Machine Learning CS 4641-7641



# **Probability and Statistics**

Mahdi Roozbahani Georgia Tech

These slides are inspired based on slides from Le Song, Sam Roweis, and Chao Zhang.

## Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

## Probability

- A sample space S is the set of all possible outcomes of a conceptual or physical, repeatable experiment. (S can be finite or infinite.)
  - E.g., S may be the set of all possible outcomes of a dice roll: S
     (1 2 3 4 5 6)
  - E.g., S may be the set of all possible nucleotides of a DNA site: S
     (A C G T)
  - E.g., S may be the set of all possible time-space positions of an aircraft on a radar screen.
- An Event A is any subset of S
  - Seeing "1" or "6" in a dice roll; observing a "G" at a site; UA007 in space-time interval

### Three Key Ingredients in Probability Theory

A **sample space** is a collection of all possible **outcomes** 

Random variables X represents **outcomes** in sample space

Probability of a random variable to happen

$$p(x) = p(X = x)$$

 $p(x) \ge 0$ 

### **Continuous variable**

Continuous probability distribution Probability density function Density or likelihood value Temperature (real number) Gaussian Distribution

$$\int_{x} p(x) dx = 1$$

#### **Discrete variable**

Discrete probability distribution Probability mass function Probability value Coin flip (integer) Bernoulli distribution

$$\sum_{x \in A} p(x) = 1$$

### **Continuous Probability Functions**

- Examples:
  - Uniform Density Function:

$$f_{x}(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b\\ 0 & \text{otherwise} \end{cases}$$

Exponential Density Function:

$$f_x(x) = \frac{1}{\mu} e^{-\frac{x}{\mu}} \qquad for \ x \ge 0$$
$$F_x(x) = 1 - e^{\frac{-x}{\mu}} \qquad for \ x \ge 0$$

Gaussian(Normal) Density Function

$$f_x(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

## **Discrete Probability Functions**

- Examples:
  - Bernoulli Distribution:

$$\begin{array}{ll}
\left\{ \begin{array}{ll}
1-p & for \ x=0\\ p & for \ x=1 \end{array} \right.
\end{array}$$

In Bernoulli, just a **single** trial is conducted

• Binomial Distribution: •  $P(X = k) = {n \choose k} p^k (1 - p)^{n-k}$ 

k is number of successes

**n-k** is number of failures

 $\binom{n}{k}$  The total number of ways of selection **k** distinct combinations of **n** trials, **irrespective of order**.

## Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

### Example



X = Throw a dice



Y = Flip a coin

X and Y are random variables

- $\mathbf{N}$  = total number of trials
- $n_{ij}$  = Number of occurrence

		Χ						C
		$x_{i=1} = 1$	$x_{i=2} = 2$	$x_{i=3} = 3$	$x_{i=4} = 4$	$x_{i=5} = 5$	$x_{i=6} = 6$	$\mathcal{L}_{j}$
Y	$y_{j=2} = tail$	$n_{ij} = 3$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 5$	$n_{ij} = 1$	$n_{ij} = 5$	20
	$y_{i=1} = head$	$n_{ij} = 2$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 2$	$n_{ij} = 4$	$n_{ij} = 1$	15
	$C_i$	5	6	6	7	5	6	N=35

$$\mathbf{Y} \begin{array}{c} y_{j=2} = tail \\ y_{j=1} = head \\ C_i \end{array} \begin{bmatrix} n_{ij} = 3 & n_{ij} = 4 & n_{ij} = 2 & n_{ij} = 5 & n_{ij} = 1 & n_{ij} = 5 & 20 \\ n_{ij} = 2 & n_{ij} = 2 & n_{ij} = 4 & n_{ij} = 2 & n_{ij} = 4 & n_{ij} = 1 & 15 \\ \hline n_{ij} = 2 & n_{ij} = 2 & n_{ij} = 4 & n_{ij} = 2 & n_{ij} = 4 & n_{ij} = 1 & 15 \\ \hline 5 & 6 & 6 & 7 & 5 & 6 & N=35 \\ \end{array}$$

### **Probability:**

Joint probability:

$$p(X = x_i) = \frac{c_i}{N}$$
$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

