

UVA CS 6316

– Fall 2015 Graduate: Machine Learning

Lecture 11: Probability Review

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

10/7/15

1

Announcements: Schedule

- Midterm – Nov. 18th or 19th / 3:30pm – 4:45pm / open notes
- HW4 is totally for sample midterm questions
- HW4 will be out next Wed, due on Nov 16th (i.e. for a good preparation for midterm. Solution will be released before due time.)

- HW3 will be out next Monday, due on Nov 6th
- Project proposal is due on Oct 11 midnight (2 pages min.)

- Grading of HW1 is available on Collab
- Solution of HW1 is available on Collab
- Grading of HW2 will be available next week
- Solution of HW2 will be available next week

10/7/15

2

Announcements: Proposal

- The following sections are expected in the proposal ,
 - Abstract / Summary (300 words limit)
 - Introduction of the target task
 - Previous solutions for the target task
 - Why is this related to machine learning ?
 - Proposed Method (optional for this phase)
 - Experimental design / where to get data / data statistics / details of the data ?
 - Why are you the right person/team for implementing this plan ?
 - References (not included in the page counts)

10/7/15

Please use the provided template (in collab proposal page);

3

Where are we ? →

Five major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
- Graphical models

10/7/15

4

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



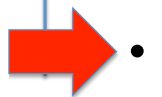
2. Generative:

- build a generative statistical model
- e.g., **naïve bayes classifier**, Bayesian networks

3. Instance based classifiers

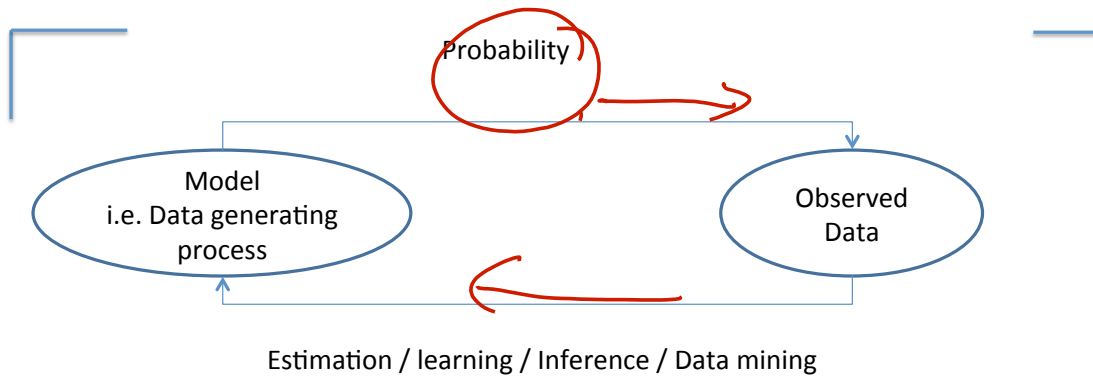
- Use observation directly (no models)
- e.g. K nearest neighbors

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.
- Structural properties
 - Independence, conditional independence

The Big Picture



But how to specify a model?

10/7/15

7

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.

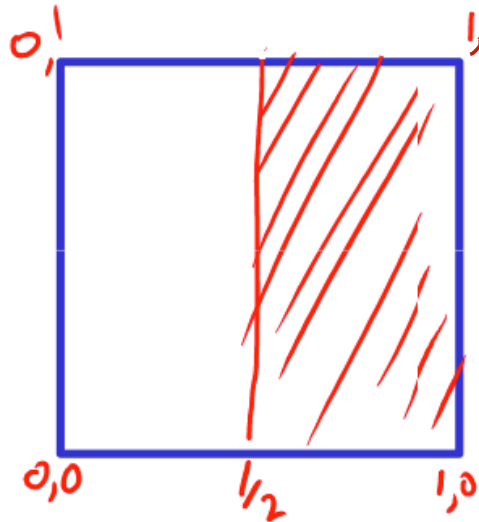
10/7/15

Adapt from Prof. Nando de Freitas's review slides

8

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?

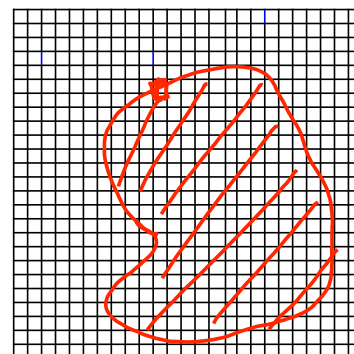


10/7/15

Adapt from Prof. Nando de Freitas's review slides ⁹

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

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

10/7/15

Adapt from Prof. Nando de Freitas's review slides

10

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.
- Structural properties
 - Independence, conditional independence

10/7/15

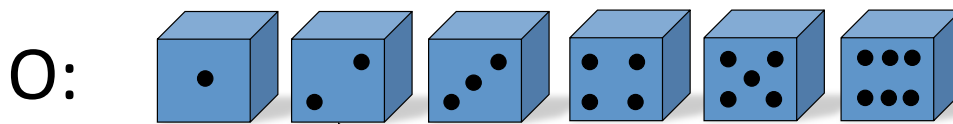
11

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:

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



Elementary Event "Throw a 2"

The elements of O are called **elementary events**.

