# Introduction to Deep Generative Modeling

# Lecture #2

**HY-673** – Computer Science Dep., University of Crete

Professors: Yannis Pantazis & Yannis Stylianou

TAs: Michail Raptakis

# What is probability?

- ## Frequentist Answer:

Let $A \subseteq \Omega$ be an event, then:

$$P(A) = \lim_{n \to \infty} \frac{n_A}{n}, \text{ where}$$

- $n$: number of experiments (independent trials),

- $n_A$: number of occurences of event A (e.g., coin or dice).

# What is probability?

- <u>Axiomatic Answer (Kolmogorov, 1933):</u>

Let $\Omega$ be a sample space, $\mathcal{F}$ be a event space (e.g., $\sigma$-algebra of $\Omega$) and $P$ be a measure. If

1. For all $A \in \mathcal{F}$: $P(A) \geq 0$.
2. $P(\Omega) = 1$.
3. $A \cap B = \emptyset \rightarrow P(A \cup B) = P(A) + P(B)$. ($\sigma$-additivity)

then, $(\Omega, \mathcal{F}, P)$ is a probability space.

Fair die: $\Omega_1 = \{\omega_1, \ldots, \omega_6\}$ with $P(\omega_i) = \frac{1}{6}$,

Fair coin: $\Omega_2 = \{\psi_1, \psi_2\}$ with $P(\psi_j) = \frac{1}{2}$,

Product spaces:

$\Omega = \Omega_1 \times \Omega_2$, or

$\Omega = \Omega_1 \times \Omega_1$

Bernoulli trials: $\underbrace{\Omega_2 \times \Omega_2 \times \cdots \times \Omega_2}_{n \text{ times}}$.

$\omega_i = $ sum of 2 independent dice,

$\Omega = \{2, 3, \ldots, 12\}$.

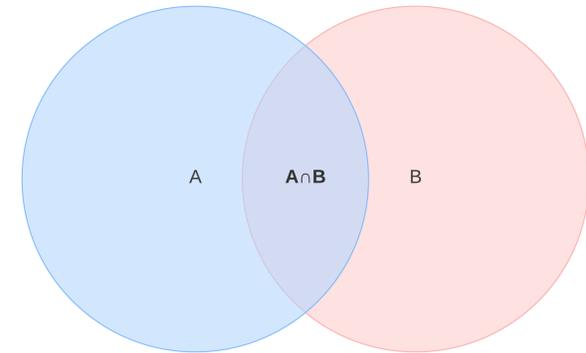## 0. Unity:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

*Venn Diagram*

## 1. Independence*:

$$A, B \text{ independent events} \to P(A \cap B) = P(A) \cdot P(B).$$

## 2. Conditional Probability:

$$P(B) \neq 0 \to P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

$$(A, B \text{ independent} \to P(A|B) = P(A)).$$

*Independence: "The occurrence of one event does not affect the probability of occurrence of the other."

# Basic properties

## 3. Chain Rule:

$$P(A_1 \cap \cdots \cap A_n) = P(A_n | A_{n-1}, \ldots, A_1) \cdot P(A_2 | A_1) \cdot P(A_1)$$

$$= \prod_{k=1}^{n} P(A_k | \cap_{j=1}^{k-1} A_j).$$

## 4. Bayes' Theorem:

If $P(B) \neq 0$ then $P(A|B) = \dfrac{P(B|A)P(A)}{P(B)}$.

# What is a Random Variable?

- A *Random Variable* (r.v. or rv) is a function from $\Omega$ to $\mathbb{R}$:

$$\omega \to x(\omega)$$

  e.g., $\omega_i$: face of a die, $x(\omega_i) = 10i$, voltage of a random source, etc.

