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

**Lecture 12:** Probability and Statistics Review

Yanjun Qi / Jane

University of Virginia

Department of Computer Science

10/02/14

1

---

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

❑ Regression (supervised)

❑ Classification (supervised)

❑ Unsupervised models

❑ Learning theory

❑ Graphical models

10/02/14

2

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

- We can divide the large variety of classification approaches into roughly three major types
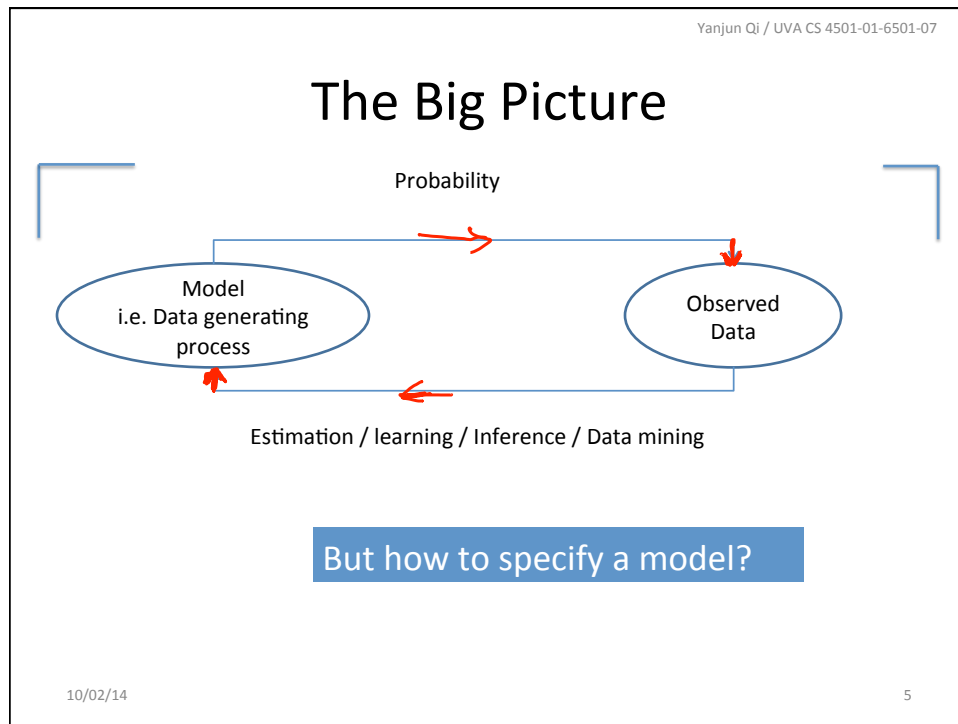
  1. Discriminative
     - directly estimate a decision rule/boundary
     - e.g., support vector machine, decision tree

  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

10/02/14                                                                 3

# Today : Probability Review

- The big picture
- Events and Event spaces
- Random variables
- Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.

10/02/14                                                                 4

# The Big Picture

Probability

Model
i.e. Data generating process

Observed
Data

Estimation / learning / Inference / Data mining

But how to specify a model?

10/02/14 5

---

# Probability as frequency
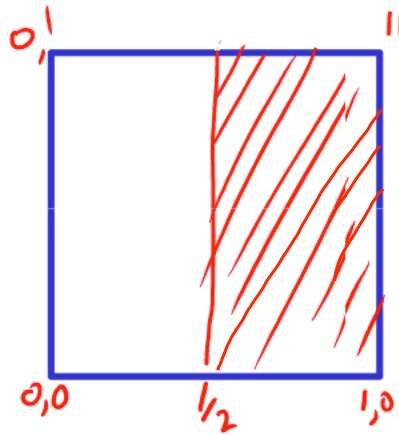
- Consider the following questions:
  - 1. What is the probability that when I flip a coin it is "heads"? We can count ➔ ~1/2
  - 2. why ?

  - 3. What is the probability of Blue Ridge Mountains to have an erupting volcano in the near future? ➔ could not count

**Message:** *The frequentist view is very useful, but it seems that we also use domain knowledge to come up with probabilities.*

10/02/14 6

Adapt from Prof. Nando de Freitas's review slides

3

---

# Probability as a measure of uncertainty

- Imagine we are throwing darts at a wall of size 1x1 and that all darts are guaranteed to fall within this 1x1 wall.

