



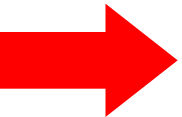
Machine Learning

Lecture 12: Probability Review

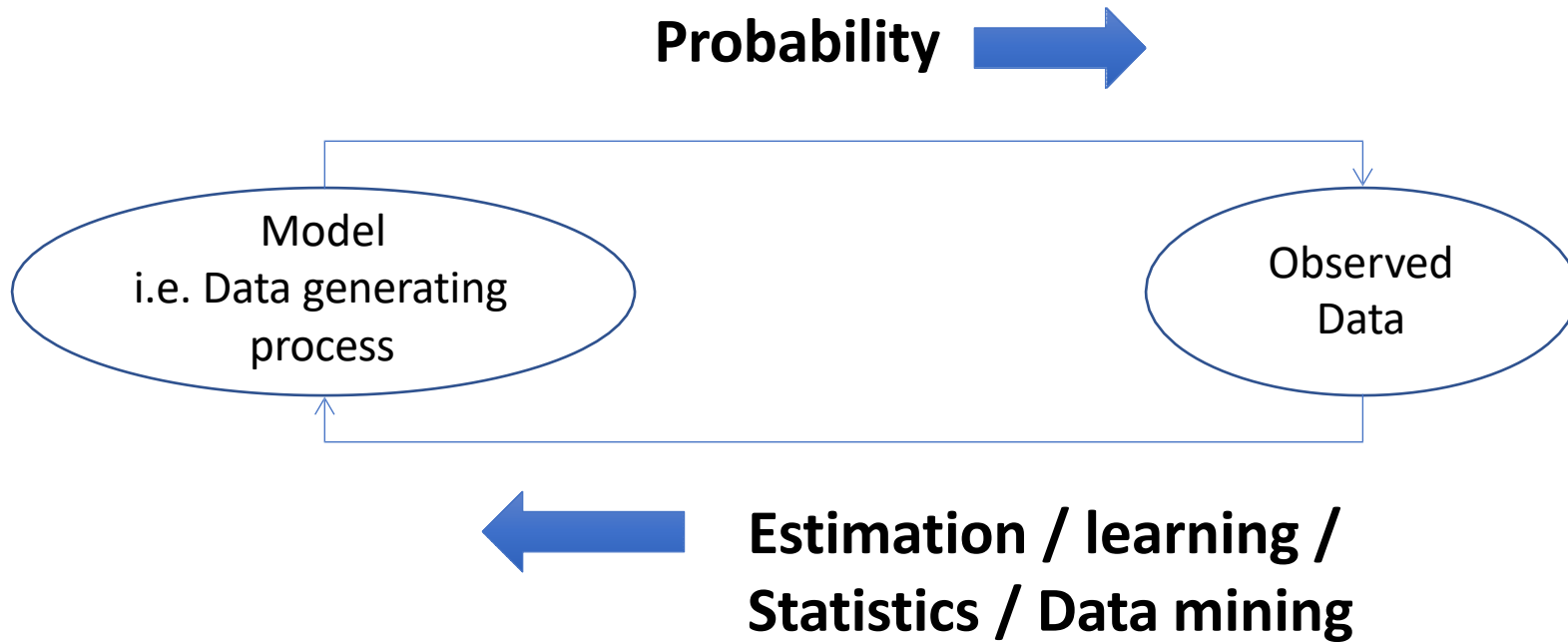
Dr. Beilun Wang

Southeast University
School of Computer Science
and Engineering

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, e.g., Independence, conditional independence
 - Maximum Likelihood Estimation

The Big Picture



Probability

- Counting
- Basics of probability
- Conditional probability
- Random variables
- Discrete and continuous distributions
- Expectation and variance
- Tail bounds and central limit theorem
-

Statistics

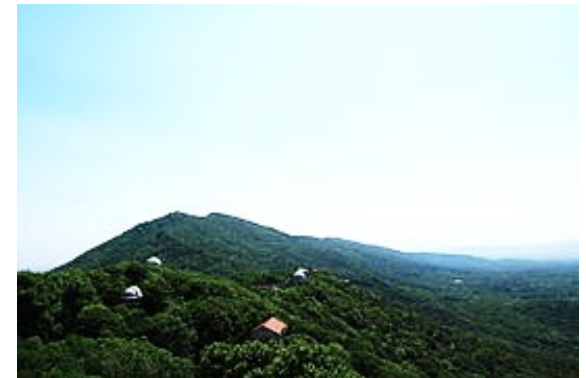
- Maximum likelihood estimation
- Bayesian estimation
- Hypothesis testing
- Linear regression
-

