

# CEE 697M: Big Data and Machine Learning for Engineers

## Lecture 1b: Probability

**Jimi Oke**

UMassAmherst  

---

College of Engineering

Feb 10, 2023

# Outline

- ① Random variables
- ② Univariate models
- ③ Multivariate models
- ④ Outlook

# Random variables

A random variable is a function that uniquely maps events in a sample space to the set of real numbers.

A random variable  $X$  may be:

- *Discrete*
- *Continuous*
- *Mixed* (probability defined over both discrete and range of continuous values)

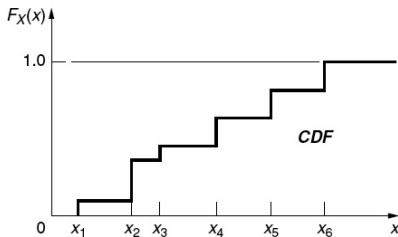
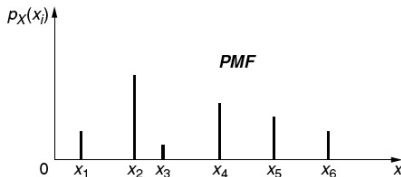
# Probability mass function (PMF)

The PMF is given by

$$p_X(x_i) \equiv P(X = x_i) \quad \forall x \quad (1)$$

CDF of discrete random variable

$$\begin{aligned} F_X(x) &= \sum_{x_i \leq x} P(X = x_i) \\ &= \sum_{x_i \leq x} p_X(x_i) \end{aligned}$$



The probability masses in a PMF sum up to 1.

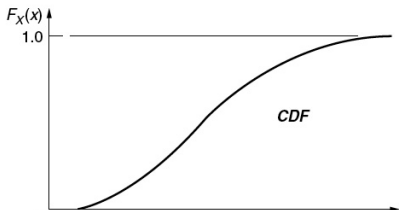
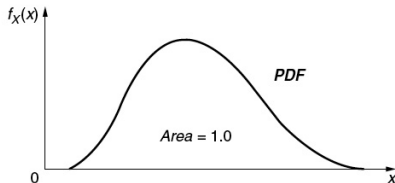
# Probability density function (PDF)

The PDF is denoted  $f_X(x)$  such that the probability of  $X$  in the interval  $(a, b]$  is:

$$P(a < X \leq b) = \int_a^b f_X(x) dx \quad (2)$$

## CDF of continuous random variable

$$\begin{aligned} F_X(x) &= P(X \leq x) \\ &= \int_{-\infty}^x f_X(\tau) d\tau \end{aligned}$$



It follows that the PDF is the derivative of the CDF:

The total area under a PDF is 1.

# Central values

These include the mean, median and mode.

- Mean: weighted average (by probability of occurrence) or expected value

$$\mathbb{E}(X) = \mu_X = \sum_i x_i p_X(x_i) \quad \text{discrete case} \quad (4)$$

$$\mathbb{E}(X) = \mu_X = \int_{-\infty}^{\infty} x f_X(x) dx \quad \text{continuous case} \quad (5)$$

## Generalized expectation

The mathematical expectation can be defined for a function  $g$  of random variable  $X$ :

$$\mathbb{E}[g(X)] = \sum_i g(x_i) p_X(x_i) \quad \text{discrete case} \quad (6)$$

$$\mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx \quad \text{continuous case} \quad (7)$$

# Measures of dispersion

## Variance

In discrete case:

$$\mathbb{V}(X) = \sum_i (x_i - \mu_X)^2 p_X(x_i) \quad (8)$$

In continuous case:

$$\mathbb{V}(X) = \int_{-\infty}^{\infty} (x - \mu_X)^2 f_X(x) dx \quad (9)$$

Expanding both equations results in:

$$\mathbb{V}(X) = \mathbb{E}(X^2) - \mu_X^2 = \mathbb{E}(X^2) - \mathbb{E}(X)^2 \quad (10)$$

# Measures of dispersion (cont.)

## Standard deviation

The standard deviation is convenient as it has the same unit as the random variable:

$$\sigma_X = \sqrt{\mathbb{V}(X)} \quad (11)$$

## Coefficient of variation

The COV gives the deviation relative to the mean. It is unitless.

