

Lecture 2: Foundation

Deep Generative Models

Sajjad Amini

Department of Electrical Engineering
Sharif University of Technology

Contents

1 Color Codes

2 Notation

3 Probability and Statistics

- Probability Mass/Density Function
- Expectation
- Distance Metrics

4 Conclusions

Section 1

Color Codes

Color Coded Blocks

Definition Block

Result Block

Note Block

Example Block

Remember Block

Section 2

Notation

Scalars, Vectors and Matrices

Type	Non-random	Random
Scalar	x or X	X
Vector	\mathbf{x}	\mathbb{X}
Matrix	\mathbf{X}	\mathbb{X}
i -th element of a vector	x_i or $[\mathbf{x}]_i$	X_i or $[\mathbb{X}]_i$
(i, j) -th element of a matrix	x_{ij} or $[\mathbf{X}]_{ij}$	X_{ij} or $[\mathbb{X}]_{ij}$
i -th row of a matrix	$\mathbf{x}_{i:}$ or $[\mathbf{X}]_{i:}$	$\mathbb{X}_{i:}$ or $[\mathbb{X}]_{i:}$
j -th column of a matrix	$\mathbf{x}_{:j}$ or $[\mathbf{X}]_{:j}$	$\mathbb{X}_{:j}$ or $[\mathbb{X}]_{:j}$

*The element index appears at the end of subscript

i -th element of vector \mathbf{x}_k : $x_{k,i}$ or $[\mathbf{x}_k]_i$

Operators

- **Element-wise product:** Assume $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^D$, then:

$$\mathbf{z} = \mathbf{x} \odot \mathbf{y} \Leftrightarrow z_i = x_i \times y_i, i = 1, \dots, D$$

Operators

- **Element-wise product:** Assume $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^D$, then:

$$\mathbf{z} = \mathbf{x} \odot \mathbf{y} \Leftrightarrow z_i = x_i \times y_i, i = 1, \dots, D$$

- **Vectorization:** Assume $\mathbf{X} \in \mathbb{R}^{m \times n}$, then:

$$\mathbf{x} = \text{vec}(\mathbf{X}) \Leftrightarrow \mathbf{x} = \begin{bmatrix} [\mathbf{X}]_{:1} \\ [\mathbf{X}]_{:2} \\ \vdots \\ [\mathbf{X}]_{:n} \end{bmatrix}$$

Operators

- **Element-wise product:** Assume $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^D$, then:

$$\mathbf{z} = \mathbf{x} \odot \mathbf{y} \Leftrightarrow z_i = x_i \times y_i, i = 1, \dots, D$$

- **Vectorization:** Assume $\mathbf{X} \in \mathbb{R}^{m \times n}$, then:

$$\mathbf{x} = \text{vec}(\mathbf{X}) \Leftrightarrow \mathbf{x} = \begin{bmatrix} [\mathbf{X}]_{:1} \\ [\mathbf{X}]_{:2} \\ \vdots \\ [\mathbf{X}]_{:n} \end{bmatrix}$$

- **Trace:** The trace of a square matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$, denoted $\text{tr}(\mathbf{X})$, is the sum of diagonal elements in the matrix as:

$$\text{tr}(\mathbf{X}) \triangleq \sum_{i=1}^n x_{ii}$$

Operators

- **Norm:** The ℓ_p norm for vector $\mathbf{x} \in \mathbb{R}^n$ is defined as:

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, p \geq 1 \Rightarrow \begin{cases} \ell_1 : \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| \\ \ell_2 : \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \\ \ell_\infty : \|\mathbf{x}\|_\infty = \max_i |x_i| \end{cases}$$

Operators

- **Norm:** The ℓ_p norm for vector $\mathbf{x} \in \mathbb{R}^n$ is defined as:

$$\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}, p \geq 1 \Rightarrow \begin{cases} \ell_1 : \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i| \\ \ell_2 : \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \\ \ell_\infty : \|\mathbf{x}\|_\infty = \max_i |x_i| \end{cases}$$

- **Transpose:** The transpose of matrix $\mathbf{X} \in \mathbb{R}^{n \times m}$, denoted by \mathbf{X}^T , is:

$$[\mathbf{X}^T]_{ji} = [\mathbf{X}]_{ij}, \begin{cases} i = 1, \dots, n \\ j = 1, \dots, m \end{cases}$$

Operators

- **Diag:**

$$\mathbf{x} \in \mathbb{R}^n \Rightarrow \mathbf{X} = \text{diag}(\mathbf{x}) = \begin{bmatrix} [\mathbf{x}]_1 & 0 & \dots & 0 & 0 \\ 0 & [\mathbf{x}]_2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & [\mathbf{x}]_{n-1} & 0 \\ 0 & 0 & \dots & 0 & [\mathbf{x}]_n \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Operators

- **Diag:**

$$\mathbf{x} \in \mathbb{R}^n \Rightarrow \mathbf{X} = \text{diag}(\mathbf{x}) = \begin{bmatrix} [\mathbf{x}]_1 & 0 & \dots & 0 & 0 \\ 0 & [\mathbf{x}]_2 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & [\mathbf{x}]_{n-1} & 0 \\ 0 & 0 & \dots & 0 & [\mathbf{x}]_n \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$\mathbf{X} \in \mathbb{R}^{n \times n} \Rightarrow \mathbf{x} = \text{diag}(\mathbf{X}) = \begin{bmatrix} [\mathbf{X}]_{11} \\ [\mathbf{X}]_{22} \\ \vdots \\ [\mathbf{X}]_{(n-1)(n-1)} \\ [\mathbf{X}]_{nn} \end{bmatrix}$$

Matrix Calculus

- **Gradient vector:** The gradient vector for $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point \mathbf{x} is:

$$\frac{\partial f}{\partial \mathbf{x}} = \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

Matrix Calculus

- **Gradient vector:** The gradient vector for $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at point \mathbf{x} is:

$$\frac{\partial f}{\partial \mathbf{x}} = \nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

- **Jacobian matrix:** The Jacobian matrix for $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at point \mathbf{x} is:

$$\mathbf{J}_f(\mathbf{x}) = \frac{\partial \mathbf{f}}{\partial \mathbf{x}^T} \triangleq \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla f_1(\mathbf{x})^T \\ \vdots \\ \nabla f_m(\mathbf{x})^T \end{bmatrix}$$