Probability as frequency

- Consider the following questions:
 - What is the probability that when I flip a coin it is “heads”?
→ 50%
 - What is the probability of Zijin Mountains to have a mudslide in the near future?
→ could not count

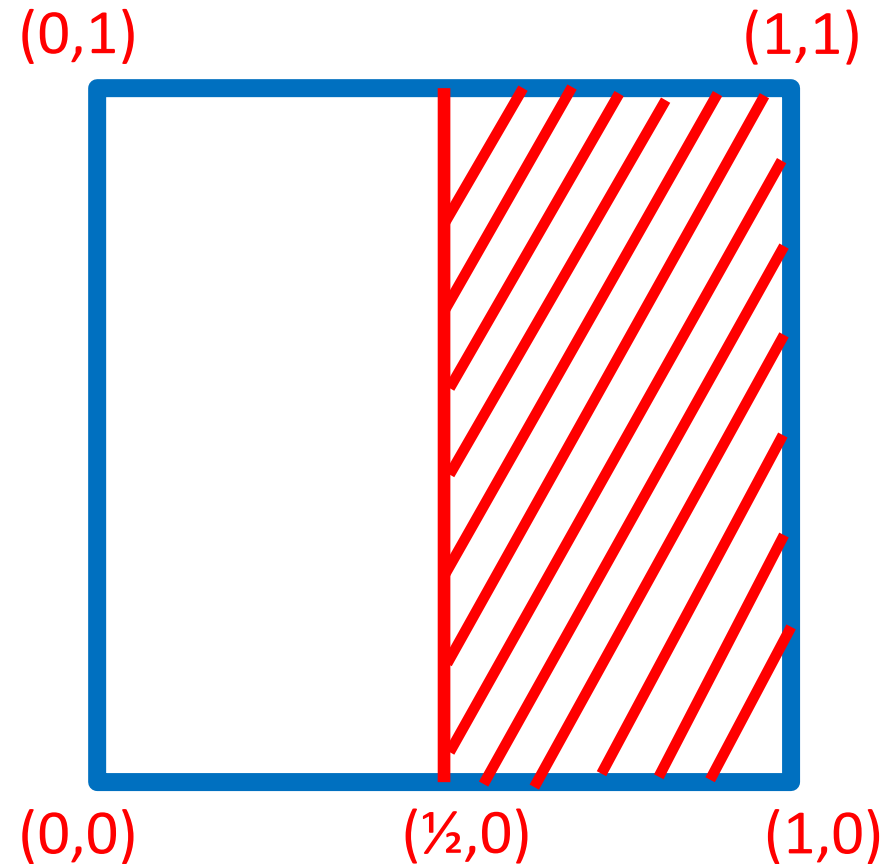


Message: The **frequentist** view is very useful, but it seems that we can also use **domain knowledge** to come up with probabilities.



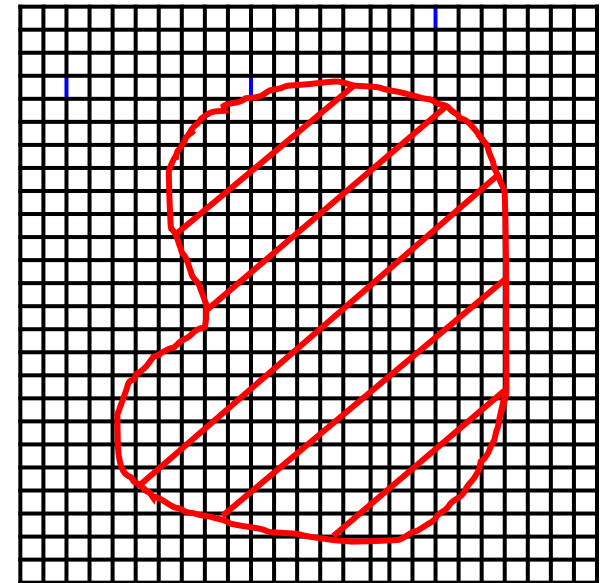
Probability as a measure of uncertainty

- Imagine we are throwing darts at a wall 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?



