

# Review on Math for AI

**Daifeng Wang**

dai feng . wang @ wisc . edu

**University of Wisconsin, Madison**

Based on slides from Xiaojin Zhu  
(<http://pages.cs.wisc.edu/~jerryzhu/cs540.html>),  
modified by Daifeng Wang



# Outline

- Probability and inference
  - Axioms of probability
  - Joint, Marginal, Conditional probability
  - Bayes rule
  - Independence, Conditional independence
  - Expected value
  - Maximum Likelihood Estimation (MLE)
  - Maximum a posteriori (MAP) estimation



## Sample space

- A space of events that we assign probabilities to
- Events can be binary, multi-valued, or continuous
- Events are mutually exclusive

## Random variable

- A variable,  $x$ , whose domain is the sample space, and whose value is somewhat uncertain

## Probability for discrete events

- Probability  $P(x=a)$  or  $P(a)$  is the fraction of times  $x$  takes value  $a$



# Probability table

- Weather

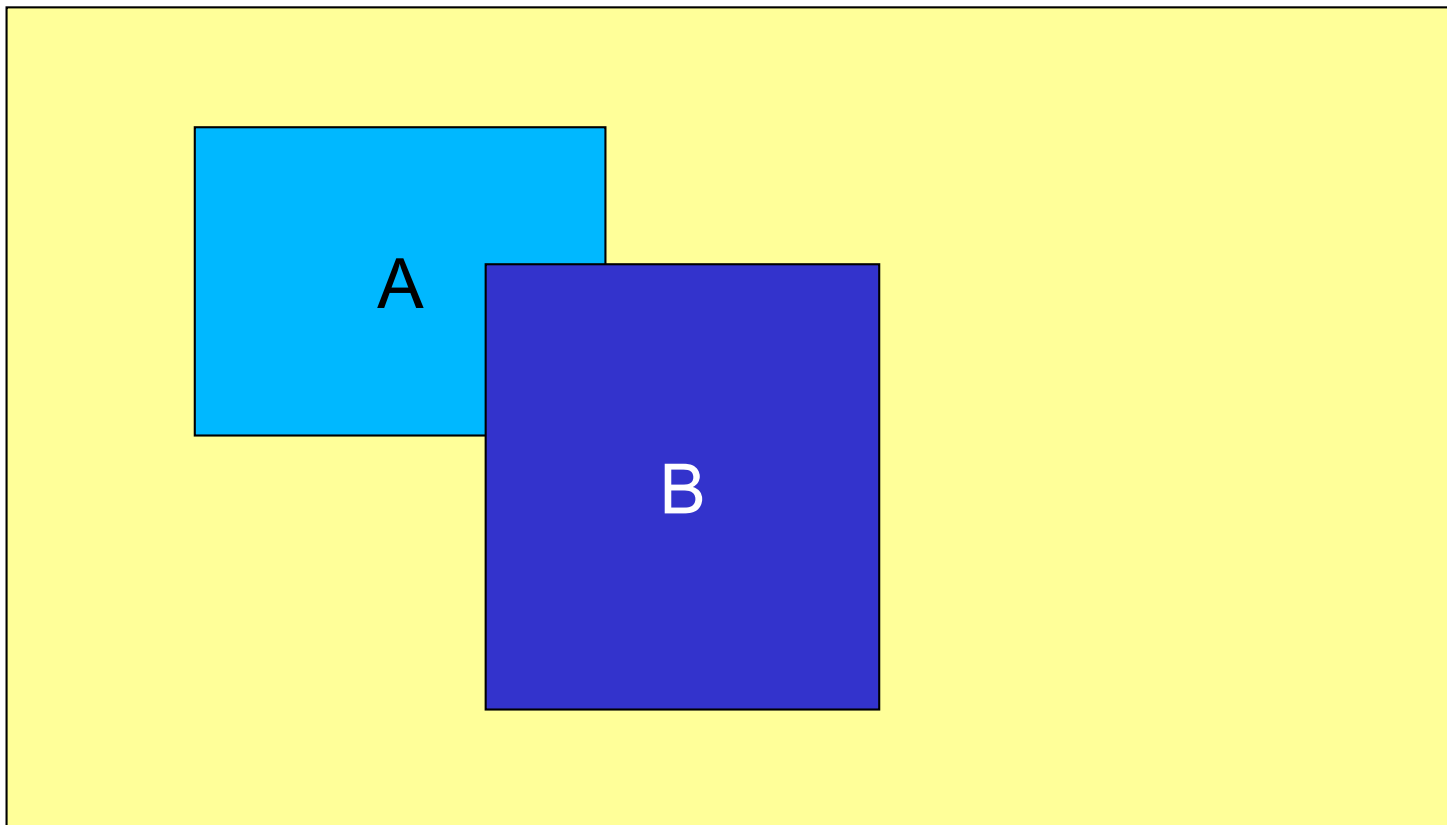
Sunny	Cloudy	Rainy
200/365	100/365	65/365

