

Review of Probability

David I. Inouye



Why probability? Probability is useful for handling uncertainty

- Inherent stochasticity
 - Quantum mechanics
 - Card games
- Incomplete observability
 - “Let’s Make a Deal” game show of three doors (called “Monty Hall” problem)
- Incomplete modeling
 - Discretization of space for object locations



Why probability? Sometimes more practical than deterministic

- “Most birds fly”
- “Birds fly, except for very young birds that have not yet learned to fly, sick or injured birds that have lost the ability to fly, flightless species of birds including the cassowary, ostrich and kiwi...”
 - (Example from Deep Learning, Goodfellow et al., 2016, Ch. 3)



Why probability? An extension of formal logic rules

- Original AI systems based on formal logic and reasoning
 - Chess
 - TurboTax
- Many AI applications based on deterministic logic were too brittle and failed often
 - Traditional linguistic approaches to natural language processing
- Modern AI systems almost always rooted in probability
 - Computer vision
 - Speech recognition
 - Natural language processing



How are these statements similar or different?

- A boardgame player: “The probability of getting a heads when flipping a fair coin is 50%.”
- The weather forecaster: “The probability of rain tomorrow is 50%.”
- Your doctor after examining your symptoms: “The probability of you having the flu is 50%.”



Frequentist and Bayesian interpretations lead to the same set of axioms

- Frequentist
 - Related to rates that events occur under repeated experimentation
- Bayesian interpretation
 - “Degree of belief”
- Pragmatic interpretation
 - They lead to the same math and are useful in similar circumstances
 - Use whichever interpretation is most useful



Big Picture Overview

- Introduction to probability
 - Motivation
 - Interpretation
- Random variables
 - Outcomes and events
- Probability distributions
 - PMF, PDF, CDF
- Multivariate distributions
 - Joint, marginal and conditional distributions
 - Chain rule and Bayes rule
 - Independence
- Expectations of random variables
 - Linearity of expectations
 - Covariance
 - Empirical expectations
- Multivariate Gaussian/Normal distribution
 - Marginal and conditional distributions
 - Affine/linear transformations



A random variable maps outcomes/events of a random/uncertain process to numbers

- Flipping a coin
 - Outcomes: {"Heads", "Tails"}
 - Possible random variable: "Heads" \rightarrow 0, "Tails" \rightarrow 1
- Flipping two coins
 - Outcomes: {(H,H), (H,T), (T, H), (T, T)}
 - Possible random variables: ## heads, ## tails, same, different
- Flipping coins until you get one tails
 - Outcomes: ?
 - Random variables: ?



A random variable maps outcomes to numbers:

Defining a random variable is the first step

- Random Tweet
 - Outcomes: ?
 - Random variables: ?
- Random Instagram image
 - Outcomes: ?
 - Random variables: ?



Random variables can be discrete or continuous

- **Discrete**

- Values are in some finite set or countably infinite set
- $\{-1, 1\}$, $\{5, 10, -20, 3\}$, $\{0, 1, 2, \dots\}$, \mathbb{Z}

- **Continuous**

- Values associated with intervals of \mathbb{R}
- $[0, 1]$, $[-1, 1]$, $[0.5, 1] \cup [-1, 0.5]$, $\mathbb{R}_+ \equiv [0, \infty)$

Note: Random variables by themselves do not provide any probability information.



An event is a set of possible outcomes

- For discrete RV such as $X \in \{0, 1, 2, \dots\}$, then events could be:
 - $E = \{0, 5, 1\}$
 - $E = \{0, 2, 4, 6, \dots\}$ (i.e., all even numbers)
- For continuous random variables $X \in \mathbb{R}$, events are sets of the real numbers:
 - $E = [0, 0.5)$
 - $E = [4, 5] \cup [8, 9]$

Big Picture Overview

- Introduction to probability
 - Motivation
 - Interpretation
- Random variables
 - Outcomes and events
- Probability distributions
 - PMF, PDF, CDF
- Multivariate distributions
 - Joint, marginal and conditional distributions
 - Chain rule and Bayes rule
 - Independence
- Expectations of random variables
 - Linearity of expectations
 - Covariance
 - Empirical expectations
- Multivariate Gaussian/Normal distribution
 - Marginal and conditional distributions
 - Affine/linear transformations



