

Foundations of Machine Learning  
African Masters in Machine Intelligence



**AIMS** | African Institute for  
Mathematical Sciences  
RWANDA

**Imperial College  
London**

# Summary Statistics

**Marc Deisenroth**

Quantum Leap Africa  
African Institute for Mathematical  
Sciences, Rwanda

Department of Computing  
Imperial College London



@mpd37

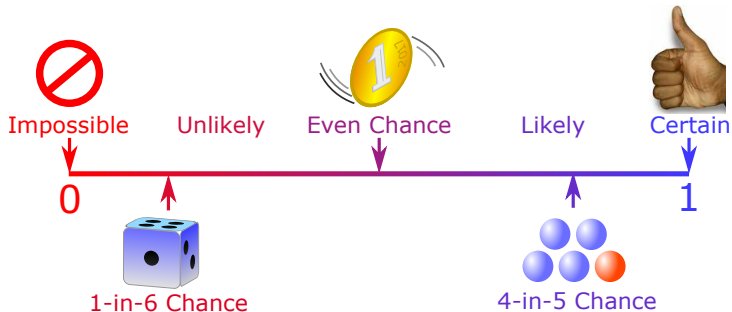
mdeisenroth@aimsammi.org

September 27, 2018

# Some Learning Material

- ▶ Book: <https://mml-book.com> (Chapter 6)
- ▶ MOOC:  
<https://www.coursera.org/learn/pca-machine-learning>  
(Week 1)

# Probabilities



- ▶ Describe a frequency ratio of events
- ▶ Express uncertainty when making predictions
- ▶ Capture a degree of belief about a hypothesis

▶ Probabilities are sufficient for reasoning under uncertainty

# Probability Distributions

- ▶ Probability density function (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

# Probability Distributions

- ▶ **Probability density function** (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

- ▶ **Probability mass function** (discrete  $x$ )

$$p(x) \geq 0, \quad \sum_x p(x) = 1$$

# Probability Distributions

- ▶ **Probability density function** (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

- ▶ **Probability mass function** (discrete  $x$ )

$$p(x) \geq 0, \quad \sum_x p(x) = 1$$

- ▶ Here: We are imprecise and call  $p(\cdot)$  a **probability distribution**

# Probability Distributions

- ▶ **Probability density function** (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

- ▶ **Probability mass function** (discrete  $x$ )

$$p(x) \geq 0, \quad \sum_x p(x) = 1$$

- ▶ Here: We are imprecise and call  $p(\cdot)$  a **probability distribution**
- ▶ If  $x$  is continuous  $p(x)$  is **not** the probability of event  $x$  happening

# Probability Distributions

- ▶ **Probability density function** (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

- ▶ **Probability mass function** (discrete  $x$ )

$$p(x) \geq 0, \quad \sum_x p(x) = 1$$

- ▶ Here: We are imprecise and call  $p(\cdot)$  a **probability distribution**
- ▶ If  $x$  is continuous  $p(x)$  is **not** the probability of event  $x$  happening
- ▶ **Cumulative distribution function**

$$p(x \leq t) = \int_{-\infty}^t p(x)dx \in [0, 1]$$



# Probability Distributions

- ▶ **Probability density function** (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

- ▶ **Probability mass function** (discrete  $x$ )

$$p(x) \geq 0, \quad \sum_x p(x) = 1$$

- ▶ Here: We are imprecise and call  $p(\cdot)$  a **probability distribution**
- ▶ If  $x$  is continuous  $p(x)$  is **not** the probability of event  $x$  happening
- ▶ **Cumulative distribution function**

$$p(x \leq t) = \int_{-\infty}^t p(x)dx \in [0, 1]$$

- ▶ **Joint distribution**  $p(x, y)$  of two random variables  $x, y$

# Probability Distributions

- ▶ **Probability density function** (continuous  $x$ )

$$p(x) \geq 0, \quad \int p(x)dx = 1$$

- ▶ **Probability mass function** (discrete  $x$ )