**Conditional probability:** 

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

Sum rule  

$$p(X = x_i) = \sum_{j=1}^{L} p(X = x_i, Y = y_j) \Rightarrow p(X) = \sum_{Y} P(X, Y)$$

### **Product rule**

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}c_i}{c_iN} = p(Y = y_j|X = x_i)p(X = x_i)$$
$$p(X, Y) = p(Y|X)p(X)$$

## **Conditional Independence**

#### • Examples:

P(Virus| DrinkBeer) = P(Virus)
iff Virus is independent of Drink Beer

P(Flu | Virus, DrinkBeer) = P(Flu |Virus)
iff Flu is independent of Drink Beer, given Virus

```
P(Headache | Flu, Virus, DrinkBeer)
= P(Headache|Flu, DrinkBeer)
```

iff Headache is independent of Virus, given Flu and Drink Beer

Assume the above independence, we obtain:

P(Headache, Flue, Virus, DrinkBeer) = P(Headache|Flu, Virus, DrinkBeer) P(Flu|Virus, DrinkBeer) P(Virus|DrinkBeer)P(DrinkBeer) = P(Headache|Flu, DrinkBeer)P(Flu|Virus)P(Virus)P(DrinkBeer)

- -

## Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

## Bayes' Rule

P(X|Y)= Fraction of the worlds in which X is true given that Y is also true.

- For example:
  - H="Having a headache"
  - F="Coming down with flu"
  - P(Headche|Flu) = fraction of flu-inflicted worlds in which you have a headache. How to calculate?
- Definition:

$$P(X|Y) = \frac{P(X,Y)}{P(Y)} = \frac{P(Y|X)P(X)}{P(Y)}$$
$$P(X,Y) = P(Y|X)P(X)$$

Corollary:

## Bayes' Rule

• 
$$P(Headache|Flu) = \frac{P(Headache,Flu)}{P(Flu)}$$
  
=  $\frac{P(Flu|Headache)P(Headache)}{P(Flu)}$ 

Other cases:

• 
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X|Y)P(Y) + P(X|\neg Y)P(\neg Y)}$$
  
•  $P(Y = y_i|X) = \frac{P(X|Y)P(Y)}{\sum_{i \in S} P(X|Y = y_i)P(Y = y_i)}$   
•  $P(Y|X,Z) = \frac{P(X|Y,Z)P(Y,Z)}{P(X,Z)} = \frac{P(X|Y,Z)P(Y,Z)}{P(X|Y,Z)P(Y,Z) + P(X|\neg Y,Z)P(\neg Y,Z)}$ 

## Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

## Mean and Variance

Expectation: The mean value, center of mass, first moment:

$$E_X[g(X)] = \int_{-\infty}^{\infty} g(x) p_X(x) dx = \mu$$

• N-th moment: 
$$g(x) = x^n$$

- N-th central moment:  $g(x) = (x \mu)^n$
- Mean:  $E_X[X] = \int_{-\infty}^{\infty} x p_X(x) dx$ 
  - $E[\alpha X] = \alpha E[X]$
  - $E[\alpha + X] = \alpha + E[X]$
- Variance(Second central moment):  $Var(x) = E_X[(X E_X[X])^2] = E_X[X^2] E_X[X]^2$ 
  - $Var(\alpha X) = \alpha^2 Var(X)$
  - $Var(\alpha + X) = Var(X)$

### Mean and average

### Variance and average:

### Covariance:

### Correlation:

## For Joint Distributions

#### Expectation and Covariance:

- E[X + Y] = E[X] + E[Y]
- $cov(X,Y) = E[(X E_X[X])(Y E_Y(Y)] = E[XY] E[X]E[Y]$

• Var(X + Y) = Var(X) + 2cov(X, Y) + Var(Y)

## Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

### **Gaussian Distribution**

Gaussian Distribution:

$$f(x|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{(x-\mu)^2}{2\sigma^2}}$$

Probability density function



### Prob vs Likelihood

### Prob vs Likelihood

### **Multivariate Gaussian Distribution**

$$p(x|\mu,\Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\{-\frac{1}{2} (x-\mu)^{\mathsf{T}} \Sigma^{-1} (x-\mu)\}$$

• Moment Parameterization  $\mu = E(X)$ 

$$\Sigma = Cov(X) = E[(X - \mu)(X - \mu)^{\mathsf{T}}]$$

- Mahalanobis Distance  $\Delta^2 = (x \mu)^T \Sigma^{-1} (x \mu)$
- Tons of applications (MoG, FA, PPCA, Kalman filter,...)