Probability as a measure of uncertainty

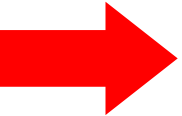
- Probability is a **measure of certainty of an event taking place.**
- i.e. in the example, we were measuring the chances of hitting the shaded area.



Its area is 1 →

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

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, e.g., Independence, conditional independence
 - Maximum Likelihood Estimation

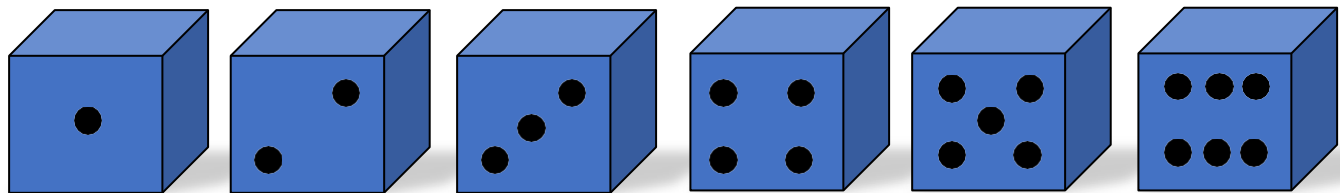
Probability

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:

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

O:



Elementary Event “Throw 2”

The elements of O are called *elementary events*.

