

# ECE 595: Machine Learning I

## Tutorial 02: Probability Review

Spring 2020

Stanley Chan

School of Electrical and Computer Engineering  
Purdue University



# Outline

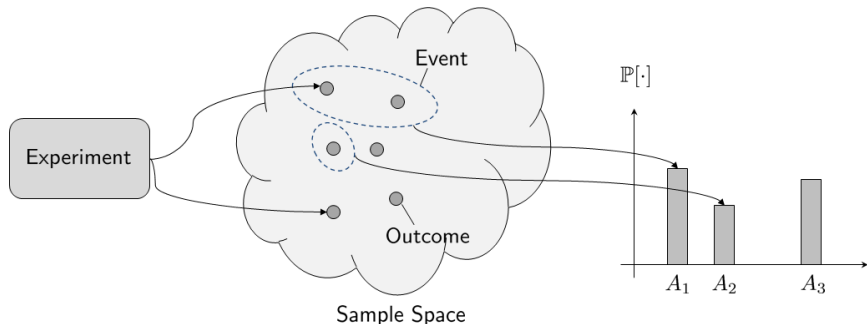
- Probability Distributions
- High-dimensional Gaussian

## Reference:

- Dimitri Bertsekas, Introduction to Probability, Athena Scientific, 2008, 2nd Edition.
- Purdue ECE 302 Course Note  
<https://engineering.purdue.edu/ChanGroup/ECE302>

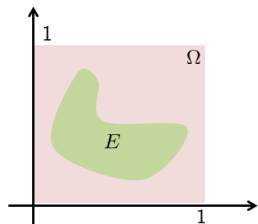
# Probability Space

- Sample space  $\Omega$  = set of all possible outcomes.
- Event Space  $\mathcal{E}$  = set of all events. Event is a subset in  $\Omega$ .
- Probability Law  $\mathbb{P}$  = a mapping from  $\mathcal{E}$  to  $[0, 1]$ .



# Interpreting Probability

$\mathbb{P}[\cdot]$  is a **measure** of the size of the event.



- $\Omega = \{1, 2, 3, 4, 5, 6\}$
- $E = \{1\}$ ,  $\mathbb{P}[E] = 1/6$
- $E = \{1, 3\}$ ,  $\mathbb{P}[E] = 2/6$
- $\mathbb{P}[E_1] \leq \mathbb{P}[E_2]$  if  $E_1 \subseteq E_2$ .

- $\Omega = [0, 1] \times [0, 1]$
- $E =$  shaded region,  $\mathbb{P}[E] =$  area.
- $E = \{(x_0, y_0)\}$ ,  $\mathbb{P}[E] = 0$ .
- $\mathbb{P}[E]$  can be 0 even if  $E \neq \emptyset$ .

# Probability Axioms

- **Non-negative:**

$$\mathbb{P}[E] \geq 0$$

- **Unity:**

$$\mathbb{P}[\Omega] = 1$$

- **Additivity:** If  $A_n$ 's are disjoint, then

$$\mathbb{P} \left[ \bigcup_{n=1}^N A_n \right] = \sum_{n=1}^N \mathbb{P}[A_n]$$

- If  $A_n$ 's are not disjoint, then Union bound holds

$$\mathbb{P} \left[ \bigcup_{n=1}^N A_n \right] \leq \sum_{n=1}^N \mathbb{P}[A_n]$$

# Conditional Probability

- Conditional probability of  $A$  given  $B$  is

$$\mathbb{P}[A | B] = \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]}$$

- Bayes Theorem:

$$\mathbb{P}[B_i | A] = \frac{\mathbb{P}[A | B_i] \mathbb{P}[B_i]}{\underbrace{\mathbb{P}[A]}_{\substack{\uparrow \\ \mathbb{P}[A] = \sum_{n=1}^N \mathbb{P}[A | B_n] \mathbb{P}[B_n]}}}$$

- Therefore,

$$\mathbb{P}[B_i | A] = \frac{\mathbb{P}[A | B_i] \mathbb{P}[B_i]}{\sum_{n=1}^N \mathbb{P}[A | B_n] \mathbb{P}[B_n]}$$

# Random Variable

- Scalar Case:  $X : \mathcal{E} \rightarrow \mathbb{R}$ 
  - **Encode** an event to a number.
  - E.g.,  $X = 1$  represents “Democrat”,  $X = 0$  represents “Republican”.
  - Each number is called a **state**.
  - $\mathbb{P}[X = x]$  = probability of  $X$  have the state  $x$ .
- Vector Case:  $\mathbf{X} : \mathcal{E}_1 \times \dots \times \mathcal{E}_N \rightarrow \mathbb{R}^N$ .
  - Encode a sequence of events to a sequence of numbers.

