# CS 581 – Advanced Artificial Intelligence

# Topic: Naïve Bayes

**Mustafa Bilgic**

🔗 http://www.cs.iit.edu/~mbilgic

# CLASSIFICATION

- Given a dataset $\mathcal{D} = \left\{\left\langle \vec{X}[m], Y[m] \right\rangle\right\}$ where

  - $\vec{X}$ is the input

  - $Y$ is the discrete-valued output

- Learn a function $f\left(\vec{X}\right) \rightarrow Y$

- We would like $f$ to **generalize** to **unseen** data

  - As opposed to memorizing the given data

# Classification Examples

- Email classification

- Medical diagnosis

- Face recognition

- Optical character/digit recognition

- Sentiment classification

- …

# CLASSIFICATION EXAMPLES

- Email classification
- Medical diagnosis
- Face recognition
- Optical character/digit recognition
- Sentiment classification
- ...

$$f(\text{Email}) \rightarrow S/\text{ys}$$

$$f(\text{patient}) \rightarrow \text{diagnos}$$

$$f(\text{mm}) \rightarrow \text{next word}$$

# ALGORITHMS

o Decision trees

o Nearest neighbor classification

o Naïve Bayes

o Logistic regression

o Support vector machines

o Neural networks

o …

# Bayes Classifier

- Input: $\vec{X} = \langle X_1, X_2, \dots, X_n \rangle$

- Output: $Y$

- Bayes classifier

$$P(Y \mid \vec{X}) = \frac{P(\vec{X} \mid Y)P(Y)}{P(\vec{X})} = \frac{P(Y)P(X_1, X_2, \dots, X_n \mid Y)}{P(X_1, X_2, \dots, X_n)}$$

$$P(X_1, X_2, \dots, X_n) = \sum_y P(Y = y)P(X_1, X_2, \dots, X_n \mid Y = y)$$

Assuming all variables are binary, how many independent parameters are needed for the Bayes classifier?

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

# BAYES CLASSIFIER

- Input: $\vec{X} = \langle X_1, X_2, \ldots, X_n \rangle$

- Output: $Y$

- Bayes classifier

$$1 + (2^n - 1) \times 2 = 2^{n+1} \sim 1$$

$$P(Y \mid \vec{X}) = \frac{P(\vec{X} \mid Y)P(Y)}{P(\vec{X})} = \frac{P(Y)P(X_1, X_2, \ldots, X_n \mid Y)}{P(X_1, X_2, \ldots, X_n)}$$

$$P(X_1, X_2, \ldots, X_n) = \sum_y P(Y = y)P(X_1, X_2, \ldots, X_n \mid Y = y)$$

Assuming all variables are binary, how many independent parameters are needed for the Bayes classifier?

7

# Naïve Bayes Assumption

$$X_i \perp X_j \mid Y$$

# Naïve Bayes

Bayes rule:

$$P(Y \mid X_1, X_2, \ldots, X_n) = \frac{P(Y)P(X_1, X_2, \ldots, X_n \mid Y)}{\sum_y P(y)P(X_1, X_2, \ldots, X_n \mid y)}$$

Assuming $X_i \perp X_j \mid Y$,
naïve Bayes:

$$P(Y \mid X_1, X_2, \ldots, X_n) = \frac{P(Y) \prod P(X_i \mid Y)}{\sum_y P(y) \prod P(X_i \mid y)}$$

**Assuming all variables are binary, how many independent parameters are needed for the naive Bayes classifier?**

# EXAMPLE

- See OneNote

# PARAMETER ESTIMATION

- Given a dataset $\mathcal{D} = \left\{ \left\langle \vec{X}[m], Y[m] \right\rangle \right\}$, how can we estimate
  - $P(Y)$
  - $P(X_i \mid Y)$
- Intuitive idea: count and normalize
  - But, why is this the right idea? Or, is it even the right idea?

# EXCURSION

- Maximum likelihood estimation

# ZERO PROBABILITIES

- Assume

  - $P(X_i = T|yes) = 0$ and $P(X_i = T|no) > 0$; and

  - $P(X_j \mid yes) > 0$ and $P(X_j \mid no) > 0$ for all other features

  - What is $P\left(yes \mid \vec{X}\right)$ if $X_i = T$?

- Assume

  - $P(X_i = T|yes) = 0$ and $P(X_i = T|no) = 0$; and

  - $P(X_j \mid yes) > 0$ and $P(X_j \mid no) > 0$ for all other features

  - What is $P\left(yes \mid \vec{X}\right)$ if $X_i = T$?

- Are these common?

- What can we do?

# Excursion

- Bayesian estimation

# MULTIPLYING SEVERAL PROBABILITY VALUES

- Assume we have 1,000 features

- What is the product of 1,001 probability values?

  - `p = np.random.random(1001)`

  - `p_c = np.clip(p, 0.01, 0.99)`

  - `print(np.product(p_c))`

- In naïve Bayes,

  - $a = P(Y = T) \prod P(X_i | Y = T)$

  - $b = P(Y = F) \prod P(X_i | Y = F)$

  - $P\left(Y = T \middle| \vec{X}\right) = \frac{a}{a+b}$

  - If $a = b = 0$ in your code, then what?

15

# Converting LogJoint to Cond Probs

- $\log\left(P\left(Y = yes, \vec{X}\right)\right) =$

  - $\log\left(P(Y = yes)\right) + \sum \log\left(P(X_i \mid Y = yes)\right)$

- $\log\left(P\left(Y = no, \vec{X}\right)\right) =$

  - $\log\left(P(Y = no)\right) + \sum \log\left(P(X_i \mid Y = no)\right)$

- What is $P\left(Y = yes \mid \vec{X}\right)$?

- Calculate it without causing an underflow.

  - You cannot use $\texttt{np.exp}\left(\log\left(P\left(Y = yes, \vec{X}\right)\right)\right)$

  - For example, try $\texttt{np.exp(-1000)}$

# Converting LogJoint to Cond Probs

- $\log\left(P\left(Y = yes, \vec{X}\right)\right) = -1001$

- $\log\left(P\left(Y = no, \vec{X}\right)\right) = -1002$

- $P\left(Y = yes \mid \vec{X}\right) = ?$

# NAÏVE BAYES IMPLEMENTATIONS

- Bernoulli / categorical naïve Bayes

  - Features are assumed to be binary / categorical

- Multinomial naïve Bayes

  - $P(\vec{X} \mid y)$ is a multinomial distribution

- Gaussian naïve Bayes

  - Each $p(x_i \mid y)$ is a Gaussian distribution

# READING

- [http://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf](http://www.cs.cmu.edu/~tom/mlbook/NBayesLogReg.pdf)

- [https://scikit-learn.org/stable/modules/naive_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)