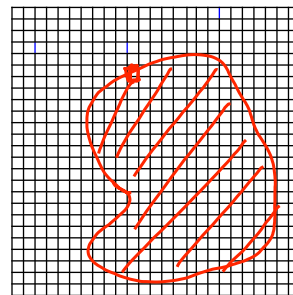- What is the probability that a dart will hit the shaded area?



10/02/14

7

Adapt from Prof. Nando de Freitas's review slides

---

# 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

$$prob = \frac{\#RedBoxes}{\#Boxes}$$

10/02/14

8

Adapt from Prof. Nando de Freitas's review slides

---

# **Today :** Probability Review

- The big picture
- Sample space, Event and Event spaces
- Random variables
- Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.
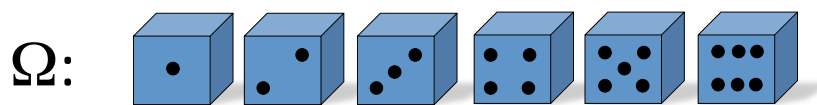
10/02/14

9

---

# Probability

***Probability*** *is the formal study of the laws of chance. Probability allows us to **manage uncertainty**.*

*The **sample space** is the set of all **outcomes**. For example, for a die we have 6 outcomes:*

$$\Omega_{die} = \{1,2,3,4,5,6\}$$

$\Omega$:

Elementary Event "Throw a 2"

The elements of $\Omega$ are called *elementary events.*

10/02/14

$$\Omega_{coin} = \{H,T\}$$

10

5

# Probability

- *Probability allows us to measure many **events**.*
- ***The events are subsets of the sample space*** $\Omega$ . *For example, for a die we may consider the following events: e.g.,*

$$\text{GREATER} = \{5, 6\}$$
$$\text{EVEN} = \{2, 4, 6\}$$

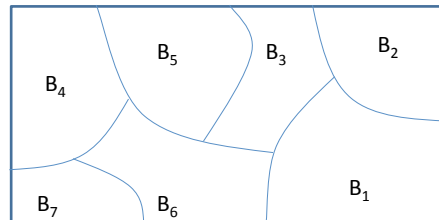- ***Assign probabilities to these events: e.g.,***

$$P(\text{EVEN}) = 1/2$$

10/02/14          Adapt from Prof. Nando de Freitas's review slides          11

# Sample space and Events

- $\Omega$ : Sample Space, result of an experiment
  - If you toss a coin twice $\Omega$ = {HH,HT,TH,TT}

$$\Omega_{toss\_once} = \{H, T\}$$

- Event: a subset of $\Omega$
  - First toss is head = {HH,HT}
- S: event space, a set of events:
  - Contains the empty event and $\Omega$

10/02/14          12

6

# Axioms for Probability

- Defined over $(\Omega, S)$ s.t.
  - $1 >= P(\alpha) >= 0$ for all $\alpha$ in S
  - $P(\Omega) = 1$
  - If $A$, $B$ are disjoint, then
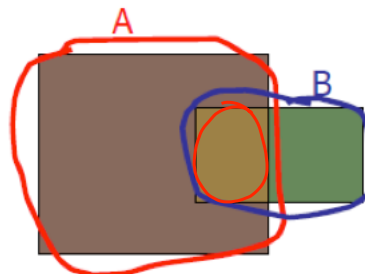    - $P(A \cup B) = p(A) + p(B)$

- $P(\Omega) = \sum P(B_i)$

$B_5$  $B_3$  $B_2$

$B_4$

$B_1$

10/02/14

$B_7$  $B_6$

13

# OR operation for Probability

- We can deduce other axioms from the above ones
  - Ex: $P(A \cup B)$ for non-disjoint events

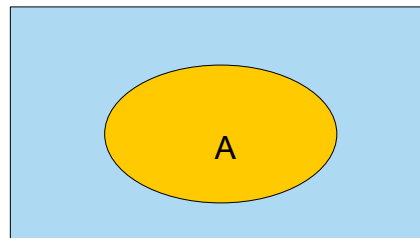$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