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_N \end{bmatrix}$$

- $\mathbb{P}[\mathbf{X} = \mathbf{x}]$  = probability of  $\{X_1 = x_1, X_2 = x_2, \dots \text{ and } X_N = x_N\}$ .

# Probability Distribution

## Cumulative Distribution Function:

- Defined as

$$F_X(x) = \mathbb{P}[X \leq x].$$

- For high-dimensional vectors:

$$F_{\mathbf{X}}(\mathbf{x}) = \mathbb{P}[\mathbf{X} \leq \mathbf{x}].$$

- This means:

$$\begin{aligned} F_{X_1, \dots, X_N}(x_1, \dots, x_N) \\ = \mathbb{P}[X_1 \leq x_1 \dots \text{and} \dots X_N \leq x_N]. \end{aligned}$$

## Probability Density Function:

- Defined as

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

- For high-dimensional vectors:

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{d}{d\mathbf{x}} F_{\mathbf{X}}(\mathbf{x})$$

- This means:

$$\begin{aligned} p_{X_1, \dots, X_N}(x_1, \dots, x_N) \\ = \frac{\partial^N}{\partial x_1 \dots \partial x_N} F(x_1, \dots, x_N). \end{aligned}$$



# Probability Distribution

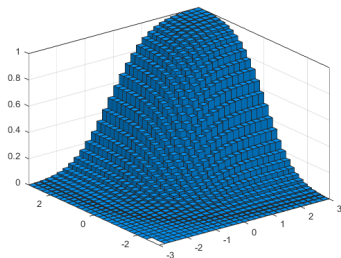
## Cumulative Distribution Function:

- Defined as

$$F_X(x) = \mathbb{P}[X \leq x].$$

- For high-dimensional vectors:

$$F_{\mathbf{X}}(\mathbf{x}) = \mathbb{P}[\mathbf{X} \leq \mathbf{x}].$$



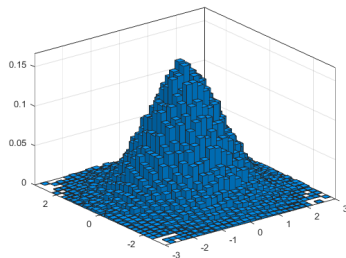
## Probability Density Function:

- Defined as

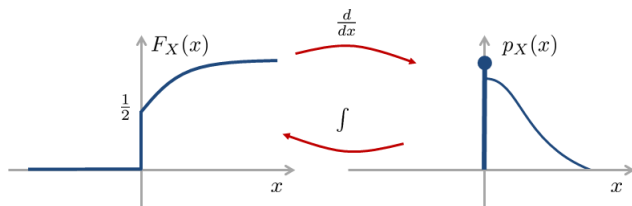
$$p_X(x) = \frac{d}{dx} F_X(x).$$

- For high-dimensional vectors:

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{d}{d\mathbf{x}} F_{\mathbf{X}}(\mathbf{x})$$



## Linking CDF and PDF



**Example.** Consider a PDF

$$p_X(x) = \begin{cases} 0, & x < 0, \\ \frac{1}{2}, & x = 0, \\ \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\}, & x > 0. \end{cases}$$

Then, the CDF of  $X$  is

$$F_X(x) = \begin{cases} 0, & x < 0, \\ \frac{1}{2} + \int_0^x \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{s^2}{2}\right\} ds, & x \geq 0. \end{cases}$$