Definitions

- **LHS and RHS:** Left Hand Side (LHS) and Right Hand Side (RHS) refer to:

$$\underbrace{z}_{\text{LHS}} = \underbrace{x \odot y}_{\text{RHS}}$$

Section 3

Probability and Statistics

Subsection 1

Probability Mass/Density Function

Discrete Random Variable

Probabilistic
Experiment



Rolling a Die

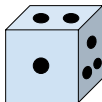


Figure: Probabilistic experiment: an experiment where the result is NOT certain a priori

Discrete Random Variable

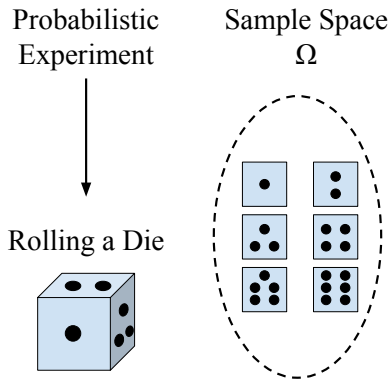


Figure: Sample space Ω : set of all possible outcomes

Discrete Random Variable

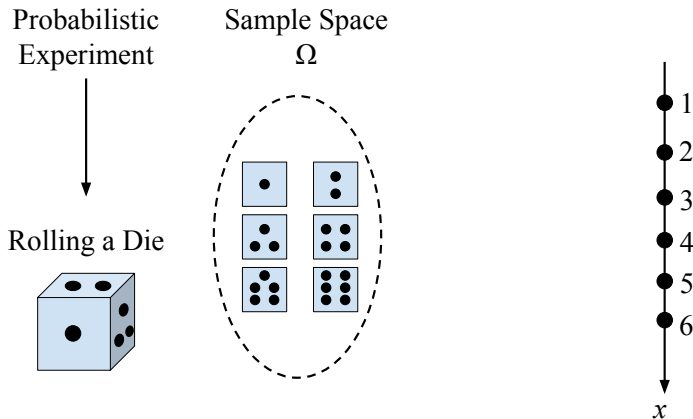


Figure: Numeric numbers x

Discrete Random Variable

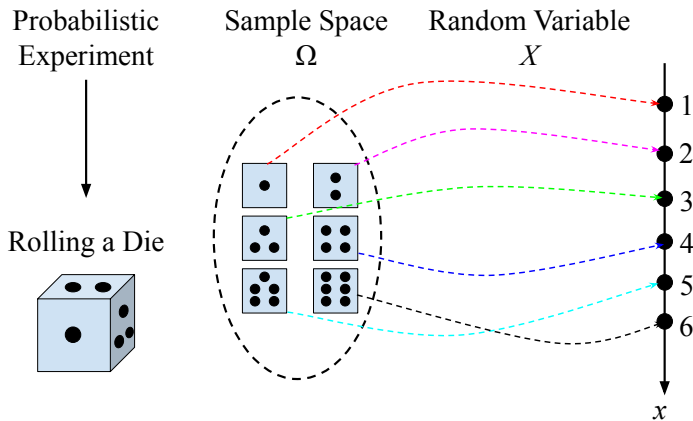


Figure: Random variable X : a function which maps every sample in Ω to a numeric number x

Event and Probability Mass Function

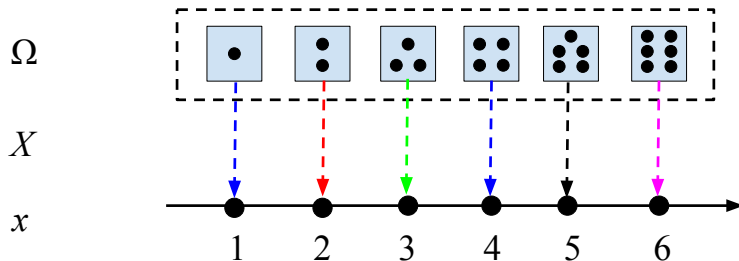


Figure: Sample space Ω , random variable X and numeric number x

Event and Probability Mass Function

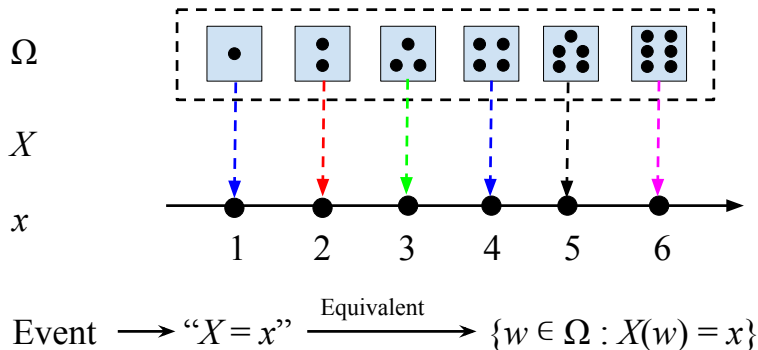


Figure: Event definition

Event and Probability Mass Function

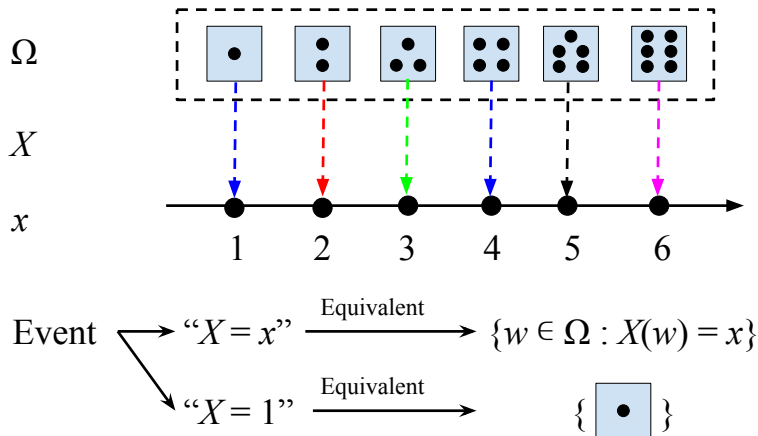


Figure: A typical event

Event and Probability Mass Function

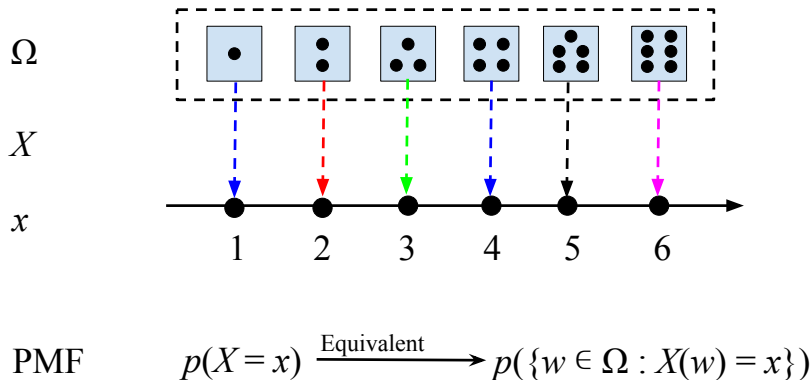


Figure: Probability mass function $p(X = x)$

Properties

PMF assigns a mass to each event corresponding to the event probability, so it must satisfy the following properties:

$$p(X = x) \geq 0$$

$$\sum_x p(X = x) = 1$$

Bernoulli

Assume $x \in \{0, 1\}$, then random variable X is Bernoulli, denoted by:

$$X \sim \text{Ber}(\theta)$$

Sample Probability Mass Function

Bernoulli

Assume $x \in \{0, 1\}$, then random variable X is Bernoulli, denoted by:

$$X \sim \text{Ber}(\theta)$$

or

$$p_{\theta}(X) = \text{Ber}(X|\theta)$$

Sample Probability Mass Function

Bernoulli

Assume $x \in \{0, 1\}$, then random variable X is Bernoulli, denoted by:

$$X \sim \text{Ber}(\theta)$$

or

$$p_{\theta}(X) = \text{Ber}(X|\theta)$$

And we have:

$$p_{\theta}(X = x) = \begin{cases} \theta & x = 1 \\ 1 - \theta & x = 0 \end{cases}$$

Note that θ must satisfy $0 \leq \theta \leq 1$.

Sample Probability Mass Function

Categorical

Assume $x \in \{1, 2, \dots, L\}$, then random variable X is Categorical, denoted by:

$$X \sim \text{Cat}(\boldsymbol{\theta})$$

Sample Probability Mass Function

Categorical

Assume $x \in \{1, 2, \dots, L\}$, then random variable X is Categorical, denoted by:

$$X \sim \text{Cat}(\boldsymbol{\theta})$$

or

$$p_{\boldsymbol{\theta}}(X) = \text{Cat}(X|\boldsymbol{\theta})$$

Sample Probability Mass Function

Categorical

Assume $x \in \{1, 2, \dots, L\}$, then random variable X is Categorical, denoted by:

$$X \sim \text{Cat}(\boldsymbol{\theta})$$

or

$$p_{\theta}(X) = \text{Cat}(X|\boldsymbol{\theta})$$

And we have:

$$p_{\theta}(X = l) = \theta_l$$

Note that θ must satisfy $0 \leq \theta \leq 1$ and $\sum_l \theta_l = 1$.

Extension to Random Vector

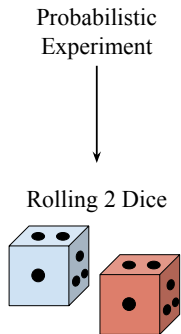


Figure: Rolling two dice experiment

Extension to Random Vector

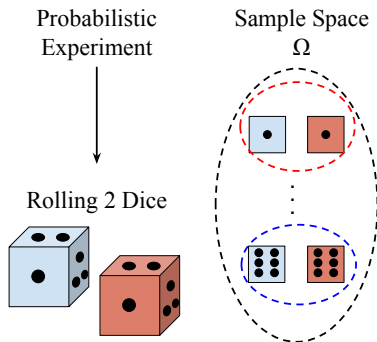


Figure: Sample space Ω

Extension to Random Vector

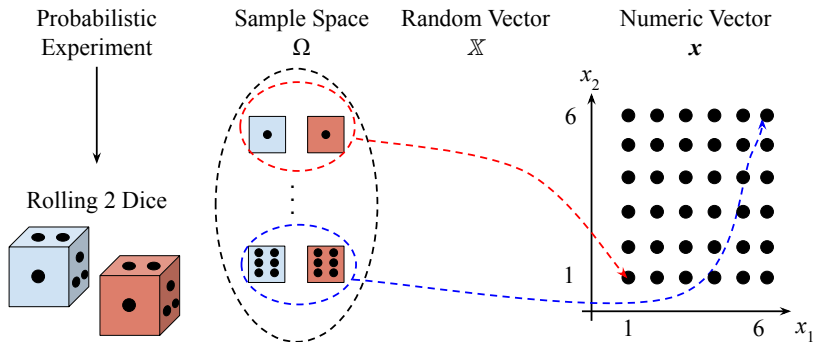


Figure: Random vector \mathbb{X} and numeric vector \mathbf{x}

Probability Mass Function

Properties

PMF properties for a random variable can be easily extended to random vectors. Assume we have:

$$\mathbb{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix},$$

then we have:

$$\begin{aligned} p(\mathbb{X} = \mathbf{x}) &= p(X_1 = x_1, X_2 = x_2) \geq 0 \\ \sum_{\mathbf{x}} p(\mathbb{X} = \mathbf{x}) &= \sum_{x_1} \sum_{x_2} p(X_1 = x_1, X_2 = x_2) = 1 \end{aligned}$$

To have PMF over only X_1 random variable, we can use *Marginalization* as:

$$p(X_1 = x_1) = \sum_{x_2} p(X_1 = x_1, X_2 = x_2)$$

Conditional Distribution

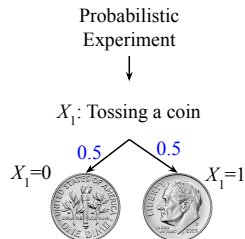


Figure: First random variable: Tossing a coin

Conditional Distribution

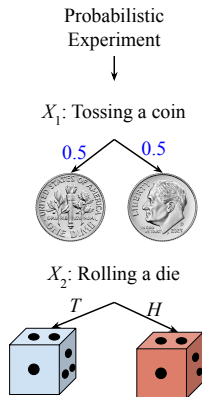


Figure: Second random variable: Rolling a die based on coin experiment

Conditional Distribution

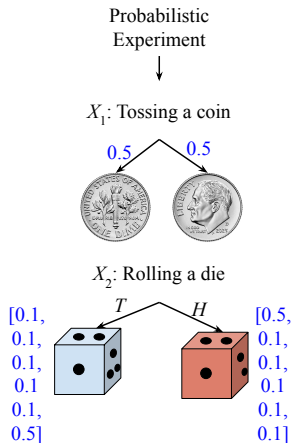


Figure: The distribution for each die (the dice are not fair)

