

Open Elective Course [OE]

Course Code: CSO507

Winter 2023-24

Lecture#

Deep Learning

Unit-1: Probability and Information Theory for Machine/Deep Learning

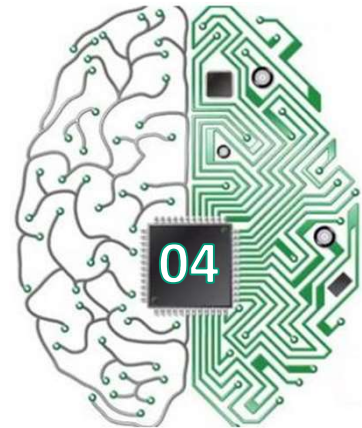
Course Instructor:

Dr. Monidipa Das

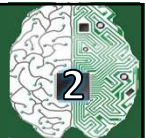
Assistant Professor

Department of Computer Science and Engineering

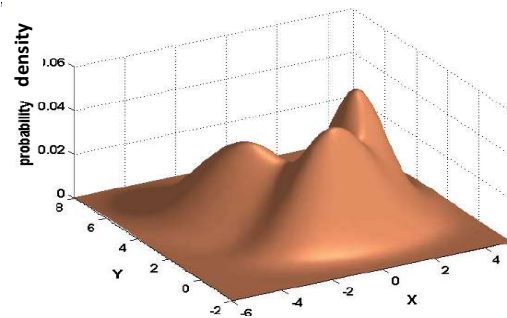
Indian Institute of Technology (Indian School of Mines) Dhanbad, Jharkhand 826004, India



Continuous Multivariate Distributions



- Same concepts of joint, marginal, and conditional probabilities apply for continuous random variables
- The probability distributions use integration of continuous random variables, instead of summation of discrete random variables

**Definition (Joint PDF for continuous random variables)**

Two random variables X and Y are *jointly continuous* if there exists a nonnegative function $f_{XY} : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, such that, for any set $A \subseteq \mathbb{R} \times \mathbb{R}$, we have

$$\mathbb{P}((X, Y) \in A) = \iint_A f_{XY}(x, y) dx dy.$$

The function $f_{XY}(x, y)$ is called the *joint probability density function* (PDF) of X and Y .

Continuous Multivariate Distributions (2)



Marginal PDFs

Suppose $f_{X,Y}(x,y)$ is a joint PDF of X and Y , then the *marginal densities* of X and of Y are given by

$$f_X(x) = \int_{-\infty}^{\infty} f_{X,Y}(x,y)dy, \quad f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y)dx.$$

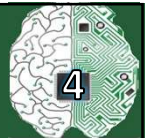
Definition (Conditional continuous random variable)

Suppose X and Y are jointly continuous, the *conditional probability density function (PDF)* of X given Y is given by

$$f_{X|Y=y}(x) = \frac{f_{X,Y}(x,y)}{f_Y(y)}.$$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Bayes' Theorem



- **Bayes' theorem** – allows to calculate conditional probabilities for one variable when conditional probabilities for another variable are known

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

- Also known as Bayes' rule
- **Multiplication rule** for the joint distribution is used: $P(X,Y) = P(Y|X)P(X)$
- By symmetry, we also have: $P(Y,X) = P(X|Y)P(Y)$
- The terms are referred to as:
 - $P(X)$, the **prior probability**, the initial degree of belief for X
 - $P(X|Y)$, the **posterior probability**, the degree of belief after incorporating the knowledge of Y
 - $P(Y|X)$, the **likelihood** of Y given X
 - $P(Y)$, the **evidence**
 - Bayes' theorem: **posterior probability** = $\frac{\text{likelihood} \times \text{prior probability}}{\text{evidence}}$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Independence



- Two random variables X and Y are **independent** if the occurrence of Y does not reveal any information about the occurrence of X
 - E.g., two successive rolls of a die are independent
- Therefore, we can write: $P(X|Y) = P(X)$
 - The following notation is used: $X \perp Y$
 - Also note that for independent random variables: $P(X, Y) = P(X)P(Y)$
- In all other cases, the random variables are **dependent**
- Two random variables X and Y are **conditionally independent** given another random variable Z if and only if $P(X, Y|Z) = P(X|Z)P(Y|Z)$
 - This is denoted as $X \perp Y|Z$