# Expectation

## Definition

The **expectation** of a random variable  $X$  is

$$\mu = \mathbb{E}[X] = \int x p(x) dx.$$

- **Second Moment:**

$$\mathbb{E}[X^2] = \int x^2 p(x) dx.$$

- **Variance:**

$$\text{Var}[X] = \mathbb{E}[(X - \mu)^2] = \int (x - \mu)^2 p(x) dx$$

# Properties of Expectation

(i) **Function.** For any function  $g$ ,

$$\mathbb{E}[g(X)] = \int g(x)p(x)dx.$$

(ii) **Linearity.** For any function  $g$  and  $h$ ,

$$\mathbb{E}[g(X) + h(X)] = \mathbb{E}[g(X)] + \mathbb{E}[h(X)].$$

(iii) **Scale.** For any constant  $c$ ,

$$\mathbb{E}[cX] = c\mathbb{E}[X].$$

(iv) **DC Shift.** For any constant  $c$ ,

$$\mathbb{E}[X + c] = \mathbb{E}[X] + c.$$

# Properties of Variance

(i) **Moment.**

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2.$$

(ii) **Scale.** For any constant  $c$ ,

$$\text{Var}[cX] = c^2 \text{Var}[X].$$

(iii) **DC Shift.** For any constant  $c$ ,

$$\text{Var}[X + c] = \text{Var}[X].$$

# Moment Generating Function

## Definition

The **moment generating function** (MGF) of a random variable  $X$  is

$$M_X(s) = \mathbb{E}[e^{sX}].$$

- MGF for Gaussian is

$$M_X(s) = \exp \left\{ \mu s + \frac{s^2 \sigma^2}{2} \right\}$$

- If  $X$  and  $Y$  are independent, then

$$\begin{aligned} M_{X+Y}(s) &= \mathbb{E}[e^{s(X+Y)}] = \mathbb{E}[e^{sX} e^{sY}] \\ &= \mathbb{E}[e^{sX}] \mathbb{E}[e^{sY}] = M_X(s) M_Y(s). \end{aligned}$$

# Gaussian Random Variable

- $X \sim \mathcal{N}(\mu, \sigma^2)$  if

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

- Transforming a Gaussian. Let  $Y = aX + b$ , then

$$Y \sim \mathcal{N}(a\mu + b, a^2\sigma^2)$$

You may check:

$$\begin{aligned}\mathbb{E}[Y] &= \mathbb{E}[aX + b] = a\mathbb{E}[X] + b = a\mu + b \\ \text{Var}[Y] &= \text{Var}[aX + b] = \text{Var}[aX] = a^2\text{Var}[X].\end{aligned}$$

# High-dimensional Gaussian

An  $d$ -dimensional **Gaussian** has a PDF

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\},$$

where  $d$  denotes the dimensionality of the vector  $\mathbf{x}$ .

- The **mean vector**  $\boldsymbol{\mu}$  is

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{X}] = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_d] \end{bmatrix}$$

- The **covariance matrix**  $\boldsymbol{\Sigma}$  is

$$\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \text{Var}[X_1] & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_d) \\ \text{Cov}(X_2, X_1) & \text{Var}[X_2] & \dots & \text{Cov}(X_2, X_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_d, X_1) & \text{Cov}(X_d, X_2) & \dots & \text{Var}[X_d] \end{bmatrix}$$

- $\boldsymbol{\Sigma}$  is always positive semi-definite. (Why?)



## Special Case: Diagonal Covariance

- Suppose that  $X_i$  and  $X_j$  are independent for all  $i \neq j$ .
- This implies  $\text{Cov}(X_i, X_j) = 0$
- Simplify  $\Sigma$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_d^2 \end{bmatrix},$$

- Then, the exponential is

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^n \frac{(x_i - \mu_i)^2}{\sigma_i^2}.$$

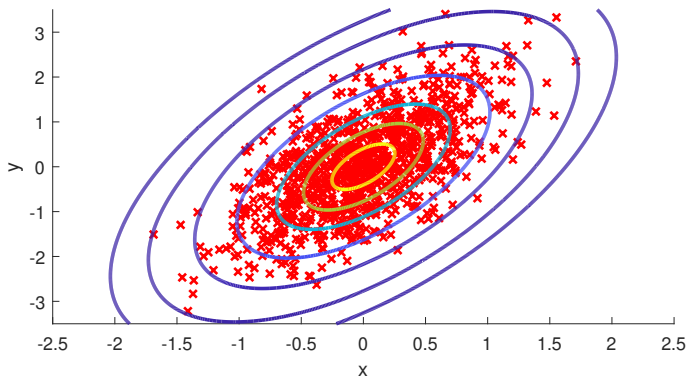
- And hence, the PDF is

$$p_{\mathbf{X}}(\mathbf{x}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp \left\{ -\frac{(x_i - \mu_i)^2}{2\sigma_i^2} \right\}.$$