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

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



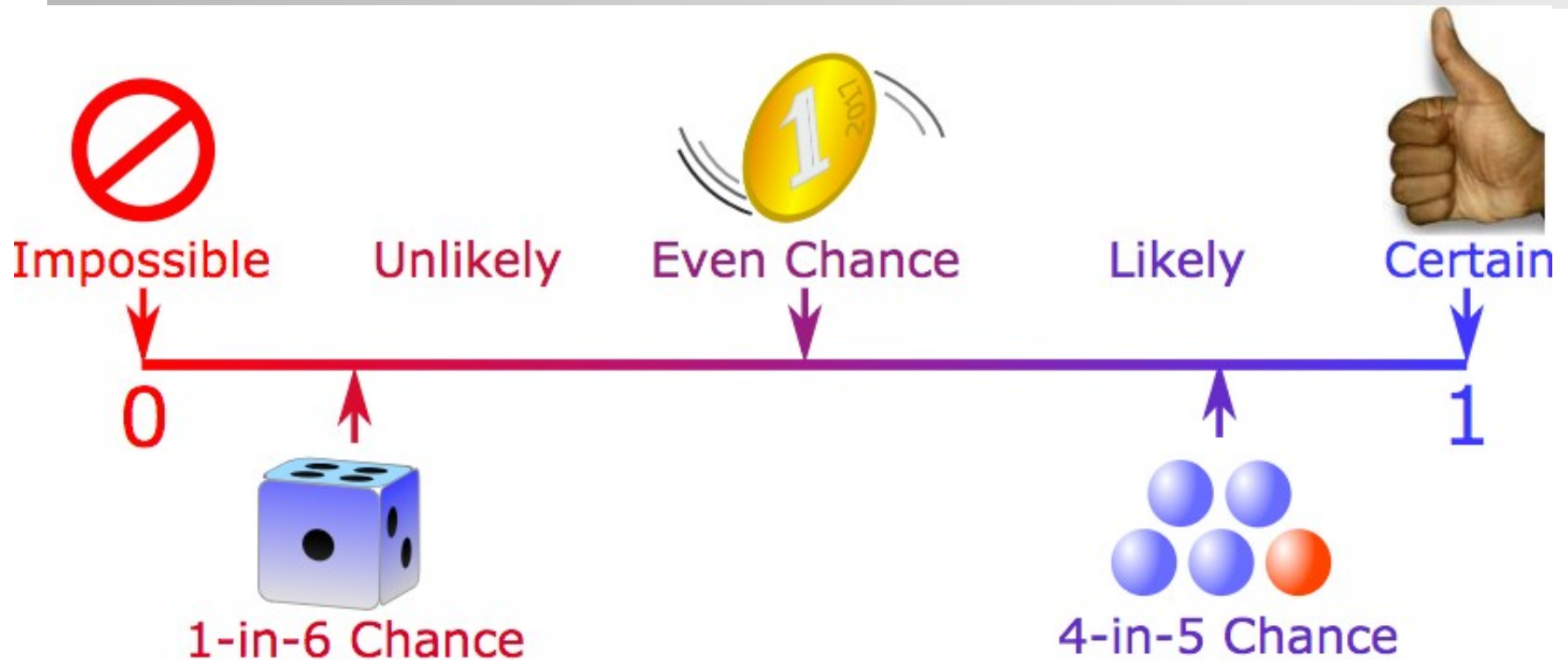


Axioms for Probability

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

Sample space
Event space

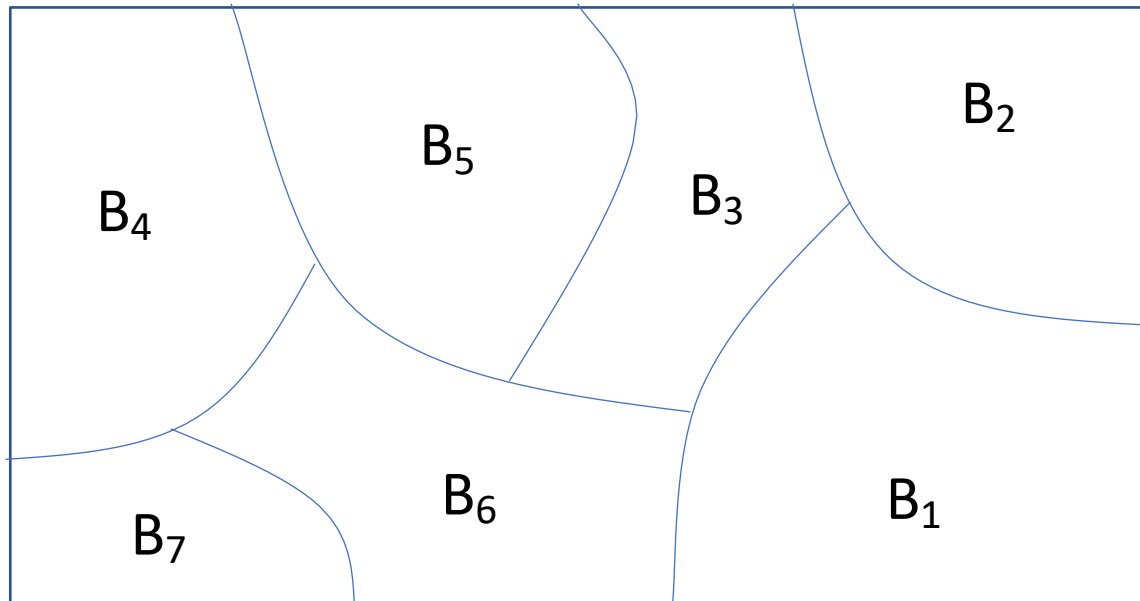
Axioms for Probability



Probability is always between 0 and 1

Axioms for Probability

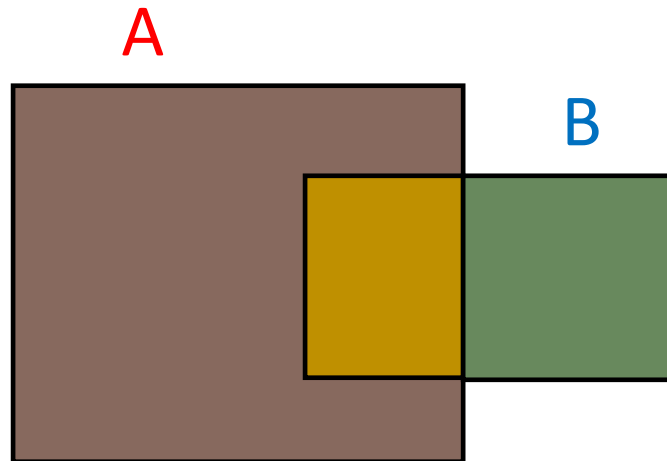
- $P(O) = \sum P(B_i) = 1$



OR & AND operation for Probability

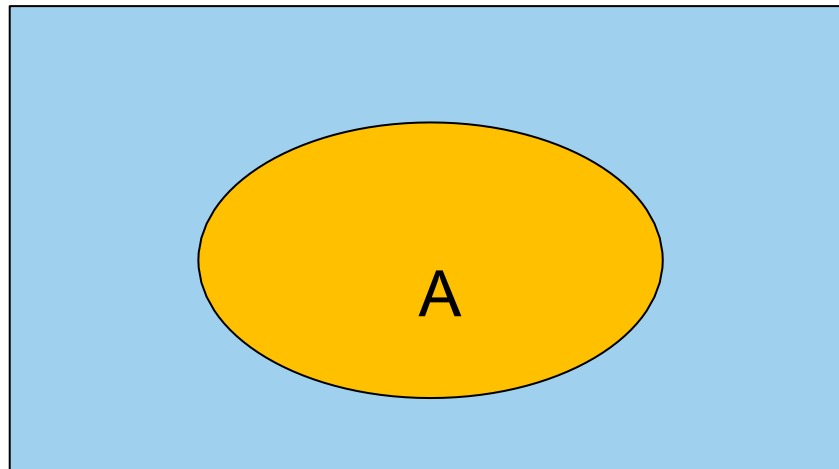
- We can deduce other axioms from the above ones
 - Ex: $P(A \cup B)$ for **non-disjoint** events

