

Introduction to Machine Learning - CS725

Instructor: Prof. Ganesh Ramakrishnan

Overview of Probability Theory¹

¹Basic notes at <https://www.cse.iitb.ac.in/~cs725/notes/classNotes/misc/BasicProbAndStats.pdf> and advanced notes at <https://www.cse.iitb.ac.in/~cs725/notes/classNotes/misc/CaseStudyWithProbabilisticModels.pdf>

A review of probability theory

δ : indicator: $\delta(x,y) = 1$ if $x=y$ & $= 0$ o/w

Sample space(S): A sample space is defined as a set of all possible outcomes of an experiment. Example of an experiment would be a coin pair toss. In this case

$S = \{ \underline{HH}, \underline{HT}, \underline{TH}, \underline{TT} \}$. $|S| = 4 = 2 \times 2$

Event (E) : An event is defined as any subset of the sample space. Total number of distinct events possible is $2^{|S|}$, where $|S|$ is the number of elements in the sample space.

Random variable (X) : A random variable is a mapping (or function) from set of events to a set of real numbers.

Continuous random variable is defined thus

$X = \sum \delta(\text{coin1}, k) + \delta(\text{coin2}, k)$
 $X \in \{0, 1, 2\}$

$X : 2^S \rightarrow \mathbb{R}$

Range : value of X defines an event

On the other hand a discrete random variable maps events to a countable set (e.g. Natural Numbers)

$X : 2^S \rightarrow$ Countable set

Range

$\{E\} : E = \{HT, TH\}$
Coin = H
 $X = \delta(\text{coin1}, H) + \delta(\text{coin2}, H)$

$k=1$
 $T=0$

Axioms of Probability

- For every event E , $0 \leq \Pr(E) \leq 1$
- $\Pr(S) = 1 \rightarrow \Pr(X \in \text{Range}) = 1$
- If E_1, E_2, \dots, E_n is a set of pairwise disjoint events, then

If X is discrete
if X is continuous

$\Pr(X=v) \in [0,1]$
(prob mass fn)
pmf

$\Pr(X \in [v, v+\delta v]) \in [0,1]$
 $\underbrace{\quad}_{v \quad v+\delta v}$

$$\Pr\left(\bigcup_{i=1}^n E_i\right) = \sum_{i=1}^n \Pr(E_i) \rightarrow$$

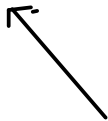
$$\Pr(X \in I) = \sum \Pr(X \in I_k)$$

$I_k \subseteq I$
(disjoint)

Eg: If X is continuous



Probability fns



Measures

Modular fns

(integrability etc)

$$f(\{e_1, \dots, e_k\}) = \sum_{i=1}^k f(\{e_i\})$$

Bayes' Theorem



$$A = \bigcup_i (A \cap B_i)$$

Let B_1, B_2, \dots, B_n be a set of mutually exclusive events that together form the sample space S . Let A be any event from the same sample space, such that $P(A) > 0$. Then,

$$Pr(A|B_i) = \frac{Pr(B_i \cap A)}{Pr(B_i)} \quad Pr(B_i/A) = \frac{Pr(B_i \cap A)}{Pr(A)}$$
$$Pr(B_i/A) = \frac{Pr(B_i \cap A)}{Pr(B_1 \cap A) + Pr(B_2 \cap A) + \dots + Pr(B_n \cap A)} \quad (1)$$

Using the relation $P(B_i \cap A) = P(B_i) \cdot P(A/B_i)$

$$Pr(B_i/A) = \frac{Pr(B_i) \cdot Pr(A/B_i)}{\sum_{j=1}^n Pr(B_j) \cdot Pr(A/B_j)} \quad (2)$$

$$\text{Basically: } Pr(B_i \cap A) = Pr(B_i/A) Pr(A) = Pr(A/B_i) Pr(B_i)$$

For random variables:

$$\begin{aligned} \Pr(X \in I_x, Y \in I_y) &= \Pr(X \in I_x | Y \in I_y) \Pr(Y \in I_y) \\ &= \Pr(Y \in I_y | X \in I_x) \Pr(X \in I_x) \end{aligned}$$