Conditional Distribution

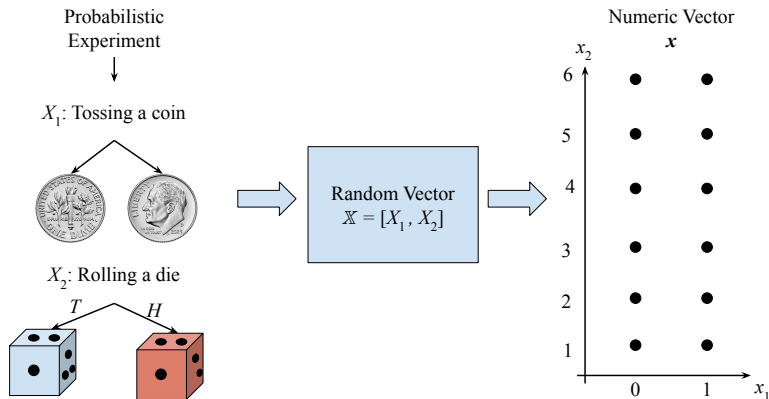


Figure: Two dimensional random variable \mathbb{X} and corresponding numeric vectors

Conditional Distribution

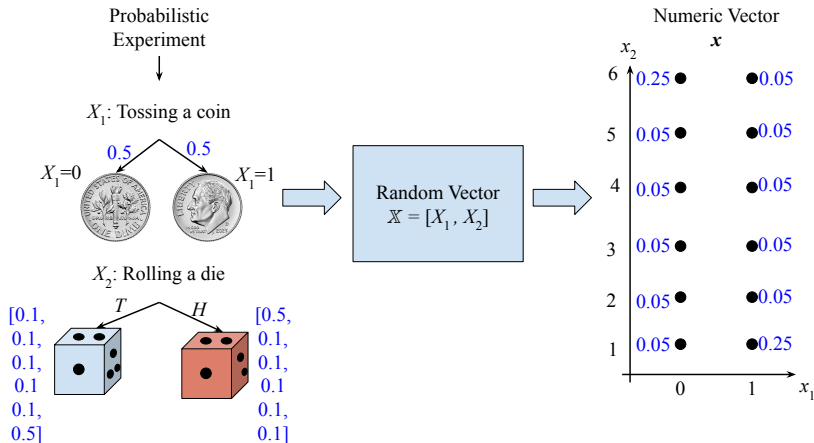
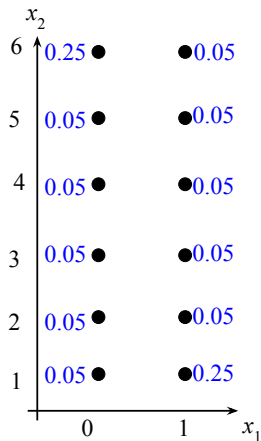
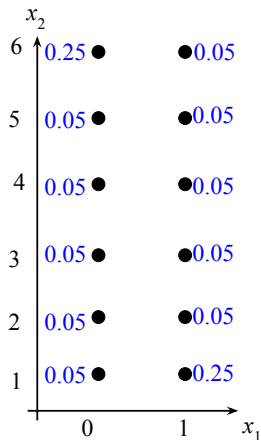


Figure: The PMF over random vector \mathbb{X}

Unconditional Event



Unconditional Event

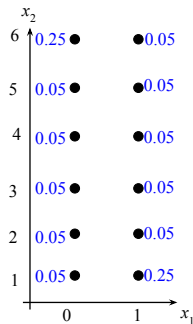


Unconditional Probability

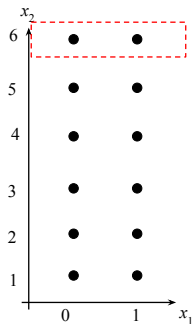
Assume we are interested in calculating the probability for event $X_1 = 0$, then using marginalization we have:

$$\begin{aligned} p(X_1 = 0) &= \sum_{x_2=1}^6 p(X_1 = 0, X_2 = x_2) \\ &= 0.05 + 0.05 + 0.05 + 0.05 + 0.05 + 0.25 = 0.5 \end{aligned}$$

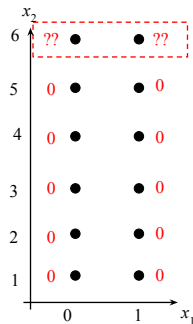
Conditional Event



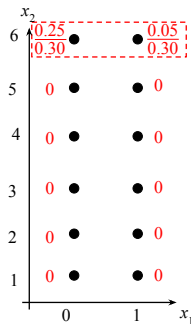
(a) $p(X_1, X_2)$



(b) Event $X_2 = 6$

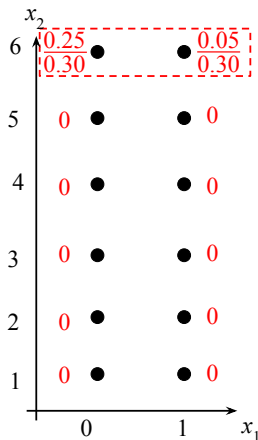


(c) Impossible events



(d) $p(X_1, X_2 | X_2 = 6)$

Conditional Event



Conditional Probability

Assume we are interested in calculating the probability for event $X_1 = 1$ conditioned on the fact that $X_2 = 6$, then using marginalization we have:

$$\begin{aligned} & p(X_1 = 0 | X_2 = 6) \\ &= \sum_{x_2=1}^6 p(X_1 = 0, X_2 = x_2 | X_2 = 6) \\ &= 0 + 0 + 0 + 0 + 0 + \frac{0.25}{0.30} \simeq 0.83 \end{aligned}$$

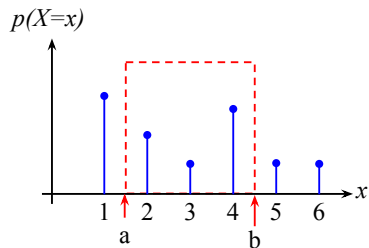
Conditional Probability Mass Function

Conditional PMF

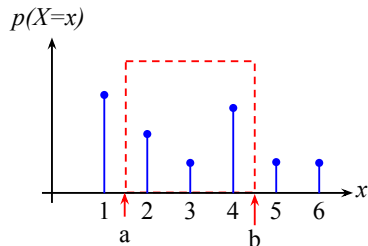
Conditional PMF is a principled way of updating a PMF given the information that some events happened. Conditional PMF is defined as:

$$p(X_1 = x_1 | X_2 = x_2) = \frac{p(X_1 = x_1, X_2 = x_2)}{p(X_2 = x_2)}$$

Extension to Continuous Random Variable [1]



Extension to Continuous Random Variable [1]

