\substack{\text{Posterior: } P(X|Y = 1) = P(X, Y = 1) / P(Y = 1) \\ P(Y = 1) = P(Y = 0) + P(Y = 1) + P(Y = 2)}} \\ &= \frac{P(X = 0|Y = 1) \times P(Y = 1)}{P(X = 0|Y = 0) \times P(Y = 0) + P(X = 0|Y = 1) \times P(Y = 1) + P(X = 0|Y = 2) \times P(Y = 2)} \\ &= \frac{0.44 \times 0.3}{0.44 \times 0.3 + 0.54 \times 0.39 + 0.49 \times 0.31} = 0.4945 \end{aligned}$$

Similarly,

$$P(Y = 0|X = 0) = 0.426, \quad P(Y = 2|X = 0) = 0.307$$

Loss or Cost

- Actions: a_{ij} , i.e., predict $y = c_j$, where $c \in \{0, \dots, C - 1\}$
- Define λ_{ij} as the loss/cost of a_{ij} when the optimal action is a_{ij} , i.e., predict $y = i$ while true class label is j
- E.g., in the previous example, $y = 0$: healthy, and $y = 1$: with covid-19 (binary classification)
- We define two corresponding actions, a_{0j} : predict $y = 0$ and a_{1j} : predict $y = 1$, and the losses as

$$\begin{aligned} \lambda_{00} &= 0 \quad \text{predict correctly} \\ \lambda_{11} &= 0 \quad \text{predict correctly} \\ \lambda_{01} &= 1 \quad \text{misclassify as 0} \quad (\text{misclassify with covid-19 as healthy}) \\ \lambda_{10} &= 1 \quad \text{misclassify as 1} \quad (\text{misclassify healthy as covid-19}) \end{aligned}$$

Community outbreak!

Decision based on Expected Risk

- Choose the action with minimum risk:

$$\text{Choose } a^* \text{ if } a^* = \arg \min_{a_c} R(a_c|\mathbf{x})$$

In the covid-19 example,

Expected risk of taking action a_0 (predict as healthy)

$$R(a_0|\mathbf{x}_c) = 1$$

Expected risk of taking action a_1 (predict with covid-19)

$$R(a_1|\mathbf{x}_c) = 0.9$$

Thus, we choose action a_1 : predict patient A is more likely with covid-19

A Special Case

- Making predictions based on maximum posterior is a special case of making decisions based on minimum expected risk

Define the losses as

$$\lambda_{ij} = \begin{cases} 0 & \text{if } i = j \\ 1 & \text{if } i \neq j \end{cases} \quad \text{All correct decisions have no loss and all errors are equally costly}$$

Known as the 0/1 loss

Marginal Probability

- Let A be a random variable (an input feature / class label in machine learning)
- Marginal probability

$$P(A = a) \quad 0 \leq P(A = a) \leq 1$$

refers to the probability that variable $A = a$

$$\sum_a P(A = a_i) = 1$$

Conditional Probability

- Conditional probability:

$$P(B = b|A = a)$$

refers to the probability that the variable B will take on the value b , given that the variable A is observed to have the value a

$$\sum_{b_i} P(B = b_i|A = a) = 1$$

$$P(Y = a) = P(Y = a|x_1, x_2, \dots, x_d) + P(Y = a|x_1, x_2, \dots, x_d)$$

all possibilities of x

Joint Probability

- Let A and B be a pair of random variables (features/labels in machine learning).
- Their joint probability

$$P(B, A = a) = P(A = a, B = b) = P(A = a)P(B = b)$$

refers to the probability that variable $A = a$ and variable $B = b$

Sum Rule

- The connection between joint probability of A and B and marginal probability of A :

$$P(A = a) = \sum_b P(A = a, B = b), \quad \text{OR} \quad P(A) = \sum_b P(A, B)$$

$$P(B = b) = \sum_a P(A = a, B = b), \quad \text{OR} \quad P(B) = \sum_a P(A, B)$$

Product Rule

- The connections between joint, conditional and marginal probabilities for A and B :

$$P(A = a, B = b) = P(B = b|A = a) \times P(A = a)$$

$$= P(A = a|B = b) \times P(B = b)$$

OR

$$P(A, B) = P(B|A) \times P(A) = P(A|B) \times P(B)$$

Bayesian Decision Theory: Summary

- If cost of misclassification on different classes is available, rather than only using posterior probabilities (usually estimated by a Bayesian classifier), Bayesian decision theory provides a way to encode the cost information into decision making.

unknown quantity

Estimate Probabilities (cont.)

- Among the 59 victories for Manchester United, 32 of them come from playing at home

$$P(X = 0|Y = 0) = \frac{32}{59} = 54\%$$

- Among the 45 games won by Manchester City, 20 of them are obtained while playing on Manchester United home ground

$$P(X = 0|Y = 1) = \frac{20}{45} = 44\%$$

- Among the 47 draw games, 23 of them were played on Manchester United home ground

$$P(X = 0|Y = 2) = \frac{23}{47} = 49\%$$

- However, the goal is to estimate:

$$P(Y = 0|X = 0) \quad \text{vs.} \quad P(Y = 1|X = 0) \quad \text{vs.} \quad P(Y = 2|X = 0)$$

Bayesian Classifiers (cont.)

- Bayesian classifiers aim to learn the mapping $f: \mathbf{x} \rightarrow y$ for supervised learning in the form of conditional probability $P(y|\mathbf{x})$ via Bayes rule

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$$

For a classification problem with C classes, given a test instance \mathbf{x}^* , a Bayesian classifier computes

$$P(y = c|\mathbf{x}^*) = \arg \max_{c \in \{0, \dots, C - 1\}} P(y = c|\mathbf{x}^*) = 1$$

Make a prediction based on the maximum posterior

$$y^* = c^* \text{ if } c^* = \arg \max_{c \in \{0, \dots, C - 1\}} P(y = c|\mathbf{x}^*)$$

where $c^* \in \{0, \dots, C - 1\}$ satisfies $P(y = c^*|\mathbf{x}^*) = \max_{c \in \{0, \dots, C - 1\}} P(y = c|\mathbf{x}^*)$

Therefore, we make a prediction based on $y^* = c^*$ if $c^* = \arg \max_{c \in \{0, \dots, C - 1\}} P(y = c|\mathbf{x}^*)$

Take-binary classification as an example: 0 vs 1

$$P(y = 0|\mathbf{x}) = \frac{P(\mathbf{x}|y = 0)P(y = 0)}{P(\mathbf{x})} \quad \text{vs.} \quad P(y = 1|\mathbf{x}) = \frac{P(\mathbf{x}|y = 1)P(y = 1)}{P(\mathbf{x})}$$

Expected Risk

- Expected risk for taking action a_i :

$$R(a_i|\mathbf{x}) = \sum_{j=0}^{C-1} \lambda_{ij} P(y = j|\mathbf{x})$$

Explanation: to estimate a risk of taking an action, one needs to consider all the possible losses

Specifically, taking action a_i predicts it belongs to class i , as the ground truth label of y can be any in the C classes, we need consider all the possible losses λ_{ij} , $j = 0, \dots, C - 1$

We can simply use the average $\frac{\sum_{j=0}^{C-1} \lambda_{ij} P(y = j|\mathbf{x})}{C}$ to estimate its risk

The possibilities of each loss occurring are different because the probability that $y = k$ is always to be different than class i are different

We use the $P(y = k)$ as a weight for each loss λ_{ik} and compute the weighted sum of all possible losses \rightarrow expected risk

$$R(a_i|\mathbf{x}) = \lambda_{i0}P(y = 0|\mathbf{x}) + \lambda_{i1}P(y = 1|\mathbf{x}) + \dots + \lambda_{i(C-1)}P(y = C - 1|\mathbf{x})$$

$$= \lambda_{i0}P(y = 0) + \lambda_{i1}(P(y = 1) + \dots + P(y = C - 1)) + P(y = C - 1|\mathbf{x})$$

$$= \sum_{j=0}^{C-1} P(y = j)P(y = j|\mathbf{x}) = P(y = i|\mathbf{x}) = 1 - P(y \neq i|\mathbf{x})$$

$$= \sum_{j=0}^{C-1} P(y = j|\mathbf{x}) = 1$$

$$P(y = i|\mathbf{x}) = 1$$

- The expected risk of taking action a_i :

$$R(a_i|\mathbf{x}) = 1 - P(y \neq i|\mathbf{x})$$

- Choose an action with minimum expected risk,

$$\text{Choose } a_i \text{ if } R(a_i|\mathbf{x}) = \min R(a_i|\mathbf{x})$$

Equivalent to

$$\text{Predict } y = c^* \text{ if } P(y = c^*|\mathbf{x}) = \max P(y = c|\mathbf{x})$$

Week 3

3ab

$$E(x \rightarrow y) : P(y|x) \quad \text{Training: } P(y|x) \\ E(x \rightarrow y) : P(y|x) \quad \text{Testing: } P(y=c|x)$$

Bayesian Classifiers: Recall

- To learn a prediction function via $P(y|x)$ using Bayes rule

$$P(y=c|x) = \frac{P(x|y=c)P(y=c)}{P(x)} = \frac{P(x|y=c)P(y=c)}{P(x)}$$

- Make predictions based on maximum posterior

$$y^* = c^* \text{ if } c^* = \arg \max_c \frac{P(x|y=c)P(y=c)}{P(x)} \quad \text{the 0/1 loss}$$

$$y^* = c^* \text{ if } c^* = \arg \max_c P(x|y=c)P(y=c)$$

Easy to estimate.

Still difficult to estimate as x contains many input variables, some are discrete, and others are continuous

How Naïve Bayes Classifier Works

$$\text{Naïve Bayes Classifier: } P(x|y=c) = \prod_{i=1}^d P(x_i|y=c) \quad \text{d easier to estimate}$$

- To classify a test record x^* , we need to compute the posteriors for each class by using

$$P(y=c|x^*) = \frac{\left(\prod_{i=1}^d P(x_i^*|y=c)\right)P(y=c)}{P(x^*)}$$

- $P(x^*)$ is constant for each class c , it is sufficient to choose the class that maximizes the numerator

$$P(y=c|x^*) = \left(\prod_{i=1}^d P(x_i^*|y=c)\right)P(y=c)$$

Example of Naïve Bayes Classifier

Naïve Bayes Classifier:

	Home Owner	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	Yes	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	93K	Yes
6	No	Married	40K	No
7	No	Divorced	220K	No
8	Yes	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

For taxable income:

If class=Yes: sample mean=100.00 sample variance=292.5

If class=No: sample mean=75.00 sample variance=87.5

PIClass = No -> 7.10
PIClass = Yes -> 3.19

	Home Owner	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	Yes	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	93K	Yes
6	No	Married	40K	No
7	No	Divorced	220K	No
8	Yes	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

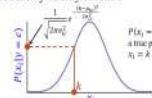
Training data

$$P(y|x)$$

Additional Notes

Probability density function $P(x_i|y=c) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i-\mu_i)^2}{2\sigma_i^2}}$

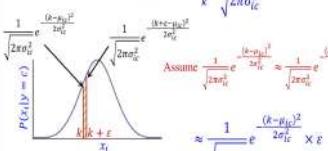
The probability density function is continuous, the probability is defined as the area under the curve of the probability density function



$P(x_i=k|y=c)$ is not a true probability that $x_i = k$ for class c

Instead, we should compute

$$\text{Small positive constant } P(k \leq x_i \leq k + \delta | y=c) = \int_{k}^{k+\delta} \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i-\mu_i)^2}{2\sigma_i^2}} dx_i$$



$$\text{Assume } \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(k+\delta-\mu_i)^2}{2\sigma_i^2}} \approx \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(k-\mu_i)^2}{2\sigma_i^2}}$$

$$\approx \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(k-\mu_i)^2}{2\sigma_i^2}} \times \delta$$

$$\approx \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(k-\mu_i)^2}{2\sigma_i^2}} \times \delta \times P(y=0)$$

Estimate Conditional Probabilities for Discrete Features

#instances in class c , whose values of feature x_i are k

$$P(x_i=k|y=c) = \frac{\#(x_i=k) \wedge (y=c)}{\#(y=c)}$$

A specific value k of the feature x_i

$$P(\text{Status} = \text{Married} | \text{Cheat} = \text{No})$$

$$= \frac{\#(\text{Status} = \text{Married} \wedge \text{Cheat} = \text{No})}{\#(\text{Cheat} = \text{No})}$$

$$P(\text{Home Owner} = \text{Yes} | \text{Cheat} = \text{Yes})$$

$$= \frac{\#(\text{Home Owner} = \text{Yes} \wedge \text{Cheat} = \text{Yes})}{\#(\text{Cheat} = \text{Yes})}$$

= 0 / 3 = 0

Independence

$$P(A|B) = P(A)$$

$P(A|B) = P(A)$

$P(B|A) = P(B)$

$P(A, B) = P(A|B) \times P(B) = P(A) \times P(B)$

$P(A, B) = P(B|A) \times P(A) = P(A) \times P(B)$

- Let A and B be two random variables.
- A is said to be **independent** of B , if the following condition holds:

$$P(A, B) = P(A|B) \times P(B) = P(A) \times P(B)$$

- E.g., let A and B denote the results of two matches in different leagues, knowing the result of one match (e.g., the value of A) does not affect the possibility of result for the other match (i.e., the value of B).

Conditional Independence (cont.)

A more general case:

- Let \mathbf{A} , \mathbf{B} , and \mathbf{C} be three sets of random variables.
- The variables in \mathbf{A} are said to be **conditionally independent** of the variables in \mathbf{B} , given the variables in \mathbf{C} observed, if the following condition holds:

$$P(\mathbf{A}, \mathbf{B}|\mathbf{C}) = P(\mathbf{A}|\mathbf{C}) \times P(\mathbf{B}|\mathbf{C})$$

- The conditional independence between \mathbf{A} and \mathbf{B} given \mathbf{C} can also be written as follows:

$$P(\mathbf{A}, \mathbf{B}|\mathbf{C}) = P(\mathbf{A}|\mathbf{C}) \times P(\mathbf{B}|\mathbf{C})$$

$$P(\mathbf{A}, \mathbf{B}|\mathbf{C}) = \frac{P(\mathbf{A}, \mathbf{B}, \mathbf{C})}{P(\mathbf{C})} \quad \text{Product rule}$$

$$P(\mathbf{A}, \mathbf{U}|\mathbf{C}) = \frac{P(\mathbf{A}, \mathbf{U}, \mathbf{C})}{P(\mathbf{C})} \quad \text{where } \mathbf{U} = \{\mathbf{B}\}$$

$$= \frac{P(\mathbf{A}, \mathbf{B}, \mathbf{C})}{P(\mathbf{B}, \mathbf{C})} \times \frac{P(\mathbf{B}, \mathbf{C})}{P(\mathbf{C})}$$

$$= P(\mathbf{A}|\mathbf{B}, \mathbf{C}) \times P(\mathbf{B}|\mathbf{C}) \quad \text{Product rule}$$

$$= P(\mathbf{A}|\mathbf{C}) \times P(\mathbf{B}|\mathbf{C}) \quad \text{Conditional independence}$$

Illustrative Example

- Consider the problem of predicting whether a loan applicant will repay his/her loan obligation (no cheat or become delinquent cheat)

Predefined categories

$$P(\text{A}|R) = P(\text{A})$$

$$P(\text{A}|B,C) = P(\text{A}|C) \times P(\text{B}|C)$$

$$P(\text{A}_B \text{~and~} \text{A}_B) < P(\text{A}_B) \times P(\text{A}_B)$$

$$(man) = \sum_{i=1}^n P(\text{A}_i)$$

$$(variable) = \sum_{i=1}^n P(\text{A}_i)$$

Estimate Conditional Probabilities for Continuous Features

For continuous features:

- Probability density estimation (more details will be introduced in the 2nd half of the semester):

- Assume the values of a feature given a class label follow a Gaussian distribution (i.e., assume $P(x_i|y=c) \sim \mathcal{N}(\mu_i, \sigma_i^2)$)

- Using training data in the class c to estimate parameters of distribution (e.g., μ_i and σ_i^2)

- Once probability density function is known, we can use it to estimate the conditional probability

- For each class c , assume values of the feature x_i follow a Gaussian distribution:

$$P(x_i|y=c) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_i-\mu_i)^2}{2\sigma_i^2}}$$

The mean of x_i of the training data in class c :

$$\mu_i = \frac{1}{N_c} \sum_{j=1}^{N_c} x_{ij} \quad \sigma_i^2 = \frac{1}{N_c-1} \sum_{j=1}^{N_c} (x_{ij} - \mu_i)^2$$

Value of feature i of the j -th instance in class c :

$$P(x_i|y=c) = \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{(x_{ij}-\mu_i)^2}{2\sigma_i^2}}$$

The estimated Gaussian distribution for $|y=c|$:

$$P(\text{Income}|\text{Cheat} = \text{No}) = \frac{1}{\sqrt{2\pi \cdot 292.54}} e^{-\frac{(120-116)^2}{2 \cdot 292.54}}$$

$P(\text{Income} = 120|\text{No}) = \frac{1}{\sqrt{2\pi \cdot 54.54}} = 0.0072$

Note: in practice, 0.0072 can be used to approximate the probability, but in theory it is much probability

Bayesian Classifiers: Recall

- Estimate $P(y|x)$ via Bayes rule:

$$P(x) = \frac{P(x,y=c)}{P(y=c)} = \frac{P(x|y=c)P(y=c)}{P(x)}$$

- Make predictions based on maximum posterior:

$$y^* = c^* \text{ if } c^* = \arg \max_c P(y=c|x) \quad \text{the 0/1 loss}$$



$$y^* = c^* \text{ if } c^* = \arg \max_c [P(x|y=c)P(y=c)]$$

Easy to estimate

Still difficult to estimate: x contains many input variables. Some are discrete, and others are continuous.

An Example

- Suppose the probability of a person having a specific disease D is 50%.
- There are two medical tests, T_1 and T_2 , of binary values (+ or -). The outcomes of T_2 are perfectly positively correlated with T_1 if a person has the disease, but are independent of T_1 if the person does not have the disease.

T_1, T_2

$+$ $-$

When a person has the disease, the outcomes of T_1 and T_2 are both positive (or negative).

- If a person has the disease, the probabilities of tests T_1 and T_2 being negative are 40%, respectively.
- If a person does not have the disease, the probabilities of T_1 and T_2 being negative are 60% and 65%, respectively.
- If the two tests T_1 and T_2 are both negative for a particular patient, diagnose whether the patient has the disease?

Variables Definition

- Let X_1 denote the outcome of T_1
 - $X_1 = 1$: positive
 - $X_1 = 0$: negative
- Let X_2 denote the outcome of T_2
 - $X_2 = 1$: positive
 - $X_2 = 0$: negative
- Let Y denote whether a person has the disease D
 - $Y = 1$: yes
 - $Y = 0$: no

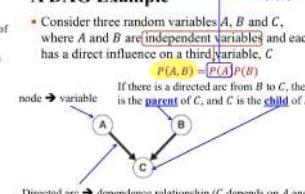
Conditioned
on Y

Bayesian Belief Networks

- A more general approach to modeling the independence and conditional independence among x and y , s.t. the computation of $P(x|y) = P(x|y)P(y)$ is tractable.
- Suppose all features are discrete (if there are both continuous and discrete features, the estimation is much more difficult).
- Representation: A Bayesian network provides a graphical representation of the probabilistic relationships among a set of random variables including features and output class.

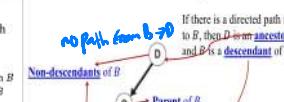
- Two key elements:
 - A directed acyclic graph (DAG) encoding the dependence relationships among a set of variables.
 - A probability table associating each node to its immediate parent nodes

A DAG Example



$$P(A|B) = P(A)$$

Another DAG Example



A Special Case: Naïve Bayes

- A Naïve Bayes classifier can be represented using a special DAG:

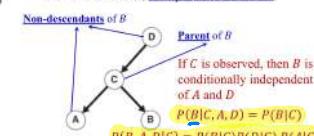


Once a class label is given, all the features are conditionally independent to each other, because each feature is a non-descendant to any other features.

$$P(x_1, x_2, \dots, x_d | y) = \prod_{i=1}^d P(x_i | y)$$

DAG: Conditional Independence

- Property (conditional independence): a node in a Bayesian network is conditionally independent of its non-descendants, if its parents are known.



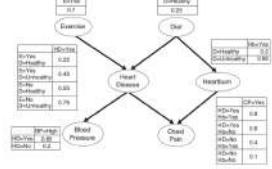
BBN Representation

- Besides the conditional independence conditions imposed by the network topology, each node is also associated with a probability table.
- If a node X does not have any parents, then the table contains only the prior probability $P(X)$
- If a node X has only one parent, Z , then the table contains the conditional probability $P(X|Z)$.
- If a node X has multiple parents, (Z_1, Z_2, \dots, Z_k) , then the table contains the conditional probability $P(X|Z_1, Z_2, \dots, Z_k)$

Suppose Heart Disease is our target variable to make prediction (i.e., output) the other variables are input features, whose values can be observed or missing.

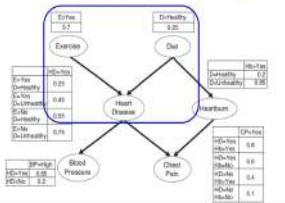
Inference: Example 1

- Without any additional information, to determine whether the person is likely to have heart disease.



$$P(HD=Yes) = \sum_{\alpha} \sum_{\beta} P(HD=Yes | E=\alpha, D=\beta)$$

Sum Rule



$$P(HD=Yes) = \sum_{\alpha} \sum_{\beta} P(HD=Yes | E=\alpha, D=\beta)$$

$$\text{Product Rule} = \sum_{\alpha} \sum_{\beta} P(HD=Yes | E=\alpha, D=\beta) P(E=\alpha) P(D=\beta)$$

Independence

$$\alpha = \{\text{Yes, No}\} \quad \beta = \{\text{Healthy, Unhealthy}\}$$

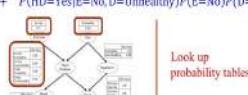
$$P(HD=Yes) = \sum_{\alpha} \sum_{\beta} P(HD=Yes | E=\alpha, D=\beta) P(E=\alpha) P(D=\beta)$$

$$= P(HD=Yes | E=Yes, D=Healthy) P(E=Yes) P(D=Healthy)$$

$$+ P(HD=Yes | E=Yes, D=Unhealthy) P(E=Yes) P(D=Unhealthy)$$

$$+ P(HD=Yes | E=No, D=Healthy) P(E=No) P(D=Healthy)$$

$$+ P(HD=Yes | E=No, D=Unhealthy) P(E=No) P(D=Unhealthy)$$



$$= 0.25 \times 0.7 \times 0.25 + 0.45 \times 0.7 \times 0.75 + 0.55 \times 0.3 \times 0.25 + 0.75 \times 0.3 \times 0.75 = 0.49$$

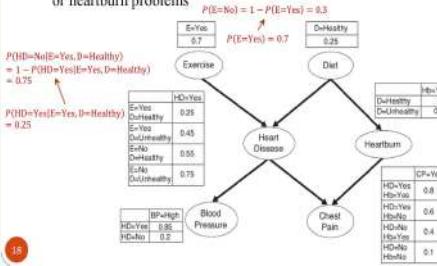
$$P(HD=Yes) = 0.49$$

$$P(HD=No) = 1 - P(HD=Yes) = 0.51$$

- Therefore, the person has a slightly higher chance of not getting the heart disease.

BBN Representation: Example

A Bayesian network for modeling patients with heart disease or heartburn problems



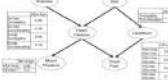
Inference: Example 2

- If the person has high blood pressure, chest pain and heartburn, but does regular exercise and eats a healthy diet, to diagnose whether the patient has heart disease:

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

VS.

$$P(HD=No | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$



$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes) / P(BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

How to estimate the joint probability in the numerator?

Keep in mind that the goal is to rewrite the joint probability into an equivalent form, such that for all the probabilities in the rewritten equivalent form, their values can be found from the tables in the BBN.

There are many ways to rewrite the above joint probability, e.g.,

$$P(D, E | HD, BP, Hb, CP) P(HD | BP, BP, CP), \text{ or } P(D, CP | E, Hb, HD) P(E, Hb, HD), \text{ many others}$$

Which one is useful for the joint probability simplification?