If Y is discrete

$$\Pr(X \in I_x) = \sum_{y \in I_y} \Pr(X \in I_x, Y \in \{y\})$$

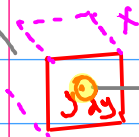
If Y is ct

$$\Pr(X \in I_x) = \int \Pr(X \in I_x, Y \in (y, y+dy)) dy$$

$(y, y+dy) \subseteq I_y$

can be imagined to be the circle

Why $\Pr(Y=y)=0$ for ct r.v. Y



if this pt has area > 0

then given there are infinite # pts, area of square $\rightarrow \infty$

Using Bayes' Theorem

$$D, \neg D \quad P, \neg P$$

$$P_r(P|D) = 0.99$$

A lab test is 99% effective in detecting a disease when in fact it is present. However, the test also yields a false positive for 0.5% of the healthy patients tested. If 1% of the population has that disease, then what is the probability that a person has the disease given that his/her test is positive?

$$P_r(D|P) = \frac{P_r(P|D)P_r(D)}{P_r(P)}$$

$$= \frac{P_r(P|D)P_r(D)}{P_r(P|D)P_r(D) + P_r(P|\neg D)P_r(\neg D)}$$

$$P_r(D) = 0.01$$

$$P_r(\neg D) = 0.99$$

$$P_r(D|P) = ?$$

$$P_r(P|\neg D) = 0.005$$

Independent Events

Two events E_1 and E_2 are called independent iff their probabilities satisfy

$$P(E_1, E_2) = P(E_1) \cdot P(E_2)$$

$$\left. \begin{aligned} &P(x \in I_x, y \in I_y) \\ &= P(x \in I_x) P(y \in I_y) \end{aligned} \right\} (3)$$

where $P(E_1, E_2)$ means $P(\underline{E_1 \cap E_2})$

In general, events belonging to a set are called as mutually independent iff, for every finite subset, E_1, \dots, E_n , of this set

$$Pr\left(\bigcap_{i=1}^n E_i\right) = \prod_{i=1}^n Pr(E_i) \quad (4)$$

3 heads
↓ in 5 tosses
→ Pr(H & S) in H5

E_1 & E_2 are independent

E_2 & E_3 are independent

↓ 2 tails in 5 tosses

} If $E_1 = E_3$ (disguised)
then E_1, E_2, E_3 are
NOT independent

Agenda

- Distribution function
- Probability Density/Mass Function
- Cumulative Distribution Function
- Statistical Measures
 - Expectation
 - Variance
 - Covariance
- Random Variables
 - Bernoulli Random Variable
 - Binomial Random Variable
 - Normal Random Variable
- Central Limit Theorem

for ct r.v.
slightly diff

→ So far

→ discrete
r.v. = dist

→ for points

→ for specific intervals
eg: $(-\infty, x]$

Uncertainty

- We are trying to build systems that understand and (possibly) interact with the real world
- We often can not prove something is true, but we can still ask how likely different outcomes are or ask for the most likely explanation

Eg: Cts v.v is Temp = 30°C more likely than Temp = 300°C

Probability theory is nothing but common sense reduced to calculation. — **Pierre Laplace, 1812**

We will restrict ourselves to a relatively informal discussion of probability theory

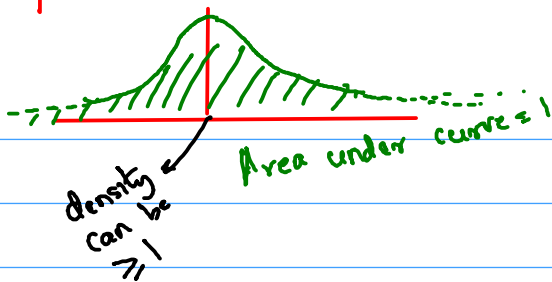
Notations

- A random variable X represents the outcome or the state of the world
- We will write $\Pr(X)$ to mean probability of event X ,
Probability($X=x$)
- **Sample space:** the space of all possible outcomes (may be discrete, continuous or mixed)