- We'll use the notation: $x(\omega), X(\omega), Y(\omega), \ldots$ or just $x, X, Y, \ldots$

• <u>Cumulative Density Function (CDF), or Distribution Function:</u>

$$F_X(x) := P(X \leq x) \equiv P(\{\omega \in \Omega : x(\omega) \leq x\})$$

- $\lim_{x \to -\infty} F_X(x) = P(\emptyset) = 0$

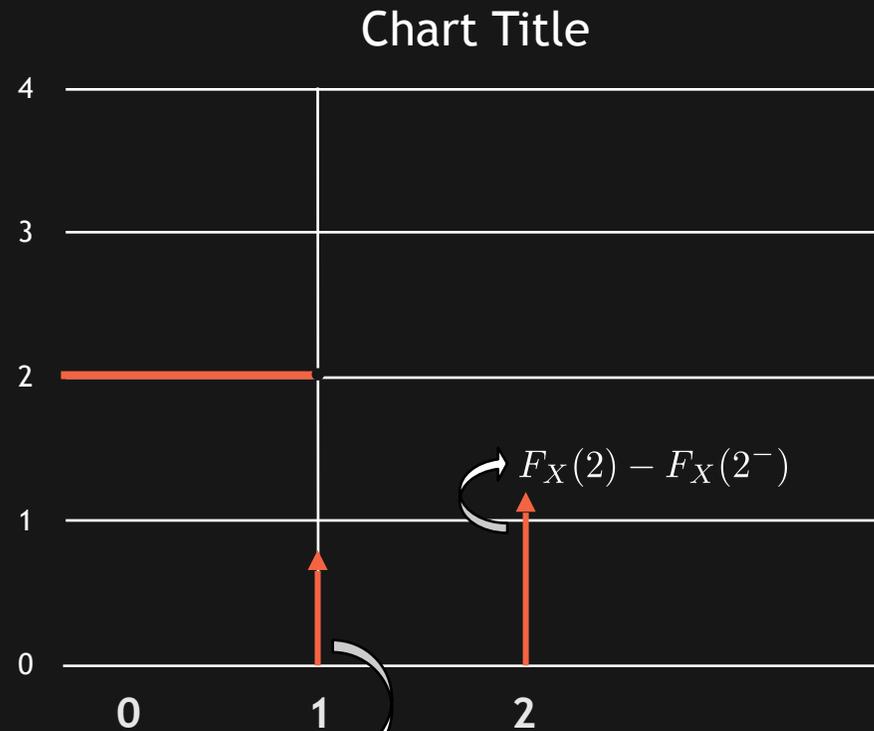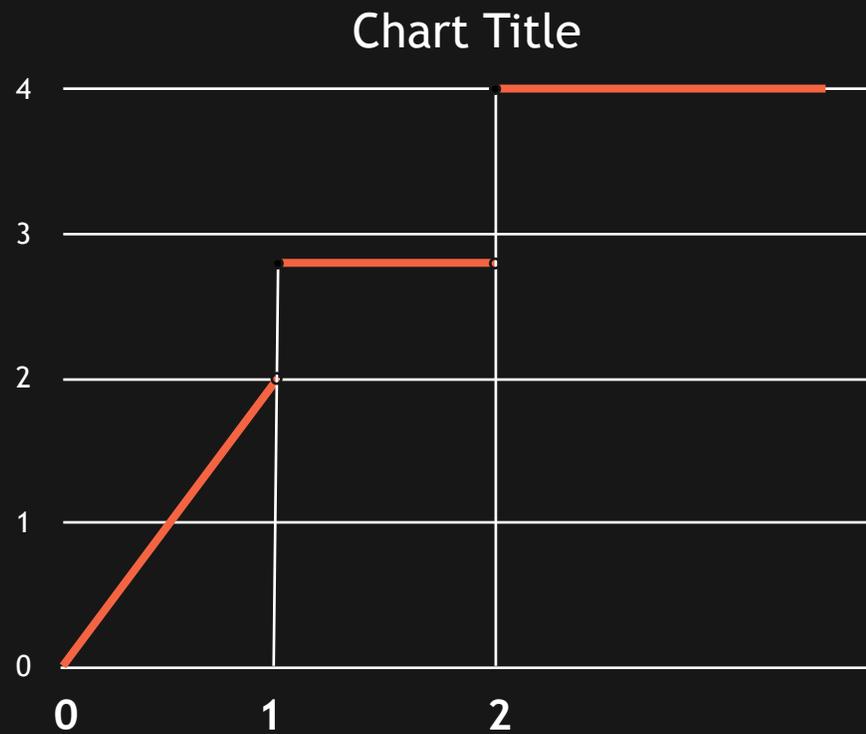- $\lim_{x \to +\infty} F_X(x) = P(\Omega) = 1$

- Probability Density Function (PDF):

$$f_X(x) = \frac{d}{dx} F_X(x).$$

- $f_X(x) \geq 0$, because $F_X(x)$ is an increasing function of $x$.

