### **CS 4100: Introduction to AI**

Wayne Snyder Northeastern University

Lecture 11: Review of Probability and Likelihood for Machine Learning



### **Random Experiments**

A Random Experiment is a process that produces uncertain outcomes from a well-defined set of possible outcomes. Usually there is some kind of physical experiment which is being modeled theoretically.



### Probability

We measure the probability of events on a real-number scale from 0 to 1:



### Finite Equiprobable Probability Spaces

For finite and equiprobable probability spaces,

it is easy to calculate the probability:

$$P(A) = \frac{|A|}{|S|}$$

Example: Roll a die, how many dots showing on the top face? Let A = "less than 5 dots are showing."

S = { 1, 2, ...., 6 }  
P = { 
$$1/6, 1/6, ...., 1/6$$
 }

P(A) = 4/6 = 2/3 = 0.6666...



#### Conditional Probability: P(A | B) = P(A) if we know B has happened

**Example:** Roll two dice. A = "the total # dots showing is > 8" and B = "the first roll was 3"

What is P(A | B)?

The key to solving such problems is to realize that there are two probability spaces:

- the one before you know whether B has happened, and 0
- the one that has been "conditioned" by knowing that B has definitely happened, so the 0 sample space has shrunk and the proportion representing event A may have changed:

Original

Conditioned by knowing B happened:

S' = B



### **Conditional Probability**

Conditioning the original sample space means changing the perspective: instead of finding the area of A inside S, we are finding the area of  $A \cap B$  inside B:

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$



# Independence and Dependence

We say that two events A and B are independent if

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A \mid B) = P(A)$$

$$P(B \mid A) = \frac{P(A \cap B)}{P(A)}$$

$$P(A \cap B) = P(A) * P(B)$$

Example:

or:

Suppose in a particular city, 40% of the population is male, and 60% female, and 20% of the population smokes. If male smokers are 8% of the population, then are smoking and gender independent? That is, are the following two events independent?

YES. Check:

A = Smoker

$$P(A \cap B) = 0.08 = 0.4 * 0.2 = P(A) * P(B)$$

B = Male

### **Bayes' Rule**

We can rearrange the conditional probability rule in a way that makes the sequence of the events irrelevant -- which happened first, A or B? Or did they happen at the same time? Does it matter?

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)} \qquad P(B \mid A) = \frac{P(B \cap A)}{P(A)}$$

We can do a little algebra to define conditional probabilities in terms of each other:

$$P(B \mid A) * P(A) = P(B \cap A) = P(A \mid B) * P(B)$$

so:

$$P(B \mid A) = \frac{P(A \mid B) * P(B)}{P(A)}$$



### Bayes' Rule

The best way to understand this is to view it with a tree diagram! P(B | A) = the probability that when A happens, it was "preceded" by B:



If A has happened, what is the probability that it did so on the path where B also occurred?

Note:

$$A = P(A \cap B) \cup P(A \cap B^{c})$$

So what percentage of A is due to  $A \cap B$ ?

Same calculation as:

$$P(B \mid A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A \cap B)}{P(A)}$$

# **Bayes' Rule**

This has an interesting flavor, because we can ask about causes of outcomes:

A Priori Reasoning -- "I randomly choose a person and observe that he is male; what the probability that it is a smoker?"

"The first toss of a pair of dice is a 5; what is the probability that the total is greater than 8?"

A Posteriori Reasoning -- "I find a cigarette butt on the ground, what is the probability that it was left by a man?"

"The total of a pair of thrown dice is greater than 8; what is the probability that the first toss was a 5?"

> This seems odd, because instead of reasoning forward from "causes to effects" we are reasoning backwards from "effects to causes" but really it is just different ways of phrasing the mathematical formulae. Time is not really relevant!

### **Discrete RandomVariables: Probability Mass Function**

The probability function of a **discrete** random variable X is a function

$$f_X$$
 = Probability Mass Function (PMF)

which assigns a probability to each real number in the range of X and follows the normal rules for a probability space:

$$f_{X} : R_{x} \to \mathcal{R}$$
$$\forall a \in R_{x} \quad f_{X}(a) \ge 0$$
$$\sum_{a \in R_{x}} f_{X}(a) = 1.0$$

### Discrete RandomVariables: Probability Distributions

We will emphasize the distributions of random variables, using graphical representations to help our intuitions.