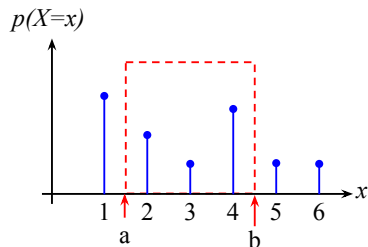


Properties

$$P(a \leq X \leq b) = \sum_{x:a \leq x \leq b} p(X = x)$$

where $p(X = x) \geq 0$ and $\sum_x p(X = x) = 1$

Extension to Continuous Random Variable [1]



Properties

$$P(a \leq X \leq b) = \sum_{x:a \leq x \leq b} p(X = x)$$

where $p(X = x) \geq 0$ and $\sum_x p(X = x) = 1$

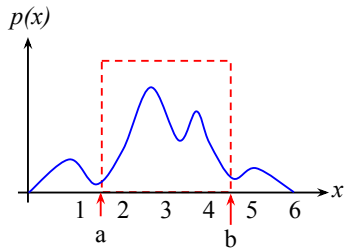
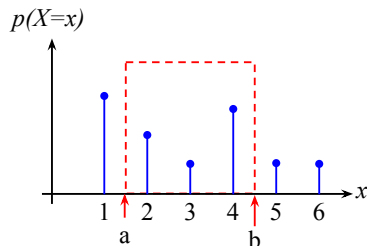


Figure: Probability Density Function

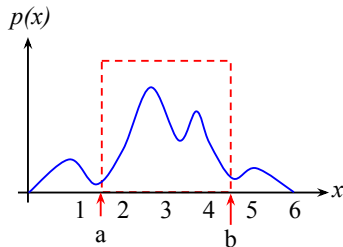
Extension to Continuous Random Variable [1]



Properties

$$P(a \leq X \leq b) = \sum_{x:a \leq x \leq b} p(X=x)$$

where $p(X=x) \geq 0$ and $\sum_x p(X=x) = 1$



Properties

$$P(a \leq X \leq b) = \int_a^b p(x)dx$$

where

$$p(x) \geq 0 \text{ and } \int_{-\infty}^{\infty} p(x)dx = 1$$

Figure: Probability Density Function

Gaussian

Gaussian random variable is an example of a continuous random variable denoted by

$$X \sim \mathcal{N}(\mu, \sigma^2) \text{ or } p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

Sample Probability Density Function

Gaussian

Gaussian random variable is an example of a continuous random variable denoted by

$$X \sim \mathcal{N}(\mu, \sigma^2) \text{ or } p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

μ is the *mean* and σ^2 is the *variance* of random variable. The probability density function (PDF) is defined as:

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

Sample Probability Density Function

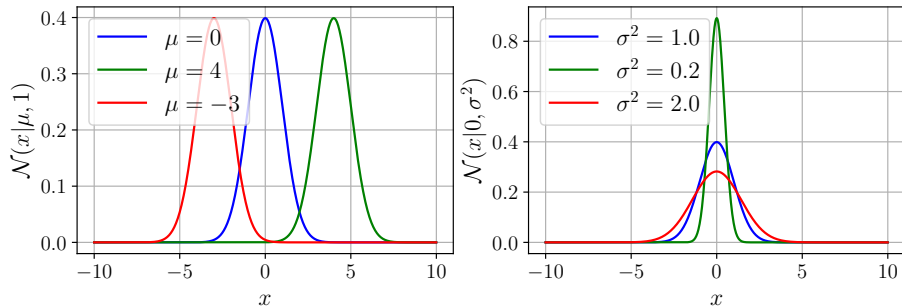


Figure: Mean μ effect (left) and Variance σ^2 (right) on the Gaussian PDF

PMF and PDF Notation

Throughout this course, we use the following notation:

Definition	Distribution	Numeric Probability
PMF/PDF	$p(X)$ or $p(\mathbb{X})$	$p(x)$ or $p(\mathbf{x})$
Conditional PMF/PDF	$p(X y)$ or $p(\mathbb{X} \mathbf{y})$	$p(x y)$ or $p(\mathbf{x} \mathbf{y})$
Model PMF/PDF	$p_{\theta}(\mathbb{X})$ or $p_{\theta}(\mathbb{X} y)$	$p_{\theta}(\mathbf{x})$ or $p_{\theta}(\mathbf{x} \mathbf{y})$
Data PMF/PDF	$p_{\text{data}}(\mathbb{X})$ or $p_{\text{data}}(\mathbb{X} y)$	$p_{\text{data}}(\mathbf{x})$ or $p_{\text{data}}(\mathbf{x} \mathbf{y})$

Subsection 2

Expectation

Expectation

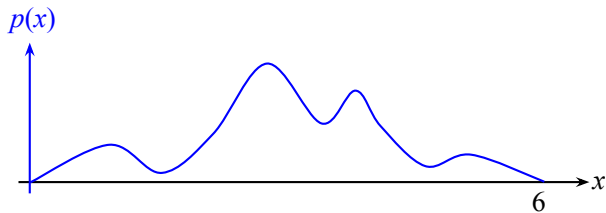


Figure: Probability density function for the time one can walk in 6 hours

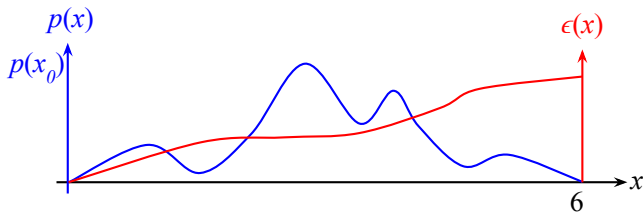


Figure: PDF with energy consumption function $\epsilon(x)$

Expectation

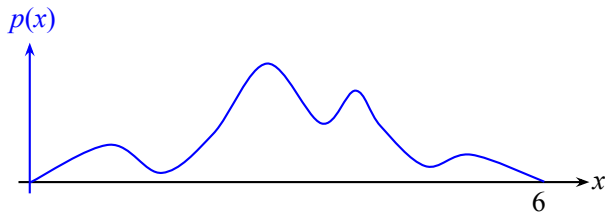


Figure: Probability density function for the time one can walk in 6 hours

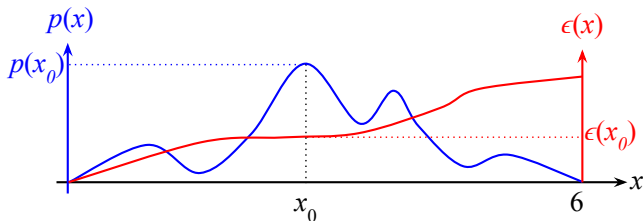


Figure: A sample high probable point x_0

Expectation

Expectation

In the general case, expectation can be interpreted as the average for the function $\epsilon(x)$ in a large number of *independent* repetitions of the experiment. This value is determined by:

$$\mathbb{E}_{x \sim p(X)} [\epsilon(x)]$$

Expectation

Expectation

In the general case, expectation can be interpreted as the average for the function $\epsilon(x)$ in a large number of *independent* repetitions of the experiment. This value is determined by:

$$\mathbb{E}_{x \sim p(X)} [\epsilon(x)]$$

and is calculated as:

$$\mathbb{E}_{x \sim p(X)} [\epsilon(x)] = \int_{-\infty}^{\infty} p(x) \epsilon(x) dx$$

