

SC4000/CZ4041/CE4041: Machine Learning

Lesson 2a: Overview of Bayesian Classifiers

Kelly KE

School of Computer Science and Engineering,
NTU, Singapore

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 $f(\mathbf{x}_i) = y_i$
- 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: which team will win?

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 $P(y = c|\mathbf{x}^*)$'s to make a prediction?
 - We first review some important probability concepts

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 $A = a$

$$\sum_{a_i} P(A = a_i) = 1$$

Joint Probability

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

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

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

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

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

Sum Rule

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

$$P(A = a) = \sum_{b_i} P(A = a, B = b_i) \quad \text{OR} \quad P(A) = \sum_B P(A, B)$$

$$P(A = a) = \sum_{c_j} \sum_{b_i} P(A = a, B = b_i, C = c_j) \quad \text{OR} \quad P(A) = \sum_C \sum_B P(A, B, C)$$

Product Rule

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

$$\begin{aligned}P(A = a, B = b) &= P(B = b|A = a) \times P(A = a) \\ &= P(A = a|B = b) \times P(B = b)\end{aligned}$$

OR

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

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 **A** and **B** are a set of variables

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

An Example



V.S.



- Suppose the next match between the two teams will be hosted by Manchester United
- Head to Head Statistics:

Team	Played	Win	Draw	Lose
Manchester United	151	59	47	45
Manchester City	151	45	47	59

- Among the 59 victories for Manchester United, 32 of them come from playing at home
- Among the games won by Manchester City, 20 of them are obtained while playing on Manchester United home ground
- Among the draw games, 23 of them were played on Manchester United home ground

What result will likely be for the match?

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 $P(Y = 2|X = 0)$

Estimate Probabilities

Team	Played	Win	Draw	Lose
Manchester United	151	59	47	45
Manchester City	151	45	47	59

- To calculate $P(Y = 0)$, $P(Y = 1)$, and $P(Y = 2)$
 - $P(Y = 0) = \frac{59}{151} \approx 39\%$
 - $P(Y = 1) = \frac{45}{151} \approx 30\%$
 - $P(Y = 2) = 1 - P(Y = 0) - P(Y = 1) = 31\%$

$$\sum_{y_i} P(Y = y_i) = 1$$

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

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



Apply Bayes Rule

- Probability that Manchester United wins: $P(Y = 0) = 0.39$
- Probability that Manchester City wins: $P(Y = 1) = 0.3$
- Probability of a draw game: $P(Y = 2) = 0.31$
- Probability that Manchester United hosted the match it won: $P(X = 0|Y = 0) = 0.54$
- Probability that Manchester United hosted the match won by Manchester City: $P(X = 0|Y = 1) = 0.44$
- Probability that Manchester United hosted the match that is a draw game: $P(X = 0|Y = 2) = 0.49$
- To use Bayes rule to compute $P(Y = 1|X = 0)$, $P(Y = 0|X = 0)$ and $P(Y = 2|X = 0)$

$$P(Y = 1|X = 0)$$

Bayes rule

$$= \frac{P(X = 0|Y = 1) \times P(Y = 1)}{P(X = 0)}$$

Sum rule: $P(X) = \sum_Y P(X, Y)$

$$P(X = 0|Y = 1) \times P(Y = 1)$$

$$= P(X = 0, Y = 1) + P(X = 0, Y = 0) + P(X = 0, Y = 2)$$

$$P(X = 0|Y = 1) \times P(Y = 1)$$

$$= P(X = 0|Y = 1) \times P(Y = 1) + P(X = 0|Y = 0) \times P(Y = 0) + P(X = 0|Y = 2) \times P(Y = 2)$$

Product rule: $P(X, Y) = P(X|Y)P(Y)$

$$= \frac{0.44 \times 0.3}{0.44 \times 0.3 + 0.54 \times 0.39 + 0.49 \times 0.31} = \frac{0.132}{0.4945} = 0.267$$

Similarly

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

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

posterior \rightarrow $P(y|\mathbf{x}) = \frac{P(y, \mathbf{x})}{P(\mathbf{x})} = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})}$ \leftarrow prior

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

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

- Make a prediction based on the maximum posterior

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

Return the value of c that
maximizes $P(y = c|\mathbf{x}^*)$

where $c \in \{0, \dots, C - 1\}$

Bayesian Classifiers (cont.)

- Based on Bayes rule

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

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

Constant w.r.t.
diff. values of c

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

- Therefore, we make a prediction based on

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

- Take binary classification as an example: 0 vs 1

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

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 $P(\mathbf{x}|y = c)$ 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 --- no need to consider all possible combinations

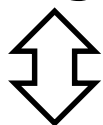
Summary on Bayesian Classifiers

- Estimate $P(y|\mathbf{x})$ via Bayes rule

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

- Make predictions based on maximum posterior

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



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

Advanced decision making: Bayesian Decision Theory (Lecture 2b)

How to estimate from training data?

- Two implementations will be introduced
 - Naïve Bayes Classifier (Lecture 3)
 - Bayesian Brief Networks or Bayesian Networks (Lecture 4)

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:

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

where $\mathbf{x} = [x_1, x_2, \dots, x_d]$

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

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)

Thank you!