- $p(x)$ is the probability mass (density) function (for a pt)
 - Assigns a number to each point in sample space
 - Non-negative, sums (integrates) to 1 (if mass ≥ 0 & volume ≥ 0 \Rightarrow density ≥ 0)
 - Intuitively: how often does x occur, how much do we believe in x .

Eg: cb r.v Y st $p(y) = \frac{5}{2}$ in $[-\frac{1}{5}, \frac{1}{5}]$ & $\int_{-\infty}^{\infty} p(y) dy = 1$
 $= 0$ o/w

Eg: Fox temperature



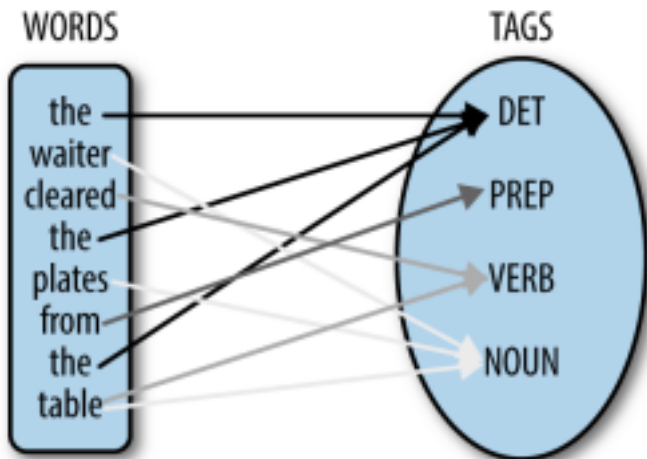
- Pr - Probability of an event in general
- F - Cumulative distribution function
- p - Probability distribution function (pdf) or probability mass function (pmf)
- pdf - pdf occurs in case of continuous random variable
- pmf - pmf occurs in case of discrete random variable

Example - Part of Speech

POS tagging is a problem of great importance in the field of Natural Language Processing, **NLP**

Input: A set of n-words

Output: POS tag for each word



Let p_k be probability of a word having pos type k .

A_k = prob of set containing pos type k .
words = m

$\Pr(A_k^c)$ = prob that set does NOT contain any pos type k

$$= \underbrace{(1-p_k)(1-p_k) \dots (1-p_k)}_{m \text{ times}} = (1-p_k)^m$$

$$\Pr(A_k) = 1 - (1-p_k)^m$$

$$\Pr(A_{\text{noun \& verb}}^c) = \Pr(\text{set contains either no noun or no verb})$$
$$= (1-p_n)^m + (1-p_v)^m - (1-p_v-p_n)^m$$

Assuming the picking of words is done independently, find probability that the set contains a 'noun' given that it contains a 'verb'.

Solution:

- Probability that a word is of part of speech type 'k' is p_k
- Let A_k be the probability that the set contains pos type 'k'

$$Pr(A_k) = 1 - (1 - p_k)^n$$

where $(1 - p_k)^n$ is that all 'n' words are not of pos of type 'k'.

$$Pr(A_{noun}/A_{verb}) = \frac{Pr(A_{noun} \cap A_{verb})}{Pr(A_{verb})}$$

$$Pr(A_{k1} \cap A_{k2}) = 1 - (1 - p_{k1})^n - (1 - p_{k2})^n + (1 - p_{k1} - p_{k2})^n$$

$$Pr(A_{noun}/A_{verb}) = \frac{1 - (1 - p_{noun})^n - (1 - p_{verb})^n + (1 - p_{noun} - p_{verb})^n}{1 - (1 - p_{verb})^n}$$