- $\int_{-\infty}^{+\infty} f_X(x) dx = 1.$

# Example of a CDF/PDF

- ## Mixed-type PDF:



Chart Title

Chart Title

$$F_X(2) - F_X(2^-)$$

**Dirac Delta Function (or impulse):**

$$\delta(x) : \int_{-\infty}^{+\infty} \phi(x)\delta(x)dx = \phi(0).$$

- Expected Value:

$$\mu_X = \mathbb{E}_X[X] = \int_{-\infty}^{+\infty} x f_X(x) dx$$

**"center of mass"**

$$\approx \frac{1}{n} \sum_{i=1}^{n} x_i, \ \ x_i \sim f_X, \ \text{i.i.d.}^*$$

**empirical**

*i.i.d: "Independent and identically distributed random variables."

- <u>Variance:</u>

$$\sigma_X^2 = \mathbb{E}_X \left[ (X - \mu_X)^2 \right]$$

$$= \int_{-\infty}^{+\infty} (x - \mu_X)^2 f_X(x) dx$$

$$= \mathbb{E}_X \left[ X^2 \right] - \mu_X^2.$$

- <u>Standard Deviation:</u>

$$\sigma_X = \sqrt{\sigma_X^2}.$$

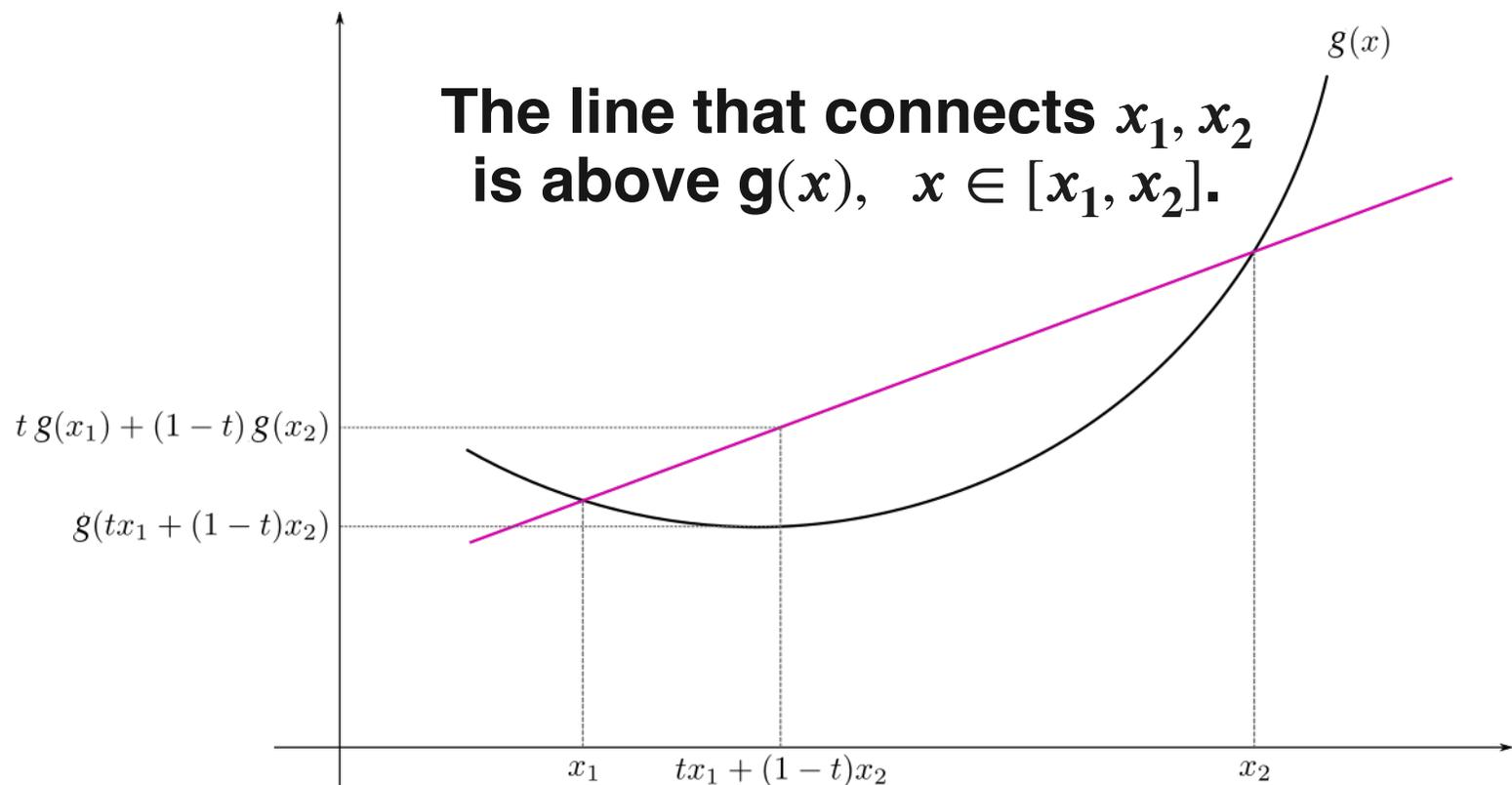$$\mathbb{E}_X[g(X)] = \int_{-\infty}^{+\infty} g(x) f_X(x) dx \neq g(\mathbb{E}_X[X]).$$