$$\delta_X = \frac{\sigma_X}{\mu_X} \quad (12)$$



# Mean of a linear function

For a continuous random variable  $X$ , the mean is defined as

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} xf_X(x)dx \quad (13)$$

Now, given that  $Z = aX + bY$ , then the mean of  $Z$  is

$$\begin{aligned} \mathbb{E}(Z) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (aX + bY)f_{X,Y}dxdy \\ &= a \int_{-\infty}^{\infty} xf_X(x)dx + b \int_{-\infty}^{\infty} yf_Y(y)dy \\ &= a\mathbb{E}(X) + b\mathbb{E}(Y) \end{aligned}$$

# Variance of a linear function

We also recall the variance of an r.v.  $X$ :

$$\mathbb{V}(X) = \mathbb{E}[(X - \mu_X)^2] \quad (14)$$

Thus, for  $Z = aX + bY$ :

$$\begin{aligned} \mathbb{V}(Z) &= \mathbb{E}[((aX + bY) - (a\mu_X + b\mu_Y))^2] \\ &= \mathbb{E}[(a(X - \mu_X) + b(Y - \mu_Y))^2] \\ &= \mathbb{E}[a^2(X - \mu_X)^2 + 2ab(X - \mu_X)(Y - \mu_Y) + b^2(Y - \mu_Y)^2] \\ &= a^2\mathbb{V}(X) + b^2\mathbb{V}(Y) + 2ab\text{Cov}(X, Y) \end{aligned}$$

# Moments

The  $m$ -th order moment of a distribution is given by:

$$\mathbb{E}(X^m) = \begin{cases} \sum_i x_i^m \cdot p_X(x_i) & \text{(discrete)} \\ \int x^m \cdot f_X(x) dx & \text{(continuous)} \end{cases} \quad (15)$$

- $m$ -th central moment:  $\mathbb{E}[(X - \mu_X)^m]$
- Normalized  $m$ -th central moment:  $\left( \frac{\mathbb{E}[(X - \mu_X)^m]}{\sigma^m} \right)$

## Examples

- **Mean:** first moment,  $\mathbb{E}(X)$
- **Variance:** second central moment,  $\mathbb{E}[(X - \mu_X)^2]$
- **Skewness:** normalized third central moment,  $\left( \frac{\mathbb{E}[(X - \mu_X)^3]}{\sigma^3} \right)$

# Bernoulli distribution

Let  $X$  be an event with only two outcomes  $\{1,0\}$ . And let the probability of the event be given by:

$$p(X) = \theta, \quad 0 \leq \theta \leq 1$$

And  $p(X = 1) = \theta$  and  $p(X = 0) = 1 - \theta$ .  $X$  is said to be Bernoulli distributed:

$$X \sim \text{Ber}(\theta) \tag{16}$$

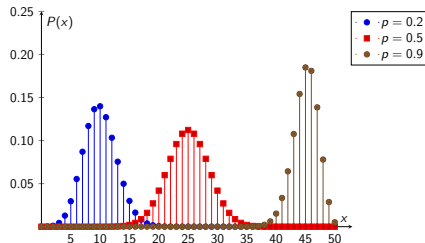
The PMF is then given by:

$$\text{Ber}(x|\theta) := \theta^x(1 - \theta)^{1-x} \tag{17}$$

# Binomial distribution

Given a Bernoulli sequence with  $X$  random number of occurrences of an event,  $N$  trials and  $\theta$  the probability of occurrence of each event:

- $X \sim \text{Bin}(N, \theta)$
- PMF:  $P(X = x) := \text{Bin}(x|N, \theta) := \binom{N}{x} p^x (1 - \theta)^{N-x}$ ,  $x = 0, 1, 2, \dots, N$
- CDF:  $F_X(x) = P(X \leq x) = \sum_{k=0}^x \binom{N}{k} \theta^k (1 - \theta)^{N-k}$
- Mean:  $\mathbb{E}(X) = N\theta$
- Variance:  $\mathbb{V}(X) = N\theta(1 - \theta)$



# Bernoulli, binomial, categorical and multinomial

- The Bernoulli distribution is a special case of the binomial distribution with  $N = 1$
- The categorical distribution is generalization of the Bernoulli to more than two outcomes for a single trial (e.g. set of labels  $x \in \{1, \dots, C\}$ ,  $C > 2$ ):

$$\text{Cat}(\mathbf{x}|\boldsymbol{\theta}) := \prod_{c=1}^C \theta_c^{x_c} \quad (18)$$

where  $\mathbf{x}$  is a one-hot vector (e.g.  $(1,0,0,0)$  for class 1 of four classes)

