# UVA CS 4501 - 001 / 6501 – 007
# Introduction to Machine Learning and Data Mining

## Lecture 13: Probability and Statistics Review (cont.) + Naïve Bayes Classifier

Yanjun Qi / Jane, PhD

University of Virginia
Department of
Computer Science

10/7/14

---

# Announcements: Schedule Change

- Midterm – rescheduled
    - Oct. 30th / 3:35pm – 4:35pm
    - Homework 4 is totally for sample midterm questions
    - HW3 will be out next Monday, due on Oct 25th
    - HW4 will be out next Tuesday, due on Oct 29th (i.e. for a good preparation for midterm. Solution will be released before due time. )

    - Grading of HW1 will be available to you this evening
    - Solution of HW1 will be available this Wed
    - Grading of HW2 will be available to you next week
    - Solution of HW2 will be available next week

10/7/14

# Where are we ? ➔
## Five major sections of this course

❑ Regression (supervised)

❑ Classification (supervised)

❑ Unsupervised models

❑ Learning theory

❑ Graphical models

10/7/14

---

# Where are we ? ➔
## Three major sections for classification

- We can divide the large variety of classification approaches into roughly three major types

1. Discriminative
   - directly estimate a decision rule/boundary
   - e.g., support vector machine, decision tree

2. Generative:
   - build a generative statistical model
   - e.g., naïve bayes classifier, Bayesian networks

3. Instance based classifiers
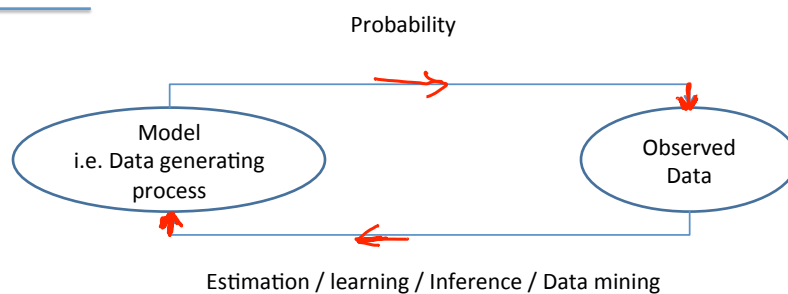   - Use observation directly (no models)
   - e.g. K nearest neighbors

10/7/14

# **Last Lecture Recap:** Probability Review

- The big picture
- Events and Event spaces
- Random variables
- Joint probability distributions, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.
- Structural properties
  - Independence, conditional independence
- Mean and Variance

10/7/14

# The Big Picture

Probability

Model
i.e. Data generating
process

Observed
Data

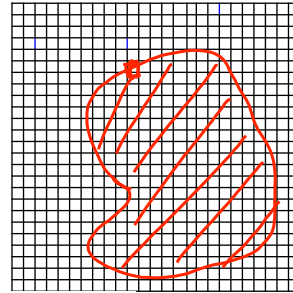Estimation / learning / Inference / Data mining

But how to specify a model?

10/7/14

# Probability as a measure of uncertainty

- *Probability is a measure of certainty of an event taking place.*

- *i.e. in the example, we were measuring the chances of hitting the shaded area.*

Its area is 1

$$prob = \frac{\#\mathrm{Re}\,dBoxes}{\#\,Boxes}$$

10/7/14

Adapt from Prof. Nando de Freitas's review slides

---

# e.g. Coin Flips

- You flip a coin
  - Head with probability 0.5

- You flip 100 coins
  - How many heads would you expect

10/7/14

# e.g. Coin Flips cont.

- You flip a coin
  - Head with probability $p$
  - Binary random variable
  - Bernoulli trial with success probability $p$
- You flip $k$ coins
  - How many heads would you expect
  - Number of heads X: discrete random variable
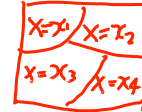  - Binomial distribution with parameters $k$ and $p$

10/7/14

# Discrete Random Variables

- Random variables (RVs) which may take on only a **countable** number of **distinct** values
  - E.g. the total number of heads X you get if you flip 100 coins

- X is a RV with arity $k$ if it can take on exactly one value out of $\{x_1, \ldots, x_k\}$
  - E.g. the possible values that X can take on are 0, 1, 2,…, 100