# Jensen's Inequality

- Convexity: $g(\cdot) \text{ convex} \iff g(tx_1 + (1-t)x_2) \leq tg(x_1) + (1-t)g(x_2)$

- Jensen's Inequality:

$$g(\cdot) \text{ convex}$$

$$\Rightarrow g\left(\mathbb{E}[X]\right) \leq \mathbb{E}\left[g(X)\right]$$

**The line that connects $x_1, x_2$ is above g$(x)$, $x \in [x_1, x_2]$.**

$g(x)$

$t\, g(x_1) + (1-t)\, g(x_2)$

$g(tx_1 + (1-t)x_2)$

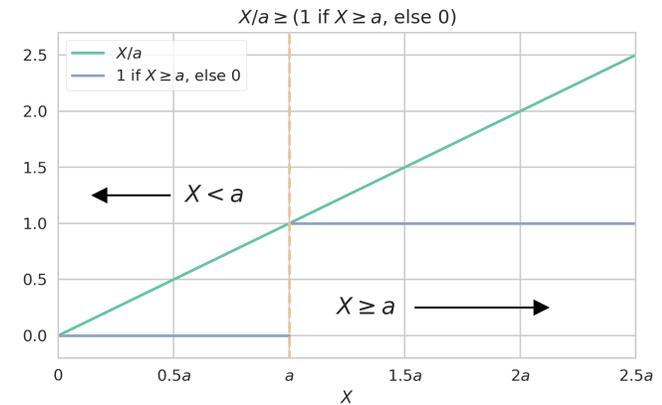$x_1 \qquad tx_1 + (1-t)x_2 \qquad\qquad x_2$

# Markov's Inequality

- Let $X$ be a non-negative r.v., and $a > 0$. Then,

$$P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.$$

- *Proof:*

$$
\begin{aligned}
\mathbb{E}[X] &= \int_{-\infty}^{+\infty} x f(x) dx \\
&= \int_0^{+\infty} x f(x) dx \\
&= \int_0^a x f(x) dx + \int_a^{+\infty} x f(x) dx \\
&\geq \int_0^a x f(x) dx + \int_a^{+\infty} a f(x) dx \\
&\geq a \int_a^{\infty} f(x) dx \\
&= a P(X \geq a).
\end{aligned}
\right\} \implies P(X \geq a) \leq \frac{\mathbb{E}[X]}{a}.
$$



$X/a \geq (1 \text{ if } X \geq a, \text{ else } 0)$

# Chebyshev's Inequality

- Let $X$ be a r.v. with finite expected value $\mu$, and finite non-zero variance $\sigma^2$. Then, for any positive real number $k > 0$:

$$P(|X - \mu| \geq k\sigma) \leq \tfrac{1}{k^2}.$$

  *Proof:* Apply Markov's inequality to $Y = (X - \mu)^2$

- E.g., for $k = 2$

$$P(|X - \mu| \leq 2\sigma) \leq \tfrac{1}{4} = \tfrac{25}{100}, \quad \forall \text{r.v. } X \text{ with } \mu, \sigma < \infty.$$

  - When $X$ is Gaussian, it holds: $P(|X - \mu| \leq 2\sigma) \leq \tfrac{5}{100}$.