The joint density function factors for independent random variables

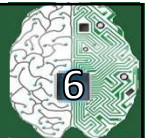
Jointly continuous random variables X and Y are independent if and only if

$$f_{X,Y}(x, y) = f_X(x)f_Y(y).$$

(f indicates PDF here)

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Expected Value



- The **expected value** or **expectation** of a function $f(x)$ with respect to a probability distribution $P(x)$ is the average (mean) when x is drawn from $P(x)$
- For a **discrete random variable** x , it is calculated as

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x)$$

- For a **continuous random variable** x , it is calculated as

$$\mathbb{E}_{x \sim P}[f(x)] = \int P(x)f(x) dx$$

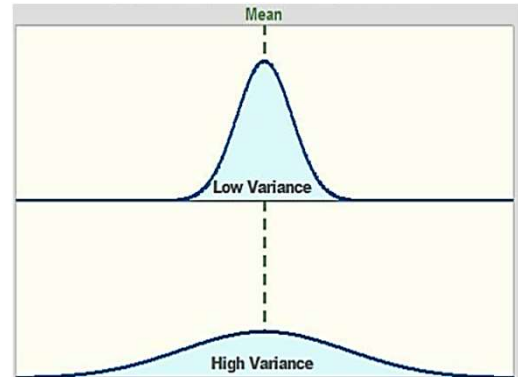
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Variance



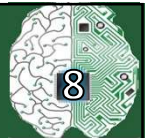
- **Variance** : how much the values of the function $f(x)$ deviate from the expected value as we sample values of x from $P(x)$

$$\text{Var}(f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$
- When the variance is low, the values of $f(x)$ cluster near the expected value
- The square root of the variance is the **standard deviation**
 - Denoted $\sigma = \sqrt{\text{Var}(f(x))}$



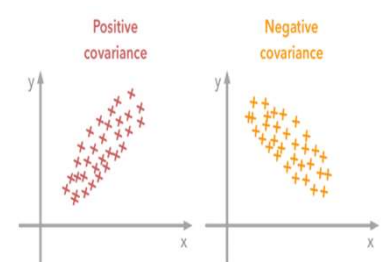
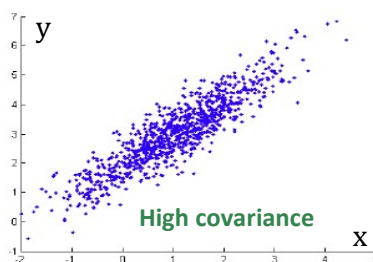
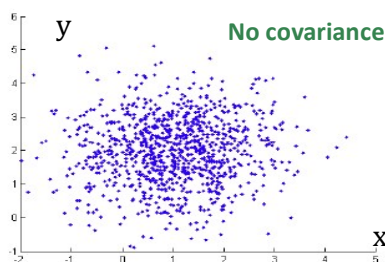
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Covariance



- **Covariance** gives the measure of how much two random variables are linearly related to each other

$$\text{Cov}(f(x), g(y)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])(g(y) - \mathbb{E}[g(y)])]$$
- The covariance measures the tendency for $f(x)$ and $g(y)$ to deviate from their means in same (or opposite) directions at same time



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Correlation



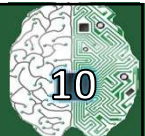
- **Correlation coefficient** is the covariance normalized by the standard deviations

$$\text{corr}(f(x), g(y)) = \frac{\text{Cov}(f(x), g(y))}{\sqrt{\text{Var}(f(x))} \cdot \sqrt{\text{Var}(g(y))}}$$

- It is also called **Pearson's correlation coefficient**
- The values are in the interval $[-1, 1]$
- It only reflects linear dependence between variables, and it does not measure non-linear dependencies between the variables

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Covariance Matrix



- **Covariance matrix** of a multivariate random variable \mathbf{x} with states $\mathbf{x} \in \mathbb{R}^n$ is an $n \times n$ matrix, such that

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j)$$

- i.e.,

$$\text{Cov}(\mathbf{x}) = \begin{bmatrix} \text{Cov}(x_1, x_1) & \text{Cov}(x_1, x_2) & \cdots & \text{Cov}(x_1, x_n) \\ \text{Cov}(x_2, x_1) & & \ddots & \text{Cov}(x_2, x_n) \\ \vdots & & & \vdots \\ \text{Cov}(x_n, x_1) & \text{Cov}(x_n, x_2) & \cdots & \text{Cov}(x_n, x_n) \end{bmatrix}$$