10/7/14

# Probability of Discrete RV

- Probability mass function (pmf): $P(X = x_i)$
- Easy facts about pmf
  - $\sum_i P(X = x_i) = 1$
  - $P(X = x_i \cap X = x_j) = 0$ if $i \neq j$
  - $P(X = x_i \cup X = x_j) = P(X = x_i) + P(X = x_j)$ if $i \neq j$
  - $P(X = x_1 \cup X = x_2 \cup \ldots \cup X = x_k) = 1$

10/7/14

# Common Distributions

- Uniform $\quad X \sim U[1,\ldots,N]$
  - X takes values 1, 2, ..., $N$
  - $P(X = i) = 1/N$
  - E.g. picking balls of different colors from a box
- Binomial $\quad X : Bin(n,p)$
  - X takes values 0, 1, ..., $n$
  - $P(X = i) = \binom{n}{i} p^i (1-p)^{n-i}$
  - E.g. coin flips

10/7/14

6

# Coin Flips of Two Persons

- Your friend and you both flip coins
  - Head with probability 0.5
  - You flip 50 times; your friend flip 100 times
  - How many heads will both of you get

10/7/14

---

# Joint Distribution

- Given two discrete RVs X and Y, their **joint distribution** is the distribution of X and Y together
  - E.g. P(You get 21 heads AND you friend get 70 heads)

- 
  - E.g. $$\sum_x \sum_y P\left(X = x \cap Y = y\right) = 1$$

$$\sum_{i=0}^{50} \sum_{j=0}^{100} P\left(\text{You get } i \text{ heads AND your friend get } j \text{ heads}\right) = 1$$

10/7/14

## Conditional Probability

- $P(X = x | Y = y)$ is the probability of $X = x$, given the occurrence of $Y = y$
  - E.g. you get 0 heads, given that your friend gets 61 heads

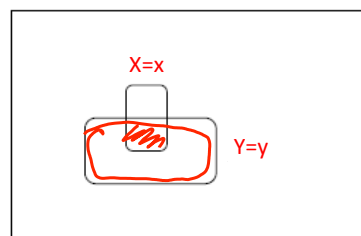- $$P(X = x | Y = y) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$

10/7/14

---

## Conditional Probability

events

$$P(X = x | Y = y) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$

But we normally write it this way:

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

X=x

Y=y

10/7/14

# Law of Total Probability

- Given two discrete RVs X and Y, which take values in $\{x_1,\ldots,x_m\}$ and $\{y_1,\ldots,y_n\}$, We have

$$P(X = x_i) = \sum_j P(X = x_i \cap Y = y_j)$$
$$= \sum_j P(X = x_i | Y = y_j)P(Y = y_j)$$

10/7/14

---

# Marginalization

$B_5$  $B_3$  $B_2$

$B_4$

A

$B_1$

$B_7$  $B_6$

Marginal Probability          Joint Probability

$$P(X = x_i) = \sum_j P(X = x_i \cap Y = y_j)$$

⇓ chain rule

$$= \sum_j P(X = x_i | Y = y_j)P(Y = y_j)$$

Conditional Probability          Marginal Probability

10/7/14

9

# Bayes Rule

- X and Y are discrete RVs…

$$P(X = x | Y = y) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$
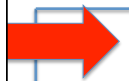
$$P(X = x_i | Y = y_j) = \frac{P(Y = y_j | X = x_i)P(X = x_i)}{\sum_k P(Y = y_j | X = x_k)P(X = x_k)}$$

10/7/14

Yanjun Qi / UVA CS 4501-01-6501-07

# **Today :** Naïve Bayes Classifier

✓ Probability review
  - Structural properties, i.e., Independence, conditional independence
✓ Naïve Bayes Classifier
  - Spam email classification

10/7/14

# Independent RVs

- Intuition: X and Y are independent means that $X = x$ **neither** makes it **more or less** probable that $Y = y$
- Definition: X and Y are independent *iff*

$$P(X = x \cap Y = y) = P(X = x)P(Y = y)$$

10/7/14

# More on Independence

- $$P(X = x \cap Y = y) = P(X = x)P(Y = y)$$

$$P(X = x \mid Y = y) = P(X = x) \qquad P(Y = y \mid X = x) = P(Y = y)$$

- E.g. no matter how many heads you get, your friend will not be affected, and vice versa

10/7/14

# More on Independence

- X is independent of Y means that knowing Y does not change our belief about X.
  - P(X|Y=y) = P(X)
  - P(X=x, Y=y) = P(X=x) P(Y=y)