# Visualization

- Generate 1000 random samples from a 2D Gaussian

- $\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ , and  $\Sigma = \begin{bmatrix} 0.25 & 0.3 \\ 0.3 & 1 \end{bmatrix}$



## MATLAB Code

```
% MATLAB code: Generate random numbers from multivariate Gaussian
mu    = [0 0];
Sigma = [.25 .3; .3 1];
x     = mvnrnd(mu,Sigma,1000);
```

```
% MATLAB code: Overlay random numbers with the Gaussian contour.
x1 = -2.5:.01:2.5;
x2 = -3.5:.01:3.5;
[X1,X2] = meshgrid(x1,x2);
F = mvnpdf([X1(:) X2(:)],mu,Sigma);
F = reshape(F,length(x2),length(x1));
figure(1);
scatter(x(:,1),x(:,2),'rx', 'LineWidth', 1.5); hold on;
contour(x1,x2,F,[.001 .01 .05:.1:.95 .99 .999], 'LineWidth', 2);
xlabel('x'); ylabel('y');
set(gcf, 'Position', [100, 100, 600, 300]);
```

# Conditional Gaussian

- Data  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .
- Class  $Y \in \{1, 2, \dots, K\}$ .
- **Likelihood:**

$p_{\mathbf{X}|Y}(\mathbf{x}|k)$  = Probability of getting  $\mathbf{X}$  given  $Y$

- **Prior:**

$p_Y(k)$  = Probability of getting  $Y$

- **Posterior:**

$p_{Y|\mathbf{X}}(k|\mathbf{x})$  = Probability of getting  $Y$  given  $\mathbf{X}$

- Related by

$$p_{Y|\mathbf{X}}(k|\mathbf{x}) = \frac{p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}{p_{\mathbf{X}}(\mathbf{x})} = \frac{p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}{\sum_k p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}$$

## Example

- Two Gaussian  $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$  and  $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ .
- **Prior** probability of getting a class is

$$p_Y(1) = \pi_1 \quad \text{and} \quad p_Y(2) = \pi_2.$$

- The **likelihood** term is

$$\begin{aligned} p_{\mathbf{X}|Y}(\mathbf{x}|k) &= \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \end{aligned}$$

- The **posterior** is

$$\begin{aligned} p_{Y|\mathbf{X}}(k|\mathbf{x}) &= \frac{p_{\mathbf{X}|Y}(\mathbf{x}|k)p_Y(k)}{p_{\mathbf{X}}(\mathbf{x})} \\ &= \frac{\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \cdot \pi_k}{\sum_{k=1}^K \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \cdot \pi_k} \end{aligned}$$

## Negative Log-Likelihood

Negative Log-Likelihood for Gaussian:

$$\begin{aligned} & -\log p_{\mathbf{X}|Y}(\mathbf{x}|k) \\ &= -\log \left( \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_k|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \right) \\ &= \underbrace{\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}_{\text{contains } \mathbf{x}} \underbrace{-\frac{n}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_k|}_{\text{no } \mathbf{x}}. \end{aligned}$$

- $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \geq 0$ , always.
- $\sqrt{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})}$  is called **Mahalanobis distance**.

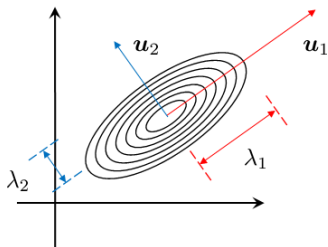
# Geometry of Gaussian

- Perform eigen-decomposition

$$\Sigma = U\Lambda U^T$$

$$= \begin{bmatrix} | & | & \dots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \dots & \mathbf{u}_n \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \lambda_n \end{bmatrix} \begin{bmatrix} - & \mathbf{u}_1^T & - \\ - & \mathbf{u}_2^T & - \\ & \vdots & \\ - & \mathbf{u}_n^T & - \end{bmatrix}.$$