From the BBN, we found that the variables BP and CP are child nodes of HD and Hb, but not a parent node for any other variables. That means we could not find any conditional probabilities in the tables, which involves BP and CP as the condition, like $P(D, E | HD, BP, Hb, CP)$. In other words, if in the rewritten equivalent form, there is a probability where BP or (and) CP is (are) in the condition, we still have to further rewrite it such that BP and CP are not in the condition of any conditional probability term.

The above analysis motivates us to transform the joint probability (numerator) using the following form based on the product rule:

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(BP=High, CP=Yes | HD=Yes, Hb=Yes, E=Yes, D=Healthy, E=Yes) P(HD=Yes | E=Yes, Hb=Yes, D=Healthy)$$

Not in the condition of the conditional probability

$$\text{Denote by } U = \{BP, CP\} \text{ and } V = \{HD, Hb, E, D\}$$

$$P(U, V) = P(U)P(V)$$

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes) / P(BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$\text{Denote by } U = \{BP, CP\} \text{ and } V = \{HD, Hb, E, D\}$$

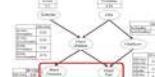
$$P(U, V) = P(U)P(V)$$

$$P(HD=Yes | BP=High, CP=Yes | HD=Yes, Hb=Yes, E=Yes, D=Healthy)$$

$$= P(HD=Yes | BP=High, CP=Yes | HD=Yes, Hb=Yes, E=Yes, D=Healthy) P(HD=Yes | E=Yes, Hb=Yes, D=Healthy)$$

$$= 0.85 \times 0.8 \times 0.7 \times 0.25 = 0.08875$$

$$\text{Slide 29}$$



$$\text{Slide 30}$$

BBN Model Building

- Two steps in the training phase:
 - Creating the structure of the network, i.e., DAG
 - Network topology can be obtained by encoding the subjective knowledge of domain experts
 - Or can be learned from data (structure learning) — still an open problem
- Estimating the probability values in the table associated with each node
 - Counting based on the definition of the corresponding probabilities
- Note: in this module we only focus on how to use a BBN to make predictions (or inference)

Recall: If A and B are conditionally independent given C, we have

$$P(A|B, C) = P(A|C)P(B|C) \quad P(A|B, C) = P(A|C)P(B|C)$$

$$1 \quad P(BP=High, CP=Yes | HD=Yes, Hb=Yes, E=Yes, D=Healthy)$$

- E and Hb are non-descendants of BP and CP
- Therefore, BP and CP are conditionally independent given HD and Hb

$$P(BP=High, CP=Yes | HD=Yes, Hb=Yes)$$

- HD is parent of BP
- BP is a non-descendant of CP
- Therefore, BP and CP are conditionally independent given HD and Hb

$$= P(BP=High, HD=Yes, Hb=Yes) P(CP=Yes | HD=Yes, Hb=Yes)$$

- HD is parent of BP
- BP is a non-descendant of CP
- Therefore, BP and CP are conditionally independent given HD

$$= P(BP=High) HD=Yes P(CP=Yes | HD=Yes, Hb=Yes)$$

$$= 0.85 \times 0.8 \times 0.68 = 0.08875$$

$$2 \quad P(HD=Yes, Hb=Yes, E=Yes, D=Healthy)$$

Denote by $U = \{HD, Hb\}$ and $V = \{D, E\}$

$$= P(HD=Yes | Hb=Yes, E=Yes, D=Healthy) P(HD=Yes | Hb=Yes, D=Healthy)$$

Gives E and Hb are conditionally independent given HD and Hb

$$P(HD=Yes | E=Yes, D=Healthy) P(HD=Yes | E=Yes, D=Healthy)$$

Gives D and Hb are conditionally independent given HD

$$P(HD=Yes | E=Yes, D=Healthy) P(HD=Yes | D=Healthy)$$

$$= 0.25 \times 0.7 \times 0.25 = 0.08875$$

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(HD=Yes | BP=High, CP=Yes | Hb=Yes, E=Yes, D=Healthy, E=Yes) P(HD=Yes | Hb=Yes, E=Yes, D=Healthy)$$

$$= 0.85 \times 0.8 \times 0.7 \times 0.25 = 0.08875$$

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(HD=Yes | BP=High, CP=Yes | Hb=Yes, E=Yes, D=Healthy, E=Yes) P(HD=Yes | Hb=Yes, E=Yes, D=Healthy)$$

$$= 0.85 \times 0.8 \times 0.7 \times 0.25 = 0.08875$$

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(HD=Yes | BP=High, CP=Yes | Hb=Yes, E=Yes, D=Healthy, E=Yes) P(HD=Yes | Hb=Yes, E=Yes, D=Healthy)$$

$$= 0.85 \times 0.8 \times 0.7 \times 0.25 = 0.08875$$

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

$$= P(HD=Yes | BP=High, CP=Yes | Hb=Yes, E=Yes, D=Healthy, E=Yes) P(HD=Yes | Hb=Yes, E=Yes, D=Healthy)$$

$$= 0.85 \times 0.8 \times 0.7 \times 0.25 = 0.08875$$

$$P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes) > P(HD=Yes | BP=High, Hb=Yes, CP=Yes, D=Healthy, E=Yes)$$

The person has a higher chance of getting the heart disease.

Inference: Example 3

- If the person has high blood pressure, but exercises regularly and eats a healthy diet, to diagnose about heart disease (estimate the probabilities).

$P(HD=Yes|BP=High, D=Healthy, E=Yes)$

vs.

$P(HD=No|BP=High, D=Healthy, E=Yes)$



Tutorial



Hint: Using BBNs for Inference

- Given a BBN, and an inference (prediction) task:
 - Translate the problem into a probabilistic language, i.e., what probabilities to be estimated?
- If the probabilities to be estimated cannot be obtained from the probability tables of the BBN, then
 - Identify a subgraph which captures the dependence between input variables (features) and the output variable (class)
- Based on the **network topology**, apply **product rule**, **sum rule** and the **properties of conditional independence and independence** to induce equivalent forms of the probabilities until all probabilities can be found from the probability tables

Bayesian Belief Networks: Summary

- BBNs use a (directed) graphical model to model dependence among variables
 - Other directed graphical models: Hidden Markov Models, Dynamic Bayesian Networks, etc.
- Undirected graphical models: e.g., Markov Random Fields, Conditional Random Fields, etc.
- Network structure construction is difficult
 - Use domain knowledge – not complete, may not accurate
 - Learn structure from data – computationally expensive, greedy algorithm → not optimal

Weeks 5 Decision Tree

Illustrative Example

- Consider the problem of predicting whether a loan applicant will repay his/her loan obligation (no/cheat) or become delinquent (cheat).

Predefined categories



High-level Algorithm

ID	Home Owner	Marital Status	Taxable Income	Cheat
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	20K	No
5	No	Divorced	95K	Yes
6	Yes	Married	60K	No
7	Yes	Divorced	20K	No
8	No	Single	65K	Yes
9	No	Divorced	20K	Yes
10	No	Single	90K	No

ID	Home Owner	Marital Status	Taxable Income	Cheat
11	Yes	Single	95K	No
12	No	Married	100K	No
13	No	Single	70K	No
14	Yes	Married	20K	No
15	No	Divorced	95K	Yes
16	Yes	Married	60K	No
17	Yes	Divorced	20K	No
18	No	Single	65K	Yes
19	No	Married	70K	No
20	No	Single	90K	No

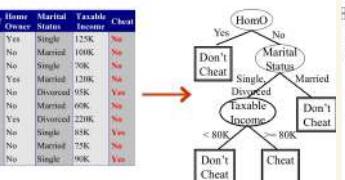
ID	Home Owner	Marital Status	Taxable Income	Cheat
21	Yes	Single	125K	No
22	No	Married	100K	No
23	No	Single	70K	No
24	Yes	Married	20K	No
25	No	Divorced	95K	Yes
26	Yes	Married	60K	No
27	Yes	Divorced	20K	No
28	No	Single	65K	Yes
29	No	Married	70K	No
30	No	Single	90K	No

ID	Home Owner	Marital Status	Taxable Income	Cheat
31	Yes	Single	125K	No
32	No	Married	100K	No
33	No	Single	70K	No
34	Yes	Married	20K	No
35	No	Divorced	95K	Yes
36	Yes	Married	60K	No
37	Yes	Divorced	20K	No
38	No	Single	65K	Yes
39	No	Married	70K	No
40	No	Single	90K	No

ID	Home Owner	Marital Status	Taxable Income	Cheat
41	Yes	Single	125K	No
42	No	Married	100K	No
43	No	Single	70K	No
44	Yes	Married	20K	No
45	No	Divorced	95K	Yes
46	Yes	Married	60K	No
47	Yes	Divorced	20K	No
48	No	Single	65K	Yes
49	No	Married	70K	No
50	No	Single	90K	No

ID	Home Owner	Marital Status	Taxable Income	Cheat
51	Yes	Single	125K	No
52	No	Married	100K	No
53	No	Single	70K	No
54	Yes	Married	20K	No
55	No	Divorced	95K	Yes
56	Yes	Married	60K	No
57	Yes	Divorced	20K	No
58	No	Single	65K	Yes
59	No	Married	70K	No
60	No	Single	90K	No

ID	Home Owner	Marital Status	Taxable Income	Cheat
61	Yes	Single	125K	No
62	No	Married	100K	No
63	No	Single	70K	No
64	Yes	Married	20K	No
65	No	Divorced	95K	Yes
66	Yes	Married	60K	No
67	Yes	Divorced	20K	No
68	No	Single	65K	Yes
69	No	Married	70K	No
70	No	Single	90K	No



Decision Tree Classifiers: Summary

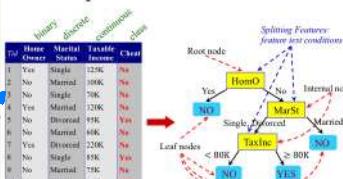
- Easy to interpret
- Efficient in both training and testing
- Effective for datasets with a lot of categorical features
- Used as a base classifier in many ensemble learning approaches (will be introduced in the 2nd half)

Motivation of Decision Trees

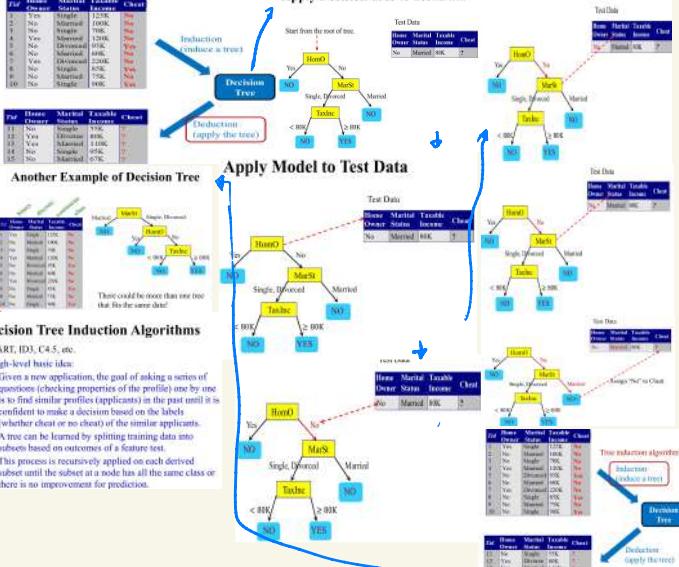
- Suppose a new applicant submits a loan application. How do we decide whether to approve or reject the application?

- To pose a series of questions about the profile of the applicant:
 - Whether the applicant is a home owner? If yes, then he/she may repay the loan obligation with a high probability.
 - After that, we may ask a follow-up question: what is the applicant's marital status? ...

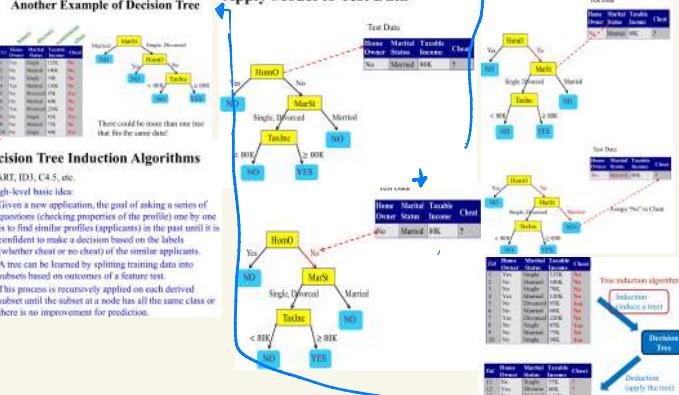
Example of a Decision Tree



Apply Decision Tree to Test Data



Another Example of Decision Tree



Decision Tree Induction Algorithms

- CART, ID3, C4.5, etc.
- High-level basic idea:
 - Given a new application, the goal of asking a series of questions (checking properties of the profile) one by one is to find similar profiles (applicants) in the past until it is confident to make a prediction based on the labels (whether cheat or no cheat) of similar applicants.
 - A tree can be learned by splitting training data into subsets based on outcome of a feature test.
 - This process is recursively applied on each derived subset until the subset at a node has all the same class or there is no improvement for prediction.

Tree Induction

- Determine how to split the data
 - How to specify the feature test condition?
 - How to determine the best split?
 - Determine when to stop splitting

Stopping Criteria for Tree Induction

- Stop expanding a node when all the data instances belong to the same class
 - Ideal case but not always possible
- Stop expanding a node when all the data instances have similar feature values
 - Early termination
 - Useful to avoid the overfitting issue (will be introduced next week)

Tree Induction

- Greedy strategy
 - Split the records based on a feature test that optimizes certain criterion
- Issues
 - Determine how to split the records
 - How to specify the feature test condition?
 - How to determine the best split?
 - Determine when to stop splitting

Splitting Based on Discrete Features (more than two distinct values)

- Multi-way split: Use as many partitions as distinct values
 - Marital Status: Single, Divorced, Married
- Binary split: Divide values into two subsets. Need to find optimal partitioning
 - Marital Status: (Married, Divorced) OR (Married, Single) OR (Divorced, Married, Single)

How to Determine the Best Split

Before Splitting:
 Before: 10 records of class 0
 After: 10 records of class 1

After Splitting:
 Left child: 4 records of class 0, 6 records of class 1
 Right child: 6 records of class 0, 4 records of class 1

Which condition test is the best?

- An intuitive idea:
 - Nodes with homogeneous class distribution are preferred
 - Need a measure of node impurity
- A split criterion is defined in terms of the difference in degree of node impurity before and after splitting
 - Non-homogeneous, the degree of impurity
 - Homogeneous, the degree of impurity

High-level Algorithm: Summary

- Let D_i be the set of training records that reach a node i .
- General Procedure:
 - If D_i contains records that belong to the same class y_i , then i is a leaf node labeled in y_i (e.g., Yes or No).
 - Otherwise if D_i is an empty set, then i is a leaf node labeled by its default class, y_d (e.g., No).
 - Otherwise D_i is partitioned into two or more smaller sets. Then a feature is selected to conduct condition test to split the data into smaller subsets.
 - A child node is created for each partition of the test condition and the records in D_i are distributed to the children based on the outcomes
 - Reursively apply the procedure to each subset

How to Specify Test Condition?

- Depends on feature types
 - Discrete
 - Continuous
- Depends on number of ways to split
 - Binary split
 - Multi-way split

Splitting Based on Binary Features

- Generate two potential outcomes
 - (a) Binary split
 - (b) Multi-way split

(a) Binary split
 Taxable Income >= 80K?
 Yes → 10 records, No → 20 records

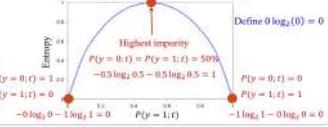
(b) Multi-way split
 Taxable Income?
 < 80K or ≥ 80K
 < 80K → 10 records, ≥ 80K → 20 records, ≥ 80K → 30 records

Discretization

- Consider all possible splits and find the best cut
- Can be very computationally intensive

Measure of Impurity: Entropy

- Entropy at a given node t :
$$\text{Entropy}(t) = -\sum_c P(y=c; t) \log_2 P(y=c; t)$$
- Entropy for a given node t , with binary classes:



Examples of Computing Entropy

$$\text{Entropy}(t) = -\sum_c P(y=c; t) \log_2 P(y=c; t)$$

Y_1	0
Y_2	1

$P(Y_1 = 0) = \frac{0}{6} = 0$ $P(Y_2 = 1) = \frac{6}{6} = 1$
 $\text{Entropy} = -0 \log_2(0) - 1 \log_2(1) = -0 - 0 = 0$

Y_1	1
Y_2	5

$P(Y_1 = 1) = \frac{1}{6}$ $P(Y_2 = 5) = \frac{5}{6}$
 $\text{Entropy} = -\left(\frac{1}{6}\right) \log_2\left(\frac{1}{6}\right) - \left(\frac{5}{6}\right) \log_2\left(\frac{5}{6}\right) = 0.65$

Y_1	2
Y_2	4

$P(Y_1 = 2) = \frac{2}{6}$ $P(Y_2 = 4) = \frac{4}{6}$
 $\text{Entropy} = -\left(\frac{2}{6}\right) \log_2\left(\frac{2}{6}\right) - \left(\frac{4}{6}\right) \log_2\left(\frac{4}{6}\right) = 0.92$



Information Gain: Motivation

Recall: a split criterion should be defined in terms of the difference in degree of node impurity before and after splitting

Information Gain:

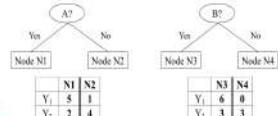
$$\Delta_{\text{info}} = \text{Entropy}(\text{parent node}) - \text{Entropy}(\text{children nodes})$$

Choose a feature whose condition test maximizes the gain (minimizes the weighted average impurity measures of the child nodes)

Information Gain: Practice

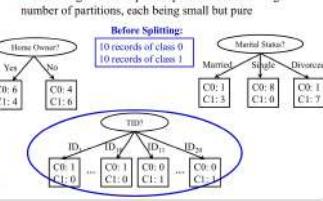
Suppose features A and B are binary

Based on Information Gain, which feature should be selected to split data?



Information Gain: Limitation

Disadvantage: tends to prefer splits that result in large number of partitions, each being small but pure



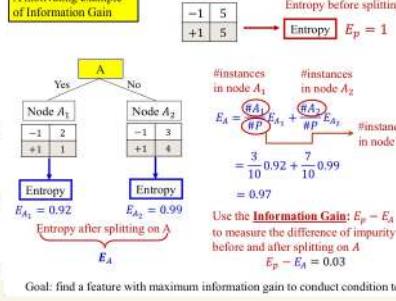
An Example

```
>>> from sklearn import tree
>>> dtC = tree.DecisionTreeClassifier()
>>> dtC.fit(X, y)
>>> predC = dtC.predict(X_t)
```

Entropy Properties

- Entropy at a given node t :
- $$\text{Entropy}(t) = -\sum_c P(y=c; t) \log_2 P(y=c; t)$$
- Total number of all possible values of y , i.e., #classes
- Maximum: $\log_2 C$ when records are equally distributed among all classes
 - Minimum: 0 when all records belong to one class

A motivating example of Information Gain



Use the **Information Gain**: $E_p - E_A$ to measure the difference of impurity before and after splitting on A

$$E_p - E_A = 0.08$$

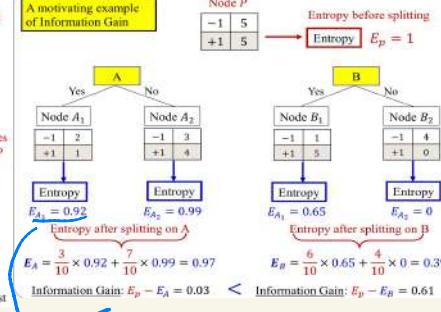
Goal: find a feature with maximum information gain to conduct condition test

Splitting Based on Entropy (cont.)

Gain Ratio:
$$\Delta_{\text{infoR}} = \frac{\Delta_{\text{info}}}{\text{SplitInfo}}$$
 where $\text{SplitInfo} = \sum_{i=1}^p \frac{n_i}{n} \log_2 \left(\frac{n_i}{n} \right)$

- Parent node t of n instances is split into p partitions (children), n_i is the number of instances in partition i
- Higher entropy partitioning (large number of small partitions) is penalized!

Refer to Question 1 in tutorial for practice



$$-\frac{2}{3} \log_2 \left(\frac{2}{3} \right) - \left(\frac{1}{3} \right) \log_2 \left(\frac{1}{3} \right)$$

Information Gain: Definition

- Suppose a parent node t is split into P partitions (children)
 - Information Gain:
- $$\Delta_{\text{info}} = \text{Entropy}(\text{parent}) - \sum_{j=1}^p \frac{n_j}{n} \text{Entropy}(j)$$
- Number of examples at node t
- To choose a feature whose condition test maximizes the gain (minimizes the weighted average impurity measures of the children nodes)

$$\text{Entropy}(\text{Parent}) = -\left(\frac{6}{12}\right) \log_2\left(\frac{6}{12}\right) - \left(\frac{6}{12}\right) \log_2\left(\frac{6}{12}\right) = 1$$

$$\text{Entropy}(N1) = -\left(\frac{5}{7}\right) \log_2\left(\frac{5}{7}\right) - \left(\frac{2}{7}\right) \log_2\left(\frac{2}{7}\right) = 0.8631$$

$$\text{Entropy}(N2) = -\left(\frac{1}{5}\right) \log_2\left(\frac{1}{5}\right) - \left(\frac{4}{5}\right) \log_2\left(\frac{4}{5}\right) = 0.7219$$

$$\text{Entropy}(\text{Split}_A) = \left(\frac{7}{12}\right) \times 0.8631 + \left(\frac{5}{12}\right) \times 0.7219 = 0.8043$$

$$\text{Entropy}(N3) = -\left(\frac{3}{9}\right) \log_2\left(\frac{3}{9}\right) - \left(\frac{6}{9}\right) \log_2\left(\frac{6}{9}\right) = 0.9183$$

$$\text{Entropy}(N4) = -\left(\frac{0}{3}\right) \log_2\left(\frac{0}{3}\right) - \left(\frac{3}{3}\right) \log_2\left(\frac{3}{3}\right) = 0$$

$$\text{Entropy}(\text{Split}_B) = \left(\frac{9}{12}\right) \times 0.9183 + \left(\frac{3}{12}\right) \times 0 = 0.6887$$

Parent
Y_1 6
Y_2 6

N1	N2
Y_1 5 1	
Y_2 2 4	

N3	N4
Y_1 6 0	
Y_2 3 3	

N1	N2
Y_1 6 0	
Y_2 3 3	

Split on A

$$\text{Entropy}(\text{Split}_A) = \left(\frac{7}{12}\right) \times 0.8631 + \left(\frac{5}{12}\right) \times 0.7219 = 0.8043$$