  - The above should hold for all $x_i$, $y_j$
  - It is symmetric and written as X ⊥ Y

  X ⊥ Y

10/7/14

---

# Conditionally Independent RVs

- Intuition: X and Y are conditionally independent given Z means that once Z is **known**, the value of X does not add any **additional** information about Y

- Definition: X and Y are conditionally independent given Z *iff*

$$P\left(X = x \cap Y = y \middle| Z = z\right) = P\left(X = x \middle| Z = z\right)P\left(Y = y \middle| Z = z\right)$$

X ⊥ Y | Z

10/7/14

12

# More on Conditional Independence

$$P\big(X = x \cap Y = y \big| Z = z\big) = P\big(X = x \big| Z = z\big) P\big(Y = y \big| Z = z\big)$$

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

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

10/7/14

# **Today :** Naïve Bayes Classifier

✓ Probability review
   ▪ Structural properties, i.e., Independence, conditional independence
✓ Naïve Bayes Classifier
   ▪ Spam email classification

10/7/14

# Where are we ? ➜
# Three major sections for classification

- We can divide the large variety of classification approaches into roughly three major types

1. Discriminative
     - directly estimate a decision rule/boundary
     - e.g., support vector machine, decision tree

2. Generative:
     - build a generative statistical model
     - e.g., naïve bayes classifier,  Bayesian networks

3. Instance based classifiers
     - Use observation directly (no models)
     - e.g. K nearest neighbors

10/7/14

---

$$X_1 \quad X_2 \quad X_3 \quad C$$

# A Dataset for
# classification

$$f : X \longrightarrow C$$

Output as Discrete Class Label
$C_1, C_2, ..., C_L$

$$P(C \mid \mathbf{X})$$

- **Data**/*points/instances/examples/samples/records*: [ rows ]
- **Features**/*attributes/dimensions/independent variables/covariates/ predictors/regressors*: [ columns, except the last]
- **Target**/*outcome/response/label/dependent variable*: special column to be predicted [ last column ]

10/7/14

# Bayes classifiers

- Treat each attribute and class label as random variables.

- Given a sample **x** with attributes ( $x_1, x_2, \ldots, x_p$ ):
  - Goal is to predict class $C$.
  - Specifically, we want to find the value of $C_i$ that maximizes $p( C_i \mid x_1, x_2, \ldots, x_p )$.

- Can we estimate $p(C_i \mid \mathbf{x}) = p( C_i \mid x_1, x_2, \ldots, x_p )$ directly from data?

10/7/14

# Bayes classifiers
# ➔ MAP classification rule

- Establishing a probabilistic model for classification
➔ **MAP** classification rule
  - **MAP**: **M**aximum **A** **P**osterior
  - Assign $x$ to $c^*$ if

$$P(C = c^* \mid \mathbf{X} = \mathbf{x}) > P(C = c \mid \mathbf{X} = \mathbf{x}), \quad c \neq c^*, \; c = c_1, \cdots, c_L$$

10/7/14

Adapt from Prof. Ke Chen NB slides

15

# Review: Bayesian Rule

- Prior, conditional and marginal probability
  - Prior probability: $P(C)$   $P(C_1), P(C_2), ..., P(C_L)$
  - Likelihood (through a generative model): $P(\mathbf{X} \,/\, C)$
  - Evidence (marginal prob. of sample ): $P(\mathbf{X})$
  - Posterior probability: $P(C \,/\, \mathbf{X})$   $P(C_1|x), P(C_2|x), ..., P(C_L|x)$
- Bayesian Rule

$$P(C \,/\, \mathbf{X}) = \frac{P(\mathbf{X} \,/\, C)P(C)}{P(\mathbf{X})} \qquad Posterior = \frac{Likelihood \times Prior}{Evidence}$$

10/7/14

Adapt from Prof. Ke Chen NB slides

---

# Bayes Classification Rule

- Establishing a probabilistic model for classification
  - **(1) Discriminative model**

$$P(C \,/\, \mathbf{X}) \quad C = c_1, \cdots, c_L, \mathbf{X} = (X_1, \cdots, X_n)$$

$P(c_1 \,|\, \mathbf{x})$   $P(c_2 \,|\, \mathbf{x})$   $\cdots$   $P(c_L \,|\, \mathbf{x})$