- $\mathbf{u}_i$  = orientation
- $\lambda_i$  = radius

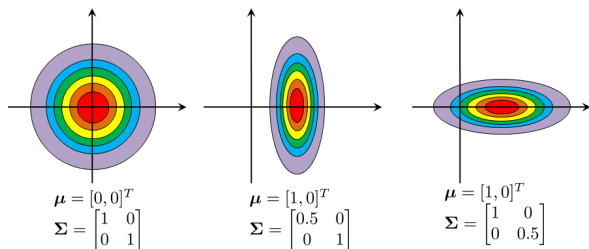


# Geometry of Gaussian

- Special Case:  $X_i$ 's are independent

$$\Sigma = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 1 \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & \lambda_n \end{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \dots & 1 \end{bmatrix}.$$

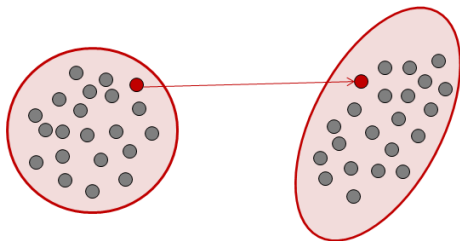
- Ellipse; Standard bases; Different radii.





# Transformation of Gaussian

- You are given  $\{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ . All  $\mathbf{X}_j$  are generated from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ .
- You want  $\{\mathbf{Y}_1, \dots, \mathbf{Y}_N\}$ , where  $\mathbf{Y}_j \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ .
- But you only have  $\{\mathbf{X}_1, \dots, \mathbf{X}_N\}$ .



# Transformation of Gaussian

How about this? Let

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$$

Can we find  $\mathbf{A}$  and  $\mathbf{b}$ ?

$$\mathbb{E}[\mathbf{Y}] = \mathbb{E}[\mathbf{A}\mathbf{X} + \mathbf{b}] = \mathbf{A}\underbrace{\mathbb{E}[\mathbf{X}]_{=0}} + \mathbf{b} = \mathbf{b}.$$

$$\begin{aligned}\text{Cov}(\mathbf{Y}) &= \mathbb{E}[(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^T] \\ &= \mathbb{E}[(\mathbf{A}\mathbf{X} + \mathbf{b} - \mathbf{b})(\mathbf{A}\mathbf{X} + \mathbf{b} - \mathbf{b})^T] \\ &= \mathbb{E}[(\mathbf{A}\mathbf{X})(\mathbf{A}\mathbf{X})^T] = \mathbb{E}[\mathbf{A}\mathbf{X}\mathbf{X}^T\mathbf{A}^T] \\ &= \mathbf{A}\mathbb{E}[\mathbf{X}\mathbf{X}^T]\mathbf{A}^T = \mathbf{A}\mathbf{A}^T = \mathbf{\Sigma}.\end{aligned}$$

So here is the choice:  $\mathbf{b} = \boldsymbol{\mu}$ , and  $\mathbf{A} = \boldsymbol{\Sigma}^{\frac{1}{2}}$ .

## Inverse Transform

If we have  $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , how to transform  $\mathbf{Y}$  to  $\mathbf{X}$  so that  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ?

The inverse transform is

$$\mathbf{X} = \boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{Y} - \boldsymbol{\mu})$$

Checking: If  $\mathbf{Y} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then

$$\mathbb{E}[\mathbf{X}] = \mathbb{E}[\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{Y} - \boldsymbol{\mu})] = \boldsymbol{\Sigma}^{-1}(\underbrace{\mathbb{E}[\mathbf{Y}] - \boldsymbol{\mu}}_{=\boldsymbol{\mu}}) = \mathbf{0}.$$

$$\begin{aligned}\text{Cov}[\mathbf{X}] &= \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T] \\ &= \mathbb{E}[(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{Y} - \boldsymbol{\mu}))(\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{Y} - \boldsymbol{\mu}))^T] \\ &= \mathbb{E}[\boldsymbol{\Sigma}^{-\frac{1}{2}}(\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-\frac{1}{2}}] \\ &= \boldsymbol{\Sigma}^{-\frac{1}{2}} \underbrace{\mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^T]}_{=\boldsymbol{\Sigma}} \boldsymbol{\Sigma}^{-\frac{1}{2}} = \mathbf{I}.\end{aligned}$$