Split on B

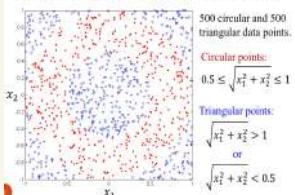
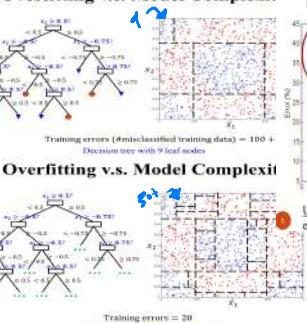
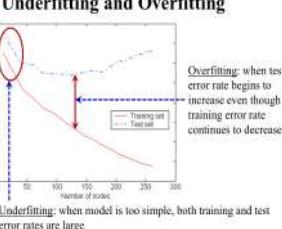
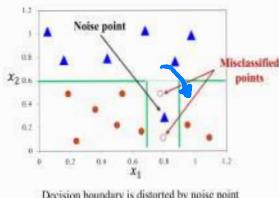
$$\text{Entropy}(\text{Split}_B) = \left(\frac{9}{12}\right) \times 0.9183 + \left(\frac{3}{12}\right) \times 0 = 0.6887$$

$$\Delta_{\text{info}}(A) = 1 - 0.8043 = 0.1957$$

$$\Delta_{\text{info}}(B) = 1 - 0.6887 = 0.3113$$

Choose B to conduct condition test to split data

$$\begin{aligned} & \text{Entropy} : \\ & P(N_1) = \frac{7}{12}, \quad P(N_2) = \frac{5}{12} \\ & E_{N1} = -\frac{7}{12} \log_2\left(\frac{7}{12}\right) - \frac{5}{12} \log_2\left(\frac{5}{12}\right) = 0.863 \\ & E_{N2} = -\frac{5}{12} \log_2\left(\frac{5}{12}\right) - \frac{7}{12} \log_2\left(\frac{7}{12}\right) = 0.7219 \\ & E_{\text{Split}_A} = \frac{7}{12} \times 0.863 + \frac{5}{12} \times 0.7219 = 0.8043 \\ & \text{Entropy} : \\ & P(N_3) = \frac{9}{12}, \quad P(N_4) = \frac{3}{12} \\ & E_{N3} = -\frac{9}{12} \log_2\left(\frac{9}{12}\right) - \frac{3}{12} \log_2\left(\frac{3}{12}\right) = 0.9183 \\ & E_{N4} = -\frac{3}{12} \log_2\left(\frac{3}{12}\right) - \frac{9}{12} \log_2\left(\frac{9}{12}\right) = 0 \\ & E_{\text{Split}_B} = \frac{9}{12} \times 0.9183 + \frac{3}{12} \times 0 = 0.6887 \\ & \Delta_{\text{info}}(A) = 1 - 0.8043 = 0.1957 \\ & \Delta_{\text{info}}(B) = 1 - 0.6887 = 0.3113 \\ & \text{Choose } B \text{ to conduct condition test to split data} \\ & \text{Entropy} : \\ & P(\text{Root}) = \frac{6}{12}, \quad P(C_0=1) = \frac{5}{12}, \quad P(C_0=0) = \frac{6}{12} \\ & E_{\text{Root}} = -\frac{6}{12} \log_2\left(\frac{6}{12}\right) - \frac{5}{12} \log_2\left(\frac{5}{12}\right) - \frac{6}{12} \log_2\left(\frac{6}{12}\right) = 0.6887 \\ & \Delta_{\text{info}}(\text{Root}) = 1 - 0.6887 = 0.3112 \quad (\text{B}) \end{aligned}$$

Underfitting and Overfitting (Example)**Overfitting v.s. Model Complexity****Underfitting and Overfitting****Overfitting due to Noise (cont.)****Estimation of Generalization Errors**

- **Training errors:** error on the training set: $e(T)$
- **Generalization errors:** error on previously unseen testing set: $e'(T)$
- Approaches to estimating generalization errors:
 - **Optimistic Estimate:** $e'(T) = e(T)$
 - Incorporating model complexity
 - Occam's Razor
 - Pessimistic Error Estimate
 - Using Validation Set

Using a Validation Set

- Divide the original training data set into two smaller subsets
- One is for training, the other (known as the validation set) is for estimating the generalization error
- The complexity of the best model can be estimated based on the performance of the model on the validation set

How to Address Overfitting**Pre-Pruning (Early Stopping Rule)**

- Stop the algorithm before it becomes a fully-grown tree
- Typical stopping conditions for a node:
 - Stop if all instances belong to the same class
 - Stop if all the feature values are the same
- Pre-Pruning (Early Stopping Rule)
 - More restrictive conditions:
 - Stop if number of instances is less than some user-specified threshold
 - Stop if expanding the current node does not improve generalization errors

Post-pruning

- Grow decision tree to its entirety
- Trim the nodes of the decision tree in a bottom-up fashion
- If generalization error improves after trimming, replace sub-tree by a new leaf node
- Class label of leaf node is determined from majority class of instances in the sub-tree



criterion : "gini"; entropy : "gini"; gini

The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

max_depth : int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves contain less than min_samples_split samples.

max_features : int, float or "auto", "sqrt", "log2", default=None

The number of features to consider when looking for the best split.

If int, then consider max_features features at each split.

If float, then max_features is a fraction and int(max_features * n_features) features are considered at each split.

If "auto", then max_features=n_features.

If "sqrt", then max_features=sqrt(n_features).

If "log2", then max_features=log2(n_features).

If None, then max_features=n_features.

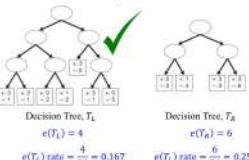
Optimistic Estimate

- Assume that the training set is a good representation of the overall data
- The training error can be used to provide an optimistic estimate for the generalization error
 - $e'(T) = e(T)$
- A decision tree induction algorithm selects the model that produces the lowest training error rate

Occam's Razor

- Definition: Given two models of similar generalization errors, one should prefer the simpler model over the more complex model.
- For complex models, there is a greater chance that it was fitted accidentally by errors in data.
- Therefore, one should include model complexity when evaluating a model.

"Everything should be made as simple as possible, but no simpler." — Albert Einstein

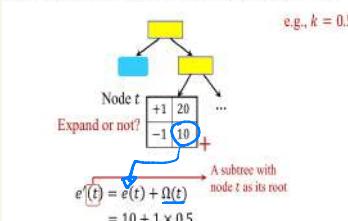
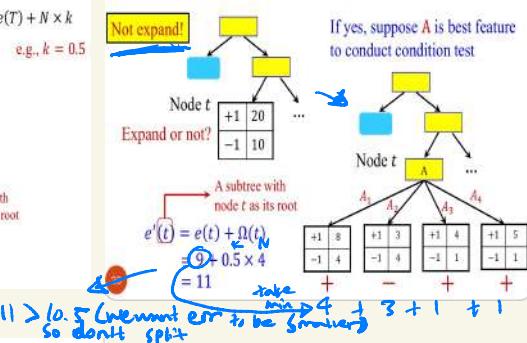
An Example of Optimistic Estimate**Pessimistic Error Estimate**

- Idea: explicitly computes generalization error as the sum of training error and a penalty term for model complexity
- $e'(T) = e(T) + \Omega(T)$
- In a decision tree, we can define a penalty term of $k > 0$ on each leaf node, i.e., $e'(t) = e(t) + \Omega(t) = e(t) + k$
- Then,

$$e'(T) = e(T) + \sum_{\text{leaf nodes}} k$$

Total number of leaf nodes

$k > 0$, e.g., $k = 0.5$

Pre-Pruning ExampleGeneralization errors: $e'(T) = e(T) + \Omega(T) = e(T) + N \times k$ Generalization errors: $e'(T) = e(T) + N \times 0.5$ **Example of Post-Pruning**

Pessimistic error

 $e'(T) = e(T) + N \times 0.5$

PRUNE?

Class = Yes	8	Class = Yes	3	Class = Yes	4	Class = Yes	5
Class = No	4	Class = No	4	Class = No	1	Class = No	1

Training errors (before pruning) = 4 + 3 + 1 + 1 = 9

Pessimistic errors (before pruning) = 9 + 4 × 0.5 = 11

Training errors (after pruning) = 10

Pessimistic errors (after pruning) = 10 + 0.5 = 10.5

PRUNE!

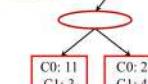
- Pessimistic error?

 $e'(T) = e(T) + N \times 0.5$

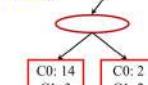
PRUNE?

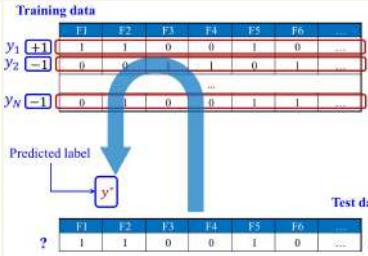
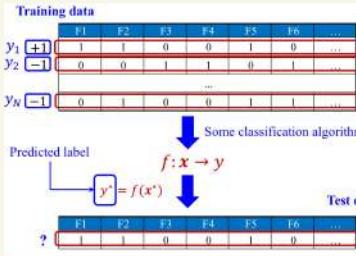
Tutorial

Case 1:



Case 2:





Summary of NN Classifier

- The K-NN classifier is a lazy learner.
 - It does not build models explicitly.
 - "Training" is very efficient.
 - Classifying unknown test instances is relatively expensive.

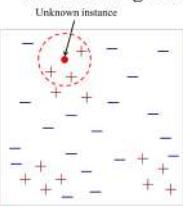
Inductive Learning

Instance Based Classifiers

K-Nearest Neighbors classifier

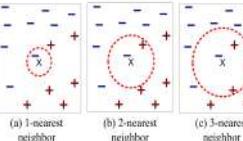
- Uses K "closest" points (nearest neighbors) for performing classification

Nearest-Neighbor Classifiers



- Requires three things
 - The set of stored labeled instances
 - Distance metric to compute distance between instances
 - The value of K , the number of nearest neighbors to retrieve
- To classify an unknown instance
 - Compute distance to all the training instances
 - Identify K nearest neighbors
 - Use class labels of nearest neighbors to determine the class label of the unknown instance (e.g., by taking majority vote)

Definition of Nearest Neighbors



K -nearest neighbors of an instance x are data points that have the K smallest distance to x .

Determine Class Label

- Determine the class from nearest neighbor list
 - Take the majority vote of class labels among the K -nearest neighbors
- Given test data x^* , majority voting
 - Indicator function that returns 1 if an input is true, otherwise 0
- $y^* = \arg \max_c \sum_{(x_i, y_i) \in N_{x^*}} I(c = y_i)$
- Every neighbor has the same impact on the classification
- This indeed makes the algorithm sensitive to the choice of K

Normalization

Min-max normalization: to $[min_{new}, max_{new}]$

- Example: To normalize income ranging from \$12,000 to \$98,000 to [0.0, 1.0], what is the value for \$73,600 after normalization?

$$x_{new} = \frac{x_{old} - min_{old}}{max_{old} - min_{old}} (max_{new} - min_{new}) + min_{new}$$

$$73,600 \rightarrow \frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

Standardization (z-score normalization) (μ : mean, σ : standard deviation):

$$v_{new} = \frac{v_{old} - \mu_{old}}{\sigma_{old}} \rightarrow \mu_{new} = 0, \text{ and } \sigma_{new} = 1$$

- Example: Let $\mu = 54,000$, $\sigma = 16,000$. What is the value for \$73,600 after standardization?

$$\frac{73,600 - 54,000}{16,000} = 1.225$$

Implementation

```
>>> from sklearn.neighbors import KNeighborsClassifier
```

...
set number of neighbors

```
>>> knnC = KNeighborsClassifier(n_neighbors=3)
>>> knnC.fit(X, y)
>>> pred = knnC.predict(X)
```

Build indices s.t. it is more efficient when making predictions on test data

Lazy Learning

Distance Metric

- Compute distance between two data points in a d -dimensional space:

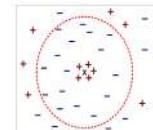
$$x_k | x_2 | \dots | x_d$$

$$x_1 | x_3 | \dots | x_d$$

Euclidean distance

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^d (x_{ik} - x_{jk})^2} = \left\| (x_i - x_j) \right\|_2$$

ℓ_2 norm of x_i : $\|x_i\|_2 = \sqrt{\sum_{k=1}^d x_{ik}^2}$



Value of K

- Choosing the value of K :
 - If K is too small, sensitive to noise points
 - If K is too large, neighborhood may include points from other classes

Other Issues

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^d (x_{ik} - x_{jk})^2}$$

- Scaling issues
 - Feature may need to be scaled to prevent distance from being dominated by some features
- Example:
 - height of a person may vary from 1.5m to 1.8m
 - weight of a person may vary from 40kg to 200kg
 - income of a person may vary from \$10K to \$1M
- Solution: normalization on features of different scales

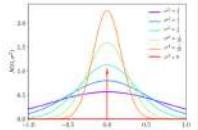
Tutorial

Purpose of Machine Learning

- Make Accurate Predictions
 - Which team will win a soccer match?
 - Which stock will see its price skyrocket?
 - Which patient is at higher risk?
- Usually it is difficult to write down rules manually
- Rather, we learn to make the predictions from paired data (x_i, y_i)

Model Uncertainty

- Germany will probably beat Japan, but what are the odds? 60/40, 70/30, or 80/20?
- I'm willing to bet more money if the odds are in my favor.
- Often translates to: what is the shape of the probability distribution?



Pattern Discovery

- We often have little prior experience, knowledge, or insight into the causal mechanisms that generated the data.
- Still, with only statistical tools, we can identify many important data characteristics
- Obviously, domain knowledge can enrich and complement statistical tools.

I have questions ...

- Full-time students: Find me after the tutorials (but not the lectures)
- Part-time students: Dedicated QA time on Thursdays
- Send questions via email boyang.li@ntu.edu.sg or via Microsoft Teams
- Make an appointment

Full-time Arrangement

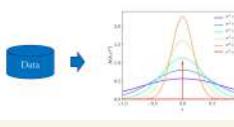
- Lectures
 - 3, 5, 6, 8, 9pm on Thursdays, LT2A
 - Weeks 9, 10, 11, 12
- Tutorials
 - 3, 5, 6, 8, 9pm on Mondays, LT2A
 - Weeks 9, 10, 11, 12
- Office Hours
 - Weeks 7, 11, 13–14, 16–17, April joint lecture room
 - Office Hours: Tuesday, 4pm–5pm
- Feedback on Municipal house w/ NLP and Deep Learning

Make Accurate Predictions

- Artificial Neural Networks (Week 7)
- Support Vector Machines (Week 8)
- Regression (Week 8)
- Ensemble Learning (Week 9)

Model Uncertainty

Density Estimation (Week 11)



Pattern Discovery

- Imagine you are an alien from another planet. You watch a soccer match. What do you observe?
- Two groups of humans. One ball.
- Behavior change when the ball goes into the net.



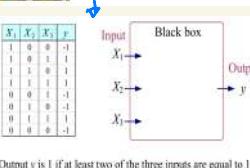
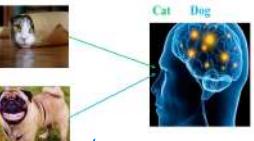
Pattern Discovery

- Clustering (Week 10)
- Dimensionality Reduction (Week 12)

Artificial Neural Networks: Perceptron

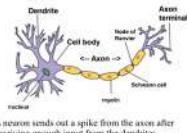
Artificial Neural Networks (ANN)

- The study of ANN was inspired by biological neural systems



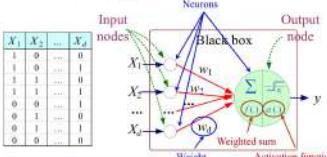
Biology

- Human brain is a densely interconnected network of neurons, connected to others via dendrites and axons.
- Dendrites receive and convert electrical signals from one neuron to another.
- The human brain learns by changing the strength of the synaptic connection between neurons.
- An ANN is composed of an interconnected assembly of nodes and directed links.



A neuron sends out a spike from the axon after receiving enough input from the dendrites.

ANN: Perceptron



Mathematically, the output of a perceptron model can be expressed in a more compact form

$$y = \text{sign} \left(\sum_{i=1}^d w_i x_i - \theta \right)$$

$$y = \text{sign}(\mathbf{w} \cdot \mathbf{x})$$

where $\mathbf{x} = (x_0, x_1, x_2, \dots, x_d)$

$$\mathbf{w} = (w_0, w_1, w_2, \dots, w_d)$$

$$w_0 = -\theta, \text{ and } x_0 = 1$$

Inner Product: Review

- Given two vectors \mathbf{x} and \mathbf{z} , which are both of dimension d , the **inner product** between \mathbf{x} and \mathbf{z} is defined as

$$\mathbf{x} \cdot \mathbf{z} = \sum_{i=1}^d x_i z_i$$

$$\mathbf{x} = (x_1, x_2, \dots, x_d)$$

$$\mathbf{z} = (z_1, z_2, \dots, z_d)$$

$$y = \text{sign}(\mathbf{w} \cdot \mathbf{x})$$

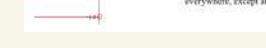
$$\mathbf{w} = (w_0, w_1, w_2, \dots, w_d)$$

$$w_0 = -\theta, \text{ and } x_0 = 1$$

ANN: Sign Function

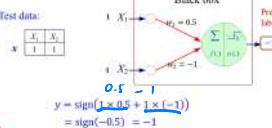
$$\text{sign}(x) = \begin{cases} 1, & x \geq 0 \\ -1, & x < 0 \end{cases}$$

Note: the $\text{sign}(x)$ function has derivative = 0 everywhere, except at $x = 0$.



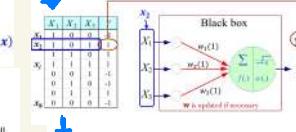
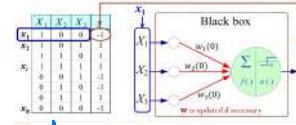
Perceptron: Making Prediction

- Given a learned perceptron with $w_1 = 0.5$, $w_2 = -1$, and $\theta = 0$

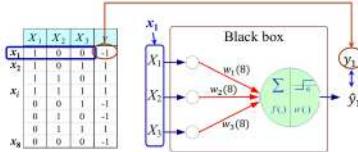


Perceptron: Learning

- During training, the weight parameters \mathbf{w} are adjusted until the outputs of the perceptron become consistent with the true outputs of training data
- The weight parameters \mathbf{w} are updated iteratively or in an online learning manner



Perceptron: Learning (cont.)



Perceptron: Learning (cont.)

- Algorithm: d dimensions

- Let $D = \{(\mathbf{x}_i, y_i) \mid i = 1, 2, \dots, N\}$ be the set of training examples, $t = 0$
- Initialize \mathbf{w} with random values \mathbf{w}_0
- Repeat
- for each training example (\mathbf{x}_i, y_i) do
- Compute the predicted output \hat{y}_i
- Update \mathbf{w}_t by $\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda(y_i - \hat{y}_i)\mathbf{x}_i$
- $t = t + 1$
- end for
- Until stopping condition is met

- Why use the following weight update rule?

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda(y_i - \hat{y}_i)\mathbf{x}_i$$

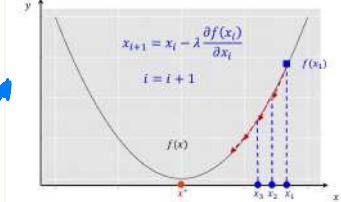
- Induced based on a gradient descent method

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}}$$

Function to be minimized, e.g., the error loss function
Learning rate $\lambda \in (0, 1]$ (not 0)

Gradient Descent

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$$



- Weight update rule

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda(y_i - \hat{y}_i)\mathbf{x}_i \quad \mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}}$$

- Consider the loss function \mathcal{L} for each training example as $e_i \triangleq y_i - \hat{y}_i$

$$\mathcal{L} = \frac{1}{2} e_i^2 = \frac{1}{2} (y_i - \hat{y}_i)^2 = \frac{1}{2} (y_i - \text{sign}(\mathbf{w} \cdot \mathbf{x}_i))^2$$

- Update the weight using a gradient descent method

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{w}_t - \lambda \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}}$$

$$\mathcal{L} = \frac{1}{2} (y - \hat{y})^2 \quad \hat{y} = \text{sign}(x) \quad z = \mathbf{w} \cdot \mathbf{x}$$

Chain rule

Chain Rule of Calculus (Review)

- Suppose that $y = g(x)$ and $z = f(y) = f(g(x))$
- Chain rule of calculus:
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$
- Generalized to the vector case: suppose $\mathbf{x} \in R^m$, $\mathbf{y} \in R^n$, $\mathbf{y} = g(\mathbf{x})$ and $\mathbf{z} = f(\mathbf{y})$
$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \sum_j \frac{\partial \mathbf{z}}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

Perceptron: Learning (cont.)

- The weight update formula for perceptron:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{w}_t - \lambda \frac{\partial \mathcal{L}(\hat{y}_i)}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{w}}$$

$$\mathcal{L} = \frac{1}{2} (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = \text{sign}(z_i)$$

$$(y_i - \hat{y}_i) \times -1$$

Note: the sign() function has derivative 0 except at $z = 0$. Therefore, to compute meaningful derivative, we remove sign() from consideration

$$\hat{y}_i = z_i$$

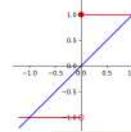
$$z_i = \mathbf{w}_i \cdot \mathbf{x}_i$$

$$\mathbf{x}_i$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda(y_i - \hat{y}_i)\mathbf{x}_i$$

Approximating the derivative?

- The equation used to compute \hat{y}_i from z_i
$$\hat{y}_i = \text{sign}(z_i)$$
- The equation used to compute $\frac{\partial \hat{y}_i}{\partial z}$
$$\hat{y}_i = z_i$$
- Why? This is approximating the step function with a linear function.



While the derivative itself is incorrect, its direction is correct.

That is, if you want to increase \hat{y}_i , you should increase z_i , and vice versa.

Perceptron Weights Update

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda(y_i - \hat{y}_i)\mathbf{x}_i$$

- If the prediction is correct, $(y - \hat{y}) = 0$, then weight remains unchanged $\mathbf{w}_{t+1} = \mathbf{w}_t$
- If $y = +1$ and $\hat{y} = -1$, then $(y - \hat{y}) = 2$
- The weights of all links with positive inputs need to be updated by increasing their values
- The weights of all links with negative inputs need to be updated by decreasing their weights

\mathbf{x}_i	x_{i1}	\dots	x_{ik}	\dots	x_{id}
----------------	----------	---------	----------	---------	----------

$x_{i1} > 0 \quad \mathbf{w}_1 \uparrow \quad \dots \quad x_{ik} = 0 \quad \mathbf{w}_k \downarrow \quad \dots \quad x_{id} < 0$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda(y_i - \hat{y}_i)\mathbf{x}_i$$

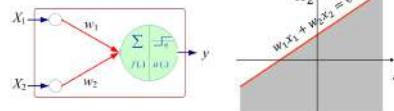
- If $y = -1$ and $\hat{y} = +1$, then $(y - \hat{y}) = -2$
- The weights of all links with positive inputs need to be updated by decreasing their values
- The weights of all links with negative inputs need to be updated by increasing their weights

\mathbf{x}_i	x_{i1}	\dots	x_{ik}	\dots	x_{id}
----------------	----------	---------	----------	---------	----------

$x_{i1} < 0 \quad \mathbf{w}_1 \downarrow \quad \dots \quad x_{ik} = 0 \quad \mathbf{w}_k \uparrow \quad \dots \quad x_{id} > 0$

Convergence

- The decision boundary of a perceptron is a linear hyperplane

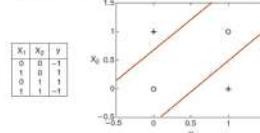


- The perceptron learning algorithm is guaranteed to converge to an optimal solution for linear classification problems

Perceptron Limitation

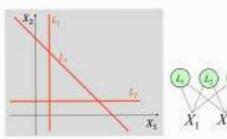
- If the problem is not linearly separable, the algorithm fails to converge

Nonlinearly separable data given by the XOR function

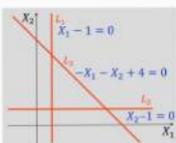


Multi-layer Perceptron

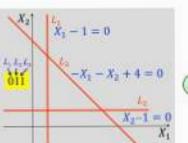
Example: Not Linearly Separable



Example: Not Linearly Separable

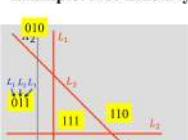


Example: Not Linearly Separable

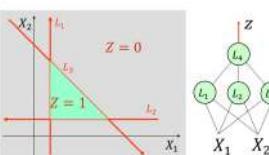
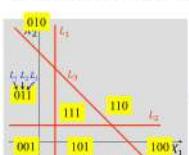


Example: Not Linearly Separable ...