**Discriminative Probabilistic Classifier**

$x_1$   $x_2$   $\cdots$   $x_n$

$$\mathbf{X} = (x_1, x_2, \cdots, x_n)$$

10/7/14

Adapt from Prof. Ke Chen NB slides

# Bayes Classification Rule

- Establishing a probabilistic model for classification (cont.)
  - **(2) Generative model**

$$P(\mathbf{X}/C) \quad C = c_1, \cdots, c_L, \ \mathbf{X} = (X_1, \cdots, X_p)$$

$P(\mathbf{x}|c_1)$  $P(\mathbf{x}|c_2)$  $P(\mathbf{x}|c_L)$

| Generative Probabilistic Model for Class *1* | Generative Probabilistic Model for Class *2* | ... | Generative Probabilistic Model for Class *L* |

$x_1 \quad x_2 \quad \cdots \quad x_p$  $x_1 \quad x_2 \quad \cdots \quad x_p$  $x_1 \quad x_2 \quad \cdots \quad x_p$

$$\mathbf{x} = (x_1, x_2, \cdots, x_p)$$

10/7/14  Adapt from Prof. Ke Chen NB slides

---

# Bayes Classification Rule

- MAP classification rule
  - **MAP**: **M**aximum **A P**osterior
  - Assign $x$ to $c^*$ if

  $$P(C = c^* \mid \mathbf{X} = \mathbf{x}) > P(C = c \mid \mathbf{X} = \mathbf{x}) \quad c \neq c^*, \ c = c_1, \cdots, c_L$$

- Generative classification with the MAP rule
  - Apply Bayesian rule to convert them into posterior probabilities

  $$P(C = c_i \mid \mathbf{X} = \mathbf{x}) = \frac{P(\mathbf{X} = \mathbf{x} \mid C = c_i)P(C = c_i)}{P(\mathbf{X} = \mathbf{x})}$$

  $$\propto P(\mathbf{X} = \mathbf{x} \mid C = c_i)P(C = c_i)$$
  $$\text{for } i = 1, 2, \cdots, L$$

  - Then apply the MAP rule

10/7/14  Adapt from Prof. Ke Chen NB slides

# Naïve Bayes

- Bayes classification

$$P(C \mid \mathbf{X}) \propto P(\mathbf{X} \mid C)P(C) = P(X_1, \cdots, X_p \mid C)P(C)$$

  Difficulty: learning the joint probability $P(X_1, \cdots, X_p \mid C)$

- Naïve Bayes classification
  - Assumption that all input attributes are conditionally independent!

$$P(X_1, X_2, \cdots, X_p \mid C) = P(X_1 \mid X_2, \cdots, X_p, C)P(X_2, \cdots, X_p \mid C)$$
$$= P(X_1 \mid C)P(X_2, \cdots, X_p \mid C)$$
$$= P(X_1 \mid C)P(X_2 \mid C) \cdots P(X_p \mid C)$$

  - MAP classification rule: for $\mathbf{x} = (x_1, x_2, \cdots, x_n)$

$$[P(x_1 \mid c^*) \cdots P(x_p \mid c^*)]P(c^*) > [P(x_1 \mid c) \cdots P(x_p \mid c)]P(c),$$
$$c \neq c^*, \ c = c_1, \cdots, c_L$$

10/7/14

Adapt from Prof. Ke Chen NB slides

---

# Naïve Bayes

- Naïve Bayes Algorithm (for discrete input attributes)
  - Learning Phase: Given a training set **S**,

    For each target value of $c_i$ $(c_i = c_1, \cdots, c_L)$
    $\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in **S**;
    For every attribute value $x_{jk}$ of each attribute $X_j$ $(j = 1, \cdots, p; \ k = 1, \cdots, K_j)$
    $\hat{P}(X_j = x_{jk} \mid C = c_i) \leftarrow$ estimate $P(X_j = x_{jk} \mid C = c_i)$ with examples in **S**;

    Output: conditional probability tables; for $X_j, K_j \times L$ elements
  - Test Phase: Given an unknown instance $\mathbf{X}' = (a_1', \cdots, a_p')$

    Look up tables to assign the label $c^*$ to $\mathbf{X}'$ if

$$[\hat{P}(a_1' \mid c^*) \cdots \hat{P}(a_p' \mid c^*)]\hat{P}(c^*) > [\hat{P}(a_1' \mid c) \cdots \hat{P}(a_p' \mid c)]\hat{P}(c),$$
$$c \neq c^*, \ c = c_1, \cdots, c_L$$