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



NOT operation for Probability

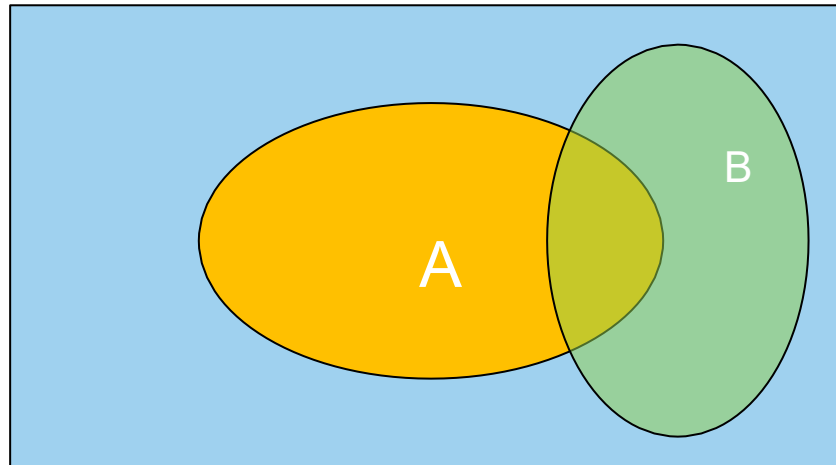
- $0 \leq P(A) \leq 1$
- $P(\textcolor{red}{A} \text{ or } \textcolor{blue}{B}) = P(\textcolor{red}{A}) + P(\textcolor{blue}{B}) - P(\textcolor{red}{A} \text{ and } \textcolor{blue}{B})$
- From these we can prove:
$$P(\text{not } A) = P(\sim A) = 1 - P(A)$$



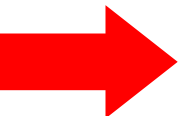
Law of Total Probability

- $0 \leq P(A) \leq 1$
- $P(\textcolor{red}{A} \text{ or } \textcolor{blue}{B}) = P(\textcolor{red}{A}) + P(\textcolor{blue}{B}) - P(\textcolor{red}{A} \text{ and } \textcolor{blue}{B})$
- From these we can prove:

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



Today: Probability Review

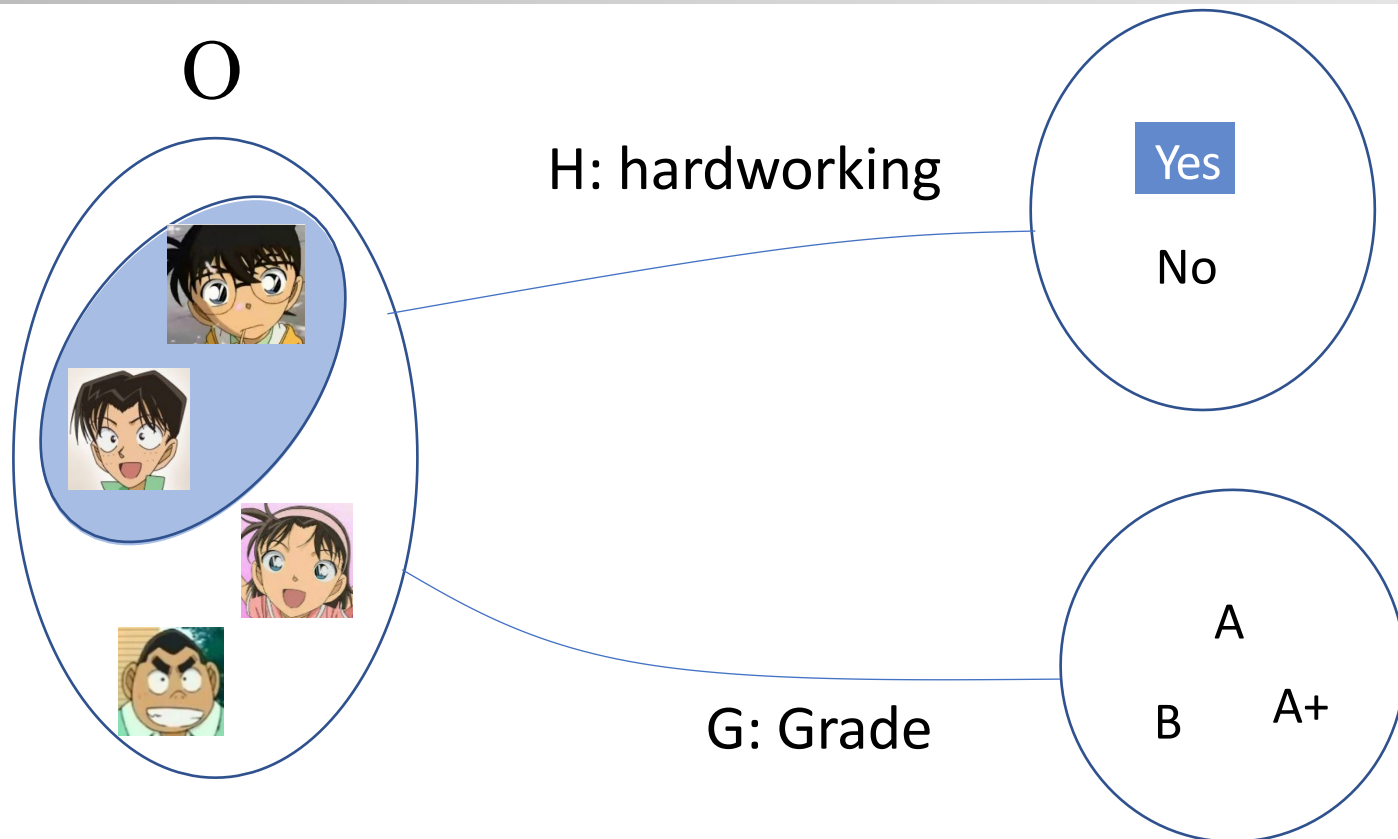
- The big picture
- Events and Event spaces
-  • Random variables
- Joint probability, Marginalization, conditioning, chain rule, Bayes Rule, law of total probability, etc.
- Structural properties, e.g., Independence, conditional independence
- Maximum Likelihood Estimation



From Events to Random Variable (RV)

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

Random Variables (RV)



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

- “functions” that maps from O to an attribute space T .

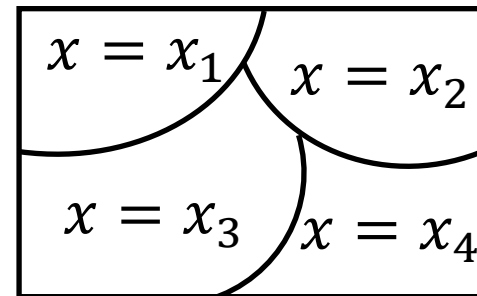
Discrete Random Variables

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

$$\sum_{i=1}^4 P(X = x_i) = 1$$



e.g. Coin Flips

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

$$\text{Binary} = \{H, T\}$$



$$\text{Integer} = \{1, 2, \dots, k\}$$

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

e.g., two common distributions

- Uniform

$$X \sim U[1, \dots, N]$$

- X takes values 1, 2, ..., N
- E.g. picking balls of different colors from a box

$$P(X = i) = \frac{1}{N}$$


- Binomial

- X takes values 0, 1, ..., k
- E.g. coin flips k times

$$X \sim B(k, p)$$

$$P(X = i) = \binom{k}{i} p^i (1 - p)^{k-i}$$

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, e.g., Independence, conditional independence
- Maximum Likelihood Estimation