Expectation

Expectation

In the general case, expectation can be interpreted as the average for the function $\epsilon(x)$ in a large number of *independent* repetitions of the experiment. This value is determined by:

$$\mathbb{E}_{x \sim p(X)} [\epsilon(x)]$$

and is calculated as:

$$\mathbb{E}_{x \sim p(X)} [\epsilon(x)] = \int_{-\infty}^{\infty} p(x) \epsilon(x) dx$$

Equivalently in the case of discrete a random variable, we have:

$$\mathbb{E}_{x \sim p(X)} [\epsilon(x)] = \sum_x p(x) \epsilon(x)$$

Sample Expectations

Sample Expectations

Consider $\mathbb{E}_{x \sim p(X)}[f(x)]$, then:

Sample Expectations

Sample Expectations

Consider $\mathbb{E}_{x \sim p(X)}[f(x)]$, then:

- If $f(x) = x$, then the resulting expectation is *mean* and denoted by μ .

Sample Expectations

Consider $\mathbb{E}_{x \sim p(X)}[f(x)]$, then:

- If $f(x) = x$, then the resulting expectation is *mean* and denoted by μ .
- If $f(x) = (x - \mu)^2$, then the resulting expectation is *variance* and denoted by σ^2 .

Sample Expectations

Consider $\mathbb{E}_{x \sim p(X)}[f(x)]$, then:

- If $f(x) = x$, then the resulting expectation is *mean* and denoted by μ .
- If $f(x) = (x - \mu)^2$, then the resulting expectation is *variance* and denoted by σ^2 .
- If $f(x) = x^n$, then the resulting expectation is *n-th raw moment* and denoted by μ'_n .

Monte Carlo Estimation

Consider random variable \mathbb{X} with distribution $p(\mathbb{X})$. The expectation of function $f(\boldsymbol{x})$ can be calculated as:

$$\mathbb{E}_{\boldsymbol{x} \sim p(\mathbb{X})} [f(\boldsymbol{x})] = \int p(\boldsymbol{x}) f(\boldsymbol{x}) d\boldsymbol{x}$$

Monte Carlo Estimation

Consider random variable \mathbb{X} with distribution $p(\mathbb{X})$. The expectation of function $f(\mathbf{x})$ can be calculated as:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

Now assume that instead of $p(\mathbb{X})$, we just have access to N independent samples of random variable \mathbb{X} as $\mathbf{x}_1, \dots, \mathbf{x}_N$.

Monte Carlo Estimation

Monte Carlo Estimation

Consider random variable \mathbb{X} with distribution $p(\mathbb{X})$. The expectation of function $f(\mathbf{x})$ can be calculated as:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

Now assume that instead of $p(\mathbb{X})$, we just have access to N independent samples of random variable \mathbb{X} as $\mathbf{x}_1, \dots, \mathbf{x}_N$. Then we define $f_N(\{\mathbf{x}_i\})$ as:

$$f_N(\{\mathbf{x}_i\}) = \frac{1}{N} \sum_n f(\mathbf{x}_n) \quad \# \text{ Monte Carlo Estimation (MCE)}$$

Monte Carlo Estimation

Consider random variable \mathbb{X} with distribution $p(\mathbb{X})$. The expectation of function $f(\mathbf{x})$ can be calculated as:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

Now assume that instead of $p(\mathbb{X})$, we just have access to N independent samples of random variable \mathbb{X} as $\mathbf{x}_1, \dots, \mathbf{x}_N$. Then we define $f_N(\{\mathbf{x}_i\})$ as:

$$f_N(\{\mathbf{x}_i\}) = \frac{1}{N} \sum_n f(\mathbf{x}_n) \quad \# \text{ Monte Carlo Estimation (MCE)}$$

Then using the *weak law of large numbers*, for arbitrary small positive ϵ :

$$\lim_{n \rightarrow \infty} P\left(\left|f_N(\{\mathbf{x}_i\}) - \mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})]\right| \geq \epsilon\right) = 0$$

Estimation Variance

Assume $X \sim N(1, 4)$, then:

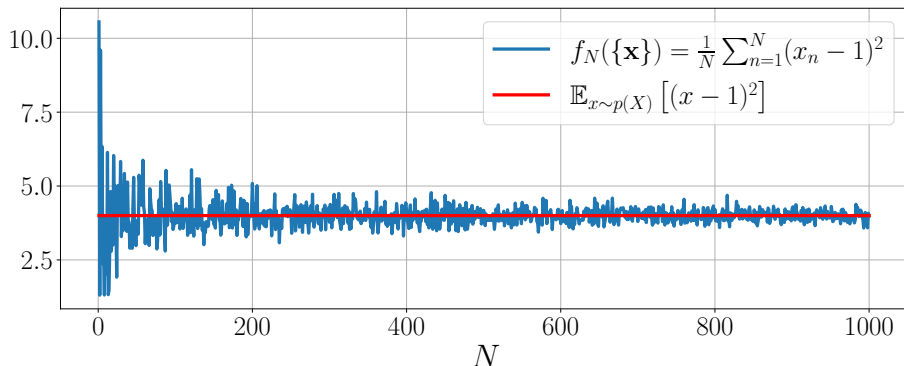
- $\mathbb{E}_{x \sim p(X)} [(x - 1)^2] = \sigma^2 = 4$

Monte Carlo Estimation

Estimation Variance

Assume $X \sim N(1, 4)$, then:

- $\mathbb{E}_{x \sim p(X)} [(x - 1)^2] = \sigma^2 = 4$
- The figure below shows the result of MCE for different values of n .