$$p(x) \geq 0, \quad \sum_x p(x) = 1$$

- ▶ Here: We are imprecise and call  $p(\cdot)$  a **probability distribution**
- ▶ If  $x$  is continuous  $p(x)$  is **not** the probability of event  $x$  happening
- ▶ **Cumulative distribution function**

$$p(x \leq t) = \int_{-\infty}^t p(x)dx \in [0, 1]$$

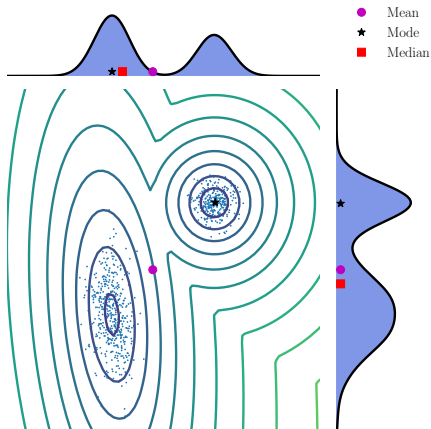
- ▶ **Joint distribution**  $p(x, y)$  of two random variables  $x, y$
- ▶ **Conditional distribution**  $p(x|y)$  of  $x$  given  $y$

# Summary Statistics

$$\mathcal{D} = \{-1, 1, 2\}$$

- ▶ Summarize datasets or random variables by describing some of their properties
- ▶ Examples:
  - ▶ **Mean/Expected value** (average):  $2/3$
  - ▶ **Variance** (related to spread of the data around the mean):  $1.56$
  - ▶ **Median** (data point “in the middle”, i.e., value so that another data point is equally likely to be greater or smaller)  $\neq$  mean:  $1$

# Mean, Mode, Median (Continuous Distributions)



## Mean (Expected Value)

- ▶ “Average”
- ▶ Does not have to be part of the dataset or a plausible realization of a random variable (▶▶ dice)

## Mean (Expected Value)

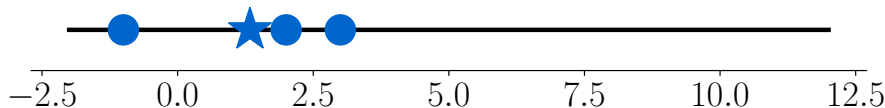
- ▶ “Average”
- ▶ Does not have to be part of the dataset or a plausible realization of a random variable (▶ dice)

$$\mathbb{E}_x[\mathbf{x}] = \int \mathbf{x}p(\mathbf{x})d\mathbf{x} =: \boldsymbol{\mu}_x \in \mathbb{R}^D \quad \text{if } \mathbf{x} \in \mathbb{R}^D \text{ is continuous}$$

$$\mathbb{E}_x[\mathbf{x}] = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n =: \boldsymbol{\mu}_x \in \mathbb{R}^D \quad \text{if } \mathbf{x} \in \mathbb{R}^D \text{ is discrete}$$

$$\mathbb{E}_x[\mathbf{x}] = \begin{bmatrix} \mathbb{E}_{x_1}[x_1] \\ \vdots \\ \mathbb{E}_{x_D}[x_D] \end{bmatrix} \in \mathbb{R}^D$$

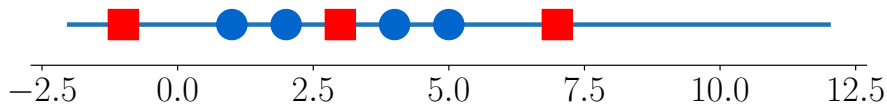
# Empirical Mean



- ▶ Random variable  $\mathbf{x} \in \mathbb{R}^D$
- ▶  $N$  concrete realizations  $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_n \in \mathbb{R}^D$
- ▶ **Empirical mean** (estimate of the true mean:

$$\frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \in \mathbb{R}^D$$

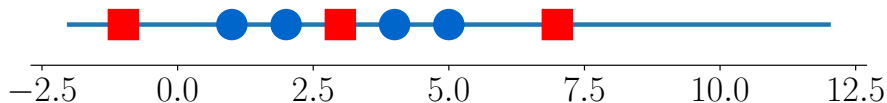
# Variance



- ▶ Both datasets have the same (empirical) mean

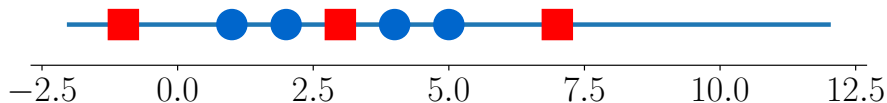


# Variance



- ▶ Both datasets have the same (empirical) mean
- ▶ Need a different quantity to describe “spread” of the data around the mean ▶ **Variance**
- ▶ Variance: Expected (squared) distance of data from the mean

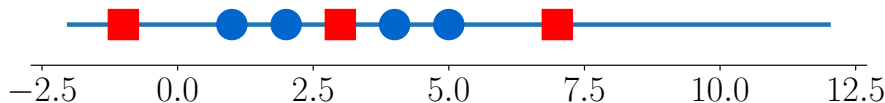
# Variance



- ▶ Both datasets have the same (empirical) mean
- ▶ Need a different quantity to describe “spread” of the data around the mean ▶ Variance
- ▶ Variance: Expected (squared) distance of data from the mean

$$\mathbb{V}[\mathbf{x}] := \mathbb{E}_x[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^\top]$$

# Variance



- ▶ Both datasets have the same (empirical) mean
- ▶ Need a different quantity to describe “spread” of the data around the mean ▶ Variance
- ▶ Variance: Expected (squared) distance of data from the mean

$$\mathbb{V}[\mathbf{x}] := \mathbb{E}_{\mathbf{x}}[(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})^{\top}]$$

$$\mathbb{V}[\mathbf{x}] := \int (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})^{\top} p(\mathbf{x}) d\mathbf{x} \in \mathbb{R}^{D \times D} \quad \text{if } \mathbf{x} \in \mathbb{R}^D \text{ is continuous}$$

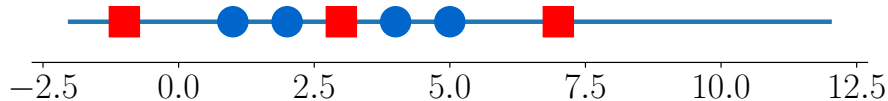
$$\mathbb{V}[\mathbf{x}] := \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{\mathbf{x}})(\mathbf{x}_n - \boldsymbol{\mu}_{\mathbf{x}})^{\top} \in \mathbb{R}^{D \times D} \quad \text{if } \mathbf{x} \in \mathbb{R}^D \text{ is discrete}$$

# Empirical Variance

- ▶ Random variable  $\mathbf{x} \in \mathbb{R}^D$
- ▶  $N$  concrete realizations  $\mathbf{x}_1, \dots, \mathbf{x}_N, \mathbf{x}_n \in \mathbb{R}^D$
- ▶ **Empirical variance** (estimate of the true variance):

$$\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^\top \in \mathbb{R}^{D \times D}$$