Probability distributions attach probabilities to all possible events of a random variable

- **Probability mass function (PMF)** is used for discrete random variables.
- A PMF P for random variable X that satisfies the following:
 1. Domain of P must include all possible states of X
 2. Unit domain: $\forall x \in X, 0 \leq P(x) \leq 1$
 3. Sum to 1: $\sum_{x \in X} P(x) = 1$



Probability distributions attach probabilities to all possible events of a random variable

- **Probability density function (PDF)** is used for continuous random variables
- A PDF p for random variable X that satisfies the following:
 1. Domain of p must include all possible states of X
 2. Non-negative: $\forall x \in X, p(x) \geq 0$
 3. Integrate to 1: $\int_X p(x) dx = 1$
- $p(x)$ is NOT a probability, rather **integrating** the PDF gives probabilities over **sets**

Are the following functions valid PDFs? Why?

- Suppose $X \in (0, 1)$ (note: 0 is not included)
- For all $x \in (0, 0.5)$, $p(x) = 2$ and for all $x \notin (0, 0.5)$, $p(x) = 0$
- $p(x) = 3x^2$
- $p(x) = -\log x$



Integrate PDF to get probabilities that random variable lies within a set (usually a range)

- The probability that X is less than q

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

- The probability that X lies between a and b

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

- The probability that X lies between (a and b) or between (c and d) where $b < c$

$$\Pr(a \leq X \leq b \text{ or } c \leq X \leq d) = \int_a^b p(x) dx + \int_c^d p(x) dx$$



Cumulative distribution function (CDF) is the integral of the PDF from the left up to query point q

- The CDF is the probability that X is less than q :

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

- What does $F(\infty)$ equal?
- The probability between a and b can be written as:

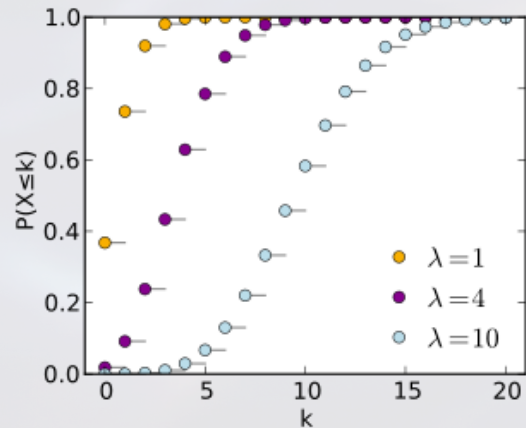
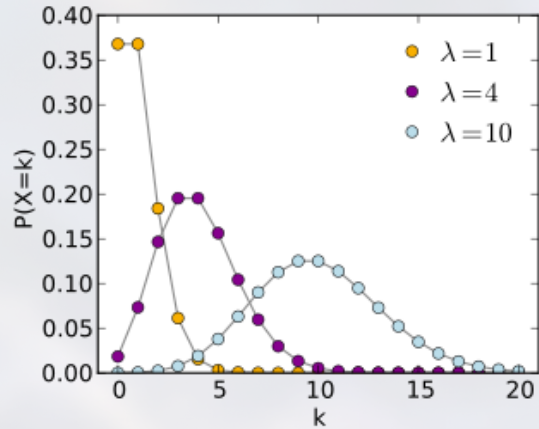
$$\Pr(a < X \leq b) = F(b) - F(a)$$

- The PDF is the derivative of the CDF:

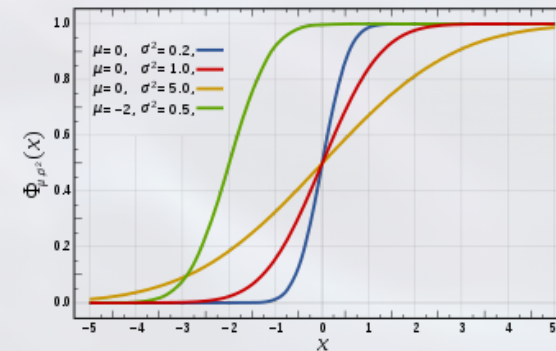
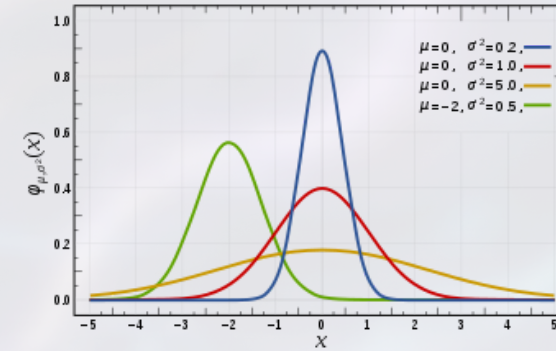
$$p(x) = \frac{dF(x)}{dx}$$

Examples of PMF/PDF and corresponding CDF

Discrete PMF/CDF



Continuous PDF/CDF



Notation: Tilde used to specify distribution of random variable (\sim in LaTeX)

- $X \sim \mathcal{N}(\mu = 0, \sigma = 1)$
 - “Random variable X is **distributed** as a normal distribution with mean of zero and standard deviation of 1.”
- $X \sim \text{Uniform}(\alpha, \beta)$
 - “Random variable X is **distributed** as a uniform distribution with parameters α and β (parameters may be unknown).”
- $X \sim P(x)$ or $X \sim \mathbb{P}(x)$
 - “Random variable X is **distributed** as the distribution represented by PMF/PDF $P(x)$ or $\mathbb{P}(x)$.”



Big Picture Overview

- Introduction to probability
 - Motivation
 - Interpretation
- Random variables
 - Outcomes and events
- Probability distributions
 - PMF, PDF, CDF
- Multivariate distributions
 - Joint, marginal and conditional distributions
 - Chain rule and Bayes rule
 - Independence
- Expectations of random variables
 - Linearity of expectations
 - Covariance
 - Empirical expectations
- Multivariate Gaussian/Normal distribution
 - Marginal and conditional distributions
 - Affine/linear transformations



Joint distribution of multiple variables

- The joint PDF/PMF is a function of two or more random variables (or a random vector).
- The joint PDF/PMF can be written as:

$$p(x, y), \quad p(x_1, x_2), \quad p(\mathbf{x})$$

- If $X \in [-1, 1]$ and $Y \in [-1, 1]$, is the following a valid PDF?

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

- If $X \in [0, 1]$ and $Y \in [0, 1]$, is the following a valid PDF?

$$p(x, y) = 4xy$$



Marginal distribution is sum/integral over other variables

- **Example:** Height and weight — “What is the distribution of height regardless of weight?”
- Given a joint distribution $P(x, y)$, the marginals are:

$$P(x) = \sum_{y \in \mathcal{Y}} P(x, y) \quad \text{and} \quad P(y) = \sum_{x \in \mathcal{X}} P(x, y)$$

- Given a joint distribution $p(x, y)$, the marginals are:

$$p(x) = \int_{y \in \mathcal{Y}} p(x, y) dy \quad \text{and} \quad p(y) = \int_{x \in \mathcal{X}} p(x, y) dx$$

- **Example:**

$$P(x, y) = \begin{bmatrix} & y = 1 & y = 0 \\ x = 0 & 0.1 & 0.3 \\ x = 1 & 0.4 & 0.2 \end{bmatrix}$$

- **Example:**

$$p(x, y) = 4xy$$



Conditional distribution is the distribution given some other event

- What is the distribution of weight given that a person is x inches tall?
- The conditional density is the joint PDF/PMF renormalized by the marginal density of the event:

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

- **Example:**

$$P(x, y) = \begin{bmatrix} & y = 1 & y = 0 \\ x = 0 & 0.1 & 0.3 \\ x = 1 & 0.4 & 0.2 \end{bmatrix}$$

- **Example:**

$$p(x, y) = 4xy$$



Note: Conditional and marginal distributions exist for **any set of variables**

- Suppose $p(\mathbf{x}) = p(x_1, x_2, x_3, x_4)$

$$p(x_1, x_3) = \int_{x_2, x_4} p(\mathbf{x}) dx_2 dx_4$$

$$p(x_1, x_2 \mid x_3) = \frac{p(x_1, x_2, x_3)}{p(x_3)} = \frac{\int_{x_4} p(\mathbf{x}) dx_4}{\int_{x_1, x_2, x_4} p(\mathbf{x}) dx_1 dx_2 dx_4}$$