- $P(\text{Weather} = \text{sunny}) = P(\text{sunny}) = 200/365$
- $P(\text{Weather}) = \{200/365, 100/365, 65/365\}$
- For now we'll be satisfied with obtaining the probabilities by counting frequency from data...



# The axioms of probability

- $P(A) \in [0, 1]$
- $P(\text{true})=1, P(\text{false})=0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



Sample  
space



## Some theorems derived from the axioms

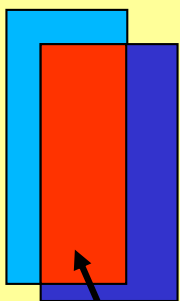
- $P(\neg A) = 1 - P(A)$  picture?
- If  $A$  can take  $k$  different values  $a_1 \dots a_k$ :  
$$P(A=a_1) + \dots P(A=a_k) = 1$$
- $P(B) = P(B \wedge \neg A) + P(B \wedge A)$ , if  $A$  is a binary event
- $P(B) = \sum_{i=1 \dots k} P(B \wedge A=a_i)$ , if  $A$  can take  $k$  values



# Joint probability

- The **joint** probability  $P(A=a, B=b)$  is a shorthand for  $P(A=a \wedge B=b)$ , the probability of both  $A=a$  and  $B=b$  happen

$P(A=a)$ , e.g.  $P(1^{\text{st}} \text{ word on a random page} = \text{"San"}) = 0.001$   
(possibly: San Francisco, San Diego, ...)



$P(B=b)$ , e.g.  $P(2^{\text{nd}} \text{ word} = \text{"Francisco"}) = 0.0008$   
(possibly: San Francisco, Don Francisco, Pablo Francisco ...)

$P(A=a, B=b)$ , e.g.  $P(1^{\text{st}} = \text{"San"}, 2^{\text{nd}} = \text{"Francisco"}) = 0.0007$



# Marginal probability

- Sum over other variables

weather

	Sunny	Cloudy	Rainy
temp			
hot	150/365	40/365	5/365
cold	50/365	60/365	60/365
$\Sigma$	200/365	100/365	65/365

$$P(\text{Weather}) = \{200/365, 100/365, 65/365\}$$

- The name comes from the old days when the sums are written on the margin of a page

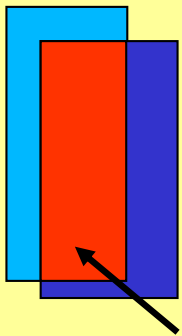




# Conditional probability

- The **conditional** probability  $P(A=a \mid B=b)$  is the fraction of times  $A=a$ , **within the region that**  $B=b$

$P(A=a)$ , e.g.  $P(1^{\text{st}} \text{ word on a random page} = \text{"San"}) = 0.001$



$P(B=b)$ , e.g.  $P(2^{\text{nd}} \text{ word} = \text{"Francisco"}) = 0.0008$

$P(A=a \mid B=b)$ , e.g.  $P(1^{\text{st}} = \text{"San"} \mid 2^{\text{nd}} = \text{"Francisco"}) = \mathbf{0.875}$   
(possibly: San, Don, Pablo ...)

Although "San" is rare and "Francisco" is rare,  
given "Francisco" then "San" is quite likely!



# The chain rule

- From the definition of conditional probability we have the chain rule

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

- It works the other way around

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

- It works with more than 2 events too

$$P(A_1, A_2, \dots, A_n) =$$

$$P(A_1) * P(A_2 | A_1) * P(A_3 | A_1, A_2) * \dots * P(A_n | A_1, A_2, \dots, A_{n-1})$$



# Inference with Bayes' rule

- $P(A|B) = P(B|A)P(A) / P(B)$  Bayes' rule
- Why do we make things this complicated?
  - Often  $P(B|A)$ ,  $P(A)$ ,  $P(B)$  are easier to get
  - Some names:
    - **Prior  $P(A)$** : probability before any evidence
    - **Likelihood  $P(B|A)$** : assuming  $A$ , how likely is the evidence
    - **Posterior  $P(A|B)$** : conditional prob. after knowing evidence
    - **Inference**: deriving unknown probability from known ones
- In general, if we have the full joint probability table, we can simply do  $P(A|B) = P(A, B) / P(B)$  – more on this later...