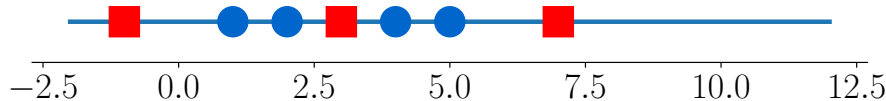
# Empirical Variance



$$\mathcal{D}_1 = \{1, 2, 4, 5\}, \quad \mathcal{D}_2 = \{-1, 3, 7\}$$

**Compute the empirical variances for both datasets**

# Empirical Variance

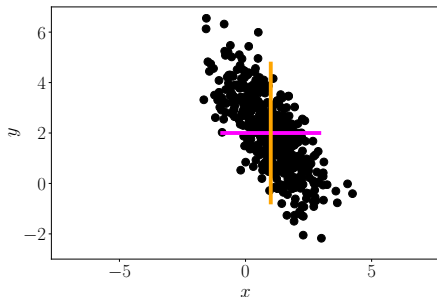


$$\mathcal{D}_1 = \{1, 2, 4, 5\}, \quad \mathcal{D}_2 = \{-1, 3, 7\}$$

**Compute the empirical variances for both datasets**

- ▶  $V[\mathcal{D}_1] = 2.5$
- ▶  $V[\mathcal{D}_2] = 10.66$  ►  $\mathcal{D}_2$  is more spread (around the mean) than  $\mathcal{D}_1$
- ▶ **Standard deviation**  $\sqrt{V[\cdot]}$  describes the spread more naturally and possesses the same units as the mean

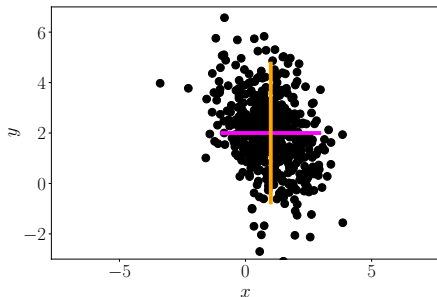
# Covariance and Cross-Covariance



- ▶ Variances along each axis remain constant, but properties of the dataset change
- ▶ Variances insufficient to characterize the relationship/correlation of two random variables

▶▶ Cross-covariance

# Covariance and Cross-Covariance

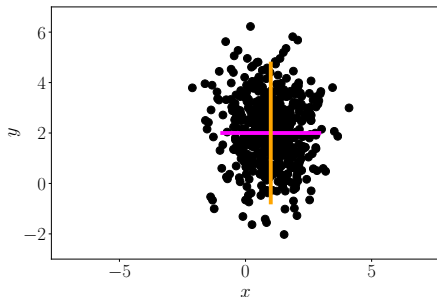


- ▶ Variances along each axis remain constant, but properties of the dataset change
- ▶ Variances insufficient to characterize the relationship/correlation of two random variables

▶▶ Cross-covariance



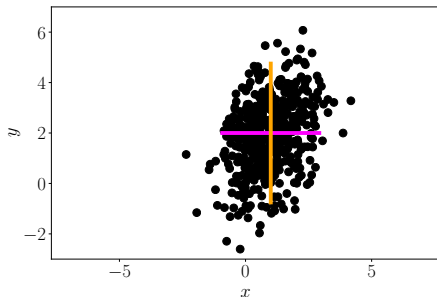
# Covariance and Cross-Covariance



- ▶ Variances along each axis remain constant, but properties of the dataset change
- ▶ Variances insufficient to characterize the relationship/correlation of two random variables

▶▶ Cross-covariance

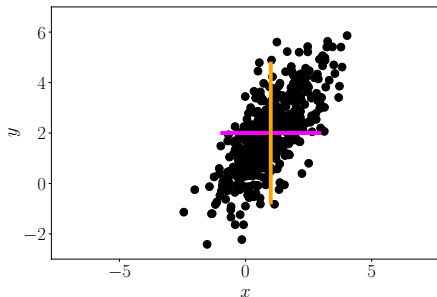
# Covariance and Cross-Covariance



- ▶ Variances along each axis remain constant, but properties of the dataset change
- ▶ Variances insufficient to characterize the relationship/correlation of two random variables

▶▶ Cross-covariance

# Covariance and Cross-Covariance



- ▶ Variances along each axis remain constant, but properties of the dataset change
- ▶ Variances insufficient to characterize the relationship/correlation of two random variables