If hard to directly estimate from data, most likely we can estimate

- Joint probability

- Use Chain Rule

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

- Marginal probability

- Use the total law of probability

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

- Conditional probability

- Use the Bayes Rule

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

(1) To calculate Joint Probability: Use Chain Rule

- Two ways to use chain rules on joint probability

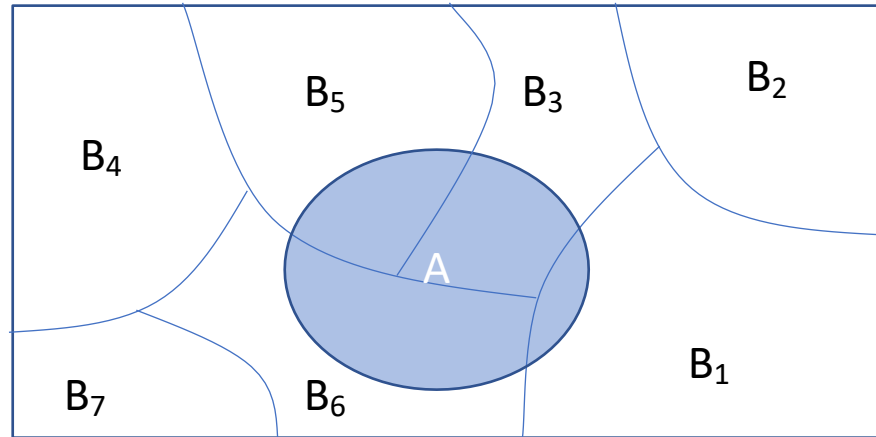
joint conditional marginal

↖ ↗ ↘

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

(2) To calculate **Marginal Probability**:

Use Rule of total probability (e.g. event version)



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

(2) To calculate **Marginal Probability**:

Use Rule of total probability (e.g. RV version)

- Given two discrete RVs X and Y , which take values in:

$$\{x_1, \dots, x_k\} \quad \{y_1, \dots, y_m\}$$

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



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



(3) To calculate **Conditional Probability**:
Use Bayes Rule (e.g. RV version)

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

One Example: Joint

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?

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



One Example: Marginal

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

$$P(B_2 = r) =$$



One Example: Conditional

*What is the probability that the 2nd ball will be red if the 1st ball drawn from the set{*r*, *r*, *r*, *b*} is red?*

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

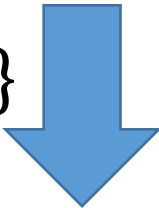


Use both Bayes Rule and Marginal

- X and Y are discrete RVs...

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

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

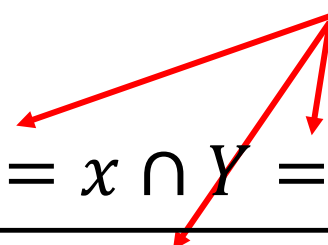


$$P(X = x_i | Y = y_j) = \frac{P(Y = y_j | X = x_i)P(X = x_i)}{\sum_k P(Y = y_j | X = x_k)P(X = x_k)}$$

Simplify Notation: Conditional Probability

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

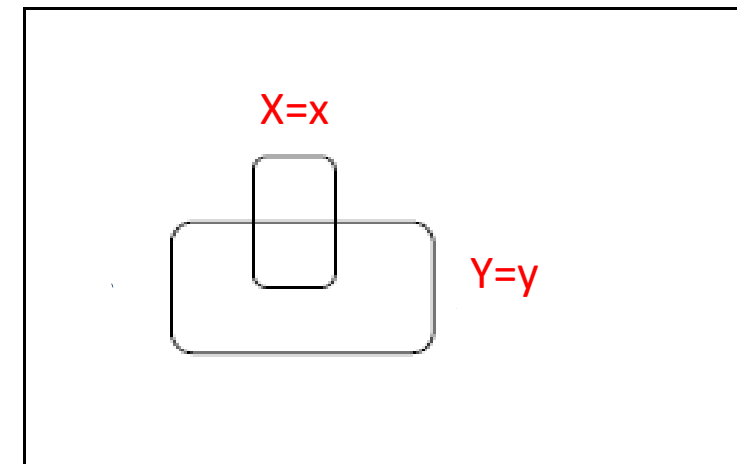
events



- But we will always write it this way:

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

- $P(X=x \text{ true}) \rightarrow P(X=x) \rightarrow P(x)$



Simplify Notation:

An Example of estimating conditional

- 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?

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



Simplify Notation:

An Example of estimating conditional

- 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?

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

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



Simplify Notation: Conditional

- Bayes Rule

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

- You can condition on **more variables**

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

Simplify Notation: Marginal

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

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

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

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

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

Simplify Notation: An Example

- 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?