Distribution Functions

- **pmf**: It is a function that gives the probability that a discrete random variable is exactly equal to some value (Src: wiki)

$$p_X(a) = Pr(X = a)$$

- **pdf**: A probability density function of a continuous random variable is a function that describes the relative likelihood for this random variable to occur at a given point in the observation space (Src: Wiki)

$$Pr(X \in D) = \int_D p(x) dx$$

\rightarrow pdf ≥ 0

Cumulative Distribution Function

Case: Discrete Random Variable

$$F(a) = \Pr(X \leq a)$$

Not often used

Case: Continuous Random Variable

$$F(a) = \Pr(X \leq a) = \int_{-\infty}^a p(x) dx$$

Note: pdf for continuous distribution can be obtained by differentiating the cdf of that random variable:

$$p(a) = \left. \frac{dF(x)}{dx} \right|_{x=a}$$

By fundamental theorem of calculus.

Joint Distribution Function

- If $p(x,y)$ is a joint pdf i.e. for continuous case:

$$F(a, b) = Pr(X \leq a, Y \leq b) = \int_{-\infty}^b \int_{-\infty}^a p(x, y) dx dy$$

$$\underline{p(a, b)} = \frac{\partial^2 F(x, y)}{\partial x \partial y} \Big|_{a, b}$$

- For discrete case i.e. $p(x,y)$ is a joint pmf:

$$F(a, b) = \sum_{x \leq a} \sum_{y \leq b} p(x, y)$$

Marginalization

- Marginal probability is the unconditional probability $P(A)$ of the event A ; that is, the probability of A , regardless of whether event B did or did not occur.
- If B can be thought of as the event of a random variable X having a given outcome, the marginal probability of A can be obtained by summing (or integrating, more generally) the joint probabilities over all outcomes for X .
- For example, if there are two possible outcomes for X with corresponding events B and B' , this means that

$$P(A) = P(A \cap B) + P(A \cap B')$$

Discrete case: $P(X = a) = \sum_y p(a, y)$

Continuous case: $P_x(a) = \int_{-\infty}^{\infty} p(a, y) dy$

Example

Let X and Y are independent continuous random variables with same density functions

$$p(x) = \begin{cases} e^{-x} & \text{if } x > 0; \\ 0 & \text{otherwise.} \end{cases}$$

Find density $\frac{X}{Y}$.

$$F_{\frac{X}{Y}}(a) = \int_{-\infty}^{\infty} \int_{-\infty}^{ay} p(x,y) dx dy = \int_{-\infty}^{\infty} \int_0^{ay} p(x)p(y) dx dy$$

\downarrow new r.v.

$$Pr\left(\frac{X}{Y} \leq a\right) = Pr(x \leq ay) \quad (y = 2a \text{ \& } x \leq ay \text{ will still mean } x/Y \leq a)$$

$$\begin{aligned}
 F_{\frac{X}{Y}}(a) &= \Pr\left(\frac{X}{Y} \leq a\right) \\
 &= \int_0^\infty \int_0^{ya} p(x,y) dx dy \\
 &= \int_0^\infty \int_0^{ya} e^{-x} e^{-y} dx dy \\
 &= 1 - \frac{1}{a+1} \\
 &= \frac{a}{a+1}
 \end{aligned}$$

Note
 $p(x)$ & $p(y)$
 are defined to
 be non-zero
 only for $x > 0$
 $y > 0$

$$\begin{aligned}
 f_{\frac{X}{Y}}(a) &= \text{derivative of } F_{\frac{X}{Y}}(a) \text{ w.r.t } a \\
 &= \frac{1}{(a+1)^2} > 0
 \end{aligned}$$

Suppose X and Y are two random variable then we can define the conditional probability density of X given Y , denoted as $X|Y$

Discrete Case

$$p_X\left(\frac{x}{Y=y}\right) = P\left(\frac{X=x}{Y=y}\right) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