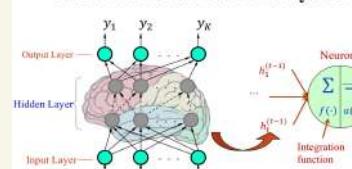
Example: Not Linearly Separable



Example: Not Linearly Separable



General Structure: Multilayer ANN



Integration Functions

- Weighted sum: $\sum_{i=1}^d w_i X_i - \theta$
- Quadratic function: $\sum_{i=1}^d w_i X_i^2 - \theta$
- Spherical function: $\sum_{i=1}^d (X_i - w_i)^2 - \theta$

Activation Functions

- Sign function (Threshold function): $a(x) = \text{sign}(x) = \begin{cases} 1 & x \geq 0 \\ -1 & x < 0 \end{cases}$
 - Unipolar sigmoid function: $a(x) = \frac{1}{1 + e^{-\lambda x}}$
- When $\lambda = 1$, it is called sigmoid function

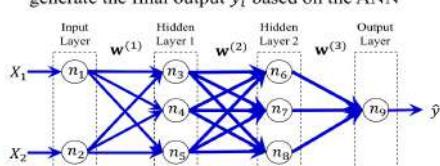
Update Weights for Multi-layer NNs

- Initialize the weights in each layer ($w^{(1)}, \dots, w^{(k)}, \dots, w^{(m)}$)
- Adjust the weights such that the output of ANN is consistent with class labels of training examples
- Loss function for each training instance:
$$L = \frac{1}{2} (y_i - \hat{y}_i)^2$$
- For each layer k , update the weights, $w^{(k)}$, by gradient descent at each iteration t :
$$w_{t+1}^{(k)} = w_t^{(k)} - \lambda \frac{\partial L}{\partial w^{(k)}}$$
- Computing the gradient w.r.t. weights in each layer is computationally expensive!

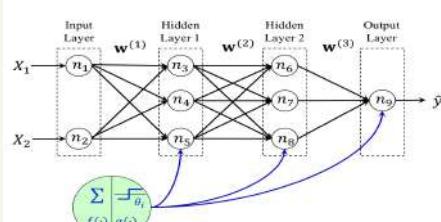
The Backpropagation Algorithm

Backpropagation: Basic Idea

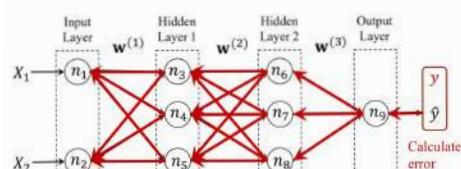
- Initialize the weights ($w^{(1)}, \dots, w^{(3)}$)
- Forward pass:** each training examples (x_i, y_i) is used to compute outputs of each hidden layer and generate the final output \hat{y}_i based on the ANN



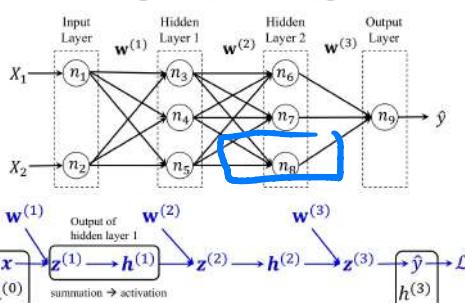
A Multi-layer Feed-forward NN



- Backpropagation:** Starting with the output layer, to propagate error back to the previous layer in order to update the weights between the two layers, until the earliest hidden layer is reached



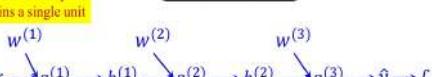
The Computational Graph



Backpropagation (BP)

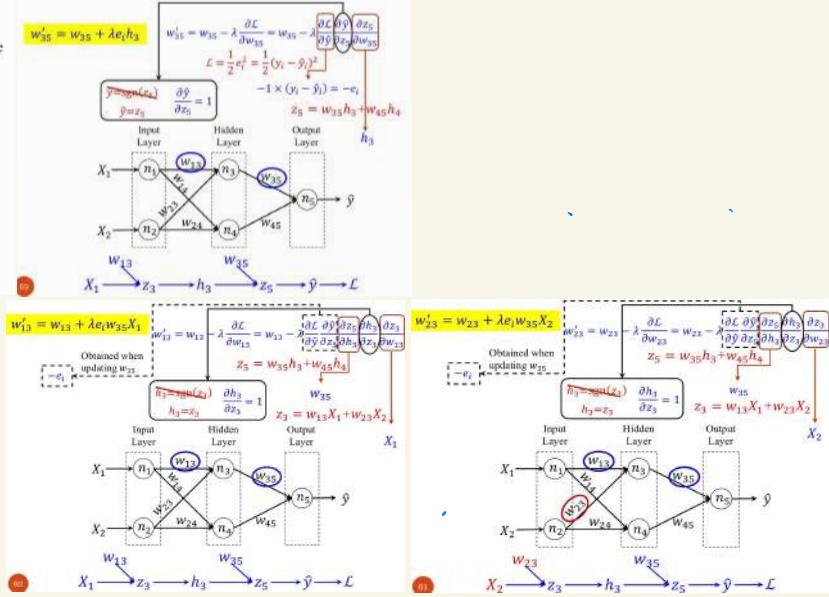
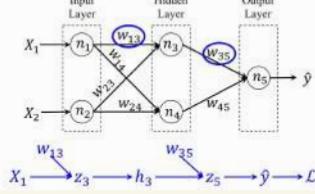
- Gradient of L w.r.t. $w^{(3)}$: $\frac{\partial L}{\partial w^{(3)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial w^{(3)}}$
- Gradient of L w.r.t. $w^{(2)}$: $\frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial w^{(2)}}$
- Gradient of L w.r.t. $w^{(1)}$: $\frac{\partial L}{\partial w^{(1)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial h^{(2)}} \frac{\partial h^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial w^{(1)}}$

Consider each layer contains a single unit



AN Example

- Consider an ANN of 1 hidden layer as follows. Suppose the sign function and the weighted sum function are used for both hidden and output nodes

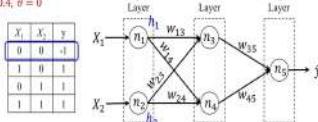


BP Algorithm: Example

Activation function: sign()

Integration function: weighted sum

$$\lambda = 0.4, \theta = 0$$



Initialization:

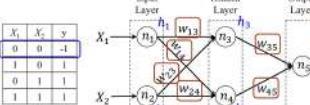
$$(w_{13} = 1, w_{14} = 1, w_{23} = 1, w_{24} = 1, w_{35} = 1, w_{45} = 1)$$

For the 1st example: $h_1 = 0$ and $h_2 = 0$

$$w'_{13} = w_{13} + \lambda e_i w_{35} X_1$$

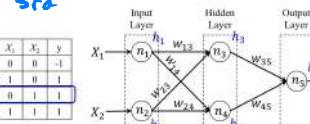
BP Algorithm: Example (cont.)

$$w'_{23} = w_{23} + \lambda e_i w_{35} X_2$$



Backpropagation:

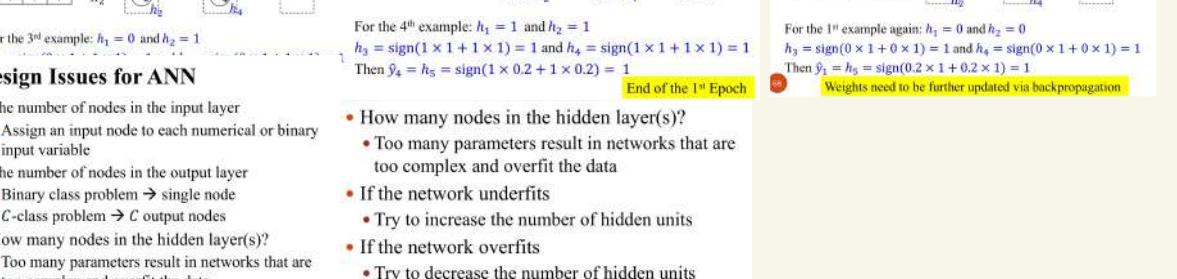
$$\begin{aligned} w_{35} &= 1 + 0.4 \times (-2) \times 1 = 0.2 & w_{45} &= 1 + 0.4 \times (-2) \times 1 = 0.2 \\ w_{13} &= 1 + 0.4 \times (-2) \times 0 \times 1 = 0 & w_{23} &= 1 + 0.4 \times (-2) \times 0 \times 1 = 0 \\ w_{23} &= 1 + 0.4 \times (-2) \times 1 \times 0 = 1 & w_{45} &= 1 + 0.4 \times (-2) \times 1 \times 0 = 1 \end{aligned}$$



For the 3rd example: $h_1 = 0$ and $h_2 = 1$

Design Issues for ANN

- The number of nodes in the input layer
- Assign an input node to each numerical or binary input variable
- The number of nodes in the output layer
 - Binary class problem \rightarrow single node
 - C-class problem $\rightarrow C$ output nodes
- How many nodes in the hidden layer(s)?
 - Too many parameters result in networks that are too complex and overfit the data



Support Vector Machines (SVMs)

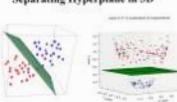
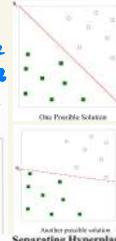
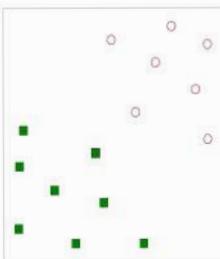
- SVMs have shown promising empirical results in many practical applications, such as computer vision, sensor networks and text mining
- The motivation behind SVMs is from the geometry perspective of linear algebra
- The objective of SVMs is to learn a maximum margin hyperplane
- Based on statistical learning theory

Separating Hyperplane

- To learn a binary classifier

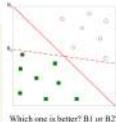


- To find a hyperplane (linear decision boundary) so that all the squares reside on one side of the hyperplane and all the circles reside on the other



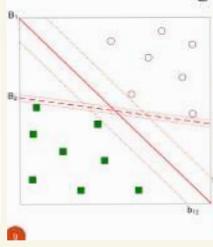
Maximum Margin

- Although all the hyperplanes shown in the figure can separate training examples perfectly, their generalization errors may differ.
- How to choose one of these hyperplanes to construct a classifier's decision boundary with small generalization errors?



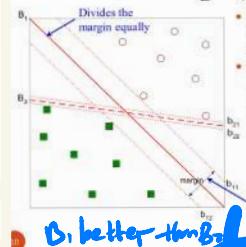
Which one is better? B_1 or B_2 ?

Maximum Margin (cont.)



- Each decision boundary B_1 is associated with a pair of parallel hyperplanes: b_{11} and b_{12} .
- b_{11} is obtained by moving the hyperplane until it touches the closest circle(s).
- b_{12} is obtained by moving a hyperplane away from the decision boundary until it touches the closest square(s).
- The distances from b_{11} and b_{12} to B_1 are the same.

Maximum Margin (cont.)



- Assumption:** larger margins imply better generalization errors.
- The margin of B_1 is much larger than that of B_2 . Therefore, B_1 is better than B_2 .

The distance between these two hyperplanes is known as the **margin** of the classifier

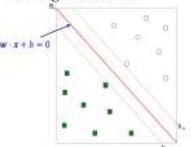
B_1 better than B_2 ↘
larger margin = better
wider = lesser error
Can tolerate disturbance

Review: Vector Addition



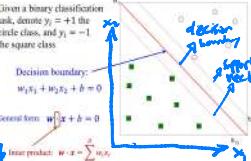
What is the direction of $a + b$?

Making Predictions



For any test example x :
 $f(x) = 1$, if $w \cdot x + b \geq 0$
 $f(x) = -1$, if $w \cdot x + b < 0$

Decision Boundary



Input product: $w \cdot x + b = 0$

Review: Inner Products of Vectors

- We use bold letters to denote vectors, such as \mathbf{a} and \mathbf{b} .
- A vector can have many dimensions: $\mathbf{a} = (a_1, a_2, \dots, a_D)$.
- The inner product of two D -dimensional vectors (D -vectors), \mathbf{a} and \mathbf{b} is defined as

$$\mathbf{a} \cdot \mathbf{b} = a_1 b_1 + a_2 b_2 + \dots + a_D b_D = \sum a_i b_i$$

Also, $\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos\theta$, where θ is the angle between \mathbf{a} and \mathbf{b} .

$\|\mathbf{a}\|$ is the Euclidean norm of \mathbf{a}

$$\|\mathbf{a}\| = \sqrt{a_1^2 + a_2^2 + \dots + a_D^2} = \sqrt{\mathbf{a} \cdot \mathbf{a}}$$

Review: Geometry of Inner Products

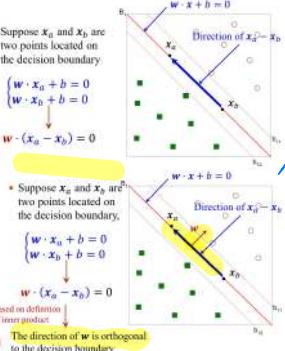
- $\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos\theta$, where θ is the angle between \mathbf{a} and \mathbf{b}
- $\|\mathbf{a} \cdot \mathbf{b}\| = \|\mathbf{a}\| \|\mathbf{b}\| \cos\theta$ is the length of the projection of \mathbf{a} on (or onto) \mathbf{b}



What is the direction of a ?

If \mathbf{a} and \mathbf{b} are orthogonal, $\theta = 90^\circ, \cos\theta = 0, \mathbf{a} \cdot \mathbf{b} = 0$

Margin – Induction



$$\begin{cases} \mathbf{w} \cdot \mathbf{x}_a + b = 0 \\ \mathbf{w} \cdot \mathbf{x}_b + b = 0 \end{cases}$$

$$\mathbf{w} \cdot (\mathbf{x}_a - \mathbf{x}_b) = 0$$

$$\begin{cases} \mathbf{w} \cdot \mathbf{x}_a + b = 0 \\ \mathbf{w} \cdot \mathbf{x}_b + b = 0 \end{cases}$$

$$\mathbf{w} \cdot (\mathbf{x}_a - \mathbf{x}_b) = 0$$

Based on definition of inner product
The direction of \mathbf{w} is orthogonal to the decision boundary

- For any circle x , located above the decision boundary:

$$\mathbf{w} \cdot \mathbf{x}_a + b \geq k, \text{ where } k > 0$$

- For any square x , located below the decision boundary:

$$\mathbf{w} \cdot \mathbf{x}_a + b \leq k', \text{ where } k' < 0$$

The two parallel hyperplanes passing the closest circle(s) and square(s) can be written as

$$\begin{cases} \mathbf{w} \cdot \mathbf{x} + b = k, \text{ where } k > 0 \\ \mathbf{w} \cdot \mathbf{x} + b = k', \text{ where } k' < 0 \end{cases}$$

It can be shown that, these two parallel hyperplanes can be further rewritten as

$$\begin{cases} \mathbf{w} \cdot \mathbf{x} + b = \bar{k} \\ \mathbf{w} \cdot \mathbf{x} + b = -\bar{k}, \text{ where } \bar{k} > 0 \end{cases}$$

Tutorial

- The two parallel hyperplanes passing the closest circle(s) and square(s) can be written as

$$\mathbf{w} \cdot \mathbf{x} + b = \bar{k} \quad \text{where } \bar{k} > 0$$

$$\mathbf{w} \cdot \mathbf{x} + b = \bar{k}' \quad \text{where } \bar{k}' < 0$$

After rescaling \mathbf{w} and b , the two parallel hyperplanes can be further rewritten as

$$\begin{cases} \mathbf{w} \cdot \mathbf{x} + b = 1 \\ \mathbf{w} \cdot \mathbf{x} + b = -1 \end{cases}$$

We should use different letters for the rescaled \mathbf{w} and b , to avoid confusion for simplicity.

$$\begin{cases} \mathbf{w} \cdot \mathbf{x} + b = 1 \\ \mathbf{w} \cdot \mathbf{x} + b = -1 \end{cases}$$

Tutorial

$$\begin{cases} \mathbf{b}_{11}: \mathbf{w} \cdot \mathbf{x}_1 + b = 1 \\ \mathbf{b}_{12}: \mathbf{w} \cdot \mathbf{x}_2 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{21}: \mathbf{w} \cdot \mathbf{x}_3 + b = 1 \\ \mathbf{b}_{22}: \mathbf{w} \cdot \mathbf{x}_4 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{31}: \mathbf{w} \cdot \mathbf{x}_5 + b = 1 \\ \mathbf{b}_{32}: \mathbf{w} \cdot \mathbf{x}_6 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{41}: \mathbf{w} \cdot \mathbf{x}_7 + b = 1 \\ \mathbf{b}_{42}: \mathbf{w} \cdot \mathbf{x}_8 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{51}: \mathbf{w} \cdot \mathbf{x}_9 + b = 1 \\ \mathbf{b}_{52}: \mathbf{w} \cdot \mathbf{x}_{10} + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{61}: \mathbf{w} \cdot \mathbf{x}_{11} + b = 1 \\ \mathbf{b}_{62}: \mathbf{w} \cdot \mathbf{x}_{12} + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{71}: \mathbf{w} \cdot \mathbf{x}_{13} + b = 1 \\ \mathbf{b}_{72}: \mathbf{w} \cdot \mathbf{x}_{14} + b = -1 \end{cases}$$

Tutorial

$$\begin{cases} \mathbf{b}_{11}: \mathbf{w} \cdot \mathbf{x}_1 + b = 1 \\ \mathbf{b}_{12}: \mathbf{w} \cdot \mathbf{x}_2 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{21}: \mathbf{w} \cdot \mathbf{x}_3 + b = 1 \\ \mathbf{b}_{22}: \mathbf{w} \cdot \mathbf{x}_4 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{31}: \mathbf{w} \cdot \mathbf{x}_5 + b = 1 \\ \mathbf{b}_{32}: \mathbf{w} \cdot \mathbf{x}_6 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{41}: \mathbf{w} \cdot \mathbf{x}_7 + b = 1 \\ \mathbf{b}_{42}: \mathbf{w} \cdot \mathbf{x}_8 + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{51}: \mathbf{w} \cdot \mathbf{x}_{10} + b = 1 \\ \mathbf{b}_{52}: \mathbf{w} \cdot \mathbf{x}_{11} + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{61}: \mathbf{w} \cdot \mathbf{x}_{12} + b = 1 \\ \mathbf{b}_{62}: \mathbf{w} \cdot \mathbf{x}_{13} + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{71}: \mathbf{w} \cdot \mathbf{x}_{14} + b = 1 \\ \mathbf{b}_{72}: \mathbf{w} \cdot \mathbf{x}_{15} + b = -1 \end{cases}$$

$$\begin{cases} \mathbf{b}_{81}: \mathbf{w} \cdot \mathbf{x}_{16} + b = 1 \\ \mathbf{b}_{82}: \mathbf{w} \cdot \mathbf{x}_{17} + b = -1 \end{cases}$$

Tutorial

high weight, robust
Optimization easy
→ deep learning if huge
small dataset can use

SVM

Margin Maximization

$$\|\mathbf{w}\|_2 \cdot d = 2 \longrightarrow d = \frac{2}{\|\mathbf{w}\|_2} \quad \|\mathbf{w}\|_2^2 = \sum_{i=1}^d (w_i \times w_i)$$

$$\text{Maximize margin } d = \frac{2}{\|\mathbf{w}\|_2} \longrightarrow \text{Minimize } \frac{\|\mathbf{w}\|_2}{2}$$

For convenience in computation

Constraints: $\mathbf{w} \cdot \mathbf{x}_i + b \geq 1$, if $y_i = 1$
 $\mathbf{w} \cdot \mathbf{x}_i + b \leq -1$, if $y_i = -1$

OR

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

Implementation Example

>>> from sklearn import svm

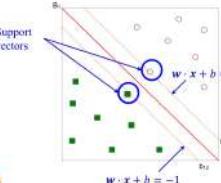
```
>>> symC = svm.LinearSVC()
>>> symC.fit(X, y)
>>> pred = symC.predict(X)
```

Optimization Problem for SVMs

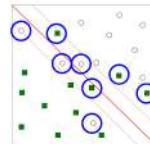
- Optimization problem of linear SVMs

$$\begin{aligned} & \min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|_2^2}{2} \\ \text{s.t. } & y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i = 1, \dots, N \end{aligned}$$

- The optimization is convex
- Many numerical approaches can be applied to find the solution
- Convex is easy. Non-convex is hard.
- The exact optimization algorithms are beyond the scope of this course.

Support Vectors**Non-separable Case**

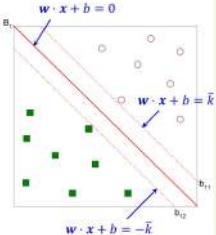
- What if data of two classes cannot be perfectly separated?



Slack variables
need to be introduced to absorb errors

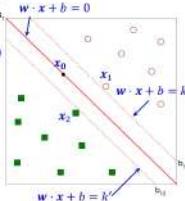
Question 2

- TUT2
- \mathbf{w} determines the orientation (slope) of the decision boundary.
 - The support vectors determine how the decision boundary moves in parallel motion.
 - Together, they determine b .

**Question 2 (cont.)**

$$\begin{aligned} b_{11}: \mathbf{w} \cdot \mathbf{x}_1 + b &= k, \text{ where } k > 0 \\ b_{12}: \mathbf{w} \cdot \mathbf{x}_2 + b &= k', \text{ where } k' < 0 \end{aligned}$$

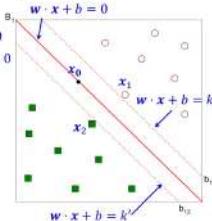
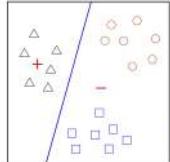
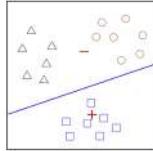
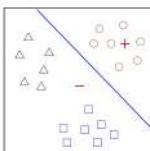
Given two support vectors (or two points on b_{11} and b_{12} respectively), I can choose b such that $k = -k'$

**Question 2 (cont.)**

$$b_{11}: \mathbf{w} \cdot \mathbf{x}_1 + b = k, \text{ where } k > 0$$

$$b_{12}: \mathbf{w} \cdot \mathbf{x}_2 + b = k', \text{ where } k' < 0$$

$$\begin{aligned} \mathbf{w} \cdot \mathbf{x}_1 + b &= -(\mathbf{w} \cdot \mathbf{x}_2 + b) \\ 2b &= -\mathbf{w} \cdot \mathbf{x}_1 - \mathbf{w} \cdot \mathbf{x}_2 \\ b &= -\frac{1}{2} \mathbf{w} \cdot (\mathbf{x}_1 + \mathbf{x}_2) \end{aligned}$$

**Multi-Class Classification**

→ Next page

Linear SVMs for Multi-Class

- Give a 3-class classification problem; C_1 : belong to C_1 , and -1 : not belong to C_1
- General approaches: 1 v.s. rest

- Binary classification 1: positive (C_1) v.s. neg

- Binary classification 2: positive (C_2) v.s. neg

- Binary classification 3: positive (C_3) v.s. neg

- For a test instance x^* , apply binary classifier

- to make predictions on x^*

- Combine predicted results of $f_1(x^*)$, $f_2(x^*)$ to make a final prediction

C_1	C_2	C_3
0	1	1
0	1	0
1	1	0

Total Votes:

not C
C2
C3
C

Linear SVMs: Separable Case

- Optimization problem of linear SVM

$$\min_{w,b} \frac{\|w\|^2}{2} \quad \text{Maximize margin}$$

$$\text{s.t. } y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, i = 1, \dots, N$$

Penalty term: $\frac{1}{2} \|w\|^2$

Maximize margin: $y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

Non training data: leave within the margin: $\mathbf{w} \cdot \mathbf{x}_i + b = 0$

Penalty term: $\frac{1}{2} \|w\|^2$

Margin: $\mathbf{w} \cdot \mathbf{x}_i + b = 1$

Margin: $\mathbf{w} \cdot \mathbf{x}_i + b = -1$

Margin: $\mathbf{w} \cdot \mathbf{x}_i + b = 0$

Linear SVMs: Nonseparable Case

- What if the problem is not separable?

Sack variables: $\xi_i \geq 0$ need to be introduced to absorb errors

For Separable Case:

$$w \cdot x_i + b \geq 1, \text{ if } y_i = 1$$

$$w \cdot x_i + b \leq -1, \text{ if } y_i = -1$$

OR

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

For Nonseparable Case:

$$w \cdot x_i + b \geq 1 - \xi_i, \text{ if } y_i = 1$$

$$w \cdot x_i + b \leq -1 + \xi_i, \text{ if } y_i = -1$$

OR

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

If $\xi_i = 0$, there is no problem with x_i

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

If $\xi_i > 1$, x_i is on the decision boundary (random guess)

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 0$

If $\xi_i > 1$, x_i is misclassified

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq -k$

$k > 0$

Can be negative

$w \cdot x_i + b \geq -k$, if $y_i = 1$

$w \cdot x_i + b \leq k$, if $y_i = -1$

Can be positive

Soft Error

- The number of misclassifications is $\#\{\xi_i > 1\}$
- The number of nonseparable points is $\#\{\xi_i > 0\}$
- Soft errors:

$\sum_i \xi_i$

Nonlinear SVMs

- What if decision boundary is not linear?