- ## Probability Integral Transform:

If $X$ is a continuous r.v. with CDF $F_X(x)$, then the r.v. $Y = F_X(X)$ has a uniform distribution in $[0, 1]$:

$$
\left.
\begin{aligned}
F_Y(y) &= P(Y \le y) \\
&= P(F_X(X) \le y) \\
&= P(X \le F_X^{-1}(y)) \\
&= F_X(F_X^{-1}(y)) \\
&= y, \ y \in [0, 1].
\end{aligned}
\right\} \implies f_Y(y) =
\begin{cases}
1, \ 0 \le y \le 1 \\
0, \ \text{elsewhere.}
\end{cases}
$$

assuming $F_X^{-1}(y)$ exists

# Sampling an R.V. given $F_X(x)$

- **Inverse version (Inverse Transform Sampling):**

    If $Y$ has a uniform distribution in $[0, 1]$ and $X$ has CDF $F_X(x)$, then the r.v. $F_X^{-1}(Y)$ has the same distribution as $X$.

    **Very important and popular example: categorical distribution or softmax.**

- **Algorithm:**

    1. Compute the inverse of $F_X(x)$, i.e., $F_X^{-1}(x)$.

    2. Generate a random number $u \sim U([0, 1])$.

    3. Compute $x = F_X^{-1}(u) \sim X$.

- Let $g(\cdot)$ be a monotonic function. The r.v. $Y = g(X)$ has PDF given by:

$$f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{d}{dy} g^{-1}(y) \right| = f_X(g^{-1}(y)) \cdot \left| \frac{1}{g'(g^{-1}(y))} \right|.$$

*Proof:*

$$F_Y(y) = P(Y \leq y)$$
$$= P(g(X) \leq y)$$
$$= P(X \leq g^{-1}(y))$$
$$= F_X(g^{-1}(y)).$$

*Examples:*    $Y = aX + b, \ a > 0,$

$$f_Y(y) = \frac{1}{a} f_X\left( \frac{y-b}{a} \right),$$

$$Y = X^3,$$

$$f_Y(y) = \ldots$$

- Joint CDF of 2 random variables:

$$P(X \leq x, Y \leq y) = F_{XY}(x, y).$$

Limits:

PDF:

$$P(x, +\infty) = F_X(x),$$

$$P(+\infty, y) = F_Y(y).$$

$$f_{XY}(x, y) = \frac{\partial^2}{\partial x \partial y} F_{XY}(x, y) \geq 0.$$

$$\int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f_{XY}(x, y) dx dy = 1.$$

# Independence

- <u>Independence of $X, Y$:</u>

$$P(X \leq x, Y \leq y) = P(X \leq x) \cdot P(Y \leq y)$$
$$\implies F_{XY}(x, y) = F_X(x) \cdot F_Y(y)$$
$$\implies f_{XY}(x, y) = f_X(x) \cdot f_Y(y).$$

- Let $X$, $Y$ be two independent random variables. The PDF of $Z = X + Y$ is given by the **convolution** of $f_x(z)$ with $f_y(z)$:

$$\left. \begin{array}{l} P(Z \leq z) = P(X + Y \leq z) \\ \qquad\quad = P(X \leq z - Y) \\ \qquad\quad = \int_{-\infty}^{+\infty} P(X \leq z - y) f_Y(y) dy \\ \qquad\quad = \int_{-\infty}^{+\infty} F_X(z - y) f_Y(y) dy \end{array} \right\} \implies$$

$$\implies f_Z(z) = \int_{-\infty}^{+\infty} f_X(z - y) f_Y(y) dy$$

$$= f_X(z) * f_Y(z).$$

_Examples:_ Sum of 2 uniform distributions, sum of 2 dice, etc.