Chain rule (or product rule) of probability

- The joint distribution can be written as a product of conditional PDFs/PMFs:

$$p(x_1, x_2) = p(x_1) p(x_2 | x_1)$$

$$p(x_1, x_2, x_3) = p(x_1) p(x_2 | x_1) p(x_3 | x_1, x_2)$$

- This can be written as:

$$p(x_1, x_2, \dots, x_d) = \prod_{i=1}^d p(x_i | x_1, \dots, x_{i-1})$$

- Consequence (order doesn't matter):

$$p(x) p(y | x) = p(y) p(x | y)$$



Bayes rule: Enables conversion between one conditional and the other (they are different)

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

- Can be derived from conditional chain rule.
- When are $p(x | y)$ and $p(y | x)$ equal?



Independence means that one variable is not affected by the other variable

- **Example:** Flip two coins, X and Y are 0 or 1.
- **Counterexample:** Roll a die for number X ; then flip that number of coins and count the number of heads Y .

- Formally, the PDF/PMF can be written as a product of functions that only involve x or y (but not both):

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

- Usually, these are the marginal densities:

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

- Equivalent definition:

$$p(x | y) = p(x) \quad \text{and} \quad p(y | x) = p(y)$$



Big Picture Overview

- Introduction to probability
 - Motivation
 - Interpretation
- Random variables
 - Outcomes and events
- Probability distributions
 - PMF, PDF, CDF
- Multivariate distributions
 - Joint, marginal and conditional distributions
 - Chain rule and Bayes rule
 - Independence
- Expectations of random variables
 - Linearity of expectations
 - Covariance
 - Empirical expectations
- Multivariate Gaussian/Normal distribution
 - Marginal and conditional distributions
 - Affine/linear transformations



An **expectation** (or **expected value**) of a function of a random variable is the average or mean value with respect to its distribution

- Formal definitions

$$\mathbb{E}_{X \sim P(x)}[f(x)] \equiv \sum_{x \in X} f(x)P(x)$$

$$\mathbb{E}_{X \sim p(x)}[f(x)] \equiv \int_{x \in X} f(x)p(x) dx$$

- Sometimes drop notation to $\mathbb{E}_X[f(x)]$ or just $\mathbb{E}[f(x)]$ if clear from context.
- Common: Mean of the distribution: $\mu = \mathbb{E}[x]$
- Examples:

$$P(x) = [0.4, 0.3, 0.1, 0.3], \quad p(x) = 3x^2$$



Expectation is a linear operator

(i.e. splits on summation and scale can come out)

- A linear operator H must satisfy two properties:

$$H(x + y) = H(x) + H(y)$$

$$H(\alpha x) = \alpha H(x)$$

- **Exercise:** Derive for expectations, i.e. $H = \mathbb{E}$

$$\mathbb{E}[af(x) + bg(x)] = a\mathbb{E}[f(x)] + b\mathbb{E}[g(x)]$$



Variance measures the “spread” of a distribution

- **Definition**

$$\begin{aligned}\text{Var}[x] &= \sigma^2 \equiv \mathbb{E}_X[(x - \mu)^2] \\ &= \mathbb{E}_X[(x - \mathbb{E}_X[x])^2]\end{aligned}$$

- Intuitively, recenter and then measure the expected value of $f(x) = x^2$.
- **Standard deviation** is the square root of variance:

$$\sigma = \sqrt{\sigma^2} = \sqrt{\mathbb{E}_X[(x - \mu)^2]}$$



Covariance and correlation measure linear relationship between two variables

- **Covariance definition**

$$\text{Cov}[x, y] \equiv \sigma_{X,Y}^2 \equiv \mathbb{E}_{(X,Y)}[(x - \mu_X)(y - \mu_Y)]$$