Kernel trick in the dual form

Nonlinear SVM: Kernel Trick

- Suppose $\varphi(\cdot)$ is given as follows, mapping an instance from 2-dimensional space to 6-dimensional space:
$$\varphi([X_1, X_2]) = [1, \sqrt{2}X_1, \sqrt{2}X_2, X_1^2, X_2^2, \sqrt{2}X_1X_2]$$
- Given two data instances: $a = [A_1, A_2]$ and $b = [B_1, B_2]$
$$\varphi(a) = [1, \sqrt{2}A_1, \sqrt{2}A_2, A_1^2, A_2^2, \sqrt{2}A_1A_2]$$

$$\varphi(b) = [1, \sqrt{2}B_1, \sqrt{2}B_2, B_1^2, B_2^2, \sqrt{2}B_1B_2]$$
- Inner product of the two instances after feature mapping:
$$\langle a, b \rangle = 1 + 2A_1B_1 + 2A_2B_2 + A_1^2B_2^2 + A_2^2B_1^2 + 2A_1A_2B_1B_2 = (1 + A_1B_1 + A_2B_2)^2$$

Lagrange Multiplier Method: Idea

- Given: an objective $f(\mathbf{w})$ to be minimized, with a set of inequality constraints to be satisfied $h_i(\mathbf{w}) \leq 0, i = 1, 2, \dots, q$
$$\min_{\mathbf{w}} f(\mathbf{w})$$

$$\text{s.t. } h_i(\mathbf{w}) \leq 0, i = 1, \dots, q$$
- The Lagrangian for the optimization problem:
$$L(\mathbf{w}, \lambda) = f(\mathbf{w}) + \sum_{i=1}^q \lambda_i h_i(\mathbf{w})$$

$$\lambda = (\lambda_1, \dots, \lambda_q)$$

The Lagrange multipliers

The Dual Form (Separable)

- By using Lagrangian Multiplier method
$$\min_{\mathbf{w}, b} \frac{\|\mathbf{w}\|^2}{2}$$

$$\text{s.t. } y_i \times (\mathbf{w} \cdot \varphi(\mathbf{x}_i) + b) \geq 1, i = 1, \dots, N$$

Primal Form

$$\max_{\mathbf{w}, b} L_D(\mathbf{w}, b) = -\left(\frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j (\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)) - \sum_{i=1}^N \lambda_i\right)$$

Dual Form

Nonlinear SVM via Kernel Trick

- Training: $\max_{\lambda} \left(\sum_{i=1}^N \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j [\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)] \right)$
- Decision boundary: $\sum_{i=1}^N \lambda_i y_i [\varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}^*)] + b = 0$
- The data points only appear as inner product
- As long as the inner product in the feature space can be calculated, no need for the explicit mapping

Linear SVMs: Nonseparable Case

- What if the problem is not separable?

Stack Variables

- For Separable Case:
$$w \cdot x_i + b \geq 1, \text{ if } y_i = 1$$

$$w \cdot x_i + b \leq -1, \text{ if } y_i = -1$$

OR

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$
- For Nonseparable Case:
$$w \cdot x_i + b \geq 1 - \xi_i, \text{ if } y_i = 1$$

$$w \cdot x_i + b \leq -1 + \xi_i, \text{ if } y_i = -1$$

OR

$$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$

If $\xi_i = 0$, there is no problem with x_i

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

If $\xi_i > 1$, x_i is on the decision boundary (random guess)

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 0$

If $\xi_i > 1$, x_i is misclassified

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

$y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \geq -k$

$k > 0$

Can be negative

$w \cdot x_i + b \geq -k$, if $y_i = 1$

$w \cdot x_i + b \leq k$, if $y_i = -1$

Can be positive

Linear SVMs: Nonseparable Case

- Linear SVMs with soft errors:

Linear SVMs: Nonseparable Case

- Linear SVMs with soft errors:

Feature Mapping

- The original input space can always be mapped to some higher-dimensional feature space where the training set is separable:

$\varphi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

$(X, X_*) \rightarrow (X_1^2, \sqrt{2}X_1, X_2^2)$

How to Apply Kernel Trick?

- Optimization problem for nonlinear SVMs

Kernel Trick: General Idea

- If $\varphi(\cdot)$ satisfies some conditions, then we can find a function $k(\cdot, \cdot)$ such that

Kernel function $\rightarrow k(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)$

Input space

Feature space

$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$

Primal Form

- Computation in the feature space can be costly because it is high dimensional
- The feature space is typically very high dimensional!
- The kernel trick comes to rescue

Dual Optimization Problem

- The dual Lagrangian involves only the Lagrange multipliers and the training data
- The negative sign in the dual Lagrangian transforms a minimization problem of the primal form to a maximization problem of the dual form
- The objective is to maximize $L_D(\lambda)$
- Can be solved using numerical techniques such as quadratic programming

Kernel Functions: Examples

- Linear kernel
$$k(x_i, x_j) = x_i \cdot x_j$$
- Radial basis function kernel with width σ
$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$
- Polynomial kernel with degree d
$$k(x_i, x_j) = (x_i \cdot x_j + 1)^d$$

Soft Margin Dual Form

- By using Lagrangian Multiplier method

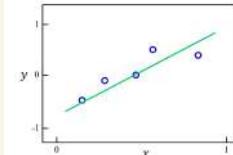
Kernel Functions: Examples

- Linear kernel
$$k(x_i, x_j) = x_i \cdot x_j$$
- Radial basis function kernel with width σ
$$k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$
- Polynomial kernel with degree d
$$k(x_i, x_j) = (x_i \cdot x_j + 1)^d$$

Regression: Example



Linear Fitting

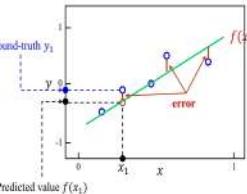


Linear Regression Model

- A special case, an instance is represented by one input feature
- To learn a linear function $f(x)$ in terms of w (drop bias term b for simplicity) from $(x_i, y_i), i = 1, \dots, N$

$$f(x) = w \cdot x$$

Such that the difference (i.e., error) between the predicted values $f(x_i)$'s and the ground-truth values y_i 's is as small as possible



- Suppose sum-of-squares error is used

$$\mathcal{L}(w) = \frac{1}{2} \sum_{i=1}^N (f(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (w \cdot x_i - y_i)^2$$

- Learn the linear model in terms of w by minimizing the error

$$w^* = \arg \min_w \mathcal{L}(w)$$

- To solve the unconstrained minimization problem, we can set the derivative of $\mathcal{L}(w)$ w.r.t. w to zero,

$$\frac{\partial \mathcal{L}(w)}{\partial w} = \frac{\partial \left(\frac{1}{2} \sum_{i=1}^N (w \cdot x_i - y_i)^2 \right)}{\partial w} = 0$$

$$\sum_{i=1}^N (w \cdot x_i - y_i) \times x_i = 0 \quad \cancel{\frac{1}{2} w \cdot x^2} - \cancel{y_i w \cdot x} = 0$$

$$w \sum_{i=1}^N x_i^2 - \sum_{i=1}^N y_i \times x_i = 0$$

$$w = \frac{\sum_{i=1}^N y_i \times x_i}{\sum_{i=1}^N x_i^2}$$

- To learn a linear function $f(x)$ in more general form in terms of w and b ,

$$f(x) = w \cdot x + b$$

- By defining $w_0 = b$, and $x_0 = 1$, w and x are of $d + 1$ dimensions

$$f(x) = w \cdot x$$

- Suppose sum-of-squares error is used

$$\mathcal{L}(w) = \frac{1}{2} \sum_{i=1}^N (f(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (w \cdot x_i - y_i)^2$$

- Learn the linear model in terms of w by minimizing the error

$$w^* = \arg \min_w \mathcal{L}(w) + \frac{\lambda}{2} \|w\|_2^2 = w \cdot w$$

Positive tradeoff parameter
A regularization term to control the complexity of the model (infinite overfitting)

Regularized Linear Regression

- To solve the unconstrained minimization problem, we can set the derivative of $\mathcal{L}(w) + \frac{\lambda}{2} \|w\|_2^2$ w.r.t. w to zero

$$\frac{\partial (\mathcal{L}(w) + \frac{\lambda}{2} \|w\|_2^2)}{\partial w} = 0$$

- We can obtain a closed-form solution for w by the above equations

Closed-Form Solution

- Denote by $X = (x_1, x_2, \dots, x_N)^T$

WxD+1

$$X = \begin{pmatrix} x_{1,0} & \cdots & x_{1,N} \\ \vdots & \ddots & \vdots \\ x_{d,0} & \cdots & x_{d,N} \end{pmatrix}^T = \begin{pmatrix} x_{1,0} & \cdots & x_{1,d} \\ x_{2,0} & \cdots & x_{2,d} \\ \vdots & \ddots & \vdots \\ x_{N,0} & \cdots & x_{N,d} \end{pmatrix}$$

- And by $y = (y_1, y_2, \dots, y_N)^T$

- The closed-form solution for w :

$$w = (X^T X + \lambda I)^{-1} X^T y$$

How to induce this closed-form solution?
Tutorial

Implementation Example

>>> from sklearn.linear_model import Ridge

...
Referred to as λ on our lecture notes

>>> rlr = Ridge(alpha=0.1)

>>> rlr.fit(X, y)

>>> pred_train_rlr = rlr.predict(X)

Evaluation

- Root Mean Square Error (RMSE)

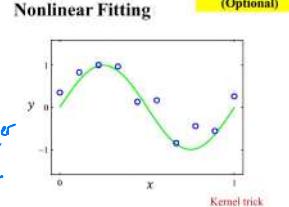
$$\sqrt{\frac{1}{N} \sum_{i=1}^N (f(x_i) - y_i)^2}$$

- Mean Absolute Error (MAE)

$$\frac{1}{N} \sum_{i=1}^N |f(x_i) - y_i|$$

Nonlinear Fitting

Additional Notes (Optional)



kernel regression notes!

Regularized Linear Regression with Kernels

- Primal objective

$$\mathcal{T}(w) = \frac{1}{2} \sum_{i=1}^N (w \cdot \phi(x_i) - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2$$

- It can be reformulated in terms of a dual form where kernel function arises naturally

$$\frac{\partial \mathcal{T}(w)}{\partial w} = 0 \rightarrow \sum_{i=1}^N (w \cdot \phi(x_i) - y_i) \phi(x_i) + \lambda w = 0$$

Closed Form Solution

- By using kernel trick $k(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$

$$f(x) = w \cdot \phi(x)$$

$$k(x) = k(x_i, x)$$

$$f(x) = k(x)^T (K + \lambda I)^{-1} y$$

$$\begin{aligned} K &= \begin{pmatrix} k(x_1, x_1) & \cdots & k(x_1, x_N) \\ k(x_N, x_1) & \cdots & k(x_N, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \cdots & k(x_N, x_N) \end{pmatrix} \\ &= \begin{pmatrix} \phi(x_1) \cdot \phi(x_1) & \cdots & \phi(x_1) \cdot \phi(x_N) \\ \phi(x_N) \cdot \phi(x_1) & \cdots & \phi(x_N) \cdot \phi(x_N) \\ \vdots & \ddots & \vdots \\ \phi(x_N) \cdot \phi(x_1) & \cdots & \phi(x_N) \cdot \phi(x_N) \end{pmatrix} \end{aligned}$$

Induction

$$\sum_{i=1}^N (w \cdot \phi(x_i) - y_i) \phi(x_i) + \lambda w = 0$$

$$w = \frac{1}{2} \sum_{i=1}^N (w \cdot \phi(x_i) - y_i) \phi(x_i)$$

$$= \sum_{i=1}^N a_i \phi(x_i)$$

$$\text{Denote by } \Phi = (\phi(x_1), \phi(x_2), \dots, \phi(x_N))^T$$

$$\text{and } a = (a_1, \dots, a_N)^T$$

$$w = \sum_{i=1}^N a_i \phi(x_i)$$

$$w = \frac{1}{2} (w \cdot \phi(x_1) - y_1) \phi(x_1)$$

$$\Rightarrow a = \frac{1}{2} (w \cdot \Phi - y) \Phi$$

$$w = \Phi^T a$$

$$w = \Phi^{-1} (y - \lambda I)$$

$$\begin{aligned} w &= \Phi^T a & n &= \frac{1}{2} (y - \Phi w) \\ w &= \frac{1}{2} (x - \Phi^T a) & \Rightarrow 2w &= y - \Phi^T w \\ & \Rightarrow 2w & \Rightarrow (I + \Phi^T \Phi) w = y \\ & \Rightarrow (I + \Phi^T \Phi)^{-1} (I + \Phi^T \Phi) w & \Rightarrow (I + \Phi^T \Phi)^{-1} y \\ & \Rightarrow w & \Rightarrow w = \Phi^{-1} (y - \lambda I) \\ & \Rightarrow w & \Rightarrow w = \Phi^{-1} (y - \lambda I) \\ & \Rightarrow w & \Rightarrow w = \Phi^{-1} (y - \lambda I) \\ & \Rightarrow w & \Rightarrow w = \Phi^{-1} (y - \lambda I) \end{aligned}$$

Lesson 9: Ensemble Learning

Ensemble Methods

- Objective:
 - To improve model performance in terms of accuracy by aggregating the predictions of multiple models
 - How to do it?
 - Construct a set of base models from the training data
 - Make predictions by combining the predicted results made by each base model
- "Two heads are better than one"

Why Ensemble Work?

- Suppose there are 3 base binary classifiers
 - Each classifier has error rate, $\epsilon = 0.35$ or accuracy $acc = 0.65$
- Given a test instance, if we choose any one of these classifiers to make prediction, the probability that the classifier makes a wrong prediction is 35%

Base classifiers: f_1, f_2, f_3, f_M

A test instance:

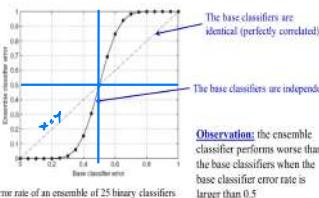
- Consider to combine the 3 base classifiers to make a prediction on a test instance using a majority vote
- The ensemble makes a wrong prediction only if more than 1 (i.e. 2 or 3) of the base classifiers predict incorrectly

Truth label: +1	f_1	f_2	f_3	f_M
+1	+1	+1	+1	+1
+1	+1	+1	+1	-1
+1	+1	+1	-1	-1
+1	+1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1
-1	-1	-1	-1	-1

The combined model makes a wrong prediction if at least two of the three base classifiers make a wrong prediction at the same time

- Assuming the three base classifiers are independent
 - The decision of one classifier does not tell us anything about another classifier.
- Each classifier has an error rate of $\epsilon = 35\%$
- What is the overall error rate of the ensemble?

Error Rate of Base Classifiers



Error rate of an ensemble of 25 binary classifiers for different base classifier error rates

Review: Independence

- If X and Y are independent, then we have
 - $P(X = x|Y = y) = P(X = x)$
 - Y provides no knowledge about X
 - E.g., no matter the outcome of the first coin toss, the second toss is still 50/50
 - Or, no matter the second toss, the first toss is still 50/50
- Equivalently, $P(X = x, Y = y) = P(X = x)P(Y = y)$
- Why is this equivalent?

$$P(X = x, Y = y) = P(X = x|Y = y)P(Y = y) \\ = P(X = x)P(Y = y)$$

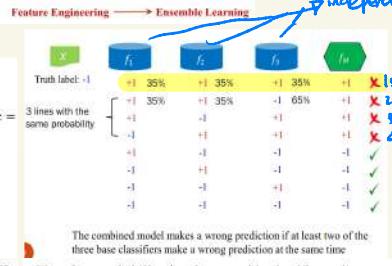
Ensemble Methods

- How to generate a set of base classifiers?
 - By manipulating the training set: multiple training sets are created by resampling the original data according to some sampling distribution. A classifier is then trained from each training set, such as Bagging, Boosting
- By manipulating the input features: a subset of input features is chosen to form each training set. A classifier is then built from each training set, such as Random Forest
- By manipulating the learning algorithm(s): applying the algorithm several times on the same training data using different parameters or applying different algorithms
- How to combine the base classifiers for predictions?

Stories of Success

- Data mining competitions on Kaggle
 - Winning teams employ ensembles of classifiers

kaggle



- Therefore, probability that the ensemble classifier makes a wrong prediction is:

$$\sum_{i=1}^3 \binom{3}{i} e^i (1-e)^{3-i} = \frac{1}{3} \times 0.35^2 \times 0.65 + 1 \times 0.35^3 \times 0.2817$$

Case 1: when there are two exact classifiers make wrong predictions, the probability is

$$\binom{3}{2} e^2 (1-e)^{3-2}$$

All possible combination $\binom{3}{2} e^2 (1-e)^{3-2}$ — The rest one makes correct prediction

Case 2: when all the three classifiers make wrong predictions, the probability is

$$\binom{3}{3} e^3 (1-e)^{3-3}$$

That is the accuracy of the ensemble classifier is 71.83%

$$\epsilon_{f_1} = 35\% \longrightarrow \epsilon_M = 28.17\%$$

- Suppose there are 25 independent base classifiers

- Therefore, probability that the ensemble classifier makes a wrong prediction is:

$$\sum_{i=1}^{25} \binom{25}{i} e^i (1-e)^{25-i} = 0.06$$

- That is the accuracy of the ensemble classifier is 94%

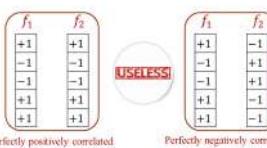
$$\epsilon_{f_1} = 35\%$$

$$\epsilon_M = 6\%$$

Correlation between Base Classifiers

- Sometimes the base classifiers are not completely independent.

- An extreme case is when they are perfectly correlated. That is, they always have the same predictions.



Review: Correlation between RVs

Difficult Input Causes Correlation

Strong positive correlation

Weak positive correlation

Strong negative correlation

Weak negative correlation

Moderate negative correlation

No correlation

One variable tells us something about the other variable.

General Procedure

- Let D denote the original training data, k denote the number of base classifiers, and T be the test dataset.
- for $i = 1$ to k do
 - Train a base classifier f_i from D
 - end for
 - for each test instance $x \in T$ do
 - Generate $f_1(x), f_2(x), \dots, f_k(x)$
 - Calculate $f_M(x) = \text{Merge}(f_1(x), f_2(x), \dots, f_k(x))$
 - end for

For example, majority voting
(can be other schemes)

For example, majority voting
(can be other schemes)

Review: Random Variables

- Consider two random variables X and Y , which may take on values from $[0, 1]$.

- Random variables are variables whose values are randomly assigned from some unspecified random experiment.

- Imagine X and Y are the results of two coin tosses.

- Head for the first toss $\Rightarrow X = 1$

- Tail for the first toss $\Rightarrow X = 0$

- Head for the second toss $\Rightarrow Y = 1$

- Tail for the second toss $\Rightarrow Y = 0$

- The probabilities are denoted as $P(X = 0)$, $P(X = 1)$, $P(Y = 0)$, $P(Y = 1)$

Necessary Conditions

- Two necessary conditions for an ensemble classifier to perform better than a single classifier:

- The base classifiers are not perfectly correlated with each other.

- In practice, the base classifiers are usually somewhat correlated.

- The base classifiers should do better than a classifier that performs random guessing (e.g., for binary classification, accuracy should be better than 0.5).

These are supposedly digits from 0 to 9.

Most models probably make incorrect predictions on these data points.

Having multiple models would help too much here.

Bagging

Known as bootstrap aggregating (bootstrapping), to repeatedly sample with replacement according to a uniform probability distribution.

Build classifier on each bootstrap sample, which is of the same size of the original data

Use majority voting to determine the class label of ensemble classifier

Bagging (cont.)

For example, majority voting
(can be other schemes)

Majority Voting

Suppose a training set D contains N examples

A training instance x has a probability of $1 - \frac{1}{N}$ of not being selected

Its probability of ending up *not* in a training set D_i is $(1 - \frac{1}{N})^N \approx \frac{1}{e} = 0.368$

A bootstrap sample D_i contains approximately 63.2% of the original training data

Optional

Implementation Example

Implementation Example

Boosting

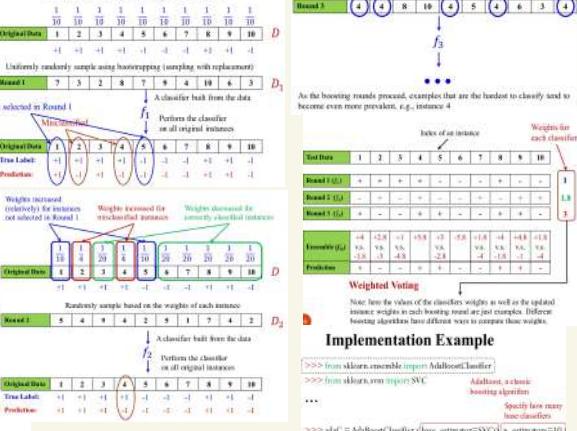
- Principles:
 - Boost a set of **weak learners** to a strong learner
 - Make instances currently misclassified more important
 - Generally,
 - To adaptively change the distribution of training data so that the base classifiers will focus more on previously misclassified records
- Specifically,
 - Initially, all N instances are assigned equal weights
 - Unlike bagging, weights may change at the end of each boosting round
 - In each boosting round, after the weights are assigned to the training instances,
 - Draw a bootstrap sample from the original data by using the weights as a sampling distribution to build a model

Boosting: Procedure

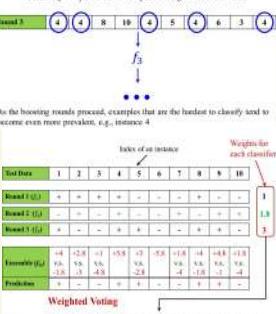
- Initially, all instances are assigned equal weights $\frac{1}{N}$, so that they are equally likely to be chosen for training. A sample is drawn uniformly to obtain a new training set.
- A classifier is induced from the training set, and used to classify all the examples in the original training set.
- The weights of the training instances are updated at the end of each boosting round
 - Instances that are wrongly classified will have their weights increased
 - Instances that are classified correctly will have their weights decreased
- Repeat Steps 2 and 3 until the stopping condition is met
- Finally, the ensemble is obtained by aggregating the base classifiers obtained from each boosting round

Boosting: Example

Initially, all the instances are assigned the same weights.