Continuous case

$$p_X\left(\frac{x}{Y=y}\right) = \frac{p_{X,Y}\left(\frac{x}{y}\right)}{p_Y(y)} = \frac{p_{X,Y}\left(\frac{x}{y}\right)}{\int_{-\infty}^{\infty} p(x,y) dx}$$

Joint Probability Distribution

- $Prob(X = x, Y = y)$
 - "Probability of $X=x$ and $Y=y$ "
 - $p(x, y)$

Conditional Probability Distribution

- $Prob(X = x|Y = y)$
 - "Probability of $X=x$ given $Y=y$ "
 - $p(x|y) = p(x, y)/p(y)$

Rules of Probability

- Sum Rule (marginalization/ summing out)

$$p(x) = \sum_y p(x, y)$$

$$p(x_1) = \sum_{x_2} \sum_{x_3} \dots \sum_{x_n} p(x_1, x_2, \dots, x_n)$$

- Product/Chain Rule

$$p(x, y) = p(y, x)p(x)$$

$$p(x_1, x_2, \dots, x_n) = p(x_1) p(x_2|x_1) \dots p(x_n|x_1, x_2, \dots, x_{n-1})$$

No independence of $x_1 \dots x_n$ was assumed... Independence will simplify further!

$$\begin{aligned} &= p(x_n|x_1, x_2, \dots, x_{n-1}) p(x_1 \dots x_{n-1}) \\ &= \dots \\ &= \dots \end{aligned}$$

$p(x_{n-1}|x_1 \dots x_{n-2})$
 $p(x_1 \dots x_{n-2})$
"
 $p(x_1, x_2)$

- One of the most important formulas in probability theory

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)} = \frac{p(y,x)p(x)}{\sum_x p(y|x)p(x)}$$

- This gives away a way of reversing conditional probabilities

Independence of Random Variables

- Two random variables are said to be **independent** iff their joint distribution factors

$$X \perp Y \iff p(x, y) = p(y|x)p(x) = p(x|y)p(y) = p(x)p(y)$$

- Two random variables are **conditionally independent** iff given a third they are independent after conditioning on the third variable

$$X \perp Y|Z \iff p(x, y|z) = p(y|x, z)p(x|z) = p(x|y, z)p(y|z) = p(x|z)p(y|z) \forall Z$$

Expectation

- **Discrete case:** Expectation is equivalent to probability weighted sums of possible values. If X is a discrete random variable

$$E(X) = \sum_i x_i Pr(x_i)$$

wt = prob
multiple values of X

If the random variable is a function of x , then

$$E(X) = \sum_i f(x_i) Pr(x_i)$$

- **Continuous case:** Expectation is equivalent to probability density weighted integral of possible values.

$$E(X) = \int_{-\infty}^{\infty} xp(x)dx$$

If the random variable is a function of x , then

$$E(X) = \int_{-\infty}^{\infty} f(x)p(x)dx$$

Properties of Expectation

① $E[X + Y] = E[X] + E[Y]$

Proof HW

② $E[(X - c)^2] \geq E[(X - \mu)^2]$

Proof HW

where $\mu = E[X]$

For any constant c and any random variable X

Expected squared deviation which is least from expected value $\mu = E(x)$

③ $E[cX] = cE[X]$

Proof HW

Variance

For any random variable X , **variance** is defined as follows:

$$\text{Var}[X] = E[(X - \mu)^2]$$

→ Expected squared deviation from $\mu = E(X)$

$$\Rightarrow \text{Var}[X] = E[X^2] - 2\mu E[X] + \mu^2$$

$$\Rightarrow \text{Var}[X] = \underline{E[X^2] - (E[X])^2}$$

$$\text{Var}[\alpha X + \beta] = \alpha^2 \text{Var}[X]$$

Covariance

For random variables X and Y , **covariance** is defined as:

$$\text{Cov}[X, Y] = E[(X - E(X))(Y - E(Y))] = E[XY] - E[X]E[Y]$$

If X and Y are independent then their covariance is 0, since in that case

$$E[XY] = E[X]E[Y]$$

Imp for course

Note: However, covariance being 0 does not necessarily imply that the variables are independent.