- The multinomial distribution generalizes the categorical distribution for multiple trials:

$$\mathcal{M}(\mathbf{x}|N, \boldsymbol{\theta}) := \binom{N}{N_1 \dots N_C} \prod_{c=1}^C \theta_c^{N_c} \quad (19)$$

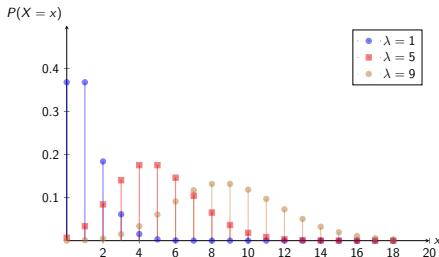
# Poisson distribution

- The Poisson distribution is used to model the probability that a number of independent events occur within a fixed time interval (or within a finite space)
- Such events are described as Poisson processes
- The PMF of a Poisson random variable with **rate parameter  $\lambda$**  is given by:

$$P(X = x) := \text{Poiss}(x|\lambda) := \frac{\lambda^x}{x!} e^{-\lambda}, \quad x \geq 0 \quad (20)$$

- The mean and variance of a Poisson random variable are equal:

$$\mathbb{E}(X) = \mathbb{V}(X) = \lambda \quad (21)$$



# Gaussian distribution

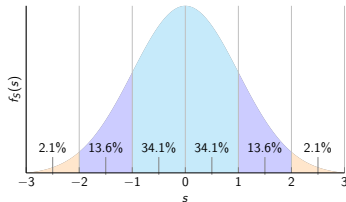
The PDF of a Gaussian (normal) distribution  $X \sim \mathcal{N}(\mu, \sigma^2)$  is given by:

$$\mathcal{N}(x|\mu, \sigma^2) := \frac{1}{\sigma\sqrt{2\pi}} \exp \left[ -\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 \right] \quad (22)$$

where  $\mu$  is the mean and  $\sigma^2$  is the variance.

$$P(a < X \leq b) = \frac{1}{\sigma\sqrt{2\pi}} \int_a^b e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2} dx = \Phi \left( \frac{b-\mu}{\sigma} \right) - \Phi \left( \frac{a-\mu}{\sigma} \right) \quad (23)$$

where  $\Phi$  is the CDF of the standard normal distribution ( $N(0, 1)$ ).

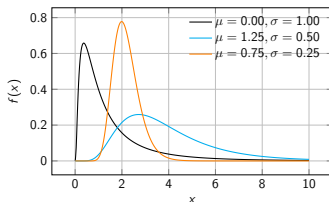




# Lognormal distribution

A random variable  $X$  that is lognormally distributed with the parameters  $\mu$  and  $\sigma^2$  (denoted  $X \sim \mathcal{LN}(\mu, \sigma^2)$ ) has the PDF:

$$\mathcal{LN}(x|\mu, \sigma^2) = \frac{1}{(\sigma x)\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{\ln(x) - \mu}{\sigma}\right)^2\right] \quad x \geq 0 \quad (24)$$



CDF:  $F_X(x) = P(X \leq x) = \Phi((\ln(x) - \mu)/\sigma)$

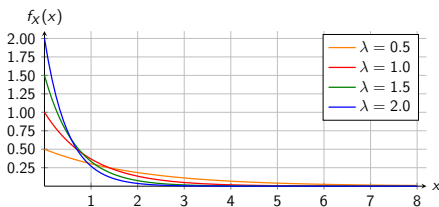
Mean:  $\mathbb{E}(X) = e^{\mu + \frac{1}{2}\sigma^2}$

Variance:  $\mathbb{V}(X) = (e^{\sigma^2} - 1)e^{2\mu + \sigma^2}$

# Exponential distribution

A random variable  $X$  exponentially distributed with parameter  $\lambda$  has the PDF:

$$\text{Exp}(x|\lambda) = \lambda e^{-\lambda x} \quad x > 0 \quad (25)$$



CDF:

$$F_X(x) = P(X \leq x) = 1 - e^{-\lambda x}, \quad x > 0 \quad (26)$$

Mean:

$$\mathbb{E}(X) = 1/\lambda \quad (27)$$

Variance:

$$\mathbb{V}(X) = 1/\lambda^2 \quad (28)$$

# Covariance and correlation

Recall that the variance of an r.v.  $X$  is given by:

$$\mathbb{V}(X) = \mathbb{E}[(X - \mu_X)^2] = \mathbb{E}(X^2) - \mathbb{E}(X)^2 \quad (29)$$

Then given two r.v.'s  $X$  and  $Y$ , the *covariance* measures the strength of the linear relationship between them.