10/7/15

12

Probability

- *Probability allows us to measure many events.*
- *The events are subsets of the sample space Ω .*
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.,*

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

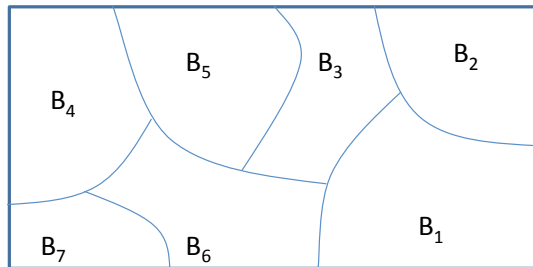
Sample space and Events

- W : **Sample Space**, result of an experiment
 - If you toss a coin twice $W = \{HH, HT, TH, TT\}$
- **Event**: a subset of W
 - First toss is head = $\{HH, HT\}$
- S : **event space, a set of events**:
 - Contains the empty event and W

Axioms for Probability

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

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

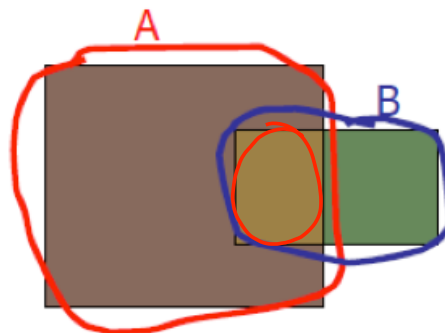


10/7/15

15

OR operation for Probability

- We can deduce other axioms from the above ones
 - Ex: $P(A \cup B)$ for **non-disjoint** events
- $$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$



10/7/15

16

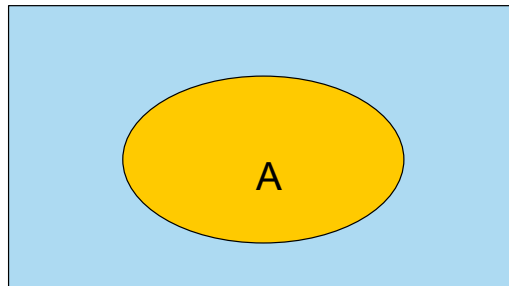
Theorems from the Axioms

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

From these we can prove:

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

Complement



10/7/15

Copyright © Andrew W. Moore

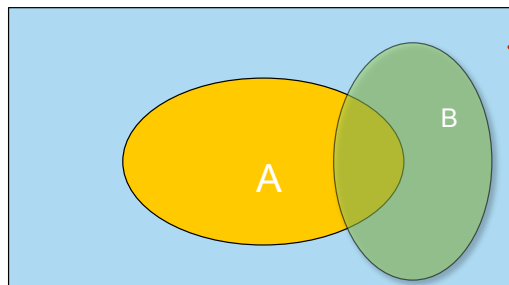
17

Another important theorem

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

From these we can prove:

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

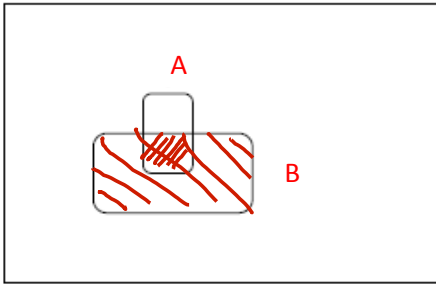


$$\begin{aligned} P(A) &= P(A \cap \Omega) \\ &= P(A \cap (B \cup \sim B)) \\ &= P((A \cap B) \cup (A \cap \sim B)) \\ &= P(A \cap B) + P(A \cap \sim B) \end{aligned}$$

10/7/15

Copyright © Andrew W. Moore

18



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
10/7/15

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

$$\begin{aligned} P(ABC) &= P(A|BC)P(B|C)P(C) \\ &= P(A|B)P(B|C)P(C) \end{aligned}$$

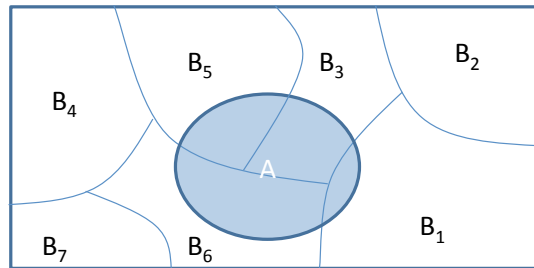
Conditional Probability / Chain Rule

- More ways to write out chain rule ...