- **Correlation is a normalized covariance**

$$\rho_{X,Y} \equiv \frac{\sigma_{X,Y}^2}{\sigma_X \sigma_Y}$$

- **Example:**

$$P(x, y) = \begin{bmatrix} y = 1 & 0.4 & 0.1 \\ y = 0 & 0.1 & 0.4 \\ & x = 0 & x = 1 \end{bmatrix}$$



Covariance and correlation example derivation

- **Example:**

$$P(x, y) = \begin{array}{c} \left[\begin{array}{cc} y = 1 & 0.4 & 0.1 \\ y = 0 & 0.1 & 0.4 \\ & x = 0 & x = 1 \end{array} \right] \end{array}$$

- Means and variances: $\mu_X = \mu_Y = 0.5$, $\sigma_X^2 = \sigma_Y^2 = 0.25$

- Covariance:

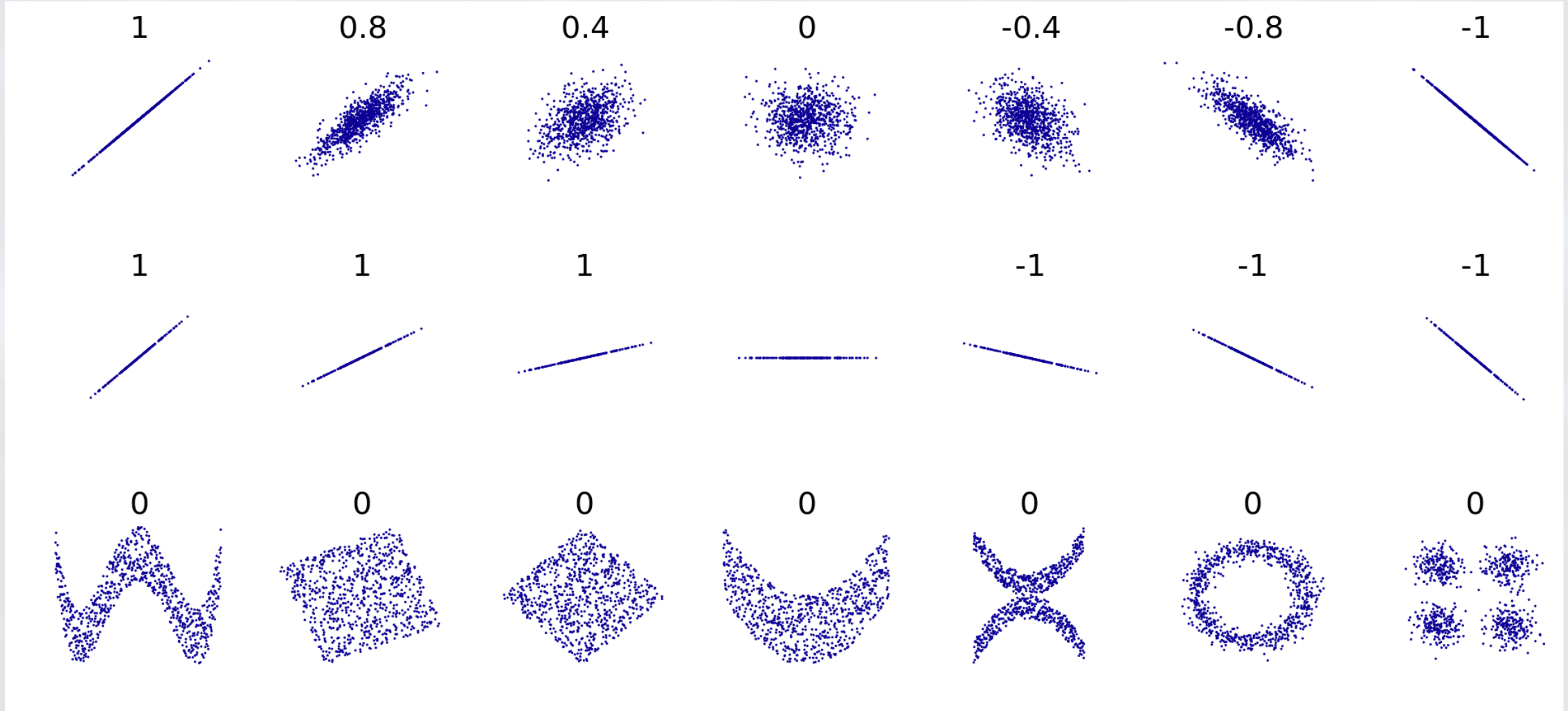
$$\begin{aligned} \sigma_{X,Y}^2 &= \mathbb{E}_{(X,Y)}[(x - \mu_X)(y - \mu_Y)] \\ &= p(0, 0)(0 - 0.5)(0 - 0.5) + p(0, 1)(0 - 0.5)(1 - 0.5) + p(1, 0)(1 - 0.5)(0 - 0.5) + p(1, 1)(1 - 0.5)(1 - 0.5) \\ &= 0.1(-0.5)(-0.5) + 0.4(-0.5)(0.5) + 0.4(0.5)(-0.5) + 0.1(0.5)(0.5) \\ &= 0.1(0.25) + 0.4(-0.25) + 0.4(-0.25) + 0.1(0.25) \\ &= 0.2(0.25) + 0.8(-0.25) \\ &= \frac{1}{5} \left(\frac{1}{4} \right) + \frac{4}{5} \left(-\frac{1}{4} \right) \\ &= -\frac{3}{20} = -0.15 \end{aligned}$$