#### Example:

Y = "The number of tosses of a fair coin until a head appears"



### **Discrete vs Continuous Distributions**

#### **Discrete Random Variables**

The Probability Mass Function (PMF) of a discrete random variable X is a function from the range of X into  $\mathcal{R}$ :

 $P_x : R_X \mapsto \mathcal{R}$ 

such that

(i) 
$$\forall y \in R_X \ P_X(y) \ge 0.0$$

(ii) 
$$\sum_{y \in R_X} P_X(y) = 1.0$$



#### **Continuous Random Variables**

The Probability Density Function (PDF) of a continuous random variable X is a function from  $\mathcal{R}$  to  $\mathcal{R}$ :

$$f_x : \mathcal{R} \mapsto \mathcal{R}$$

such that

i) 
$$\forall y \ f_X(y) \ge 0.0$$





I am going to show examples from the Brown Corpus, an early (1961) and still-used corpus of various English texts:

Category	Genre (Code)	# of	Total	%
		texts	Tokens	
INFORMATIVE	Learned (J)	80	160,000	16.0%
INFORMATIVE	Belles Lettres,	75	150,000	15.0%
	Biography, Memoirs, etc (G)			
INFORMATIVE	Popular Lore (F)	48	96,000	9.6%
INFORMATIVE	Press: Reportage (A)	44	88,000	8.8%
INFORMATIVE	Skills and Hobbies (E)	36	72,000	7.2%
INFORMATIVE	Miscellaneous (H)	30	60,000	6.0%
IMAGINATIVE	General Fiction (K)	29	58,000	5.8%
IMAGINATIVE	Adventure and	29	58,000	5.8%
	Western Fiction (N)			
IMAGINATIVE	Romance and Love	29	58,000	5.8%
	Story (P)			
INFORMATIVE	Press: Editorial (B)	27	54,000	5.4%
IMAGINATIVE	Mystery and Detective	24	48,000	4.8%
	Fiction (L)			
INFORMATIVE	Press: Reviews	17	34,000	3.4%
	(theatre, books, music,			
	dance) (C)			
INFORMATIVE	Religion (D)	17	34,000	3.4%
IMAGINATIVE	Humor (R)	9	18,000	1.8%
IMAGINATIVE	Science Fiction (M)	6	12,000	1.2%
	TOTAL	500	1,000,000	100.0%

Brown Corpus Sample (untagged)		
A01 0010 The Fulton County Grand Jury said		
Friday an investigation		
A01 0020 of Atlanta's recent primary election		
produced "no evidence" that		
A01 0030 any irregularities took place. The jury		
further said in term-end		
A01 0040 presentments that the City Executive		
Committee, which had over-all		
A01 0050 charge of the election, "deserves the		
praise and thanks of the		
A01 0060 City of Atlanta" for the manner in		
which the election was conducted.		
		Ł

#### Brown Corpus Sample (tagged)

A01\_FO 0010\_MC The\_AT Fulton\_NP1 County\_NN1 Grand\_JJ Jury\_NN1 said\_VVD Friday\_NPD1 an\_AT1 investigation NN1 A01 FO 0020 MC of IO Atlanta 93's 03 recent JJ primary\_JJ election\_NN1 produced\_VVD "\_" no\_AT evidence\_NN1 "\_" that\_CST A01\_FO 0030\_MC any\_DD irregularities\_NN2 took\_VVD place\_NN1 .\_. The\_AT jury\_NN1 further\_RRR said\_VVD in\_II termend NN1 A01\_FO 0040\_MC presentments\_NN2 that\_CST the\_AT City\_NN1 Executive\_NN1 Committee\_NN1 ,\_, which\_DDQ had\_VHD over-all\_RR A01\_FO 0050\_MC charge\_NN1 of\_IO the\_AT election\_NN1 ,\_, "\_" deserves\_VVZ the\_AT praise\_NN1 and\_CC thanks\_NN2 of\_IO the\_AT A01 FO 0060 MC City NN1 of IO Atlanta NP1" " for\_IF the\_AT manner\_NN1 in\_II which\_DDQ the\_AT election\_NN1 was\_VBDZ conducted\_VVN .\_.



The composition of the Brown Corpus



First of all, be sure you understand the difference between a histogram and a probability distribution!







There are 1,161,192 occurrences of words in the Brown Corpus.

There are 47,451 unique words; 20,478 of these words are *hapax legomena*, i.e., they occur only once in the corpus.



#### The first 100 words are:

['the', 'of', 'and', 'to', 'a', 'in', 'that', 'is', 'was', 'he', 'for', 'it', 'with', 'as', 'his', 'on', 'be', 'at', 'by', 'i', 'this', 'had', 'not', 'are', 'but', 'from', 'or', 'have', 'an', 'they', 'which', 'one', 'you', 'were', 'he r', 'all', 'she', 'there', 'would', 'their', 'we', 'him', 'been', 'has', 'when', 'who', 'will', 'more', 'if', 'no']





['the', 'of', 'and', 'to', 'a', 'in', 'that', 'is', 'was', 'he', 'for', 'it', 'with', 'as', 'his', 'on', 'be', 'at', 'by', 'i', 'this', 'had', 'not', 'are', 'but', 'from', 'or', 'have', 'an', 'they', 'which', 'one', 'you', 'were', 'he r', 'all', 'she', 'there', 'would', 'their', 'we', 'him', 'been', 'has', 'when', 'who', 'will', 'more', 'if', 'no']

### Bayesian Reasoning and Maximum Likelihood Estimates (MLE)

A typical machine learning workflow is the following:



Typically:

Pick initial values for the parameters  $\Theta$ , based on some initial knowledge (e.g., from a textbook, a similar experiment, a "hunch," etc.),

Repeat

Perform experiment and find new  $\mu$ 

Based on knowledge gained from last run revise  $\Theta$ 

Until .....