A

B

10/02/14

14

7

# Theorems from the Axioms

- $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

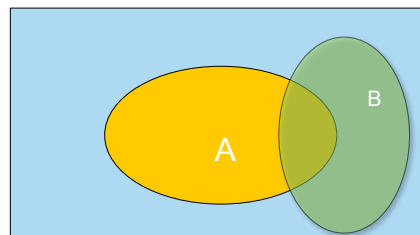From these we can prove:

$$P(\text{not } A) = P(\sim A) = 1 - P(A)$$

# Another important theorem

- $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

From these we can prove:

$$P(A) = P(A \wedge B) + P(A \wedge \sim B)$$

$$P(A) = P(A \cap \Omega)$$
$$= P(A \cap (B \cup \sim B))$$
$$= P(A \cap B) + P(A \cap \sim B)$$

# Conditional Probability



$$P(A \text{ given } B) = P(A \text{ and } B) / P(B)$$

That is, in the frequentist interpretation, we calculate the ratio of the number of times both A and B occurred and divide it by the number of times B occurred.

For short we write: $P(A|B) = P(AB)/P(B)$; or $P(AB)=P(A|B)P(B)$, where $P(A|B)$ is the *conditional* probability, $P(AB)$ is the *joint*, and $P(B)$ is the *marginal*.

If we have more events, we use the chain rule:

from Prof. Nando de
Freitas's review

10/02/14

17

$$P(ABC) = P(A|BC)\ P(B|C)\ P(C)$$

# Conditional Probability / Chain Rule

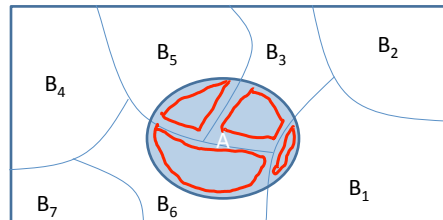- More ways to write out chain rule …

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

10/02/14

18

9

## Rule of total probability => Marginalization

Diagram with rectangle divided into regions $B_1$ through $B_7$ with circle $A$ in the middle.

$$p(A) = \sum P(B_i) P(A \mid B_i)$$

WHY ???

$P(A) = P(A \cap \Omega) = P(A \wedge (B_1 \cup B_2 \cdots \cup B_k))$

$= P((A \wedge B_1) \cup (A \wedge B_2) \cup (A \wedge B_3) \cdots \cup (A \wedge B_k))$

$= P(A \wedge B_1) + P(A \wedge B_2) + \cdots + P(A \wedge B_k)$

$= P(B_1) P(A \mid B_1) + P(B_2) P(A \mid B_2) + \cdots + P(B_k) P(A \mid B_k)$

10/02/14

---

# **Today :** Probability Review

- The big picture
- Events and Event spaces
- Random variables
- Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.

10/02/14

20

## Slide 21

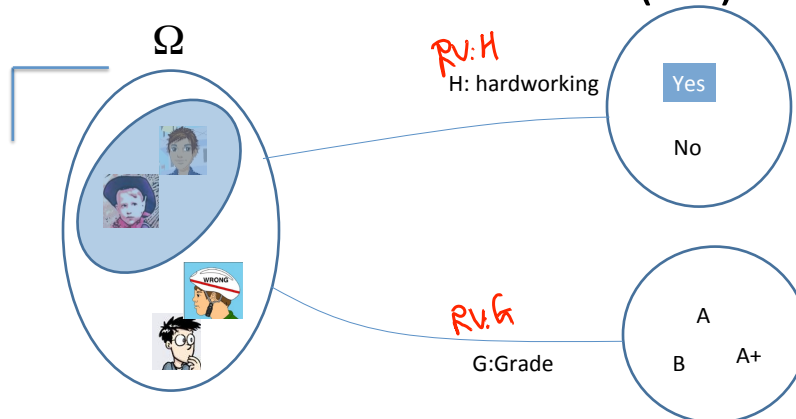# From Events to Random Variable

- Concise way of specifying attributes of outcomes
- Modeling students (Grade and Intelligence):
  - $\Omega$ = all possible students (sample space)
  - What are events (subset of sample space)
    - Grade_A = all students with grade A
    - Grade_B = all students with grade B
    - HardWorking_Yes = … who works hard
  - Very cumbersome