Randomly sample based on the updated weights of each instance



Implementation Example

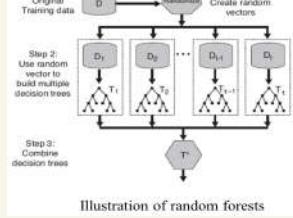
```
>>> from sklearn.ensemble import AdaBoostClassifier
>>> from sklearn.ensemble import SVC
...<br/>
>>> adaC = AdaBoostClassifier(base_estimator=SVC(), n_estimators=10)
>>> adaC.fit(X,y)
>>> pred = adaC.predict(X)
```

Addition: a classic boosting algorithm

Specify how many base classifiers

Random Forests

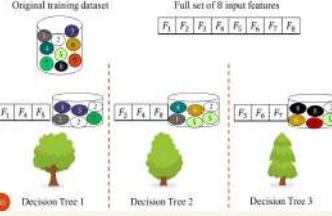
- A class of ensemble methods specifically designed for decision tree classifiers
- Random Forests grow many trees
- Each tree is generated based on a random subset of features
- Final result on classifying a new instance – voting
- Forrest chooses the classification result having the most votes (over all the trees in the forest)



Random Forests: Algorithm

- Choose T : number of trees to grow
- Choose $m' < m$ (m is the number of total features): number of features used to calculate the best split at each node (typically 20%)
- For each tree
 - Choose a training set via bootstrapping
 - For each node, randomly choose m' features and calculate the best split
 - Fully grown and no pruned
- Use majority vote among all the trees

Example



Random Forests: Discussions

- Bagging + random features
- Improve Accuracy
 - Incorporate more diversity
- Improve Efficiency
 - Searching among subsets of features is much faster than searching among the complete set

Implementation Example

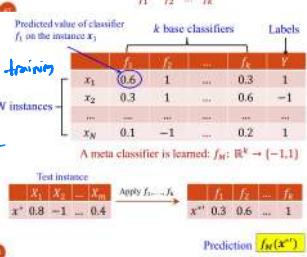
```
>>> from sklearn.ensemble import RandomForestClassifier
...<br/>
>>> rfc = RandomForestClassifier(n_estimators=20, max_depth=3)
>>> rfc.fit(X,y)
>>> pred = rfc.predict(X)
```

Set parameter for the base tree

Specify how many base classifiers

Combining by Learning

- Stacking:
 - A general procedure where a learner is trained to combine the individual learners
 - Individual learners: first-level learners
 - Combiner: second-level learner or meta-learner
 - Suppose given a binary classification problem,



Combination Methods

- Voting
 - Majority voting
 - Weighted voting
- Average
 - Simple average
 - Weighted average
- Combining by learning

Voting

- Majority voting:
 - Takes the class label that receives the largest number of votes as the final winner
- Weighted voting:
 - A generalized version of majority voting by introducing weights for each classifier

Average

- Simple average:

$$f_M(x) = \frac{1}{T} \sum_{t=1}^T f_t(x)$$
- Weighted average:

$$f_M(x) = \sum_{t=1}^T w_t f_t(x)$$

where $w_t \geq 0$, and $\sum_{t=1}^T w_t = 1$
- Prediction:

$$\hat{y}(x) = \begin{cases} 1, & f_M(x) \geq 0 \\ -1, & \text{otherwise} \end{cases}$$

Classification
for binary

Combining by Learning: Summary

- Represent each training instances using classifier-generated outputs

$$\mathbf{x}'_i = (f_1(x_i), f_2(x_i), \dots, f_k(x_i))$$
- Learn a meta classifier f_M from $\{\mathbf{x}'_i, y_i\}, i = 1, \dots, N$.
- For a test instance \mathbf{x}^*
 - Represent it by $\mathbf{x}'^* = (f_1(\mathbf{x}^*), f_2(\mathbf{x}^*), \dots, f_k(\mathbf{x}^*))$
 - Use f_M to make a prediction $f_M(\mathbf{x}^*)$

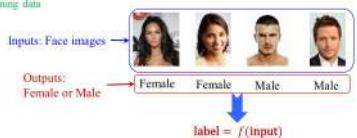
Similar to multi layer perception

Lesson 10: Clustering

Supervised Learning

- The examples presented to computers are pairs of **inputs** and the corresponding **outputs**, the goal is to "learn" a **mapping or model** from inputs to labels

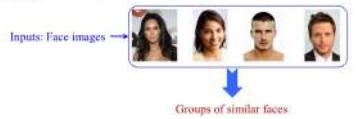
Labeled training data



Unsupervised Learning

- The examples presented to computers are a set of **inputs without any outputs**, the goal is to "learn" an **intrinsic structure** of the examples, e.g., clusters of examples, density of the examples

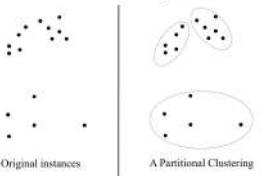
Unlabeled training data



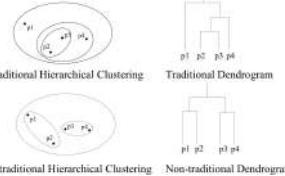
Types of Clusterings

- A clustering is a set of clusters
- Important distinction between **hierarchical** and **partitional** sets of clusters
- Partitional clustering
- Divide data instances into **non-overlapping** subsets (clusters) so that each data instance is in exactly one subset
- Hierarchical clustering
- A set of **nested** clusters organized as a hierarchical tree

Partitional Clustering



Hierarchical Clustering



Other Distinctions Between Clusterings

- Exclusive** versus **non-exclusive**
 - In non-exclusive clustering, instances may belong to multiple clusters
- Fuzzy** versus **non-fuzzy**
 - In fuzzy clustering, a point belongs to every cluster with some weight between 0 and 1
 - Weights must sum to 1
- Partial** versus **complete**
 - In partial clustering, only some of the instances are clustered

Clustering Algorithms

- K**-means and its variants
- Hierarchical clustering

K-means Clustering

- Partitional clustering approach
- Number of clusters, K , must be specified
- Each cluster is associated with a **centroid** (center point)
- Each point is assigned to the cluster with the closest centroid
- Basic algorithm:
 - Select K data instances as the initial centroids
 - Repeat**
 - Form K clusters by assigning all data instances to the closest centroid
 - Recompute the centroid of each cluster
 - Until** The centroids do not change

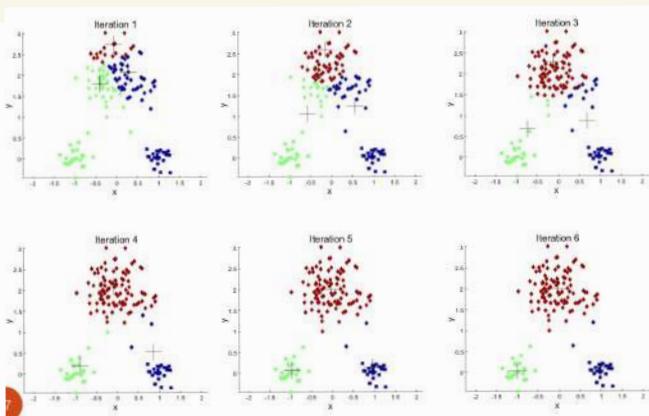
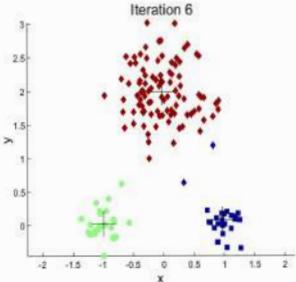
k = number of classes

close to
calculate average

K-means Clustering – Details

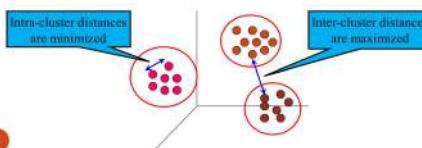
- Initial centroids are often chosen randomly
- The centroid is (typically) the mean of the data instances in the cluster
- 'Closeness' is measured by a proximity
 - Distance: Euclidean distance, etc
 - Similarity: cosine similarity, correlation, etc
- K**-means will converge for common similarity measures mentioned above
 - In practice, it converges in the first few iterations
 - Often the stopping condition is changed to 'Until relatively few instances change clusters'

K-means Illustration



What is Cluster Analysis?

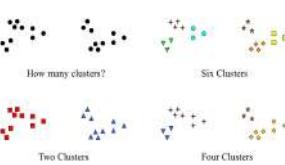
- Finding groups of data instances such that the data instances in a group are
 - similar to one another
 - different from data instances in other groups



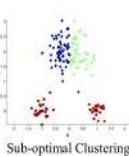
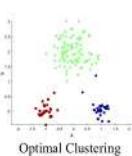
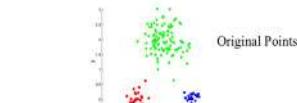
What is not Cluster Analysis?

- Supervised classification
 - Have class label information
- Simple segmentation
- Dividing customers into different groups alphabetically, by last name

Notion of a Cluster can be Ambiguous



Two different K-means Clusterings



Evaluating K-means Clusters

- Most common measure is Sum of Squared Error (SSE)
- For each data instance, the "error" is the distance to the nearest cluster that is represented by a centroid
- To get an overall SSE, we square these errors and sum them

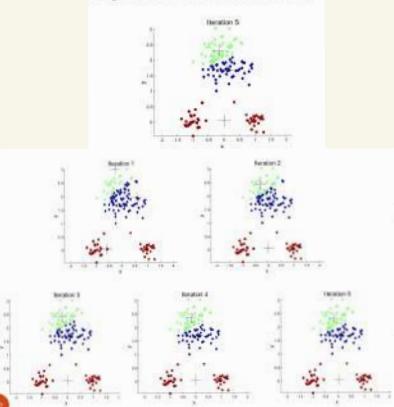
$$\text{Total SSE} \quad \text{SSE} = \sum_{i=1}^K \sum_{x \in C_i} \text{dist}(c_i, x)^2$$

Centroid of the cluster C_i

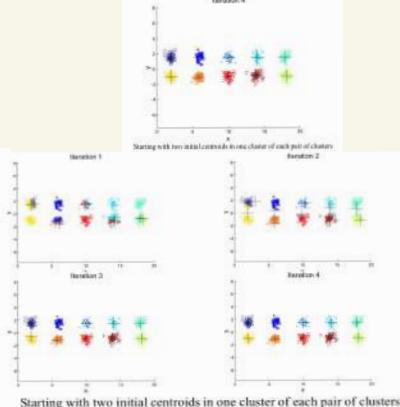
$$\text{Cluster SSE for cluster } C_i = \sum_{x \in C_i} \text{dist}(c_i, x)^2$$

- Given two different runs of K-means, we can choose the one with the smallest **Total SSE**
- One easy way to reduce SSE is to increase K , the number of clusters
- However, a good clustering with smaller K can have a lower SSE than a poor clustering with higher K

Importance of Initial Centroids



10 Clusters Example



Initial Centroids Issue

- Multiple runs
- Postprocessing
- Decompose a cluster with high Cluster SSE,
- Merge clusters with low Cluster SSE, which are close to each other
- Bisecting K-means

Pre-processing and Post-processing

- Pre-processing
 - Normalize the data
 - Eliminate outliers
- Post-processing
 - Eliminate small clusters that may represent outliers
 - Split 'loose' clusters, i.e., clusters with relatively high SSE
 - Merge clusters that are 'close' and that have relatively low SSE

Empty Clusters Issue

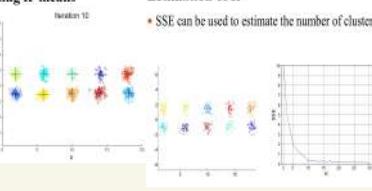
- Basic K-means algorithm can yield empty clusters
- Several strategies to choose a replacement centroid
 - Choose the point that contributes most to SSE
 - Choose a point from the cluster with the highest Cluster SSE
 - If there are several empty clusters, the above can be repeated several times

Bisecting K-means

- Basic algorithm:
 - Initialize the list of clusters to be a single cluster that contains all points.
 - Repeat
 - Select a cluster with highest SSE from the list of clusters
 - For $i = 1$ to D do
 - Bisect the selected cluster using basic K-means
 - Add the two clusters from the bisection with lowest SSE to the list of clusters
 - Until the list of clusters contains K clusters
 - (After bisecting) run K-means with K clusters found so far

Bisecting K-means

Iteration 10



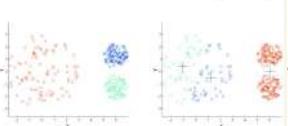
Estimation of K

- SSE can be used to estimate the number of clusters

Limitations of K-means

- K-means has problems when clusters are of differing
 - Sizes
 - Densities
 - Non-globular shapes
- K-means has problems when the data contains outliers

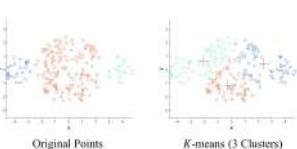
Limitations of K-means: Differing Density



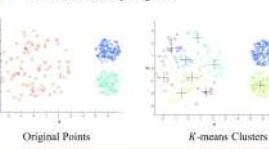
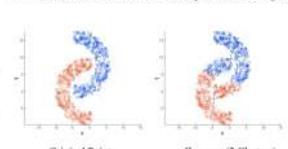
Solution

One solution is to use many clusters. Find parts of clusters, but need to put together.

Limitations of K-means: Differing Sizes



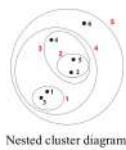
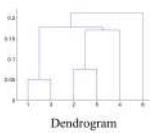
Limitations of K-means: Non-globular Shapes



Hierarchical Clustering

Don't know how to select

- Produces a set of nested clusters organized as a hierarchical tree
- Can be visualized as a dendrogram
- A tree-like diagram that records the sequences of merges or splits

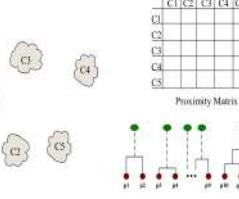


Agglomerative Clustering Algorithm

- Basic algorithm:
 1. Compute the proximity matrix
 2. Let each data instance be a cluster
 3. Repeat
 4. Merge the two closest clusters
 5. Update the proximity matrix
 6. Until only a single cluster remains
- Key operation is to compute the proximity of two clusters
- Different approaches to defining the proximity between clusters lead to different clustering results

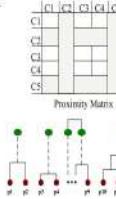
Intermediate Situation

- After some merging steps, we have some clusters

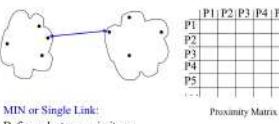


Intermediate Situation...

- We want to merge the two closest clusters (C2 and C5) and update the proximity matrix.



Inter-Cluster Similarity (I)



MIN or Single Link:

- Defines cluster proximity as
- the proximity between the closest two data points that are in different clusters
- or the shortest edge (single link) between two nodes in different subsets (using graph terms)

Clustering with MIN (cont.)

- Step 1: Merge the two closest clusters (largest similarity)

	P1UP2	P3	P4	P5
P1UP2	1.00	0.70	0.65	0.50
P3	0.70	1.00	0.40	0.30
P4	0.65	0.40	1.00	0.80
P5	0.50	0.30	0.80	1.00

Similarity matrix



Clustering with MIN (cont.)

- Step 2: Update proximity matrix based on MIN: proximity of two clusters is based on the two closest points in different clusters (largest similarity)

	P1UP2UP3	P4UP5
P1UP2UP3	1.00 0.70 0.65	0.65 0.70 0.90
P4UP5	0.65 1.00	1.00

Similarity matrix



Strengths of Hierarchical Clustering

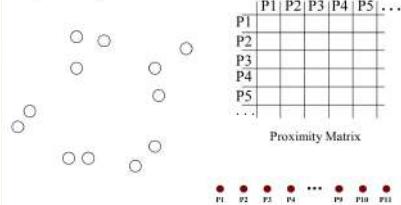
- Do not have to assume any particular number of clusters
- Any desired number of clusters can be obtained by 'cutting' the dendrogram at the proper level
- They may correspond to meaningful taxonomies
- Examples in document organization, biological sciences

Hierarchical Clustering (cont.)

- Two main types of hierarchical clustering
 - Agglomerative (bottom-up):
 - Start with the instances as individual clusters
 - At each step, merge the closest pair of clusters until only one cluster (or K clusters) left
 - Divisive (top-down):
 - Start with one, all-inclusive cluster
 - At each step, split a cluster until each cluster contains a point (or there are K clusters)
- Traditional hierarchical algorithms use a proximity matrix (similarity or distance) to merge or split one cluster at a time

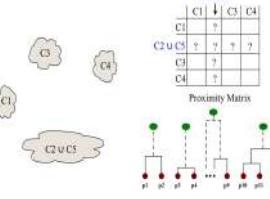
Starting Situation

- Start with clusters of individual instances and a proximity matrix

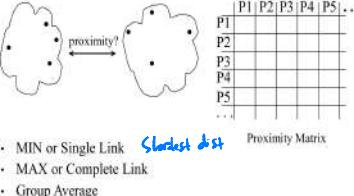


After Merging

- The question is "How do we update the proximity matrix?"

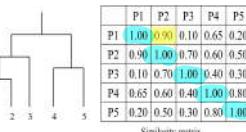


Define Inter-Cluster Proximity



Clustering with MIN (cont.)

- Step 1: Merge the two closest clusters (largest similarity)



Clustering with MIN (cont.)

- Step 1: Merge the two closest clusters (largest similarity)

	P1UP2	P3	P4UP5
P1UP2	1.00	0.70	0.65
P3	0.70	1.00	0.40
P4UP5	0.65	0.40	1.00

Similarity matrix



attention → **PICK min**
distance/similarity matrix
Color on Join 3
 Single = smallest distance / max similarity
 Complete = max distance / min similarity

Clustering with MIN (cont.)

- Step 2: Update proximity matrix based on MIN: proximity of two clusters is based on the two closest points in different clusters (largest similarity)

	P1UP2UP3	P4UP5
P1UP2UP3	1.00 0.70 0.65	0.65 0.70 0.90
P4UP5	0.65 1.00	1.00

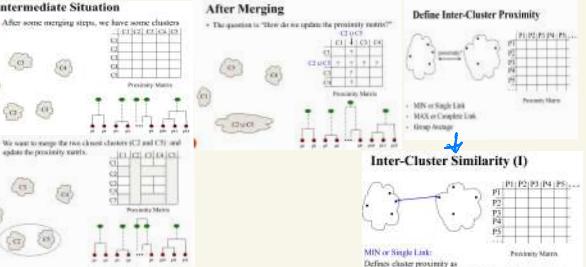
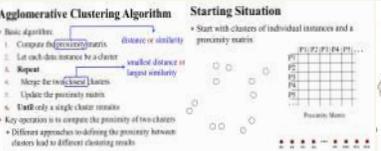
Similarity matrix

What is the proximity matrix? It is a distance matrix?

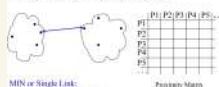


Strengths of Hierarchical Clustering

- Do not have to assume any particular number of clusters
- Any desired number of clusters can be obtained by 'cutting' the dendrogram at the proper level
- They may correspond to meaningful taxonomies
- Examples in document organization, biological sciences

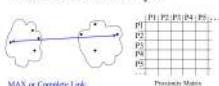


Inter-Cluster Similarity (I)



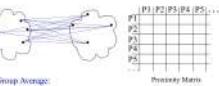
MIN or Single Link:
Defines cluster proximity as
- the proximity between the closest two data points that are in different clusters
- or the shortest edge (single link) between two nodes in different subsets (using graph terms)

Inter-Cluster Similarity (II)



MAX or Complete Link:
Defines cluster proximity as
- the proximity between the farthest two points that are in different clusters
- or the longest edge (complete link) between two nodes in different subsets (using graph terms)

Inter-Cluster Similarity (III)



Group Average:
Defines cluster proximity as
- the average pairwise proximities of all pairs of points from different clusters
- Or average length of edges between nodes in different subsets (using graph terms)

MIN or Single Link

- Similarity of two clusters is based on the two closest points (most similar) in the different clusters

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00



Clustering with MIN

- Step 1: Merge the two closest clusters (largest similarity)

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00

Similarity matrix

1 2 3 4 5

- Step 1: Merge the two closest clusters (largest similarity)

	P1UP2	P3	P4UP5
P1UP2	1.00	0.70	0.65
P3	0.70	1.00	0.40
P4UP5	0.65	0.40	1.00

Similarity matrix

1 2 3 4 5

- Step 1: Merge the two closest clusters (largest similarity)

	P1UP2UP3	P4UP5	
P1UP2UP3	1.00	0.70	0.65
P4UP5	0.65	1.00	

Similarity matrix

1 2 3 4 5

- Step 1: Merge the two closest clusters (largest similarity)

	P1UP2UP3	P4UP5	
P1UP2UP3	1.00	0.70	0.65
P4UP5	0.65	1.00	

Similarity matrix

1 2 3 4 5

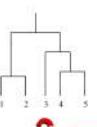
What about the proximity matrix is a distance matrix?



MAX or Complete Link

- Similarity of two clusters is based on the two most similar (most distant) points in the different clusters

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00



Tutorial

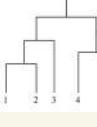
Group Average

- Proximity of two clusters is the average of pairwise proximity between points in the two clusters

$$\text{Proximity}(\mathcal{C}_i, \mathcal{C}_j) = \frac{\sum_{x_i \in \mathcal{C}_i, x_j \in \mathcal{C}_j} \text{Proximity}(x_i, x_j)}{|\mathcal{C}_i| \times |\mathcal{C}_j|}$$

- Need to use average connectivity for scalability since total proximity favors large clusters

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00



Tutorial

Agglomerative Clustering: Limitations

- Once a decision is made to combine two clusters, it cannot be undone
- No objective function is directly minimized
- Different schemes have problems with one or more of the following:
 - Sensitivity to noise and outliers
 - Difficulty in handling different-sized clusters

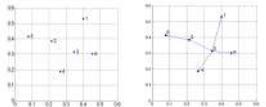
Divisive Hierarchical Clustering

- Basic algorithm:

- Generate a minimum spanning tree to connect all data instances as a single cluster
- Repeat
 - Create a new cluster by breaking the link corresponding to the largest distance (smallest similarity)
 - Until only singleton clusters remain

Minimum Spanning Tree (MST)

- Start with a tree that consists of any point
- In successive steps, look for the closest pair of points (x_i, x_j) such that one point (x_i) is in the current tree but the other (x_j) is not
- Add (x_j) to the tree and put an edge between x_i and x_j



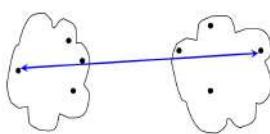
An Example of DHC

The distance matrix between 5 points

	P1	P2	P3	P4	P5
P1	0	90	10	65	20
P2	90	0	70	69	56
P3	10	70	0	40	30
P4	65	60	40	0	80
P5	20	50	30	80	0

Suppose $K = 2$, and P3 is chosen at the beginning for constructing the MST

Inter-Cluster Similarity II



	P1	P2	P3	P4	P5	...
P1	1.00	0.90	0.10	0.65	0.20	
P2	0.90	1.00	0.70	0.60	0.50	
P3	0.10	0.70	1.00	0.40	0.30	
P4	0.65	0.60	0.40	1.00	0.80	
P5	0.20	0.50	0.30	0.80	1.00	

MAX or Complete Link:

Defines cluster proximity as

- the proximity between the farthest two points that are in different clusters
- or the longest edge (complete link) between two nodes in different subsets (using graph terms)

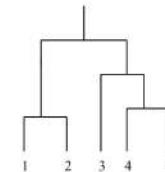
Proximity Matrix

MAX or Complete Link

- Similarity of two clusters is based on the two least similar (most distant) points in the different clusters

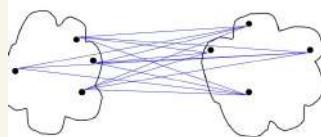
	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00

Similarity matrix



Tutorial

Inter-Cluster Similarity (III)



	P1	P2	P3	P4	P5	...
P1	1.00	0.90	0.10	0.65	0.20	
P2	0.90	1.00	0.70	0.60	0.50	
P3	0.10	0.70	1.00	0.40	0.30	
P4	0.65	0.60	0.40	1.00	0.80	
P5	0.20	0.50	0.30	0.80	1.00	

Proximity Matrix

Group Average:

Defines cluster proximity as

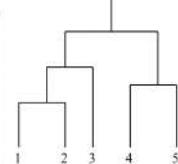
- the average pairwise proximities of all pairs of points from different clusters
- Or average length of edges between nodes in different subsets (using graph terms)

↙
due to
outlier
bias

Group Average

- Proximity of two clusters is the average of pairwise proximity between points in the two clusters
- $$\text{Proximity}(C_i, C_j) = \frac{\sum_{x_i \in C_i, x_j \in C_j} \text{Proximity}(x_i, x_j)}{|C_i| \times |C_j|}$$
- Need to use average connectivity for scalability since total proximity favors large clusters

	P1	P2	P3	P4	P5
P1	1.00	0.90	0.10	0.65	0.20
P2	0.90	1.00	0.70	0.60	0.50
P3	0.10	0.70	1.00	0.40	0.30
P4	0.65	0.60	0.40	1.00	0.80
P5	0.20	0.50	0.30	0.80	1.00

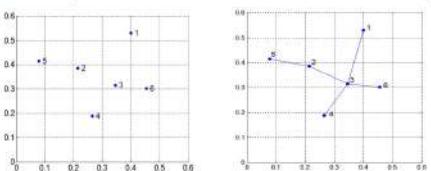


Agglomerative Clustering: Limitations

- Once a decision is made to combine two clusters, it cannot be undone
- No objective function is directly minimized
- Different schemes have problems with one or more of the following:
 - Sensitivity to noise and outliers
 - Difficulty in handling different-sized clusters

Divisive Hierarchical Clustering (cont.)