- ## Conditional Probability: *Example:* Die with even/odd events as conditions

$$f_{Y|X}(y|X=x) = f_{Y|X}(y|x) = \frac{f_{XY}(x,y)}{f_X(x)} = \frac{f_{XY}(x,y)}{\int_{-\infty}^{+\infty} f_{XY}(x,y)dy}$$

- ## Marginal Probability:

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

- ## Bayes Theorem:

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

Both a consequence of:

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

- <u>Notation:</u>

$$x = (x_1, \ldots, x_d)^T, \quad x \sim \mathcal{N}(\mu, \Sigma).$$

- <u>PDF General Form:</u>

$$\mathcal{N}(x|\mu, \Sigma) = f_X(x_1, \ldots, x_d)$$

$$= \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}, \quad |\Sigma| \triangleq \det(\Sigma).$$

$$|\Sigma| \neq 0 \rightarrow \mathcal{N}(x|\mu, \Sigma) \text{ non-degenerate.}$$

- ## Mean Vector:

$$\mu = \mathbb{E}[X] = (\mathbb{E}[X_1], \mathbb{E}[X_2], \dots, \mathbb{E}[X_d])^T$$

$$= \int_{-\infty}^{+\infty} x f_X(x|\mu, \Sigma) dx \in \mathbb{R}^d.$$

- ## Covariance Matrix:

$$\Sigma_{ij} = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)] = \mathrm{Cov}[X_i, X_j] \in \mathbb{R}^{d \times d}.$$
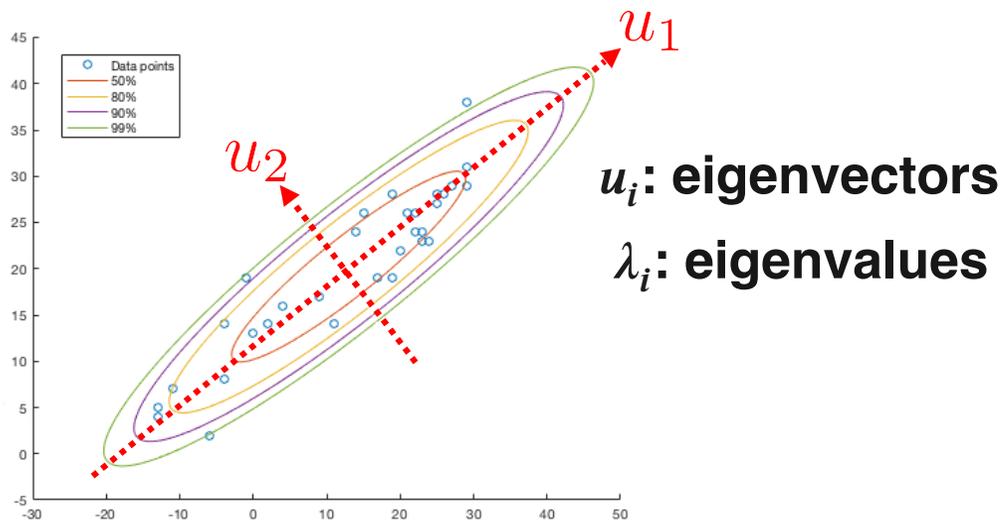
# Multivariate Gaussian Distribution

- **Marginals (2D example):**

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}.$$

$$p(x_1) = f_{X_1}(x_1) = \int f_X(X|\mu, \Sigma) dx_2 = \mathcal{N}(x_1|\mu_1, \sigma_1^2).$$

$$p(x_2) = \ldots = \mathcal{N}(x_2|\mu_2, \sigma_2^2).$$

- <u>Geometric interpretation:</u>

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}, \quad \rho = \frac{\sigma_{12}}{\sigma_1\sigma_2}.$$

**Correlation coefficient**



$u_i$: **eigenvectors**

$\lambda_i$: **eigenvalues**

$$\Delta^2 = (x - \mu)^T \Sigma^{-1} (x - \mu)$$

$$\Sigma^{-1} = \sum_{i=1}^{d} \frac{1}{\lambda_i} u_i u_i^T \Bigg\} \implies$$

$$\Delta^2 = \sum_{i=1}^{d} \frac{y_i^2}{\lambda_i}, \quad y_i = u_i^T (x - \mu).$$