  - Need "functions" that maps from $\Omega$ to an attribute space T.
  - P(H = YES) = P({student ∈ $\Omega$ : H(student) = YES})

*I → RV*

21

## Slide 22

# Random Variables (RV)



$\Omega$

RV: H
H: hardworking — Yes / No

RV: G
G: Grade — A, B, A+

P(H = Yes) = P( {all students who is working hard on the course})

22

11

# Notation Digression

- P(A) is shorthand for P(A=true)

- P(~A) is shorthand for P(A=false)

- Same notation applies to other binary RVs: P(Gender=M), P(Gender=F)

- Same notation applies to *multivalued* RVs: P(Major=history), P(Age=19), P(Q=c)

- Note: upper case letters/names for *variables*, lower case letters/names for *values*

10/02/14     23

# Discrete Random Variables

- Random variables (RVs) which may take on only a **countable** number of **distinct** values

- X is a RV with arity $k$ if it can take on exactly one value out of $\{x_1, \ldots, x_k\}$

10/02/14     24

# Probability of Discrete RV

- Probability mass function (pmf): $P(X = x_i)$

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

10/02/14                                                        25

# **Today :** Probability Review

- The big picture
- Events and Event spaces
- Random variables
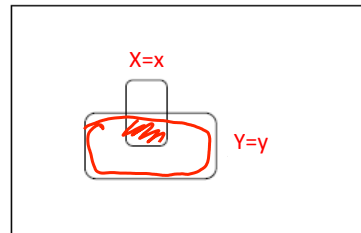- Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.

10/02/14                                                        26

13

# Conditional Probability

events

$$P\left(X = \underline{x}\middle|Y = \underline{y}\right) = \frac{P\left(X = x \cap Y = y\right)}{P\left(Y = y\right)}$$

But we will always write it this way:

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

X=x

Y=y

10/02/14　　　　　　　　　　　　　　　　　　27

---

# Marginalization

- We know p(X, Y), what is P(X=x)?
- We can use the law of total probability, why?

total prob. law

$$\underline{p(x)} = \sum_y P(x, y)$$

margin Prob.

⇓ chain rule

$$= \sum_y P(y)P(x \mid y)$$

$B_5$　$B_3$　$B_2$

$B_4$

A

$B_1$

$B_7$　$B_6$

10/02/14　　　　　　　　　　　　　　　　　　28

14

# Marginalization Cont.

- Another example

$$p(x) = \sum_{y,z} P(x,y,z)$$

↓ Chain Rule

$$= \sum_{z,y} P(y,z)P(x \mid y,z)$$

10/02/14                                                                                                29

---

# Bayes Rule

- We know that P(rain) = 0.5
  - If we also know that the grass is wet, then how this affects our belief about whether it rains or not?

$$P(rain \mid wet) = \frac{P(rain)P(wet \mid rain)}{P(wet)}$$

$$P(x \mid y) = \frac{P(x)P(y \mid x)}{P(y)}$$

10/02/14                                                                                                30

15

# What we just did...

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

This is Bayes Rule

**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London,* **53:370-418**

31

---

# More General Forms of Bayes Rule

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

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

$$P(A = a_1 \mid B) = \frac{P(B|A = a_1)P(A = a_1)}{\sum_i P(B|A = a_i)P(A = a_i)}$$

32

16

---

# Bayes Rule cont.

- You can condition on more variables

$$P(x \mid y, z) = \frac{P(x \mid z)P(y \mid x, z)}{P(y \mid z)}$$

---

# Conditional Probability Example

*Assume we have a dark box with 3 red balls and 1 blue ball. That is, we have the set {r,r,r,b}. What is the probability of drawing 2 red balls in the first 2 tries?*

