

Naïve Bayes Classifier

Il-Chul Moon
Dept. of Industrial and Systems Engineering
KAIST

icmoon@kaist.ac.kr

Weekly Objectives

- Learn the optimal classification concept
 - Know the optimal predictor
 - Know the concept of Bayes risk
 - Know the concept of decision boundary
- Learn the naïve Bayes classifier
 - Understand the classifier
 - Understand the Bayesian version of linear classifier
 - Understand the conditional independence
 - Understand the naïve assumption
- Apply the naïve Bayes classifier to a case study of a text mining
 - Learn the bag-of-words concepts
 - How to apply the classifier to document classifications

NAÏVE BAYES CLASSIFIER

Dataset for Optimal Classifier Learning

| Sky | Temp | Humid | Wind | Water | Forecst | EnjoySpt |
|-------|------|--------|--------|-------|---------|----------|
| Sunny | Warm | Normal | Strong | Warm | Same | Yes |
| Sunny | Warm | High | Strong | Warm | Same | Yes |
| Rainy | Cold | High | Strong | Warm | Change | No |
| Sunny | Warm | High | Strong | Cool | Change | Yes |

- $f^*(x) = \operatorname{argmax}_{Y=y} P(X = x|Y = y)P(Y = y)$
 - $P(X=x|Y=y)$
 $= P(x_1=\text{sunny}, x_2=\text{warm}, x_3=\text{normal}, x_4=\text{strong}, x_5=\text{warm}, x_6=\text{same}|y=\text{Yes})$
 - $P(Y=y)=(y=\text{Yes})$
- How many parameters are needed? How many observations are needed?
 - $P(X=x|Y=y)$ for all x, y $\rightarrow (2^d-1)k$
 - $P(Y=y)$ for all y $\rightarrow k-1$

Often, what happens is $N \gg (2^d-1)k \gg |D|$
- Remember that we are not living in the perfect world!
 - Noise exists, so need to model it as a random variable with a distribution
 - Replications are needed!

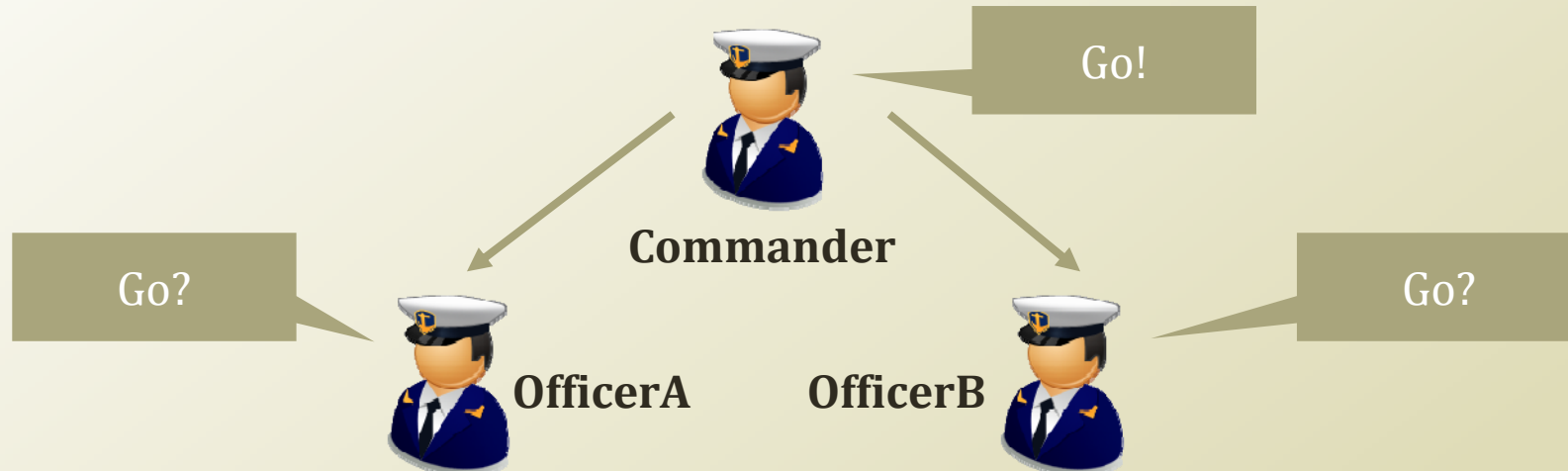
Why need an additional assumption?

- $f^*(x) = \operatorname{argmax}_{Y=y} P(X = x|Y = y)P(Y = y)$
 - To learn the above model, we need a very large dataset that is impossible to get
- The model has relaxed unrealistic assumptions, but now the model has become impossible to learn.
 - Time to add a different assumption
 - An assumption that is not so significant like the ones being relaxed
- What are the major sources of the dataset demand?
 - $P(X=x|Y=y)$ for all $x, y \rightarrow (2^d-1)k$
 - x is a vector value, and the length of the vector is d
 - d is the source of the demand
 - Then, reduce d ?
 - Or, ????

Conditional Independence

- A passing-by statistician tells us
 - Hey, what if?
 - $P(X = \langle x_1, \dots, x_i \rangle | Y = y) \rightarrow \prod_i P(X_i = x_i | Y = y)$
 - Your response: Is it possible?
 - Statistician: Yes! If x_1, \dots, x_i are conditionally independence given y
- Conditional Independence
 - x_1 is conditionally independent of x_2 given y
 - $(\forall x_1, x_2, y) \quad P(x_1 | x_2, y) = P(x_1 | y)$
 - Consequently, the above asserts
 - $P(x_1, x_2 | y) = P(x_1 | y)P(x_2 | y)$
 - Example,
 - $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightening})$
 - If there is a **lightening**, there will be a **thunder** with a prob. **p** regardless of **raining**

Conditional vs. Marginal Independence



- Marginal independence
 - $P(\text{OfficerA}=\text{Go}|\text{OfficerB}=\text{Go}) > P(\text{OfficerA}=\text{Go})$
 - **This is not marginally independent!**
 - X and Y are independent if and only if $P(X)=P(X|Y)$
 - Consequently, $P(X,Y)=P(X)P(Y)$
- Conditional independence
 - $P(\text{OfficerA}=\text{Go}|\text{OfficerB}=\text{Go},\text{Commander}=\text{Go}) = P(\text{OfficerA}=\text{Go}|\text{Commander}=\text{Go})$
 - **This is conditionally independent!**