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

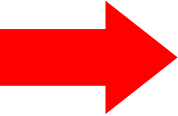
↓
↓
↓

weather
grass
 $P(\text{rain})P(\text{rain}|\text{wet})$

{rain, sunny}
{wet, dry}
+ $P(\text{sunny})P(\text{sunny}|\text{wet})$

0.5
 1

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, e.g., Independence, conditional independence
- Maximum Likelihood Estimation (next)

Independent RVs

- Definition: X and Y are independent *iff*

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

More on Independence

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



$$P(X = x | Y = y) = P(X = x)$$



$$P(Y = y | X = x) = P(Y = y)$$

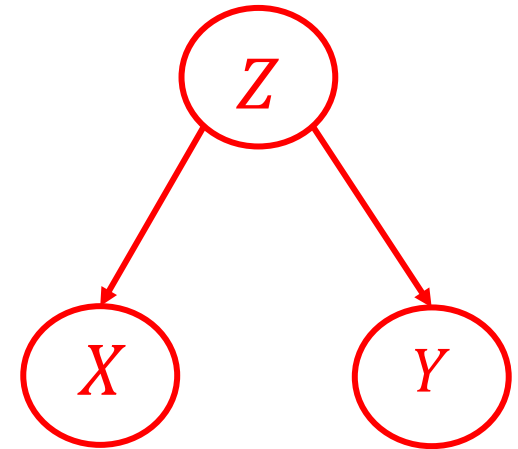
- **E.g.** no matter how many heads you get, your friend will not be affected, and vice versa

More on Independence

- X is independent of Y means that knowing Y does not change our belief about X .
- The following forms are **equivalent**:
 - $P(X=x, Y=y) = P(X=x) P(Y=y)$
 - $P(X=x | Y=y) = P(X=x)$
- **The above should hold for all x_i, y_j**
- It is symmetric and **written as** $X \perp Y$

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*



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

If holding for all x_i, y_j, z_k $X \perp Y | Z$

More on Conditional Independence

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



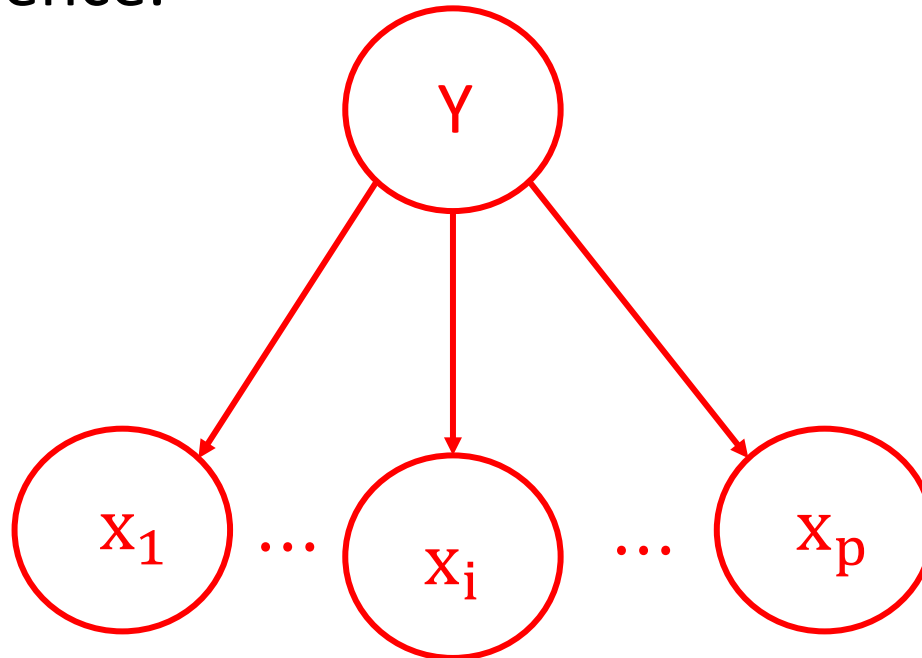
$$P(X = x | Y = y, Z = z) = P(X = x | Z = z)$$



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

Independence and Conditional Independence

- Independence does not imply conditional independence.
- Conditional independence does not imply independence.



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, e.g., Independence, conditional independence
- **Maximum Likelihood Estimation (next)**

References

- <https://qiyanjun.github.io/2019f-UVA-CS6316-MachineLearning/>
- Prof. Andrew Moore's review tutorial
- Prof. Nando de Freitas's review slides
- Prof. Carlos Guestrin recitation slides



Thanks for listening