Subsection 3

Distance Metrics

Kullback-Leibler Divergence

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X)||q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

Kullback-Leibler Divergence

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X) \| q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- $\text{KL} (p(X) \| q(X)) \neq \text{KL} (q(X) \| p(X)) \Rightarrow$ KLD is not symmetric

Kullback-Leibler Divergence

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X) \| q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- $\text{KL} (p(X) \| q(X)) \neq \text{KL} (q(X) \| p(X)) \Rightarrow$ KLD is not symmetric
- $\text{KL} (p(X) \| q(X)) \geq 0 \Rightarrow$ KLD is non-negative

Kullback-Leibler Divergence

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X) \| q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- $\text{KL} (p(X) \| q(X)) \neq \text{KL} (q(X) \| p(X)) \Rightarrow$ KLD is not symmetric
- $\text{KL} (p(X) \| q(X)) \geq 0 \Rightarrow$ KLD is non-negative
- $\text{KL} (p(X) \| q(X)) = 0 \Leftrightarrow p(x) = q(x) \ \forall x$

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X) \| q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- $\text{KL} (p(X) \| q(X)) \neq \text{KL} (q(X) \| p(X)) \Rightarrow$ KLD is not symmetric
- $\text{KL} (p(X) \| q(X)) \geq 0 \Rightarrow$ KLD is non-negative
- $\text{KL} (p(X) \| q(X)) = 0 \Leftrightarrow p(x) = q(x) \ \forall x$
- KLD is not upper-bounded in general.

Kullback-Leibler Divergence

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X)||q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- $\text{KL} (p(X)||q(X)) \neq \text{KL} (q(X)||p(X)) \Rightarrow$ KLD is not symmetric
- $\text{KL} (p(X)||q(X)) \geq 0 \Rightarrow$ KLD is non-negative
- $\text{KL} (p(X)||q(X)) = 0 \Leftrightarrow p(x) = q(x) \forall x$
- KLD is not upper-bounded in general.
- KLD is not a distance metric.

Kullback-Leibler Divergence

Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X)||q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

- $\text{KL} (p(X)||q(X)) \neq \text{KL} (q(X)||p(X)) \Rightarrow$ KLD is not symmetric
- $\text{KL} (p(X)||q(X)) \geq 0 \Rightarrow$ KLD is non-negative
- $\text{KL} (p(X)||q(X)) = 0 \Leftrightarrow p(x) = q(x) \forall x$
- KLD is not upper-bounded in general.
- KLD is not a distance metric.
 - It is not symmetric.

Kullback-Leibler Divergence

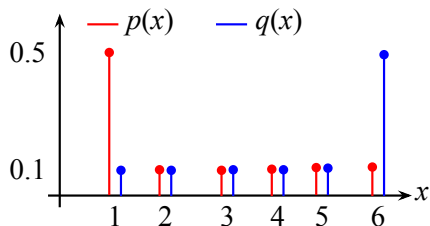
Kullback-Leibler Divergence

Kullback-Leibler divergence (KLD) is a metric to calculate the distance between two distributions. KLD for two distributions p and q defined over discrete random variable X is:

$$\text{KL} (p(X)||q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

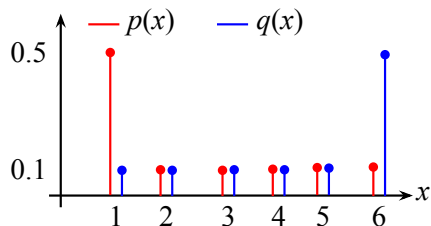
- $\text{KL} (p(X)||q(X)) \neq \text{KL} (q(X)||p(X)) \Rightarrow$ KLD is not symmetric
- $\text{KL} (p(X)||q(X)) \geq 0 \Rightarrow$ KLD is non-negative
- $\text{KL} (p(X)||q(X)) = 0 \Leftrightarrow p(x) = q(x) \forall x$
- KLD is not upper-bounded in general.
- KLD is not a distance metric.
 - It is not symmetric.
 - It does not satisfy the triangle inequality.

Kullback-Leibler divergence

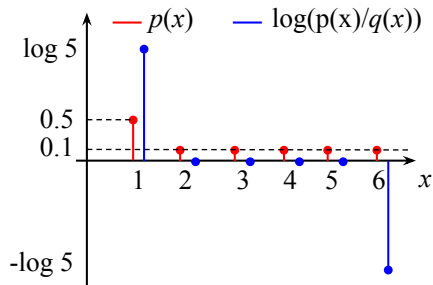


(a) Two distributions $p(X)$ and $q(X)$

Kullback-Leibler divergence

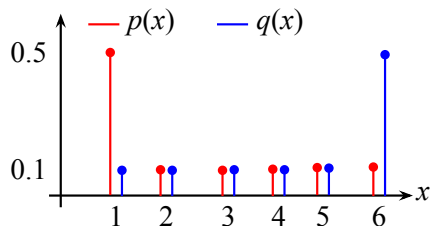


(a) Two distributions $p(X)$ and $q(X)$

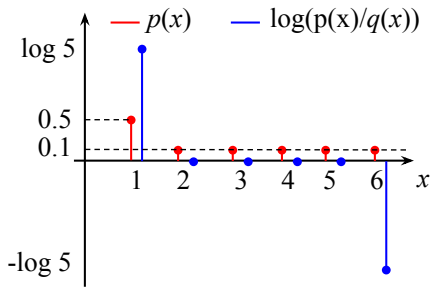


(b) $p(X)$ and log likelihood ratio

Kullback-Leibler divergence



(a) Two distributions $p(X)$ and $q(X)$



(b) $p(X)$ and log likelihood ratio

KLD Between Two Categorical Distributions

$$\begin{aligned}\text{KL}(p(X)||q(X)) &= 0.5 \times \log 5 + 0.1 \times 0 + 0.1 \times 0 + 0.1 \times 0 + 0.1 \times 0 + 0.1 \times \log \frac{1}{5} \\ &= 0.5 \times \log 5 - 0.1 \times \log 5 = 0.4 \times \log 5 \simeq 0.64\end{aligned}$$

Kullback-Leibler divergence

Kullback-Leibler divergence

KLD is similarly defined for continuous random variables as:

$$\text{KL} (p(X) \| q(X)) \triangleq \mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{q(x)} \right] = \int_x p(x) \log \frac{p(x)}{q(x)} dx$$

Jensen-Shanon Divergence

Jensen-Shanon Divergence (JSD) is a symmetric and smoothed version of the KLD. To define JSD for two distributions p and q , first we should define new distribution $m(X)$ as:

$$m(X) \triangleq \frac{p(X) + q(X)}{2} \Leftrightarrow m(X = x) = \frac{p(X = x) + q(X = x)}{2} \quad \forall x$$