# Independence

- Two events A, B are **independent**, if (the following are equivalent)
  - $P(A, B) = P(A) * P(B)$
  - $P(A | B) = P(A)$
  - $P(B | A) = P(B)$

## Conditional independence

- In general A, B are conditionally independent given C
  - if  $P(A | B, C) = P(A | C)$ , or
  - $P(B | A, C) = P(B | C)$ , or
  - $P(A, B | C) = P(A | C) * P(B | C)$

# Expected values

- The *expected value* of a random variable that takes on numerical values is defined as:

$$\mathbf{E}[X] = \sum_x xP(x)$$

This is the same thing as the *mean*

- We can also talk about the expected value of a function of a random variable

$$\mathbf{E}[g(X)] = \sum_x g(x)P(x)$$

Suppose we have data  $\mathcal{D} = \{x^{(i)}\}_{i=1}^N$

## Maximum Likelihood Estimation (MLE)

Find optimal  $\theta^*$  to maximize the likelihood given the data

$$\theta_{MLE}^* = \operatorname{argmax}_{\theta} P(\mathcal{D}|\theta)$$

## Maximum a posteriori (MAP) estimation


$$\begin{aligned}\theta_{MAP}^* &= \operatorname{argmax}_{\theta} P(\theta|\mathcal{D}) && \leftarrow \text{Posterior} \\ &= \operatorname{argmax}_{\theta} \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})} && \left. \begin{array}{l} \text{Bayes rule} \\ \text{Posterior} \end{array} \right\} \\ &= \operatorname{argmax}_{\theta} P(\mathcal{D}|\theta)P(\theta)\end{aligned}$$

# Play outside or not?

- If weather is sunny, would you like to play outside?  
Posterior probability  $P(\text{Yes}|\text{Sunny})$  vs  $P(\text{No}|\text{Sunny})$

- Weather = {Sunny, Rainy, Overcast}
- Play = {Yes, No}
- Observed data {Weather, Play on Day  $m$ },  
 $m=1,2,\dots,N$

How can we calculate posterior probabilities?

$$P(\text{Play}|\text{Weather}) = \frac{P(\text{Weather}|\text{Play})P(\text{Play})}{P(\text{Weather})}$$


Bayes rule

# Play outside or not?

- **Step 1:** Convert the data to a frequency table of Weather and Play

- **Step 2:** Based on the frequency table, calculate likelihoods  $P(\text{Weather}|\text{Play})$  and priors  $P(\text{Play})$

Weather	Play
Sunny	No
Overcast	Yes
Rainy	Yes
Sunny	Yes
Sunny	Yes
Overcast	Yes
Rainy	No
Rainy	No
Sunny	Yes
Rainy	Yes
Sunny	No
Overcast	Yes
Overcast	Yes
Rainy	No



Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9



Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

$$e. g., P(\text{Play} = \text{Yes}) = \frac{9}{14} = 0.64$$

$$P(\text{Sunny}|\text{Yes}) = \frac{3}{9} = 0.33$$



# Play outside or not?

Frequency Table		
Weather	No	Yes
Overcast		4
Rainy	3	2
Sunny	2	3
Grand Total	5	9

Likelihood table				
Weather	No	Yes		
Overcast		4	=4/14	0.29
Rainy	3	2	=5/14	0.36
Sunny	2	3	=5/14	0.36
All	5	9		
	=5/14	=9/14		
	0.36	0.64		

- Step 3:** Based on the likelihoods and priors, calculate posteriors  $P(\text{Play} | \text{Weather})$
- $P(\text{No} | \text{Sunny})$   
 $= P(\text{Sunny} | \text{No}) * P(\text{No}) / P(\text{Sunny})$   
 $= 0.4 * 0.36 / 0.36$   
 $= 0.4$
- $P(\text{Yes} | \text{Sunny})$   
 $= P(\text{Sunny} | \text{Yes}) * P(\text{Yes}) / P(\text{Sunny})$   
 $= 0.33 * 0.64 / 0.36$   
 $= 0.6$
- $P(\text{Yes} | \text{Sunny}) > P(\text{No} | \text{Sunny})$ , you should go outside and play!