- Correlation:

$$\rho_{X,Y} = \frac{\sigma_{X,Y}^2}{\sigma_X \sigma_Y} = \frac{-\frac{3}{20}}{\sqrt{0.25}\sqrt{0.25}} = \frac{-\frac{3}{20}}{\frac{1}{4}} = -\frac{3}{5}$$



Uncorrelated ($\rho_{X,Y} = 0$) is **NOT** the same as independence
(because it only measures **linear** relationships)



Covariance and correlation matrices are generalizations for vectors

- Covariance matrix has covariance of every pair of random variables

$$\underline{\Sigma} = \begin{bmatrix} \sigma_{X_1 X_1}^2 & \cdots & \sigma_{X_1 X_d} \\ \vdots & \vdots & \vdots \\ \sigma_{X_d X_1}^2 & \cdots & \sigma_{X_d X_d} \end{bmatrix}$$

- Matrix has variance along diagonal $\sigma_{X_i, X_i}^2 = \sigma_{X_i}^2$
- Correlation matrix is similar but with 1s on diagonal

$$R = \begin{bmatrix} 1 & \cdots & \rho_{X_1, X_d} \\ \vdots & \vdots & \vdots \\ \rho_{X_d X_1} & \cdots & 1 \end{bmatrix}$$

- Both matrices are symmetric $\Sigma = \Sigma^T$ and $R = R^T$



The **empirical expectation** is a sample version of the population-level expectation

- **Empirical expectation** is the average over n samples

$$\mathbb{E}_n[f(X)] = \frac{1}{n} \sum_{i=1}^n f(x_i)$$

where x_1, x_2, \dots, x_n are i.i.d. samples from the true distribution $p(x)$.

- **Law of large numbers** ensures this approaches the population expectation as the number of samples grows

$$\lim_{n \rightarrow \infty} \hat{\mathbb{E}}_n[f(X)] \rightarrow \mathbb{E}_{p(x)}[f(X)]$$



Big Picture Overview

- Introduction to probability
 - Motivation
 - Interpretation
- Random variables
 - Outcomes and events
- Probability distributions
 - PMF, PDF, CDF
- Multivariate distributions
 - Joint, marginal and conditional distributions
 - Chain rule and Bayes rule
 - Independence
- Expectations of random variables
 - Linearity of expectations
 - Covariance
 - Empirical expectations
- Multivariate Gaussian/Normal distribution
 - Marginal and conditional distributions
 - Affine/linear transformations



The Most Ubiquitous Multivariate Distribution Is the **Multivariate Gaussian/Normal Distribution**

- Compare univariate to multivariate:

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

- $p(x_1, \dots, x_d) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$

- μ is mean and Σ is covariance.
- $\Theta = \Sigma^{-1}$ is called the **precision matrix** (or **inverse covariance**).
- Σ and Θ must be positive definite $\Sigma > 0$.
- (Suppose $\Sigma = I$, suppose $\mu = 0$)



Why Are Multivariate Gaussian Distributions So Ubiquitous?