10/7/14

Adapt from Prof. Ke Chen NB slides

# Example

$X_1 X_2 X_3 C$

- Example: Play Tennis

*PlayTennis*: training examples

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1  | Sunny   | Hot         | High     | Weak   | No  |
| D2  | Sunny   | Hot         | High     | Strong | No  |
| D3  | Overcast | Hot        | High     | Weak   | Yes |
| D4  | Rain    | Mild        | High     | Weak   | Yes |
| D5  | Rain    | Cool        | Normal   | Weak   | Yes |
| D6  | Rain    | Cool        | Normal   | Strong | No  |
| D7  | Overcast | Cool       | Normal   | Strong | Yes |
| D8  | Sunny   | Mild        | High     | Weak   | No  |
| D9  | Sunny   | Cool        | Normal   | Weak   | Yes |
| D10 | Rain    | Mild        | Normal   | Weak   | Yes |
| D11 | Sunny   | Mild        | Normal   | Strong | Yes |
| D12 | Overcast | Mild       | High     | Strong | Yes |
| D13 | Overcast | Hot        | Normal   | Weak   | Yes |
| D14 | Rain    | Mild        | High     | Strong | No  |

10/7/14

---

# Example

- Learning Phase

$P(X_2|C_1), P(X_2|C_2)$

| Outlook | Play=*Yes* | Play=*No* |
|---------|------------|-----------|
| *Sunny* | 2/9 | 3/5 |
| *Overcast* | 4/9 | 0/5 |
| *Rain* | 3/9 | 2/5 |

| Temperature | Play=*Yes* | Play=*No* |
|-------------|------------|-----------|
| *Hot*  | 2/9 | 2/5 |
| *Mild* | 4/9 | 2/5 |
| *Cool* | 3/9 | 1/5 |

$P(X_4|C_1), P(X_4|C_2)$

| Humidity | Play=*Yes* | Play=*No* |
|----------|------------|-----------|
| *High*   | 3/9 | 4/5 |
| *Normal* | 6/9 | 1/5 |

| Wind | Play=*Yes* | Play=*No* |
|------|------------|-----------|
| *Strong* | 3/9 | 3/5 |
| *Weak*   | 6/9 | 2/5 |

$P(\text{Play}=Yes) = 9/14$    $P(\text{Play}=No) = 5/14$    $P(C_1), P(C_2), ..., P(C_L)$

10/7/14

---

# Example

- Test Phase
    - Given a new instance,

      $\mathbf{x}'$ =(Outlook=*Sunny*, Temperature=*Cool*, Humidity=*High*, Wind=*Strong*)
    - Look up tables

      | |
      |---|
      | P(Outlook=*Sunny*|Play=*Yes*) = 2/9 |
      | P(Temperature=*Cool*|Play=*Yes*) = 3/9 |
      | P(Huminity=*High*|Play=*Yes*) = 3/9 |
      | P(Wind=*Strong*|Play=*Yes*) = 3/9 |
      | P(Play=*Yes*) = 9/14 |

      P(Outlook=S*unny*|Play=*No*) = 3/5
      P(Temperature=*Cool*|Play==*No*) = 1/5
      P(Huminity=*High*|Play=*No*) = 4/5
      P(Wind=*Strong*|Play=*No*) = 3/5
      P(Play=*No*) = 5/14
    - MAP rule

      P(*Yes*|$\mathbf{x}'$): [P(*Sunny*|*Yes*)P(*Cool*|*Yes*)P(*High*|*Yes*)P(*Strong*|*Yes*)]P(Play=*Yes*) = 0.0053

      P(*No*|$\mathbf{x}'$): [P(*Sunny*|N*o*) P(*Cool*|*No*)P(*High*|*No*)P(*Strong*|*No*)]P(Play=*No*) = 0.0206

      Given the fact P(*Yes*|$\mathbf{x}'$) < P(*No*|$\mathbf{x}'$), we label $\mathbf{x}'$ to be "*No*".

      Adapt from Prof. Ke Chen NB slides

---

# **Next:** Naïve Bayes Classifier

- ✓ Probability review
    - ▪ Structural properties, i.e., Independence, conditional independence
- ✓ Naïve Bayes Classifier
    - ▪ Text article classification

# References

❑ Prof. Andrew Moore's review tutorial

❑ Prof. Ke Chen NB slides

❑ Prof. Carlos Guestrin recitation slides

10/7/14