### **Properties of Gaussian Distribution**

The linear transform of a Gaussian r.v. is a Gaussian. Remember that no matter how x is distributed

E(AX + b) = AE(X) + b  $Cov(AX + b) = ACov(X)A^{T}$ this means that for Gaussian distributed quantities:

$$X \sim N(\mu, \Sigma) \rightarrow AX + b \sim N(A\mu + b, A\Sigma A^{\top})$$

The sum of two independent Gaussian r.v. is a Gaussian

$$Y = X_1 + X_2, X_1 \perp X_2 \rightarrow \mu_y = \mu_1 + \mu_2, \Sigma_y = \Sigma_1 + \Sigma_2$$

 The multiplication of two Gaussian functions is another Gaussian function (although no longer normalized)

$$N(a,A)N(b,B) \propto N(c,C),$$
  
where  $C = (A^{-1} + B^{-1})^{-1}, c = CA^{-1}a + CB^{-1}b$ 

### **Central Limit Theorem**



Probability mass function of a **biased** dice

Let's say, I am going to get a sample from this pmf having a size of n = 4

$$S_1 = \{1, 1, 1, 6\} \Rightarrow E(S_1) = 2.25$$

$$S_2 = \{1, 1, 3, 6\} \Rightarrow E(S_2) = 2.75$$

:  
$$S_m = \{1,4,6,6\} \Rightarrow E(S_m) = 4.25$$



According to CLT, it will follow a bell curve distribution (normal distribution)

## Outline

- Probability Distributions
- Joint and Conditional Probability Distributions
- Bayes' Rule
- Mean and Variance
- Properties of Gaussian Distribution
- Maximum Likelihood Estimation

## Maximum Likelihood Estimation

- Probability: inferring probabilistic quantities for data given fixed models (e.g. prob. of events, marginals, conditionals, etc).
- Statistics: inferring a model given fixed data observations (e.g. clustering, classification, regression).

### Main assumption:

Independent and identically distributed random variables i.i.d

### Maximum Likelihood Estimation

For Bernoulli (i.e. flip a coin):

Objective function:  $P(x_i|\theta) = \theta^{x_i}(1-\theta)^{1-x_i}$   $x_i \in \{0,1\} \text{ or } \{head, tail\}$ 

$$L(\theta|X) = L(\theta|X = x_1, X = x_2, X = x_3, \dots, X = x_n)$$
  
i.i.d assumption  
$$L(\theta|X) = \prod_{i=1}^n P(x_i|\theta)$$
$$L(\theta|X) = \prod_{i=1}^n P(x_i|\theta) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{1-x_i}$$

$$\begin{split} L(\theta|X) &= \theta^{x_1}(1-\theta)^{1-x_1} \times \theta^{x_2}(1-\theta)^{1-x_2} \dots \times \theta^{x_n}(1-\theta)^{1-x_n} \\ &= \theta^{\sum x_i}(1-\theta)^{\sum (1-x_i)} \end{split}$$

### We don't like multiplication, let's convert it into summation

What's the trick? Take the log  $L(\theta|X) = \theta^{\sum x_i} (1-\theta)^{\sum(1-x_i)}$   $logL(\theta|X) = l(\theta|X) = log(\theta) \sum_{i=1}^n x_i + log(1-\theta) \sum_{i=1}^n (1-x_i)$ 

How to optimize  $\theta$ ?

$$\frac{\partial l(\theta | X)}{\partial \theta} = 0 \qquad \qquad \frac{\sum_{i=1}^{n} x_i}{\theta} - \frac{\sum_{i=1}^{n} (1 - x_i)}{1 - \theta} = 0$$

$$\theta = \frac{1}{n} \sum_{i=1}^{n} x_i$$