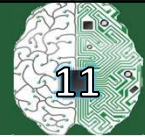
- The diagonal elements of the covariance matrix are the variances of the elements of the vector

$$\text{Cov}(x_i, x_i) = \text{Var}(x_i)$$

- Also note that the covariance matrix is symmetric, since $\text{Cov}(x_i, x_j) = \text{Cov}(x_j, x_i)$

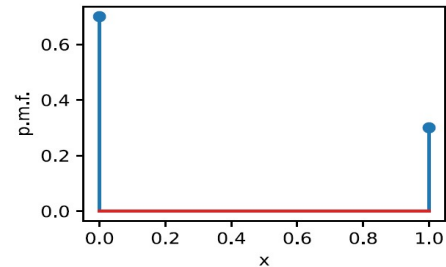
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Probability Distributions



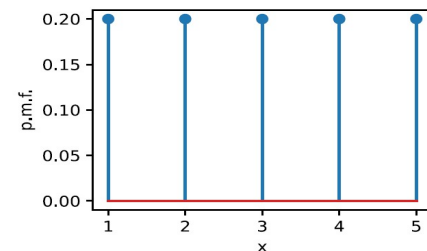
- **Bernoulli distribution**

- Binary random variable x with states $\{0, 1\}$
- The random variable can encode a coin flip which comes up 1 with probability p and 0 with probability $1 - p$
- Notation: $x \sim \text{Bernoulli}(p)$



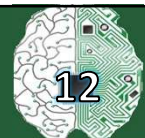
- **Uniform distribution**

- The probability of each value $i \in \{1, 2, \dots, n\}$ is $p_i = \frac{1}{n}$
- Notation: $x \sim U(n)$
- Figure: $n = 5, p = 0.2$



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Probability Distributions

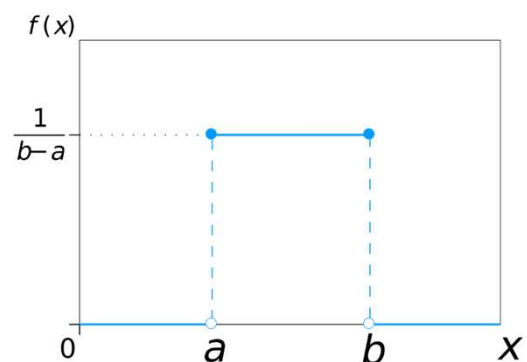


- **Uniform distribution [continuous]**

The pdf of a uniform random variable with domain $[a, b]$, where $b > a$ are real numbers, is given by

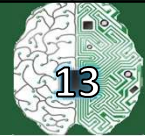
$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a \leq x \leq b, \\ 0, & \text{otherwise.} \end{cases}$$

Notation: $x \sim U(a, b)$



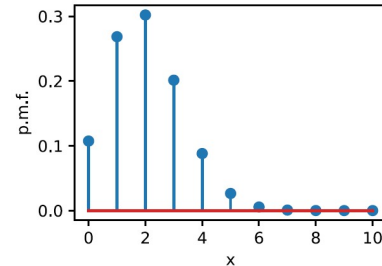
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Probability Distributions



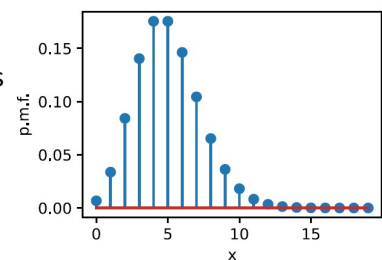
Binomial distribution

- Performing a sequence of n independent experiments, each of which has probability p of succeeding, where $p \in \{0, 1\}$
- The probability of getting k successes in n trials is $P(x = k) = \binom{n}{k} p^k (1-p)^{n-k}$
- Notation: $x \sim \text{Binomial}(n, p)$



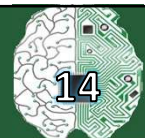
Poisson distribution

- A number of events occurring independently in a fixed interval of time with a known rate λ
- A discrete random variable x with states $k \in \{0, 1, 2, \dots\}$ has probability $P(x = k) = \frac{\lambda^k \cdot e^{-\lambda}}{k!}$
- The rate λ is the average number of occurrences of the event
- Notation: $x \sim \text{Poisson}(\lambda)$