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

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

Rule of total probability => Marginalization



$$\Rightarrow P(A) = \sum P(B_i)P(A|B_i)$$

WHY ???

$$\begin{aligned} P(A) &= P(A \cap \Omega) \\ &= P(A \cap (B_1 \cup B_2 \dots \cup B_k)) \\ &= \sum P(A \cap B_i) \end{aligned}$$

10/7/15

21

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.
- Structural properties
 - Independence, conditional independence

10/7/15

22

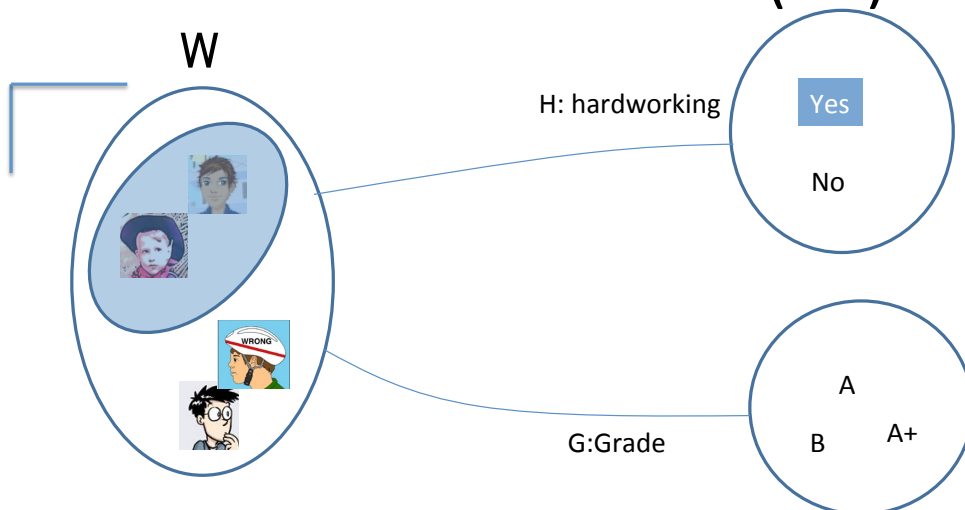
From Events to Random Variable

- Concise way of specifying attributes of outcomes
- Modeling students (Grade and Intelligence):
 - W = 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 W to an attribute space T .
- $P(H = \text{YES}) = P(\{\text{student} \in W : H(\text{student}) = \text{YES}\})$

10/7/15

23

Random Variables (RV)



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

10/7/15

24

Notation Digression

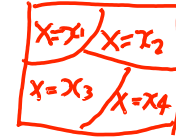
- $P(A)$ is shorthand for $P(A=\text{true})$
- $P(\sim A)$ is shorthand for $P(A=\text{false})$
- Same notation applies to other binary RVs:
 $P(\text{Gender}=\text{M})$, $P(\text{Gender}=\text{F})$
- Same notation applies to *multivalued* RVs:
 $P(\text{Major}=\text{history})$, $P(\text{Age}=19)$, $P(Q=c)$
- Note: upper case letters/names for *variables*,
lower case letters/names for *values*

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 $\{x_1, \dots, x_k\}$

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 $i \neq j$
 - $P(X = x_1 \cup X = x_2 \cup \dots \cup X = x_k) = 1$



10/7/14

e.g. Coin Flips

- You flip a coin
 - Head with probability 0.5
- You flip 100 coins
 - How many heads would you expect

10/7/15

28

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

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 $\{x_1, \dots, x_k\}$
 - **E.g.** the possible values that X can take on are 0, 1, 2, ..., 100

e.g., two Common Distributions


- Uniform $X \sim U[1, \dots, N]$
 - X takes values 1, 2, ..., N
 - $P(X=i) = 1/N$
 - E.g. picking balls of different colors from a box

- Binomial $X \sim Bin(k, p)$
 - X takes values 0, 1, ..., k
 - $P(X=i) = \binom{k}{i} p^i (1-p)^{k-i}$
 - E.g. coin flips k times

10/7/15

31

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.
- Structural properties
 - Independence, conditional independence

10/7/15

32

e.g., Coin Flips by 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

Joint Distribution

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

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

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

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

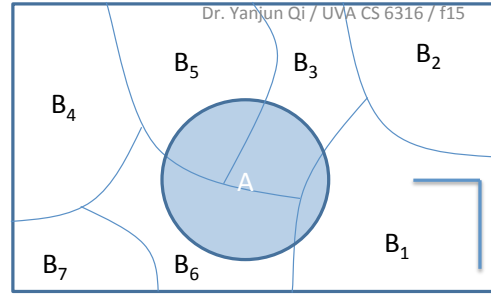
- $$P(\underbrace{X = x}|Y = y) = \frac{P(X = x \cap Y = y)}{P(Y = y)}$$