*↓
If X & Y are ind $\text{cov}(X, Y) = 0$
But converse does not hold !!*

Properties:

- 1 $\text{Cov}[X + Z, Y] = \text{Cov}[X, Y] + \text{Cov}[Z, Y]$
- 2 $\text{Cov}[\sum_i X_i, Y] = \sum_i \text{Cov}[X_i, Y]$
- 3 $\text{Cov}[X, X] = \text{Var}[X]$

Q: Why?

Chebyshev's Inequality

Chebyshev's inequality states that if X is any random variable with mean μ and variance σ then $\forall k > 0$

$$Pr[|X - \mu| \geq k] \leq \frac{\sigma^2}{k^2}$$

Implications:

- If n tends to infinity, then the data mean tends to converge to μ , giving rise to the *weak law of large numbers*.
- If X_i are independent and identically distributed random variables,

$Pr[|\frac{X_1+X_2+\dots+X_n}{n} - \mu| \geq k]$ tends to 0 as n tends to ∞

Important Random Variables

Bernoulli Random Variable: It is a *discrete* random variable taking values 0,1

Say, $Pr[X_i = 0] = 1 - q$ where $q \in [0, 1]$

Then $Pr[X_i = 1] = q$

- $E[X] = (1 - q) * 0 + q * 1 = q$
- $Var[X] = q - q^2 = q(1 - q)$

Note: It represents the probability of success in a random event. For example: Coin toss experiment can be modeled as a *Bernoulli random variable* with $Pr[Head] = Pr[X_i = 1] = q$

Binomial Random Variable It is a *discrete* variable where the distribution is of number of 1's in a series of n experiments with $\{0,1\}$ value, with the probability that the outcome of a particular experiment is 1 being q .

A binomial distribution is the distribution of n -times repeated bernoulli trials.

① $Pr[X = k] = \binom{n}{k} q^k (1 - q)^{n-k}$

② $E[X] = \sum_i E[Y_i]$ where Y_i is a bernoulli random variable

$$E[X] = nq$$

③ $Var[X] = \sum_i Var[Y_i]$ (since Y_i 's are independent)

$$Var[X] = nq(1 - q)$$

Example:

An example of Binomial distribution is the distribution of number of heads when a coin is tossed n times.

Normal (Gaussian) Distribution

- It is a continuous distribution

$$p(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$

- μ is the mean
 - σ^2 is the variance
-
- Exercise: Verify there mean and variance. For e.g.

$$E(X) = \mu = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} dx$$

- Multivariate Gaussian

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

- x is now a vector
- μ is the mean vector
- Σ is the co-variance matrix

Properties of Normal Distribution

- All marginals of a Gaussian are again Gaussian
- Any conditional of a Gaussian is Gaussian
- The product of two Gaussians is again Gaussian
- Even the sum of two independent Gaussian RVs is a Gaussian

Note: Many of the standard distributions belong to the family of **exponential distributions**

- Bernoulli, binomial/multinomial, Poisson, Normal (Gaussian), beta/Dirichlet ...
- Share many important properties - e.g. They have a conjugate prior. (We will discuss this in next lecture)

Central Limit Theorem

If X_1, X_2, \dots, X_m is a sequence of i.i.d. random variables each having mean μ and variance σ^2

Then for large m , $X_1 + X_2 + \dots + X_m$ is approximately normally distributed with mean $m\mu$ and variance $m\sigma^2$

If $X \sim N(\mu, \sigma^2)$

Then $P[x] = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

It can be shown by CLT

- $\frac{X_1 + X_2 + \dots + X_n - n\mu}{\sigma\sqrt{n}} \sim N(0, 1)$
- Sample Mean: $\hat{\mu} \sim N(\mu, \frac{\sigma^2}{m})$