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

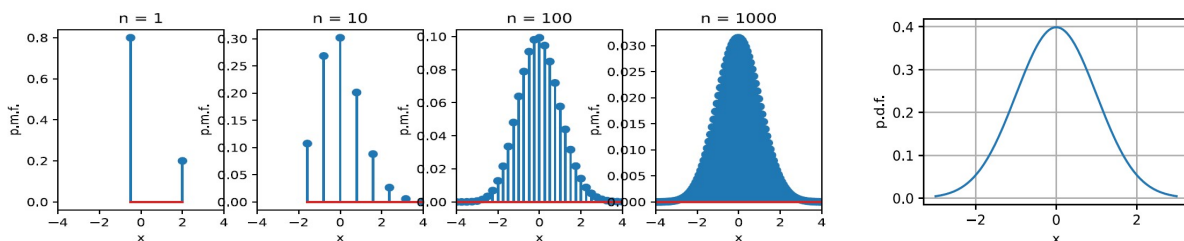
Probability Distributions



Gaussian distribution

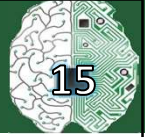
- The most well-studied distribution
 - Referred to as **normal distribution** or informally **bell-shaped distribution**
- Defined with the mean μ and variance σ^2
- Notation: $x \sim \mathcal{N}(\mu, \sigma^2)$
- For a random variable x with n independent measurements, the density is

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



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

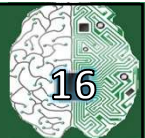
Probability Distributions



- **Multinoulli distribution**

- It is an extension of the Bernoulli distribution, from binary class to multi-class
- Multinoulli distribution is also called **categorical distribution** or **generalized Bernoulli distribution**
- For example, in multi-class classification in machine learning, we have a set of data examples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, and corresponding to the data example \mathbf{x}_i is a k -class label $\mathbf{y}_i = \{y_{i1}, y_{i2}, \dots, y_{ik}\}$ representing **one-hot encoding**
 - Let's denote the probabilities for assigning the class labels to a data example by $\{p_1, p_2, \dots, p_k\}$
 - The multinoulli probability of the data example \mathbf{x}_i is $P(\mathbf{x}_i) = p_1^{y_{i1}} \cdot p_2^{y_{i2}} \dots p_k^{y_{ik}} = \prod_j p_j^{y_{ij}}$
 - Similarly, we can calculate the probability of all data examples as $\prod_i \prod_j p_j^{y_{ij}}$

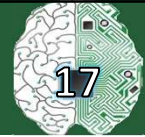
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



Information Theory

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

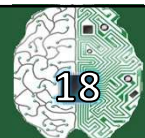
Information Theory



- **Information theory** studies encoding, decoding, transmitting, and manipulating information
- Father of information theory: *Claude Elwood Shannon*
- As such, information theory provides fundamental language for discussing the information processing in computer systems
 - E.g., machine learning applications use the cross-entropy loss [to be discussed in the next lecture], derived from information theoretic considerations

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Self-information



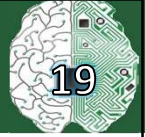
- Shannon defined the **self-information** of an event X as

$$I(X) = -\log(P(X))$$
 - $I(X)$ is the self-information, and $P(X)$ is the probability of the event X
- The self-information outputs the bits of information received for the event X
 - For example, if we want to send the code "0010" over a channel
 - The event "0010" is a series of codes of length n (in this case, the length is $n = 4$)
 - Each code is a **bit** (0 or 1), and occurs with probability of $\frac{1}{2}$; for this event $P = \frac{1}{2^n}$
$$I("0010") = -\log(P("0010")) = -\log\left(\frac{1}{2^4}\right) = -\log_2(1) + \log_2(2^4) = 0 + 4 = 4 \text{ bits}$$

[Log base e => unit is nats Log base 2 => unit is bits]