- **Reason from nature**

- The sum of independent random variables approaches a Gaussian distribution.
- **Central limit theorem!**

- **Math reason**

- Closed-form marginal and conditionals!

(Usually, very difficult to compute because sum/integral!)

- Affine/linear transformations of Gaussians are Gaussians.



Marginal and Conditional Distributions Are Gaussian and Can Be Computed in Closed-Form

- **2D case:**

- $x = [x_1, x_2] \sim \mathcal{N} \left(\mu = [\mu_1, \mu_2], \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} \right)$

- **Marginal distributions:**

- $x_1 \sim \mathcal{N}(\mu = \mu_1, \sigma^2 = \sigma_1^2)$

- $x_2 \sim \mathcal{N}(\mu = \mu_2, \sigma^2 = \sigma_2^2)$

- **Conditional distributions:**

- $x_1 | x_2 = a \sim \mathcal{N} \left(\mu = \mu_1 + \frac{\sigma_{12}}{\sigma_2^2} (a - \mu_2), \sigma^2 = \sigma_1^2 - \frac{\sigma_{21}^2}{\sigma_2^2} \right)$



Marginal and Conditional Distributions Are Gaussian and Can Be Computed in Closed-Form

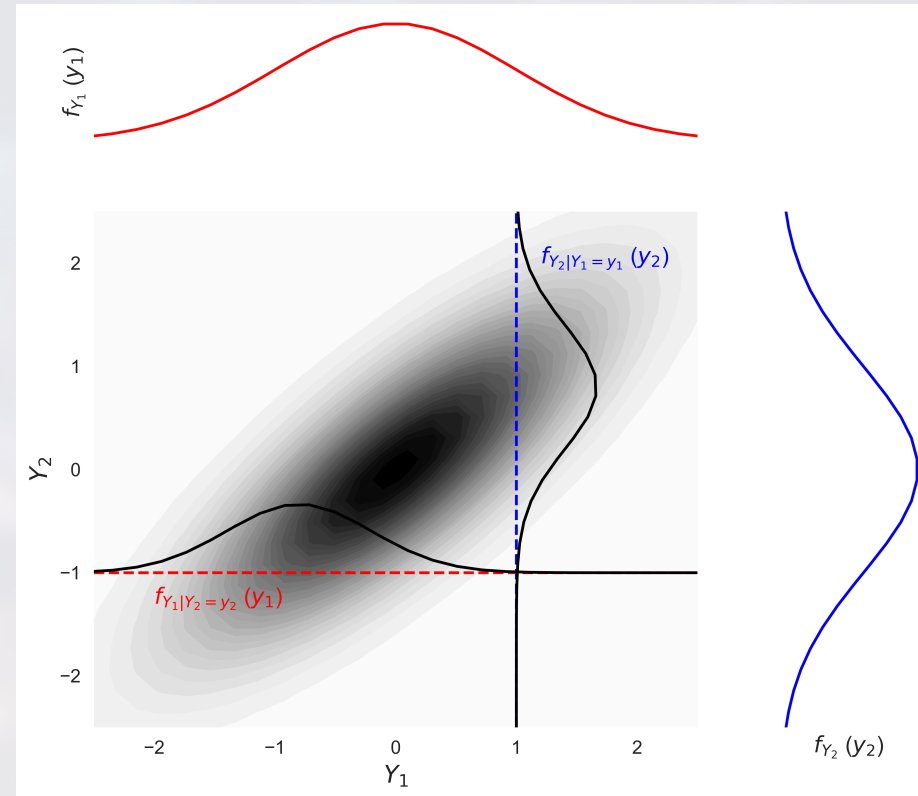
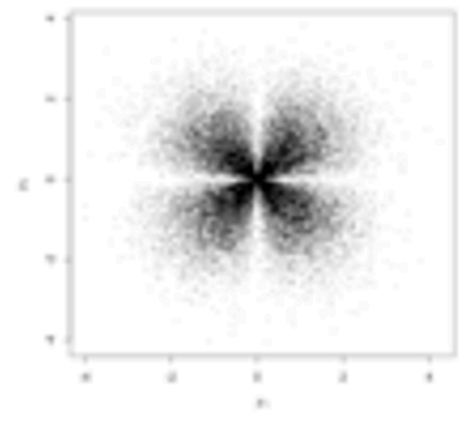
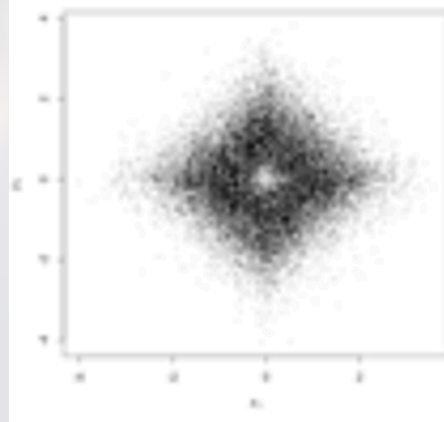
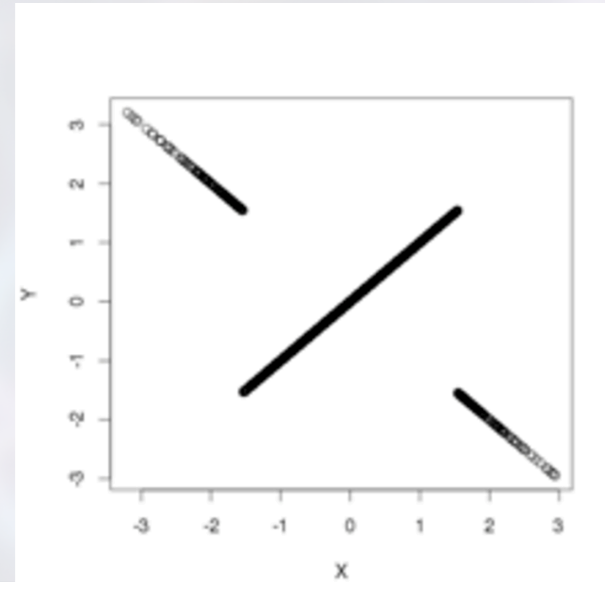
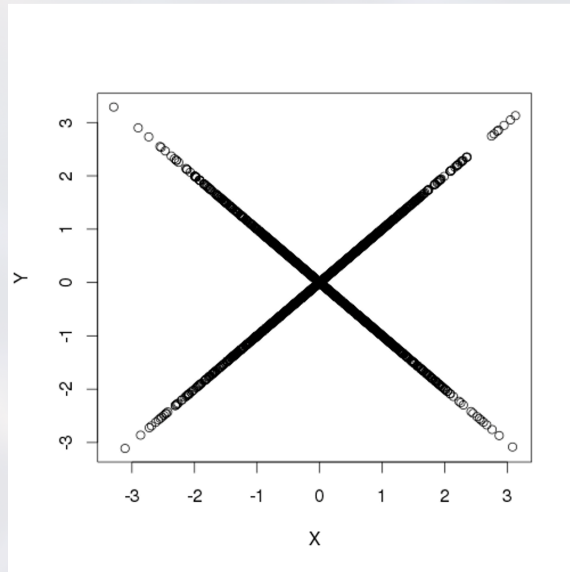


Image from <https://geostatisticslessons.com/lessons/multigaussian>

Caveat: Gaussian Marginals Does **Not** Imply Jointly Multivariate Gaussian (Converse **Not** Generally True)



Affine / Linear Transformations of Multivariate Gaussian Vector Are Also Multivariate Gaussian

- If $x \sim \mathcal{N}(\mu, \Sigma)$ and $y = Ax + b$, then $y \sim \mathcal{N}(A\mu + b, A\Sigma A^T)$.
- Special case: Marginal distribution when A is:

$$A_i = \begin{cases} 1, & \text{if } i = k \\ 0, & \text{otherwise} \end{cases} \implies y = x_k \sim p(x_k)$$

- **Key point:** Marginals, conditionals and affine functions known in **closed-form**.
- **Consequence 1:** Easy to manipulate.
- **Consequence 2:** Gaussians and linear ideas play nicely with each other.