But notice this is a bit... weird! We are used to thinking of programs as starting with inputs and producing outputs, but in ML the emphasis is on using the outputs to determine the inputs!





Bayesian Reasoning is a way of thinking about probability which emphasizes this point of view, and has become standard in NLP and ML.

Formally, we have Bayes Rule and some terminology:





But notice this is a bit... weird! We are used to thinking of programs as starting with inputs and producing outputs, but in ML the emphasis is on using the outputs to determine the inputs!



This is generally called the Bayesian interpretation of the experiment: Given some experimental result, how likely is the parameter  $\theta$  to be some particular value? Note: Probability refers to outcomes, and likelihood refers to parameters.

Let's consider a simple random experiment:

We flip a (possibly) unfair coin n times and record the outcomes.

Frequentist: The data will converge in the limit to a specific probability  $\theta$  which existed before we started to flip:

```
[ ] from numpy import mean
theta = 0.65
for n in range(7):
    N = 10**n
    print("Over",N," iterations, the mean is", mean([X(theta) for k in range(N)]))
Over 1 iterations, the mean is 1.0
Over 10 iterations, the mean is 0.7
Over 100 iterations, the mean is 0.67
Over 1000 iterations, the mean is 0.648
Over 10000 iterations, the mean is 0.657
Over 100000 iterations, the mean is 0.65094
Over 1000000 iterations, the mean is 0.650158
```

Bayesian: Given some experimental result, what is the MOST LIKELY value for the parameter  $\theta$ ?

In our example we would naturally assume that  $\theta$  is simply the mean of the outcomes observed so far, for example, if we get outcomes

```
[] theta = 0.25
n = 100
outcomes = [X(theta) for k in range(n) ]
print(outcomes)
[0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0]
```

we would assume that  $\theta$  is the mean of the *n* outcomes, i.e.,

 $\hat{\theta} = \frac{count(outcomes = 1)}{count(all outcomes)}$ 

```
[ ] theta_hat = mean(outcomes)
```

```
print("theta_hat =", theta_hat)
```

```
theta_hat = 0.24
```

This is the most likely estimate for  $\theta$ , since we have no other information.

The more outcomes we see, the more accurate an estimate of  $\theta$  we would get. However, notice that in general, we will not get  $\theta$  exactly right, but only an *estimate* based on our experiment. This is appropriate, because we did not know  $\theta$  to start with, and we want the *most likely* value.

Now we can introduce the general notion of the Maximum Likelihood Estimate (MLE), which applies to our simple experiment and to many more complex situations which occur often in NLP and ML.

If we express the situation mathematically, keeping the parameter  $\theta$  in the foreground, let us consider the *n* outcomes of the previous experiment as a sequence of *n* independent and identically-distributed random variables, X<sub>1</sub>, X<sub>2</sub>, ..., X<sub>n</sub>, where

 $X_i(\theta) \approx Bernoulli(\theta).$ 

and the PDF (which assigns a probability to each outcome x) is

$$f(x,\theta) = \begin{cases} \theta & \text{if } x = 1\\ 1 - \theta & \text{if } x = 0 \end{cases}$$

In the frequentist interpretation, we would ask: Given the parameter  $\theta$ , how likely are the outcomes  $X_1, \ldots, X_n$ , and get the answer

$$f(X_1,\theta) * f(X_2,\theta) * \cdots * f(X_n,\theta) = \prod_{i=1}^n f(X_i,\theta).$$

```
[ ] import numpy as np
def f(x,theta):
    if(x==1):
        return theta
    else:
        return 1-theta
outcome_probs = [f(x,theta) for x in outcomes]
print("Outcomes", outcomes)
print("Outcome probabilities:",outcome_probs)
print("Outcome probabilities:",outcome_probs)
print("With theta =",theta,"the probability of this particular list of outcomes is ",
        np.prod( outcome_probs))
```

In the Bayesian interpretation, where we don't assume we know  $\theta$ , we would ask a more useful question: Given the set of outcomes, what is the *most likely* value for  $\theta$ ? Alternately, what estimate of  $\theta$  is most likely given our outcomes? This is formalized as the *likelihood* function for *n* outcomes with parameter  $\theta$ :

$$\mathcal{L}_n(\theta) = \prod_{i=1}^n f(X_i, \theta).$$

The maximum likelihood estimator (MLE) is then

$$\hat{\theta} = \arg \max_{\theta} \mathcal{L}_n(\theta) = \text{the value of } \theta \text{ that maximizes } \mathcal{L}_n(\theta)$$

If we try many values of the parameter, we would see that it reaches a maximum at the value found by simple counting....



The MLE formalizes the idea that we are searching for the best possible (= most likely) set of parameters; we shall see (e.g., in Logistic Regression) that when there is an explicit formula for  $\mathcal{L}_n$ , we can calculate the MLE directly using the derivative of  $\mathcal{L}_n$ .