- Minimum Spanning Tree (MST)
 - Start with a tree that consists of any point
 - In successive steps, look for the closest pair of points (x_i, x_j) such that one point (x_i) is in the current tree but the other (x_j) is not
 - Add (x_j) to the tree and put an edge between x_i and x_j



Divisive Hierarchical Clustering

- Basic algorithm:
 - Generate a minimum spanning tree to connect all data instances as a single cluster
 - Repeat
 - Create a new cluster by breaking the link corresponding to the largest distance (smallest similarity)
 - Until only singleton clusters remain

An Example of DHC

no cycle

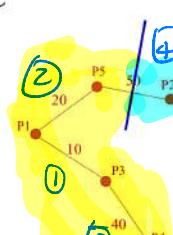
The distance matrix between 5 points

	P1	P2	P3	P4	P5
P1	—	90	—	—	—
P2	90	—	70	60	50
P3	—	70	—	—	30
P4	—	60	—	—	80
P5	—	50	—	—	—

See column cross all p when chose
Suppose $K = 2$, and P3 is chosen at the beginning for constructing the MST

pick longest distance

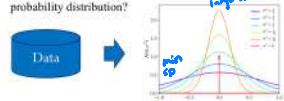
$K=2$
(2 classes)



Lesson 11: Density Estimation

Purpose of Machine Learning

- > Model Uncertainty
- > Germany will probably beat Japan, but what are the odds? 60/40, 70/30, or 80/20?
- > I'm willing to bet more money if the odds are in my favor.
- > Often translates to: what is the shape of the probability distribution?



Density Estimation

- Density estimation aims to estimate an unobservable underlying probability density function based on observed data
 - Denote by $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ the set of observed data points, drawn from an unknown $p(x)$, $x_i \sim p(x)$, for $i = 1, 2, \dots, N$
- The goal is to estimate the probability density function $p(x)$

Density Estimation Approaches

Parametric Density Estimation

- Assume that $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ are drawn from some known probability density family, $P(x|\theta)$, defined up to parameters, θ
- E.g., Gaussian distribution $N(\mu, \sigma^2)$, $\theta = [\mu, \sigma^2]$
- Estimate θ from the observed data points

- Maximum Likelihood Estimation

- Nonparametric density estimation

The General Principle

- The observed data points $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ are assumed to be a sample of N random variables independent and identically distributed (i.i.d.)
- **Identically distributed:** for any $x_i \in \mathcal{D}$, it is sampled from the same probability distribution
- **Independent:** all the data points $x_i \in \mathcal{D}$ are independent events

Why Estimating Distributions

- Recall Naïve Bayes Classifiers

$$y_i = \arg \max_c P(x|y=c)P(y=c)$$

$$= \arg \max_c \prod_{i=1}^d P(x_i|y=c)P(y=c)$$

A native, simplifying assumption, which we may remove if we can estimate the distribution properly.

State-of-the-Art Density Estimation

- Generative Adversarial Networks
 - The competition drives the counterfeiter to learn the distribution of real data.



Parametric Density Estimation

- Assume that $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ are drawn from some known probability density family, $P(x|\theta)$, defined up to parameters, θ
- $x_i \sim P(x|\theta)$
- We seek θ that makes x_i as likely as possible under $P(x|\theta)$
- Approach: maximum likelihood estimation

Maximum Likelihood Estimation

- Likelihood of parameter θ given sample \mathcal{D} :

$$l(\mathcal{D}; \theta) \triangleq P(\mathcal{D}|\theta)$$

- As $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ are i.i.d., the above likelihood is the product of the likelihoods of the individual data points
- $l(\mathcal{D}; \theta) \triangleq P(\mathcal{D}|\theta) = \prod_{i=1}^N P(x_i|\theta)$
- In MLE, we aim to find θ that makes \mathcal{D} the most likely to be drawn from. Mathematically, we aim to search for θ such that

$$\hat{\theta} = \arg \max_{\theta} l(\mathcal{D}; \theta)$$

- Typically, we maximize the log-likelihood:

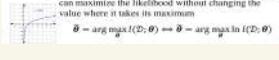
$$\mathcal{L}(\mathcal{D}; \theta) \triangleq \ln(l(\mathcal{D}; \theta))$$

- Why?
 - The $\ln(\cdot)$ function converts the product into a sum

$$\ln(l(\mathcal{D}; \theta)) = \ln P(\mathcal{D}|\theta) = \ln \left(\prod_{i=1}^N P(x_i|\theta) \right) = \sum_{i=1}^N \ln P(x_i|\theta)$$

- The $\ln(\cdot)$ function is a strictly increasing function, one can maximize the likelihood without changing the value where it takes its maximum

$$\hat{\theta} = \arg \max_{\theta} l(\mathcal{D}; \theta) \implies \hat{\theta} = \arg \max_{\theta} \ln(l(\mathcal{D}; \theta))$$



Discrete vs Continuous Probability Distributions

- Discrete Probability Distribution
 - Usually a finite number of outcomes
 - A 6-sided die \Rightarrow 6 possible outcomes
 - The distribution can be described with six numbers
 - Non-negative
 - Sum up to 1

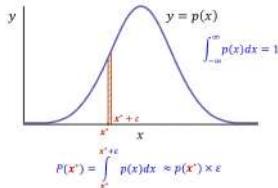
Outcomes	1	2	3	4	5	6
Probabilities	0.1	0.2	0.1	0.1	0.3	0.2

Continuous Probability Distribution

- The outcome is a real number in a range
- For example, $[0, 1], (-\infty, \infty)$
- An infinite number of outcomes between any two real numbers.
- This is really tricky.
- There is a formal branch of mathematics (measure theory) that deal with this, though we will not introduce it here.

Probability Density Function

Key message: Probability is area under the curve



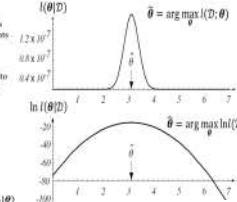
Reference: 18.05

Stability AI, the startup behind Stable Diffusion, raises \$101M

New Report: 2023 Q3 Report | 1-100M USD Fundraising in 2023



MLE Illustration



Solution of MLE

- Suppose θ contains p parameters,

$$\theta = [\theta_1, \theta_2, \dots, \theta_p]^T$$

- $\max_{\theta} \ln(l(\mathcal{D}; \theta))$ is an unconstrained optimization problem.

To solve it, we first set the derivative of $\ln(l(\mathcal{D}; \theta))$ w.r.t. θ to zero

$$\boxed{\text{Gradient operator}} \quad \nabla_{\theta} \ln(l(\mathcal{D}; \theta)) = \nabla_{\theta} \left(\sum_{i=1}^N \ln P(x_i|\theta) \right) = \sum_{i=1}^N \nabla_{\theta} \ln P(x_i|\theta) = 0,$$

$$\nabla_{\theta} = \begin{bmatrix} \frac{\partial}{\partial \theta_1} & \frac{\partial}{\partial \theta_2} & \dots & \frac{\partial}{\partial \theta_p} \end{bmatrix}^T$$

- We then obtain a solution $\hat{\theta}$ by solving the above system of equations

Univariate Gaussian

- Suppose $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$. Each data point x_i is a scalar, and is drawn from a Gaussian distribution with unknown μ and σ^2 :

$$P(x_i | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\mu = \boxed{\mathbb{E}[x]} = \int P(x; \mu, \sigma^2) dx$$

$$\sigma^2 = \text{Var}(x) = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^T]$$

Expectation, considering all possible values (infinite)

- The log-likelihood is

$$\ln l(\mathcal{D}; \theta) = \sum_{i=1}^N \ln P(x_i | \theta)$$

$$= \sum_{i=1}^N \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right)$$

$$= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{N}{2} \sum_{i=1}^N \ln(x_i - \mu)^2$$

$$= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln \sigma^2$$

Multivariate Gaussian

- Suppose $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$, and each data point x_i has d dimensions, and is drawn from a Gaussian distribution with unknown μ and Σ

$$d\text{-dimensional mean vector} \quad \mu = \mathbb{E}[x] = \int P(x | \mu, \Sigma) x dx$$

$$d \times d \text{ covariance matrix} \quad \Sigma = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^T]$$

$$P(x | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

$$|\Sigma|: \text{the determinant of } \Sigma$$

- The log-likelihood of multivariate Gaussian is

$$\ln l(\mathcal{D}; \theta) = \sum_{i=1}^N \ln P(x_i | \mu, \Sigma) = -\frac{Nd}{2} \ln(2\pi) - \frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)$$

- The derivative of the log-likelihood w.r.t. μ is

$$\frac{\partial}{\partial \mu} \ln l(\mathcal{D}; \theta) = \sum_{i=1}^N \Sigma^{-1} (x_i - \mu)$$

- By setting the derivative to zero, we have

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\begin{aligned} \text{Unbiased estimation:} \\ \mathbb{E}[\hat{\mu}] = \mu \end{aligned}$$

Summary: Estimating Multivariate Gaussian

- Suppose $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$. Each data point x_i has d dimensions, and is drawn from a Gaussian distribution with unknown μ and Σ :

$$P(x | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

- The unbiased estimators for μ and Σ are

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

Note: in this module, whenever you are asked to estimate the mean vector and the covariance matrix of data samples, use the unbiased estimators.

Is MLE always the best choice?

- What if you only have a single data point x_1 ?

The MLE $\mu = x_1$

Estimating the entire distribution based on a single data point seems unwise.



Unbiased Estimator

- An estimator is a rule for estimating a quantity based on observations.
- The MLE estimator for the Gaussian mean is $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$
- An estimator is unbiased if its expectation is the same as the true quantity.

$$\begin{aligned} \text{Unbiased estimation} \\ \mathbb{E}[\hat{\mu}] = \boxed{\mu} = \text{True} \\ \mathbb{E}[\hat{\mu}] = \frac{N-1}{N} \mu \end{aligned}$$

$$\begin{aligned} \text{To correct bias} \\ \hat{\sigma}^2 = \frac{N-1}{N} \sigma^2 \\ = \frac{1}{N-1} \sum_{i=1}^N (x_i - \mu)^2 \end{aligned}$$

Wait... Expectation of an Estimator?

- How can you talk about an expectation without a probability distribution?
- Imagine you want to estimate the average height of Singaporeans:
- Conceptually simple: measure everyone!
- Actually feasible: measure 100 randomly selected Singaporeans.
- However, the estimate will change depending on who are selected.
- These different estimators form a distribution.
- The calculation of this distribution is out of the scope of this course.

$$\begin{aligned} E[F(x)] &= \int f(x) p(x) dx \\ \prod_{i=1}^n ab &= \sum_{i=1}^n (\log a + \log b) \\ \ln e^x &= X \\ e^{bx} &= x \end{aligned}$$

- By computing the derivative of the log-likelihood w.r.t. Σ (each Σ_{ij} , i.e., the entry of the i -th row and the j -th column), and setting to be zero, we have

$$\hat{\Sigma}_{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T \quad \text{MLE estimator}$$

- The above estimation is biased!

$$\mathbb{E}[\hat{\Sigma}] = \frac{N-1}{N} \Sigma$$

- We can correct the bias by timing $\hat{\Sigma}$ with $\frac{N}{N-1}$.

$$\hat{\Sigma}_{\text{Unbiased}} = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T \quad \text{Unbiased estimator (not MLE)}$$

Covariance Matrix Σ

$$\begin{array}{ll} \text{Variance of the first dimension} & \text{Covariance of the first and second dimensions} \\ \Sigma = \begin{bmatrix} \text{Var}(x^{(1)}) & \text{Cov}(x^{(1)}, x^{(2)}) & \dots & \text{Cov}(x^{(1)}, x^{(d)}) \\ \text{Cov}(x^{(2)}, x^{(1)}) & \text{Var}(x^{(2)}) & \dots & \text{Cov}(x^{(2)}, x^{(d)}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x^{(d)}, x^{(1)}) & \text{Cov}(x^{(d)}, x^{(2)}) & \dots & \text{Var}(x^{(d)}) \end{bmatrix} \end{array}$$

$$\begin{aligned} \text{Var}(x^{(1)}) &= \mathbb{E}[(x^{(1)} - \mathbb{E}[x^{(1)}])^2] \\ \text{Cov}(x^{(1)}, x^{(j)}) &= \mathbb{E}[(x^{(1)} - \mathbb{E}[x^{(1)}])(x^{(j)} - \mathbb{E}[x^{(j)}])] \\ &= \int (x^{(1)} - \mathbb{E}[x^{(1)}])(x^{(j)} - \mathbb{E}[x^{(j)}]) P(x^{(1)}) P(x^{(j)}) dx^{(1)} dx^{(j)} \end{aligned}$$

Exam hint year

Summary: MLE for General Distributions

- Suppose $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$. Each data point x_i is of d dimensions, and drawn from a distribution with unknown parameter θ

$$x_i \sim P(x | \theta)$$

- Step 1: Compute the likelihood of θ given sample \mathcal{D} :

$$l(\mathcal{D}; \theta) \triangleq P(\mathcal{D} | \theta) = \prod_{i=1}^N P(x_i | \theta)$$

- Step 2: Compute the log-likelihood of θ given \mathcal{D} : $\ln l(\mathcal{D}; \theta)$

- Step 3: Compute the derivative of $\ln l(\theta | \mathcal{D})$ w.r.t. θ and set it to zero:

$$\nabla_{\theta} \ln l(\theta | \mathcal{D}) = 0$$

- Step 4: Solve the above system of equations to obtain the maximum likelihood estimate $\hat{\theta}$

$$\hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^N \tilde{x}_i \tilde{x}_i^T = \frac{1}{N-1} \sum_{i=1}^N \tilde{x}_i \tilde{x}_i^T = \frac{1}{N-1} \tilde{X}^T \tilde{X}$$

The Bayesian Approach

- We may introduce a prior distribution $P(\theta)$ that represents our belief about θ in the absence of any data.

$$\begin{aligned} \text{Likelihood} &= \frac{P(D|\theta)P(\theta)}{P(D)} \\ \text{Posterior} &= \frac{P(D|\theta)P(\theta)}{\int P(D|\theta)P(\theta)d\theta} \\ &\quad \text{Normalization constant} \end{aligned}$$

If the likelihood is overly concentrated, the prior can provide some smoothing effect.

$$\begin{aligned} \text{Likelihood} &= \frac{P(D|\theta)P(\theta)}{P(D)} \\ \text{Posterior} &= \frac{P(D|\theta)P(\theta)}{\int P(D|\theta)P(\theta)d\theta} \\ &\quad \text{Normalization constant} \end{aligned}$$

Density Estimation Approaches

- Parametric density estimation
- Nonparametric density estimation
 - Without assuming any forms for the underlying density
 - Assume that similar inputs have similar outputs: if x_i and x_j are similar, then $P(x_i)$ and $P(x_j)$ are similar
 - Approaches
 - Histogram Estimator
 - Naïve Estimator / Parzen Windows / Kernel Estimator
 - K-NN Estimator

Non-Parametric Density Estimation

- Assume that $\mathcal{D} = \{x_1, x_2, \dots, x_N\}$ are drawn from some unknown probability density $P(x)$
- To learn the estimator $\hat{P}(x)$ for $P(x)$
- We first focus on the univariate case, i.e., x_i is scalar
- Note that the introduced approaches can be generalized to the multivariate case easily

Look-up table for Volume of an n -ball

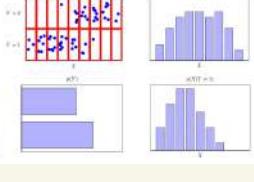
https://en.wikipedia.org/wiki/Volume_of_an_n-ball

	Dimension	Volume or a ratio of radiuses	Radius of a ball of volume 1
0	1	$\frac{1}{2}$	0.5 = $\sqrt[3]{V}$
1	$2R$	$\frac{1}{2}\pi R^2$	$\sqrt{\frac{1}{2}} = 0.707 = \sqrt[3]{V^{\frac{2}{3}}}$
2	πR^2	$\pi R^2 = 3.1415 \times R^2$	$\sqrt{\frac{\pi}{4}} = 0.364 = \sqrt[3]{V^{\frac{1}{2}}}$
3	$\frac{4\pi}{3} R^3$	$4.339 \times R^3$	$(\frac{4\pi}{3})^{\frac{1}{3}} = 0.620 = \sqrt[3]{V^{\frac{1}{3}}}$
4	$\frac{8\pi}{15} R^4$	$8.305 \times R^4$	$(\frac{8\pi}{15})^{\frac{1}{4}} = 0.670 = \sqrt[3]{V^{\frac{1}{4}}}$
5	$\frac{16\pi}{105} R^5$	$16.264 \times R^5$	$(\frac{16\pi}{105})^{\frac{1}{5}} = 0.717 = \sqrt[3]{V^{\frac{1}{5}}}$
6	$\frac{32\pi}{315} R^6$	$31.508 \times R^6$	$(\frac{32\pi}{315})^{\frac{1}{6}} = 0.763 = \sqrt[3]{V^{\frac{1}{6}}}$
7	$\frac{64\pi}{2520} R^7$	$47.125 \times R^7$	$(\frac{64\pi}{2520})^{\frac{1}{7}} = 0.801 = \sqrt[3]{V^{\frac{1}{7}}}$
8	$\frac{128\pi}{20160} R^8$	$67.652 \times R^8$	$(\frac{128\pi}{20160})^{\frac{1}{8}} = 0.839 = \sqrt[3]{V^{\frac{1}{8}}}$
9	$\frac{256\pi}{120960} R^9$	$91.088 \times R^9$	$(\frac{256\pi}{120960})^{\frac{1}{9}} = 0.876 = \sqrt[3]{V^{\frac{1}{9}}}$
10	$\frac{512\pi}{725760} R^{10}$	$125.424 \times R^{10}$	$(\frac{512\pi}{725760})^{\frac{1}{10}} = 0.913 = \sqrt[3]{V^{\frac{1}{10}}}$
11	$\frac{1024\pi}{3628800} R^{11}$	$169.760 \times R^{11}$	$(\frac{1024\pi}{3628800})^{\frac{1}{11}} = 0.944 = \sqrt[3]{V^{\frac{1}{11}}}$
12	$\frac{2048\pi}{1814400} R^{12}$	$224.096 \times R^{12}$	$(\frac{2048\pi}{1814400})^{\frac{1}{12}} = 0.970 = \sqrt[3]{V^{\frac{1}{12}}}$
13	$\frac{4096\pi}{9037440} R^{13}$	$288.432 \times R^{13}$	$(\frac{4096\pi}{9037440})^{\frac{1}{13}} = 0.994 = \sqrt[3]{V^{\frac{1}{13}}}$
14	$\frac{8192\pi}{4011520} R^{14}$	$362.768 \times R^{14}$	$(\frac{8192\pi}{4011520})^{\frac{1}{14}} = 1.017 = \sqrt[3]{V^{\frac{1}{14}}}$
15	$\frac{16384\pi}{18042240} R^{15}$	$447.104 \times R^{15}$	$(\frac{16384\pi}{18042240})^{\frac{1}{15}} = 1.037 = \sqrt[3]{V^{\frac{1}{15}}}$
16	$\frac{32768\pi}{8920800} R^{16}$	$541.440 \times R^{16}$	$(\frac{32768\pi}{8920800})^{\frac{1}{16}} = 1.066 = \sqrt[3]{V^{\frac{1}{16}}}$
17	$\frac{65536\pi}{4460400} R^{17}$	$645.776 \times R^{17}$	$(\frac{65536\pi}{4460400})^{\frac{1}{17}} = 1.094 = \sqrt[3]{V^{\frac{1}{17}}}$
18	$\frac{131072\pi}{2230200} R^{18}$	$750.112 \times R^{18}$	$(\frac{131072\pi}{2230200})^{\frac{1}{18}} = 1.123 = \sqrt[3]{V^{\frac{1}{18}}}$
19	$\frac{262144\pi}{1115100} R^{19}$	$864.448 \times R^{19}$	$(\frac{262144\pi}{1115100})^{\frac{1}{19}} = 1.152 = \sqrt[3]{V^{\frac{1}{19}}}$
20	$\frac{524288\pi}{557550} R^{20}$	$988.784 \times R^{20}$	$(\frac{524288\pi}{557550})^{\frac{1}{20}} = 1.181 = \sqrt[3]{V^{\frac{1}{20}}}$

Analysis on Histogram Estimator

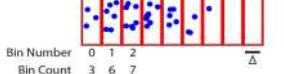
- Advantages:
 - Simple to evaluate and simple to use
 - One can throw away \mathcal{D} once the histogram is computed
 - Can be updated incrementally
- Disadvantages:
 - The estimated density has discontinuities due to the bin edges rather than any property of the underlying density
 - Scales poorly to multiple variables: we would have m^d bins (hypercubes) if we divided each feature (dimension) in a d -dimensional space into m bins

Histogram Estimator



- Simply partition x into distinct bins of a fixed width Δ
 - Count the number N_t of data points falling into bin t
 - Turn this count into a normalized probability density via dividing by the total number of observed data points N and by the width of the bins:
- $$p_t = \frac{N_t}{N\Delta}$$
- Why divide by Δ ?
- The model for the density $p(x)$ is constant over the width of each bin: find the bin where x is in (e.g., bin t), then
- $$p_t = \frac{\#\{x_i \mid x_i \text{ in the same bin as } x\}}{N\Delta} = p_t(x)$$
- For a bin t , given a density function, the probability that a data instance falling into the bin t is
- $$p_t = \int_{\Delta t} p(x)dx = p_t(x)\Delta$$
- On the other hand,
- $$p_t = \frac{\#\{x_i \mid x_i \text{ in bin } t\}}{N} = \frac{N_t}{N}$$
- Therefore:
- $$p_t(x)\Delta = \frac{N_t}{N} \quad \Rightarrow \quad p_t(x) = \frac{N_t}{N\Delta}$$

- Note: different bins may have different widths Δ_t in general, but in practice, we use the same width Δ



- Histogram density as a function of bin width Δ
- The true density is in the underlying true density from which the data points were drawn.
- When Δ is very small, the resulting density is quite noisy and fails to capture a lot of structure not present in the true density.
- When Δ is very big, the resulting density is quite smooth and consequently fails to capture the bimodality of the true density.