Joint

chain Rule

$$P(B_1 = r, B_2 = r) = P(B_1 = r) P(B_2 = r \mid B_1 = r)$$

$$= \frac{3}{4} , \quad \frac{2}{3}$$

$$= \frac{1}{2}$$

Adapt from Prof. Nando de Freitas's review slides

## Slide 1

# Conditional Probability Example

*What is the probability that the $2^{nd}$ ball drawn from the set $\{r,r,r,b\}$ will be red?*

total prob. law

*Using marginalization,* $P(B_2 = r) = P(B_2=r, B_1=r) + P(B_2=r, B_1=b)$

$$= P(B_1=r)\, P(B_2=r \mid B_1=r) + P(B_1=b)\, P(B_2=r \mid B_1=b)$$

chain Rule

$$= \frac{3}{4} \cdot \frac{2}{3} + \frac{1}{4} \times 1$$

$$= \frac{3}{4}$$

35

## Slide 2

# Conditional Probability Example
# ➔ Matrix Notation

- $X\_1$: random variable representing first draw
- $X\_2$: random variable representing second draw
- X ==1 means "red ball", 0 mean "blue ball"

Blue → Red

*use the math notation:* $X \in \{0,1\}$

*drawn from the set* $\{r,r,r,b\}$

36

18

# Conditional Probability Example
## ➔ Matrix Notation

- $P(X_1=0) =$
- $P(X_1=1) =$
- $P(X_2=0|X_1=0) =$
- $P(X_2=1|X_1=0) =$
- $P(X_2=0|X_1=1) =$
- $P(X_2=1|X_1=1) =$

- ➔ $P(X_2=0)$
- ➔ $P(X_2=1)$



$$\pi_2 = \begin{bmatrix} P(X_2=1) \\ P(X_2=0) \end{bmatrix} \quad \text{2×1 vector}$$

$$= \begin{bmatrix} P(X_2=1, X_1=0) + P(X_2=1, X_1=1) \\ P(X_2=0, X_1=0) + P(X_2=0, X_1=1) \end{bmatrix}$$

$$= \begin{bmatrix} P(X_2=1|X_1=0)P(X_1=0) + P(X_2=1|X_1=1)P(X_1=1) \\ P(X_2=0|X_1=0)P(X_1=0) + P(X_2=0|X_1=1)P(X_1=1) \end{bmatrix}$$

$$= \underbrace{\begin{bmatrix} P(X_2=1|X_1=1) & P(X_2=1|X_1=0) \\ P(X_2=0|X_1=1) & P(X_2=0|X_1=0) \end{bmatrix}}_{G^T} \underbrace{\begin{bmatrix} P(X_1=1) \\ P(X_1=0) \end{bmatrix}}_{\pi_1}$$

10/02/14                                                                 37

---

# Conditional Probability Example
## ➔ Matrix Notation

$\{0,1\} \begin{cases} \to \text{Blue} \\ \to \text{Red} \end{cases}$

*We can obtain an expression for $P(X_2)$ easily using matrix notation:*



$$\underbrace{\begin{bmatrix} 3/4 & 1/4 \end{bmatrix}}_{\pi_2^T} = \underbrace{\begin{bmatrix} 3/4 & 1/4 \end{bmatrix}}_{\pi_1^T} \underbrace{\begin{bmatrix} 2/3 & 1/3 \\ 1 & 0 \end{bmatrix}}_{G}$$

$$\pi_2 = G^T \pi_1 \quad \Longleftrightarrow \quad \pi_2^T = \pi_1^T G$$

10/02/14                                                                 38

19

# Conditional Probability Example
# ➔ Matrix Notation

*We can obtain an expression for $P(X_2)$ easily using matrix notation:*

$$P(X_2) = \sum_{X_1 \in \{0,1\}} P(X_1)\, P(X_2|X_1)$$

*For short, we write this using vectors and a **stochastic matrix**:*

$$\underset{1 \times 2}{\Pi_1^T}\ \underset{2 \times 2}{G} = \underset{1 \times 2}{\Pi_2^T} \quad \equiv \quad \Pi_2(j) = \sum_{i=0}^{1} \Pi_1(i)\, G(i,j)$$

10/02/14

Adapt from Prof. Nando de Freitas's review slides    39

---

# **Today :** Probability Review

- The big picture
- Sample space, Event and Event spaces
- Random variables
- Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.

10/02/14    40

# References

❑ Prof. Andrew Moore's review tutorial
❑ Prof. Nando de Freitas's review slides
❑ Prof. Carlos Guestrin recitation slides