# Multivariate Gaussian Distribution

- ## How to sample:

$$x = \mu + \sigma z \sim \mathcal{N}(\mu, \sigma^2) \, , \quad z \sim \mathcal{N}(0, 1)$$

**Reparametrization trick**

**Cholesky decomposition**

$$x = \mu + Lz \sim \mathcal{N}(\mu, \Sigma) \, , \quad z \sim \mathcal{N}(0, I_d) \quad \text{and} \ \Sigma = LL^T.$$

*Sampling from 1d normal*
import numpy as np
mu, sigma = 0, 0.1 # mean and standard deviation
x = np.random.normal(mu, sigma, 1000)
x.shape →   (1000,)
type(x) →   <class 'numpy.ndarray'>

*Sampling from multivariate normal*
mu = [1, 2]
Sigma = [[1, 2], [2, 4]]
x = np.random.multivariate_normal(mu, Sigma, 1000)
x.shape →   (1000,2)

- ## Conditional probability:

$$x = \begin{bmatrix} x_A \\ x_B \end{bmatrix} \qquad x_A = \begin{bmatrix} x_1 \\ \vdots \\ x_{d_1} \end{bmatrix}$$

$$x_B = \begin{bmatrix} x_{d_1+1} \\ \vdots \\ x_d \end{bmatrix}$$

$$\mu = \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}.$$

**Schur complement**

$$p(x_A|x_B) = \mathcal{N}(x_A|\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(x_B - \mu_B), \overbrace{\Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}})$$

$$p(x_B|x_A) = \mathcal{N}(x_B|\mu_B + \Sigma_{BA}\Sigma_{AA}^{-1}(x_A - \mu_A), \Sigma_{BB} - \Sigma_{BA}\Sigma_{AA}^{-1}\Sigma_{AB})$$

# Multivariate Gaussian Distribution

- <u>Conditional probability example:</u>

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix}$$



$$p(x_2|x_1) = \mathcal{N}(x_2|\mu_2 + \sigma_{21}\sigma_1^{-2}(x_1 - \mu_1), \ \sigma_2^2 - \sigma_{21}^2\sigma_1^{-2})$$

− Variance (a.k.a. uncertainty) is reduced whenever there is correlation!

# Asymptotics

- *Law of Large Numbers:*

Let $\{X_k\}$ be a sequence of i.i.d. r.v.s and $S_n = \sum_{k=1}^{n} X_k$ be the sum r.v. If $\mu = \mathbb{E}[X_k]$ exists, then $\forall \epsilon > 0$:

$$\lim_{n \to \infty} P\left( \left| \frac{1}{n} S_n - \mu \right| > \epsilon \right) = 0.$$

E.g., a fair coin, fair dice, etc.

*Proof:* Application of Chebyshev's inequality.

- *Central Limit Theorem:*

Let $\{X_k\}$ be a sequence of i.i.d. r.v.s and $S_n = \sum_{k=1}^{n} X_k$ be the sum r.v. Suppose that $\mu = \mathbb{E}[X]$, and $\sigma^2 = \text{var}(X)$ exist. Then, $\forall \beta$ fixed:

$$\lim_{n \to \infty} P\left(\frac{S_n - n\mu}{\sigma\sqrt{n}} < \beta\right) = \frac{1}{2}\left(1 + \text{erf}\left(\frac{\beta}{\sqrt{2}}\right)\right),$$

where $\text{erf}(\cdot)$ is the error function defined by: $\text{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$.

- In other words: $\frac{S_n - n\mu}{\sigma\sqrt{n}} \sim \mathcal{N}(0, 1)$, as $n \to \infty$.

1. All of statistics: A Concise Course in Statistical Inference (*Chapters 1–4*)
   Larry Wasserman, Springer (2004)

2. Probabilistic Machine Learning: An Introduction (*Chapters 2–3*)
   Kevin P Murphy, The MIT Press (2022)

# Introduction to Deep Generative Modeling

# Lecture #2

**HY-673** – Computer Science Dep., University of Crete

Professors: Yannis Pantazis & Yannis Stylianou

TAs: Michail Raptakis