Law of Total Probability

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

$$\begin{aligned}
 &= P(X = x_i \cap \Omega) \\
 P(\underbrace{X = x_i}) &= \sum_j P(\underbrace{X = x_i \cap Y = y_j}) \\
 &= \sum_j P(X = x_i | Y = y_j) P(Y = y_j)
 \end{aligned}$$

Marginalization



Marginal Probability

Joint Probability

$$\begin{aligned}
 P(X = x_i) &= \sum_j P(X = x_i \cap Y = y_j) \\
 &= \sum_j P(X = x_i | Y = y_j) P(Y = y_j)
 \end{aligned}$$

Conditional Probability
Marginal Probability

chain rule

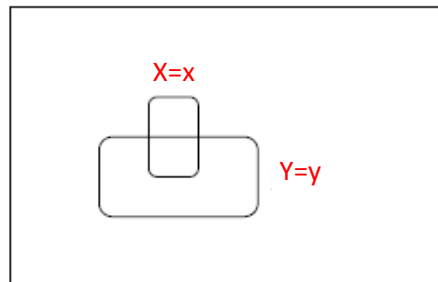
Simplify Notation: Conditional Probability

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

events

But we will always write it this way:

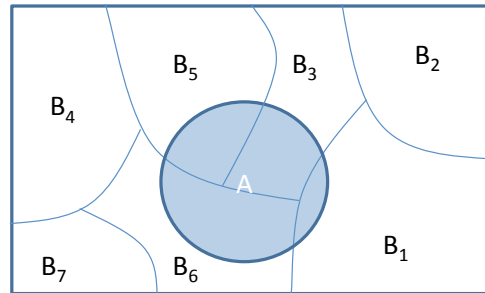
$$P(x | y) = \frac{p(x, y)}{p(y)}$$



Simplify Notation: Marginalization

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

$$p(x) = \sum_y P(x, y)$$
$$= \sum_y P(y)P(x|y)$$



Marginalization Cont.

- Another example

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

$\sum_y \sum_z p(y, z) = 1$

Bayes Rule

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

$$P(\text{rain} | \text{wet}) = \frac{P(\text{rain})P(\text{wet} | \text{rain})}{P(\text{wet})}$$

$$P(W=S | \text{wet})$$

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

Bayes Rule

- X and Y are discrete RVs...

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



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

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



10/7/15

Copyright © Andrew W. Moore

43

More General Forms of Bayes Rule

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

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

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

10/7/15

Copyright © Andrew W. Moore

44

Bayes Rule cont.

- You can condition on more variables

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

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?

$$\begin{aligned} P(B_1 = r, B_2 = r) &= P(B_1 = r) P(B_2 = r | B_1 = r) \\ &= \frac{3}{4} \times \frac{2}{3} = \frac{1}{2} \end{aligned}$$

Conditional Probability Example

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

$$\begin{aligned} \text{Using marginalization, } P(B_2 = r) &= P(B_2 = r \wedge B_1 = r) \\ &\quad + P(B_2 = r \wedge B_1 = b) \\ &= P(B_1 = r) P(B_2 = r | B_1 = r) + P(B_1 = b) P(B_2 = r | B_1 = b) \end{aligned}$$

10/7/15

47

$$\begin{bmatrix} P(B_2 = r) \\ P(B_2 = b) \end{bmatrix} = \begin{bmatrix} P(B_2 = r | B_1 = r) P(B_1 = r) + P(B_2 = r | B_1 = b) P(B_1 = b) \\ P(B_2 = b | B_1 = r) P(B_1 = r) + P(B_2 = b | B_1 = b) P(B_1 = b) \end{bmatrix}$$

Dr. Yanjun Qi / UVA CS 6316 / f15

→ Matrix Notation

$$\begin{bmatrix} P(B_2 = r) \\ P(B_2 = b) \end{bmatrix} = \begin{bmatrix} P(B_2 = r | B_1 = r), & P(B_2 = r | B_1 = b) \\ P(B_2 = b | B_1 = r), & P(B_2 = b | B_1 = b) \end{bmatrix} \begin{bmatrix} P(B_1 = r) \\ P(B_1 = b) \end{bmatrix}$$

$$\Rightarrow \text{matrix notation form } \Pi_2 = G^T \Pi_1$$

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

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.
- Structural properties
 - Independence, conditional independence



Today Recap : Probability Review

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

References

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