Then JSD is defined as:

$$\begin{aligned} \text{JS} (p(X) \| q(X)) &\triangleq \frac{1}{2} \left(\text{KL} (p(X) \| m(X)) + \text{KL} (q(X) \| m(X)) \right) \\ &= \frac{1}{2} \left(\mathbb{E}_{x \sim p(X)} \left[\log \frac{p(x)}{m(x)} \right] + \mathbb{E}_{x \sim q(X)} \left[\log \frac{q(x)}{m(x)} \right] \right) \end{aligned}$$

Jensen-Shanon Divergence

$$\text{JS} (p(X) \| q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X) \| m(X)) + \text{KL} (q(X) \| m(X)) \right)$$

Jensen-Shanon Divergence

$$\text{JS} (p(X) \| q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X) \| m(X)) + \text{KL} (q(X) \| m(X)) \right)$$

- $\text{JS} (p(X) \| q(X)) = \text{JS} (q(X) \| p(X)) \Rightarrow \text{JSD is symmetric}$

Jensen-Shanon Divergence

$$\text{JS} (p(X) \| q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X) \| m(X)) + \text{KL} (q(X) \| m(X)) \right)$$

- $\text{JS} (p(X) \| q(X)) = \text{JS} (q(X) \| p(X)) \Rightarrow \text{JSD is symmetric}$
- $0 \leq \text{JS} (p(X) \| q(X)) \leq 1 \Rightarrow \text{JSD is non-negative and upper-bounded.}$

Jensen-Shanon Divergence

$$\text{JS} (p(X) \| q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X) \| m(X)) + \text{KL} (q(X) \| m(X)) \right)$$

- $\text{JS} (p(X) \| q(X)) = \text{JS} (q(X) \| p(X)) \Rightarrow \text{JSD is symmetric}$
- $0 \leq \text{JS} (p(X) \| q(X)) \leq 1 \Rightarrow \text{JSD is non-negative and upper-bounded.}$
- $\text{JS} (p(X) \| q(X)) = 0 \Leftrightarrow p(x) = q(x) \ \forall x$

Jensen-Shanon Divergence

$$\text{JS} (p(X) \| q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X) \| m(X)) + \text{KL} (q(X) \| m(X)) \right)$$

- $\text{JS} (p(X) \| q(X)) = \text{JS} (q(X) \| p(X)) \Rightarrow \text{JSD is symmetric}$
- $0 \leq \text{JS} (p(X) \| q(X)) \leq 1 \Rightarrow \text{JSD is non-negative and upper-bounded.}$
- $\text{JS} (p(X) \| q(X)) = 0 \Leftrightarrow p(x) = q(x) \ \forall x$
- JSD is not a distance metric.

Jensen-Shanon Divergence

$$\text{JS} (p(X)\|q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X)\|m(X)) + \text{KL} (q(X)\|m(X)) \right)$$

- $\text{JS} (p(X)\|q(X)) = \text{JS} (q(X)\|p(X)) \Rightarrow \text{JSD}$ is symmetric
- $0 \leq \text{JS} (p(X)\|q(X)) \leq 1 \Rightarrow \text{JSD}$ is non-negative and upper-bounded.
- $\text{JS} (p(X)\|q(X)) = 0 \Leftrightarrow p(x) = q(x) \forall x$
- JSD is not a distance metric.
 - It does not satisfy the triangle inequality.

Jensen-Shanon Divergence

Jensen-Shanon Divergence

$$\text{JS} (p(X)\|q(X)) \triangleq \frac{1}{2} \left(\text{KL} (p(X)\|m(X)) + \text{KL} (q(X)\|m(X)) \right)$$

- $\text{JS} (p(X)\|q(X)) = \text{JS} (q(X)\|p(X)) \Rightarrow \text{JSD}$ is symmetric
- $0 \leq \text{JS} (p(X)\|q(X)) \leq 1 \Rightarrow \text{JSD}$ is non-negative and upper-bounded.
- $\text{JS} (p(X)\|q(X)) = 0 \Leftrightarrow p(x) = q(x) \forall x$
- JSD is not a distance metric.
 - It does not satisfy the triangle inequality.

Square Root of JSD

The square root of JSD $\sqrt{\text{JS} (p(X)\|q(X))}$ is a distance metric.

Comparing Divergences

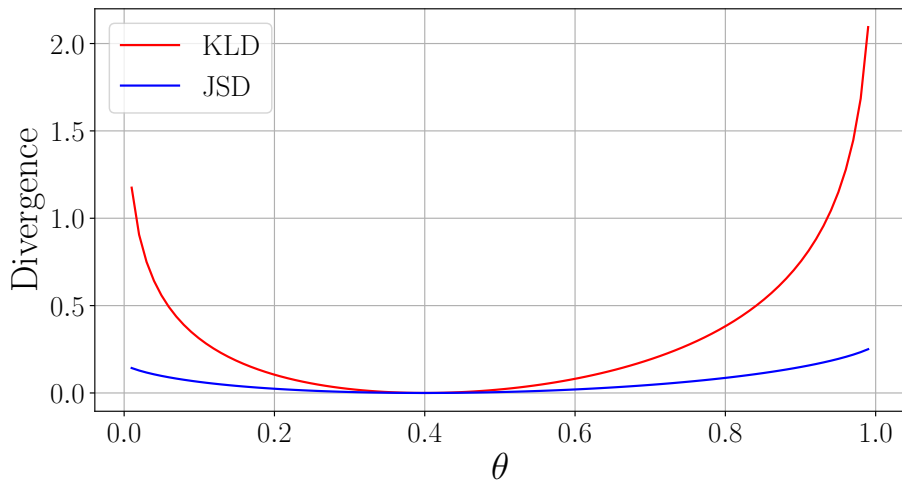


Figure: Distance between $p(X) = \text{Ber}(X|0.4)$ and $q(X) = \text{Ber}(X|\theta)$ as a function of θ

Section 4

Conclusions

Our Foundation

- Notation
- Probability and Statistics
 - Probability Mass/Density Function
 - expectation
 - Distance measurement between densities

List of Abbreviations

Complete	Abbreviation
Jensen-Shanon Divergenc	JSD
Kullback-Leibler divergence	KLD
Left Hand Side	LHS
Monte Carlo Estimation	MCE
Probability Density Function	PDF
Probability Mass Function	PMF
Right Hand Side	RHS

References I



John Tsitsiklis and Patrick Jaillet,

“Mit res.6-012 introduction to probability, spring 2018,” Spring 2018.