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Entropy



- For a discrete random variable X that follows a probability distribution P with a probability mass function $P(X)$, the expected amount of self-information is **entropy** (or **Shannon entropy**):

$$H(X) = \mathbb{E}_{X \sim P}[I(X)] = -\mathbb{E}_{X \sim P}[\log P(X)]$$

- Based on the expectation definition $\mathbb{E}_{X \sim P}[f(X)] = \sum_X P(X)f(X)$, we can rewrite the entropy as

$$H(X) = -\sum_X P(X) \log P(X)$$

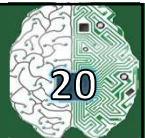
- If X is a continuous random variable that follows a probability distribution P with a probability density function $P(X)$, the entropy is

$$H(X) = -\int_X P(X) \log P(X) dX$$

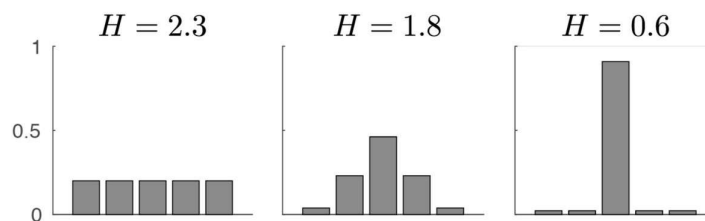
- For continuous random variables, the entropy is also called **differential entropy**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Entropy

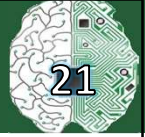


- Distributions that are closer to a uniform distribution have high entropy
- Because there is little surprise when we draw samples from a uniform distribution, since all samples have similar values



Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Kullback–Leibler Divergence



- **Kullback-Leibler (KL) divergence** (or **relative entropy**)
- provides a measure of how different two probability distribution are
- For two probability distributions $P(X)$ and $Q(X)$ over the same random variable X , the KL divergence is

$$D_{KL}(P||Q) = \mathbb{E}_{X \sim P} \left[\log \frac{P(X)}{Q(X)} \right] = \mathbb{E}_{X \sim P} [\log P(X) - \log Q(X)] = -\mathbb{E}_{X \sim P} \left[\log \frac{Q(X)}{P(X)} \right]$$

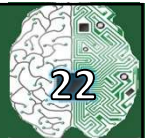
- For discrete random variables, this formula is equivalent to

$$D_{KL}(P||Q) = \sum_X P(X) \log \frac{P(X)}{Q(X)} = -\sum_X P(X) \log \frac{Q(X)}{P(X)}$$

- KL divergence is **non-negative**

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Jensen-Shannon divergence



- KL divergence is not a true distance metric, because it is **not symmetric**

$$D_{KL}(P||Q) \neq D_{KL}(Q||P)$$

- An alternative divergence which is non-negative and symmetric is the **Jensen-Shannon divergence**, defined as

$$D_{JS}(P||Q) = \frac{1}{2} D_{KL}(P||M) + \frac{1}{2} D_{KL}(Q||M)$$

- In the above, M is the average of the two distributions, $M = \frac{1}{2}(P + Q)$

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Exercise



- What is the gradient ($\nabla_x \sigma(x)$) of the sigmoid function $\sigma(x)$ as defined below?

$$\sigma(x) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$$

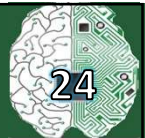
- What is the SVD decomposition of the following matrix:

$$\mathbf{A} = \begin{bmatrix} 3 & 2 \\ 2 & 4 \\ 1 & 2 \end{bmatrix}$$

(write using the following form : $U\Sigma V^T$ where U and V are orthogonal, and Σ is diagonal)

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Exercise



- Compute expectation and variance of the random variable X following:

- $P(X = k) = \phi(1 - \phi)^{k-1}$ for $k = 1, 2, \dots$ $\phi \in [0, 1]$

- $P(X = k) = e^{-\lambda} \frac{\lambda^k}{k!}$ for $k = 1, 2, \dots$ $\lambda > 0$

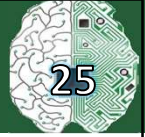
- Compute expectation and variance of the random variable X following:

- $p(X = x) = \frac{1}{b - a}$ $\forall x \in (a, b)$

- $p(X = x) = \lambda e^{-\lambda x}$ $x \geq 0, \lambda > 0$

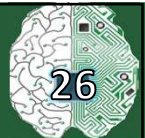
Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad

Exercise



- What can be the maximum value of D_{KL} ?
- What does $D_{KL} = 0$ indicate?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad



Questions?

Prof. Monidipa Das, Department of CSE, IIT (ISM) Dhanbad