Naïve Estimator: An Alternative

- In Histogram Estimator, besides Δ , we have to choose an origin x_0 as well, the bins are the intervals defined as

$$[x_0 + m\Delta, x_0 + (m+1)\Delta) \quad m \text{ is an integer}$$

- The Naïve Estimator does not need to set an origin

$$\hat{p}(x) = \frac{\#\{x_i \mid x_i \text{ in the same bin as } x\}}{N\Delta}$$

Given x , use x as a center to create a bin with a length of Δ

$$\hat{p}(x) = \frac{\#\{x_i \mid x - \frac{\Delta}{2} \leq x_i < x + \frac{\Delta}{2}\}}{N\Delta}$$

$$\hat{p}(x) = \frac{\#\{x_i \mid x - \frac{\Delta}{2} \leq x_i < x + \frac{\Delta}{2}\}}{N\Delta}$$

The Naïve Estimator can also be written as

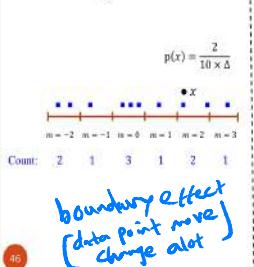
$$\hat{p}(x) = \frac{1}{N\Delta} \sum_{i=1}^N w\left(\frac{x_i - x}{\Delta}\right)$$

Windowing function $w(u) = \begin{cases} 1 & \text{if } -\frac{1}{2} \leq u \leq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$

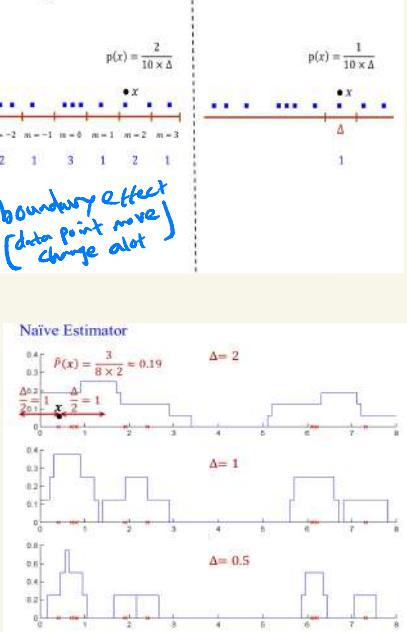
Parzen windows

Histogram v.s. Naïve Estimator

Histogram Estimator



Naïve Estimator



Generalization to Multivariate

- Suppose the observed data points $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ are d -dimensional
- In a d -dimensional space, we define \mathcal{R} is a d -dimensional hypercube with h being the length of each edge. Then the volume of the hypercube is given by

$$V = h^d$$

- The windowing function can be defined as

$$w(\mathbf{u}) = \begin{cases} 1 & \text{if } -\frac{1}{2} \leq u_j < \frac{1}{2} \text{ for all } j \in \{1, 2, \dots, d\} \\ 0 & \text{otherwise} \end{cases}$$

Defines a hypercube of length h centered at \mathbf{x} , if x_i falls in the cube, then the count is increased by 1

E.g., $d = 2$

- Hence, $w\left(\frac{\mathbf{x}_i - \mathbf{x}}{h}\right)$ is equal to unity if x_i falls within the hypercube of volume V centered at \mathbf{x} , and is zero otherwise

The density estimator can be written as

$$\hat{P}(\mathbf{x}) = \frac{\#\{\mathbf{x}_i | \mathbf{x}_i \text{ in the same hypercube as } \mathbf{x}\}}{NV}$$

Parzen window function can be a kernel function – Kernel estimator (Appendix, optional)

$$\hat{P}(\mathbf{x}) = \frac{1}{NV} \sum_{i=1}^N w\left(\frac{\mathbf{x}_i - \mathbf{x}}{h}\right)$$

$$\downarrow$$

$$\hat{P}(\mathbf{x}) = \frac{3}{16h^2}$$

K-NN Estimator

- The K-NN Estimator *adapts* the amount of smoothing to the *local* density of data, and the degree of smoothing is controlled by K , the number of neighbors

Consider K nearest neighbors of \mathbf{x}

$$p(\mathbf{x}) = \frac{K}{NV_x}$$

The volume of the space centered at \mathbf{x} that exactly contains K nearest neighbors of \mathbf{x}

- For the multivariate case

$$\hat{p}(\mathbf{x}) = \frac{K}{NV_{d_K(\mathbf{x})}}$$

The volume of the d -ball of the radius $d_K(\mathbf{x})$ centered at \mathbf{x} . And $d_K(\mathbf{x})$ is the distance of \mathbf{x} to the K -th nearest observed instance

- For the univariate case

$$\hat{p}(\mathbf{x}) = \frac{K}{N2d_K(\mathbf{x})}$$

Appendix: Kernel Estimator Optional

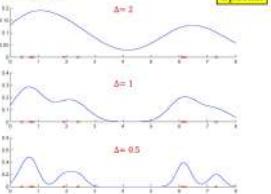
- To get a smooth estimate, we use a smooth weight function, *kernel function*, e.g., the Gaussian kernel
- For the univariate case, i.e., each data point is 1-dimensional, the Gaussian kernel is defined as

$$k(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right)$$

- The Kernel Estimator is computed via

$$\hat{P}(\mathbf{x}) = \frac{1}{N\Delta} \sum_{i=1}^N k\left(\frac{\mathbf{x}_i - \mathbf{x}}{\Delta}\right)$$

Kernel Estimator



K-Nearest Neighbor Estimator

- Recall

$$P(\mathbf{x}) = \frac{\#\{\mathbf{x}_i | \mathbf{x}_i \text{ in the same hypercube as } \mathbf{x}\}}{VN}$$

- In the previous approaches, V or Δ (or for univariate) is fixed for different queries \mathbf{x} 's
- The K-NN Estimator *adapts* the amount of smoothing to the *local* density of data, and the degree of smoothing is controlled by K , the number of neighbors

Consider K nearest neighbors of \mathbf{x}

$$P(\mathbf{x}) = \frac{K}{NV_x}$$

The volume of the space centered at \mathbf{x} that exactly contains K nearest neighbors of \mathbf{x}

- Note that the number of nearest neighbors is typically much smaller than the number of training data points, i.e., $K \ll N$

- With a predefined K , for a particular data point \mathbf{x} ,

- Compute distance between \mathbf{x} and all the observed data, e.g., Euclidean distance $\|\mathbf{x} - \mathbf{x}_i\|_2$
- Sort the observed data points based on the distances in ascending order:

$$d_1(\mathbf{x}) \leq d_2(\mathbf{x}) \leq \dots \leq d_j(\mathbf{x}) \leq \dots \leq d_N(\mathbf{x})$$

$d_1(\mathbf{x})$ is the distance of \mathbf{x} to the nearest observed instance

$d_j(\mathbf{x})$ is the distance of \mathbf{x} to the j -th nearest observed instance

- The K -NN density estimate is

$$\hat{P}(\mathbf{x}) = \frac{K}{NV_{d_K(\mathbf{x})}}$$

A ball in d -dimensional Euclidean space

The volume of the d -ball of the radius $d_K(\mathbf{x})$ centered at \mathbf{x} . And $d_K(\mathbf{x})$ is the distance of \mathbf{x} to the K -th nearest observed instance

Naïve Estimator v.s. K-NN Estimator

Univariate Case

Naïve Estimator

$$\hat{p}(\mathbf{x}) = \frac{\#\{\mathbf{x}_i | \mathbf{x}_i - \frac{\Delta}{2} \leq \mathbf{x}_i < \mathbf{x} + \frac{\Delta}{2}\}}{N\Delta}$$

K-NN Estimator

$$\hat{p}(\mathbf{x}) = \frac{K}{N(2d_K(\mathbf{x}))}$$

2-NN

$$\hat{p}(\mathbf{x}) = \frac{2}{7 \times (2 \times d_2(\mathbf{x}))}$$

Fix K , the number of observed data points to fall in the bin, and compute the bin size

K-NN Estimator

$k = 5$

$k = 3$

$k = 1$

Kernel Estimator (cont.)

Optional

- For the multivariate case, i.e., each data point is d -dimensional, the Gaussian kernel is defined as

$$k(\mathbf{u}) = \frac{1}{(2\pi)^{d/2}} \exp\left(-\frac{\|\mathbf{u}\|^2}{2}\right)$$

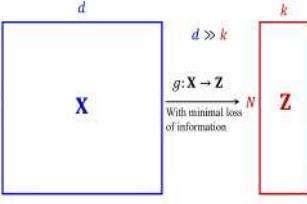
- The Kernel Estimator is computed via

$$\hat{P}(\mathbf{x}) = \frac{1}{N\Delta} \sum_{i=1}^N k\left(\frac{\mathbf{x}_i - \mathbf{x}}{\Delta}\right)$$

Lesson 12: Dimensionality Reduction

High-level Idea

- To summarize observed high-dimensional data points with low-dimensional vectors



Why Dimensionality Reduction

- To avoid curse of dimensionality
 - Decision boundary in SVM: $\mathbf{w} \cdot \mathbf{x} + b = 0$
 - Linear Regression: $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$
 - Perceptron: $f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x})$
 - One parameter to learn for every input dimension
 - Difficult to accurately estimate the best parameters when the number of data points is small
 - To identify the features, or the transformations of features that capture the most important data characteristics
 - All measurements contain error or noise.
- Mathematically we may write
- $$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{e}, \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$
- All measurements contain error or noise
- $$\mathbf{x} = \bar{\mathbf{x}} + \mathbf{e}, \mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$
- Valid hypothesis: Data reside on a straight line, but the noise is isotropic (equal amount of variation) in all three dimensions.



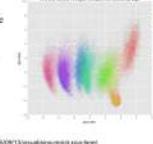
What is noise?

- These photos depict the same person
- Different lighting, makeup, facial hair, expressions, etc., create visual differences
- The data variation that we care about is far smaller than the pixel-level differences.



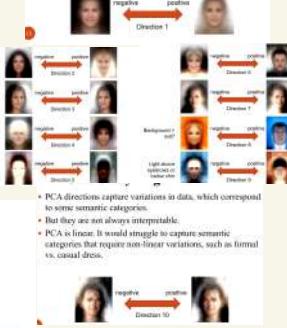
A Case Study: MNIST

- We use PCA to reduce the number of dimensions to 2 and visualize the results.
- Some clustering structure



A Case Study: Eigenface

- We use PCA to identify 10 directions that capture most variations in the data and discard the noise.
- Benefits
 - Reduces storage requirements
 - Allows visualization in 2D or 3D
 - Reduces noise and improve the performance of machine learning



Dimensionality Reduction Approaches

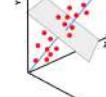
- Feature Selection
 - To select a subset of k features from the original d features to represent each data instance
 - Brute-force approach
 - Greedy search
- Feature Extraction
 - To learn k new features from the original d features to represent each data instance
 - Linear combination of original features
 - Principal component analysis
 - Nonlinear combination of original features

Principal Component Analysis

- One of the most widely-used (unsupervised) dimensionality reduction methods
- Takes a data matrix of N data points by d features, and summarizes it by principal components that are linear combinations of the original d variables
- The first k components display as much as possible of the variation among data instances

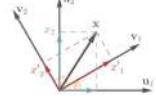
PCA: Geometric Rationale

- Data may only occupy a small subspace of the high-dimensional \mathbb{R}^d space.
- We want to pick a few ($< d$) basis vectors that data are projected onto.
- Which basis vectors would you pick?
- Goal: to find a projection or rotation of the original d -dimensional coordinate system to capture the largest amount of variance in the data
 - Ordered s.t. the 1st principal component has the highest variance, the 2nd component has the next highest variance, ..., the d -th component has the lowest variance
 - Principal components are orthogonal to each other



Linear Algebra: Change of Basis

- Vector $x = (x_1, x_2) = x_1 \mathbf{u}_1 + x_2 \mathbf{u}_2$ where $\mathbf{u}_1 = (1, 0)$ and $\mathbf{u}_2 = (0, 1)$ are the basis vectors.
- Can we change the basis vectors to \mathbf{v}_1 and \mathbf{v}_2 ?
- $x = x'_1 \mathbf{v}_1 + x'_2 \mathbf{v}_2$. We just need to find x'_1 and x'_2 by projecting x onto \mathbf{v}_1 and \mathbf{v}_2 .

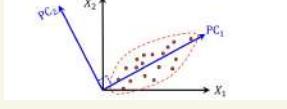


Consider a projection of a data point x onto a vector v

The projection of x onto v is

$$\frac{\mathbf{v}^T \mathbf{x}}{\|\mathbf{v}\|_2} = \frac{\|\mathbf{v}\|_2 \|\mathbf{x}\|_2 \cos(\theta)}{\|\mathbf{v}\|_2} = \|\mathbf{x}\|_2 \cos(\theta)$$

For simplicity, consider v with unit length, i.e., $\|\mathbf{v}\|_2 = 1$



Eigenvalues & Eigenvectors

- Given a d -by- d square matrix A , if there exists a non-zero d -dimensional vector \mathbf{u} , s.t.

$$A\mathbf{u} = \lambda \mathbf{u} \quad \text{scalar}$$

then \mathbf{u} is an eigenvector of A , and λ is called the corresponding eigenvalue

- Notes:

- There are d eigenvectors and eigenvalues
- An eigenvalue can be positive, negative or zero
- An eigenvector cannot be a zero vector
- For symmetric matrices, eigenvectors are orthogonal to each other

Properties of Eigenvalues

Given a square matrix A (d -by- d)

- A is invertible ($A^{-1}A = I$ or $AA^{-1} = I$) if all the eigenvalues of A are non-zero (positive or negative)
- If all the eigenvalues of A are non-negative, then A is a positive semi-definite matrix:
 - For any non-zero vector $x \in \mathbb{R}^{d \times 1}$, we have $x^T A x \geq 0$
- If all the eigenvalues of A are positive, then A is a positive definite matrix:
 - For any non-zero vector $x \in \mathbb{R}^{d \times 1}$, we have $x^T A x > 0$
- Recall: when inducing a closed form solution of regularized linear regression model, we mentioned that if a matrix A can be written as

$$A = \mathbf{X}^T \mathbf{X} + \lambda I, \text{ where } \mathbf{X} \in \mathbb{R}^{N \times d}, I \in \mathbb{R}^{d \times d} \text{ and } \lambda > 0$$

then A is always invertible:

$$\exists A^{-1}, \text{ s.t. } A^{-1}A = I \text{ and } AA^{-1} = I$$

- We first prove A is positive definite

$$\begin{aligned} \text{For any non-zero vector } \mathbf{x} \in \mathbb{R}^{d \times 1} \\ x^T A x = x^T (\mathbf{X}^T \mathbf{X} + \lambda I) x \\ = x^T (\mathbf{X}^T \mathbf{X}) x + x^T (\lambda I) x \end{aligned}$$

$$\begin{aligned} \text{Denote } \mathbf{z} = \mathbf{X} \mathbf{x} \\ = \mathbf{z}^T \mathbf{z} + \lambda \mathbf{x}^T \mathbf{x} \\ = \|\mathbf{z}\|_2^2 + \lambda \|\mathbf{x}\|_2^2 \\ \|\mathbf{z}\|_2^2 \geq 0 \text{ and } \|\mathbf{x}\|_2^2 = 0 \text{ if and only if } \mathbf{z} = \mathbf{0} \\ \|\mathbf{z}\|_2^2 > 0 \text{ because } \mathbf{x} \neq \mathbf{0} \Rightarrow \lambda \|\mathbf{x}\|_2^2 > 0 \text{ as long as } \lambda > 0 \\ x^T A x > 0 \end{aligned}$$

- As A is positive definite, all of its eigenvalues are positive, i.e., non-zero
- Recall: A is invertible if all the eigenvalues of A are non-zero (either positive or negative)
- Therefore, if a matrix A can be written as $A = \mathbf{X}^T \mathbf{X} + \lambda I$, where $\mathbf{X} \in \mathbb{R}^{N \times d}$, $I \in \mathbb{R}^{d \times d}$ and $\lambda > 0$ then A is invertible!

PCA: Variance Preservation

- The first k components display as much as possible of the variation among data instances
- Consider a projection of a data point \mathbf{x} onto a vector going through the origin, represented by \mathbf{u}
- The projection of \mathbf{x} onto \mathbf{u} is

$$\frac{\mathbf{u}^T \mathbf{x}}{\|\mathbf{u}\|_2} = \frac{\|\mathbf{u}\|_2 \|\mathbf{x}\|_2 \cos(\theta)}{\|\mathbf{u}\|_2} = \|\mathbf{x}\|_2 \cos(\theta)$$

- For simplicity, consider \mathbf{u} with unit length, i.e., $\|\mathbf{u}\|_2 = 1$
- The projected instances $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ onto \mathbf{u} are

$$\{\mathbf{u}^T \mathbf{x}_1, \mathbf{u}^T \mathbf{x}_2, \dots, \mathbf{u}^T \mathbf{x}_N\}$$

- In PCA, data points are centered at the beginning

$$\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$$

- After projection onto \mathbf{u} , the mean of data points is still 0

$$\frac{1}{N} \sum_{i=1}^N \mathbf{u}^T \mathbf{x}_i = \mathbf{u}^T \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \mathbf{0}$$

- The variance of the data points projected onto \mathbf{u} is

$$\frac{1}{N-1} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i - 0)^2 = \frac{1}{N-1} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)^2 = \mathbf{u}^T \Sigma \mathbf{u}$$

Each row is a data instance
↑

Determine Value of k

- Wrapper approaches
- Dimensionality reduction is usually an intermediate step for some downstream tasks, such as classification, regression, clustering
- Use cross-validation based on the performance of the final task to tune the value of k
- Based on the percentage of variance preserved

$$p_{var} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} \times 100$$

- All the λ_i 's are nonnegative
- Predfine a value for the percentage of variance to determine the value of k

Compute Eigenvalues and Eigenvectors

- How to compute eigenvalues and eigenvectors of $\Sigma = \frac{1}{N-1} \mathbf{X}^T \mathbf{X}$?
- In a general case, if a d -by- d square matrix \mathbf{A} can be written as

$$\mathbf{A} = \mathbf{X}^T \mathbf{X}, \text{ where } \mathbf{X} \in \mathbb{R}^{N \times d}$$

then eigenvectors and eigenvalues of $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ can be computed by performing Singular Value Decomposition (SVD) on \mathbf{X}

Orthogonal Vectors

- Two vectors \mathbf{v}_1 and \mathbf{v}_2 are said to be orthogonal if they are perpendicular to each other, i.e., the inner or dot product of two vectors is 0
- $\mathbf{v}_1 \cdot \mathbf{v}_2 = 0$
- A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ are mutually orthogonal if every pair of vectors are orthogonal
- $\mathbf{v}_1 \cdot \mathbf{v}_j = 0$, for any $i \neq j$

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \quad \mathbf{v}_2 = \begin{pmatrix} 1/\sqrt{2} \\ 1 \end{pmatrix}, \quad \mathbf{v}_3 = \begin{pmatrix} 1/\sqrt{2} \\ -1 \end{pmatrix}, \quad \mathbf{v}_4 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$\mathbf{v}_1 \cdot \mathbf{v}_2 = \mathbf{v}_1 \cdot \mathbf{v}_3 = \mathbf{v}_2 \cdot \mathbf{v}_3 = 0$$

A set of vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ are mutually orthonormal if every pair of vectors are orthogonal, and the L_2 norm of each vector is 1

- $\mathbf{v}_1 \cdot \mathbf{v}_j = 0$, for any $i \neq j$
- $\|\mathbf{v}_i\|_2 = \sqrt{\mathbf{v}_i^T \mathbf{v}_i} = 1$

A set of orthogonal vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_d\}$ can be normalized to orthonormal via $\left(\frac{\mathbf{v}_1}{\|\mathbf{v}_1\|_2}, \dots, \frac{\mathbf{v}_d}{\|\mathbf{v}_d\|_2} \right)$

$$\mathbf{v}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \quad \|\mathbf{v}_1\|_2 = \sqrt{2}, \quad \mathbf{v}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \|\mathbf{v}_2\|_2 = \sqrt{2}, \quad \mathbf{v}_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \|\mathbf{v}_3\|_2 = 2$$

$\mathbf{v}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad \|\mathbf{v}_1\|_2 = 1, \quad \mathbf{v}_2 = \begin{pmatrix} 1/\sqrt{2} \\ 1 \end{pmatrix}, \quad \|\mathbf{v}_2\|_2 = 1, \quad \mathbf{v}_3 = \begin{pmatrix} 1/\sqrt{2} \\ -1 \end{pmatrix}, \quad \|\mathbf{v}_3\|_2 = 1$

Given a matrix $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_d)$, where \mathbf{v}_i is an N -dimensional column vector, and $N \geq d$

If the columns of \mathbf{V} are mutually orthonormal, then we have

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}_d$$

$$\mathbf{I}_d = \begin{pmatrix} 1 & & & & \\ 0 & 1 & \dots & 0 & \\ \vdots & \vdots & \ddots & \vdots & \\ 0 & 0 & \dots & 1 & \end{pmatrix} \xrightarrow{d}$$

Obtain Eigenvectors via SVD

- Perform SVD on \mathbf{X} to get $\mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{U}^T$
- Then \mathbf{A} can be rewritten as $\mathbf{A} = \mathbf{X}^T \mathbf{X} = (\mathbf{V} \mathbf{D} \mathbf{U}^T)^T \mathbf{V} \mathbf{D} \mathbf{U}^T = \mathbf{U} \mathbf{D}^T \mathbf{V}^T \mathbf{V} \mathbf{D} \mathbf{U}^T$

$$\mathbf{A} = \mathbf{U} \mathbf{D}^T \mathbf{V}^T$$

Denote $\tilde{\mathbf{D}} = \mathbf{D}^T \mathbf{D} = \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T$

$\tilde{\mathbf{D}}$ is a d -by- d diagonal matrix with diagonal elements $\tilde{\lambda}_1^2 \geq \tilde{\lambda}_2^2 \geq \dots \geq \tilde{\lambda}_d^2 \geq 0$



Suppose $d < N$. All elements are zero except for the first d -by- d submatrix diagonal being $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$

$$\mathbf{A} = \mathbf{U} \tilde{\mathbf{D}} \mathbf{U}^T$$

$\mathbf{U}^T \mathbf{U} = \mathbf{I}_d$

$$(\mathbf{A} \times \mathbf{u}_1, \mathbf{A} \times \mathbf{u}_2, \dots, \mathbf{A} \times \mathbf{u}_d) = [\lambda_1^2 \mathbf{u}_1, \lambda_2^2 \mathbf{u}_2, \dots, \lambda_d^2 \mathbf{u}_d]$$

$$\mathbf{A} \mathbf{u}_i = \lambda_i^2 \mathbf{u}_i, i = 1, \dots, d$$

Each column \mathbf{u}_i of \mathbf{U} is an eigenvector of \mathbf{A} with the eigenvalue λ_i^2

Reference (Optional)

- For feature subset selection:
 - An Introduction to Variable and Feature Selection, Isabelle Guyon, Andre Elisseeff, in JMLR 2003
- For dimensionality reduction:
 - Dimensionality Reduction: A Comparative Review, L.J.P. van der Maaten and E. O. Postma and H. J. van den Herik, Technical Report, 2008
 - <https://lvdmaaten.github.io/drtoolbox/>

Appendix (optional)

Derivation of PCA

* The variance preservation view

* The first k components display as much as possible of the variation among data instances

* The minimum reconstruction view

* The first k components convey maximum useful information of original data instances

Appendix (optional)

Minimum Reconstruction Error

* Given any orthonormal basis $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_d$, if data point \mathbf{x}_i has been centered, we can write as

$$\mathbf{x}_i = \sum_{j=1}^d \langle \mathbf{x}_i, \mathbf{v}_j \rangle \mathbf{v}_j = \sum_{j=1}^d \langle \mathbf{x}_i, \mathbf{v}_j \rangle v_j = \sum_{j=1}^d a_j \mathbf{v}_j$$

* Consider the k -term approximation of \mathbf{x}_i

$$\mathbf{x}_i \approx \sum_{j=1}^k \langle \mathbf{x}_i, \mathbf{v}_j \rangle \mathbf{v}_j$$

* The error of the approximate over all data points is

$$\frac{1}{N} \sum_{i=1}^N \| \mathbf{x}_i - \hat{\mathbf{x}}_i \|^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^k \sum_{l=1}^d a_{ij} v_{jl} \hat{v}_{il}$$

Appendix (optional)

Minimum Reconstruction Error (cont.)

* The error of the approximate over all data points

$$E = \frac{1}{N} \sum_{i=1}^N \| \mathbf{x}_i - \hat{\mathbf{x}}_i \|^2 = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d a_{ij}^2 v_{j1}^2 \hat{v}_{1i}$$

$$+ \frac{1}{N} \sum_{i=1}^N \sum_{j=2}^d a_{ij}^2 v_{j2}^2 \hat{v}_{2i} + \dots + \frac{1}{N} \sum_{i=1}^N \sum_{j=d}^d a_{ij}^2 v_{jd}^2 \hat{v}_{di}$$

* Suppose $k = d-1$, i.e., we aim to remove a single dimension, then residual optimization problem is

$$\min_{\mathbf{u} \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^d a_{ij}^2 \| \mathbf{x}_i - \mathbf{u} \|^2$$

Appendix (optional)

Minimum Reconstruction Error (cont.)

* By setting the gradient of the Lagrangian w.r.t. \mathbf{u} to zero, we get

$$\hat{\mathbf{x}}_{d-1} - 2\lambda \mathbf{u} = 0 \longrightarrow \hat{\mathbf{x}}_{d-1} = \lambda \mathbf{u}$$

The desired direction \mathbf{u} is an eigenvector of Σ .

If local regression, which is:

$$\hat{\mathbf{x}}_{d-1} = \mathbf{u} + \lambda \mathbf{v}_{d-1}$$

* Our goal is to minimize the reconstruction error $\mathbf{v}_d^T \mathbf{v}_d$

$$\frac{\partial}{\partial \lambda} \mathbf{v}_d^T \mathbf{v}_d = -2\lambda \mathbf{v}_d^T \mathbf{v}_d = -2\lambda = 0 \longrightarrow \lambda = 0$$

* Therefore, \mathbf{v}_d would be the eigenvector \mathbf{u}_d of Σ with the smallest eigenvalue because $\lambda_d = \lambda_{d-1}$

* Similarly, the other dimensions to remove are subsequently the eigenvectors corresponding to the least eigenvalues

$$e^{du} = q^1 e^u$$

$$x^2 = 2x$$