## Covariance

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)] = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) \quad (30)$$

## Correlation coefficient

This is the normalized covariance

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} := \frac{\text{Cov}(X, Y)}{\sqrt{\mathbb{V}(X)\mathbb{V}(Y)}} \quad (31)$$

# Useful results

If  $\mathbf{x}$  is a  $D$ -dimensional r.v., then its **covariance matrix** is defined as:

$$\text{Cov}[\mathbf{x}] := \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] := \Sigma \quad (32)$$

This implies that:

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \Sigma + \mu\mu^T \quad (33)$$

The covariance of a linear transformation is defined as:

$$\text{Cov}[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\text{Cov}[\mathbf{x}]\mathbf{A}^T \quad (34)$$

The cross-covariance between two random vectors  $\mathbf{x}$  and  $\mathbf{y}$  is:

$$\text{Cov}[\mathbf{x}\mathbf{y}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{y} - \mathbb{E}[\mathbf{y}])^T] \quad (35)$$

# Uncorrelated does not imply independent

The following scatterplots indicate pairs of variables with various correlation values.

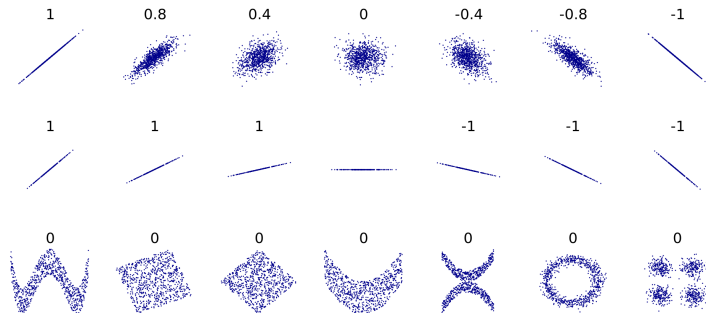
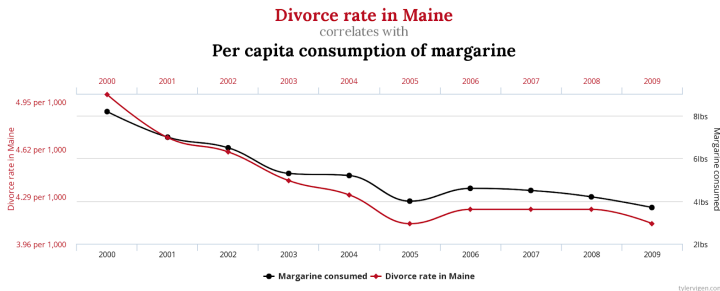


Figure: Source:

[https://en.wikipedia.org/wiki/Pearson\\_correlation\\_coefficient](https://en.wikipedia.org/wiki/Pearson_correlation_coefficient)

Note that some with 0 correlation still have functional dependence (but non-linear).

# Correlation does not imply causation



**Figure:** Source: <https://sitn.hms.harvard.edu/flash/2021/when-correlation-does-not-imply-causation-why-your-gut-microbes-may-not-yet-be-a-silver-bullet-to-all-your-problems/>

Visit <https://www.tylervigen.com/spurious-correlations> for more examples.

# Simpson's paradox

Trends appearing in different groups may be reversed or disappear when groups are combined

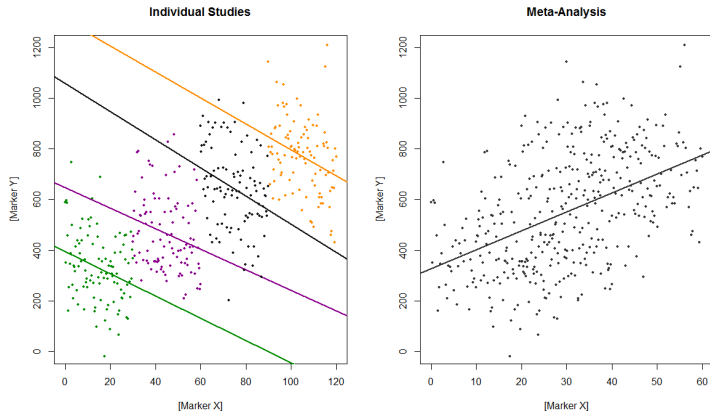


Figure: Source:

[https://rinterested.github.io/statistics/simpsons\\_paradox.html](https://rinterested.github.io/statistics/simpsons_paradox.html)

# Joint distributions

Given two random variables  $X$  and  $Y$ :

## Discrete case

The joint PMF is:

$$p_{X,Y}(x_i, y_j) = P(X = x_i, Y = y_j) \quad (36)$$

The CDF is:

$$F_{X,Y}(x, y) = \sum_{x_i \leq x} \sum_{y_j \leq y} p_{X,Y}(x_i, y_j) \quad (37)$$

## Continuous case

The joint probability is given by:

$$P(a < X \leq b, c < Y \leq d) = \int_a^b \int_c^d f_{X,Y}(x, y) dy dx \quad (38)$$



# Conditional distributions of continuous random variables

Recall the definition of conditional probability (multiplication rule):

$$P(A|B) = \frac{P(AB)}{P(B)} \quad (39)$$

$$P(AB) = P(A|B)P(B) = P(B|A)P(A) \quad (40)$$

Similarly, for two continuous r.v.'s, the conditional PDF of  $X$  given  $Y$  is:

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} \quad (41)$$

## Joint PDF and CDF of two variables

The joint PDF is given by:

$$f_{X,Y}(x,y) = f_{X|Y}(x|y)f_Y(y) = f_{Y|X}(y|x)f_X(x) \quad (42)$$

While the joint CDF is given by:

$$F_{X,Y}(a,b) = P(X \leq a, Y \leq b) = \int_{-\infty}^a \int_{-\infty}^b f_{X,Y}(x,y) dy dx \quad (43)$$

# Marginal distributions of continuous random variables

Recall the theorem of total probability:

$$P(A) = \sum_{i=1}^n P(A|E_i)P(E_i) \quad (44)$$

Similarly, the marginal PDFs from a joint distribution of two continuous r.v.'s  $X$  and  $Y$  is given as:

$$f_X(x) = \int_{-\infty}^{\infty} f_{X|Y}(x|y)f_Y(y)dy = \int_{-\infty}^{\infty} f_{X,Y}(x,y)dy \quad (45)$$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{Y|X}(y|x)f_X(x)dx = \int_{-\infty}^{\infty} f_{X,Y}(x,y)dx \quad (46)$$

# Multivariate normal distribution (MVN)

The MVN PDF is given by:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) := \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (47)$$

where:

- $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}] \in \mathbb{R}^D$  is the mean vector
- $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{x}]$  is the  $D \times D$  covariance matrix:

$$\text{Cov}[\mathbf{x}] := \mathbb{E} [(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] \quad (48)$$

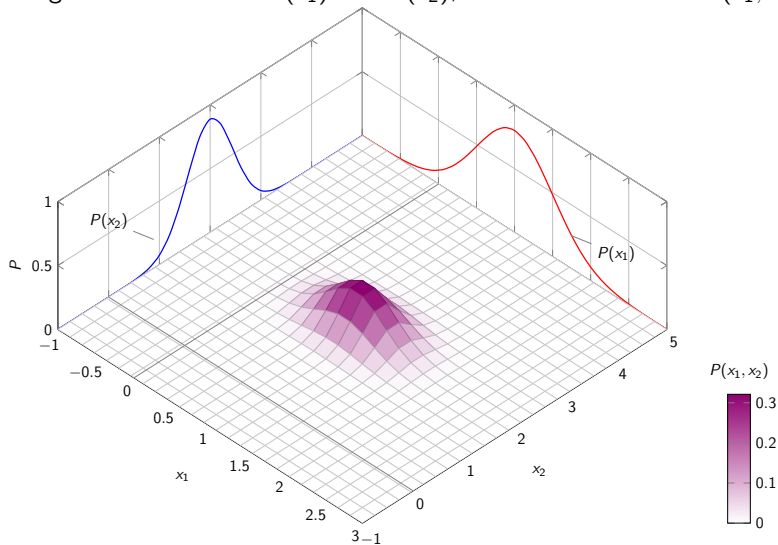
In 2D:

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_1^2 & \sigma_{12}^2 \\ \sigma_{21}^2 & \sigma_2^2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{pmatrix} \quad (49)$$

where  $\rho$  is the correlation coefficient.

# Bivariate MVN

Marginal distributions:  $P(x_1)$  and  $P(x_2)$ ; Joint distribution:  $P(x_1, x_2)$ .



# Reading

- PMLI 1, 2, 3
- PMLCE 1, 3, 4