▶▶ Cross-covariance

## Cross-Covariance (2)

$x \in \mathbb{R}^D, y \in \mathbb{R}^E$ . Then

- ▶ Cross-covariance:

$$\text{Cov}[x, y] := \mathbb{E}_{x,y}[(x - \mu_x)(y - \mu_y)^\top]$$

## Cross-Covariance (2)

$\mathbf{x} \in \mathbb{R}^D, \mathbf{y} \in \mathbb{R}^E$ . Then

► Cross-covariance:

$$\text{Cov}[\mathbf{x}, \mathbf{y}] := \mathbb{E}_{\mathbf{x}, \mathbf{y}}[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^\top]$$

$$\text{Cov}[\mathbf{x}, \mathbf{y}] := \iint (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^\top p(\mathbf{x})p(\mathbf{y})d\mathbf{x}d\mathbf{y} \in \mathbb{R}^{D \times E}$$

if  $\mathbf{x}, \mathbf{y}$  are continuous

$$\text{Cov}[\mathbf{x}, \mathbf{y}] := \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_x)(\mathbf{y}_n - \boldsymbol{\mu}_y)^\top \in \mathbb{R}^{D \times E}$$

if  $\mathbf{x}, \mathbf{y}$  are discrete

## Cross-Covariance (2)

$\mathbf{x} \in \mathbb{R}^D, \mathbf{y} \in \mathbb{R}^E$ . Then

- ▶ Cross-covariance:

$$\text{Cov}[\mathbf{x}, \mathbf{y}] := \mathbb{E}_{\mathbf{x}, \mathbf{y}}[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^\top]$$

$$\text{Cov}[\mathbf{x}, \mathbf{y}] := \iint (\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^\top p(\mathbf{x})p(\mathbf{y})d\mathbf{x}d\mathbf{y} \in \mathbb{R}^{D \times E}$$

if  $\mathbf{x}, \mathbf{y}$  are continuous

$$\text{Cov}[\mathbf{x}, \mathbf{y}] := \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_x)(\mathbf{y}_n - \boldsymbol{\mu}_y)^\top \in \mathbb{R}^{D \times E}$$

if  $\mathbf{x}, \mathbf{y}$  are discrete

- ▶  $\text{V}[\mathbf{x}] = \text{Cov}[\mathbf{x}, \mathbf{x}] \in \mathbb{R}^{D \times D}$
- ▶  $\text{Cov}[\mathbf{x}, \mathbf{y}] = \text{Cov}[\mathbf{y}, \mathbf{x}]^\top \in \mathbb{R}^{D \times E}$

# Covariance Matrix

- ▶ Random variable

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix} \in \mathbb{R}^D$$

- ▶ Variance of this  $D$ -dimensional random variable is given by a **covariance matrix**

$$\mathbb{V}_x[\mathbf{x}] = \begin{bmatrix} \mathbb{V}[x_1] & \text{Cov}[x_1, x_2] & \cdots & \text{Cov}[x_1, x_D] \\ \text{Cov}[x_2, x_1] & \mathbb{V}[x_2] & \cdots & \text{Cov}[x_2, x_D] \\ \vdots & & \ddots & \vdots \\ \text{Cov}[x_D, x_1] & \cdots & & \mathbb{V}[x_D] \end{bmatrix} \in \mathbb{R}^{D \times D}$$

- ▶ Covariance matrix is **symmetric, positive (semi-)definite**

# Mean and (Co)Variance

Mean and (co)variance are often useful to describe properties of data distributions (expected values and spread).

## Summary

$$\mathbb{E}_x[\mathbf{x}] = \int \mathbf{x}p(\mathbf{x})d\mathbf{x} =: \boldsymbol{\mu} \quad \left( \frac{1}{N} \sum_{n=1}^N x_n \right)$$

$$\mathbb{V}_x[\mathbf{x}] = \mathbb{E}_x[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] = \mathbb{E}_x[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top =: \boldsymbol{\Sigma}$$

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{x,y}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}_x[\mathbf{x}]\mathbb{E}_y[\mathbf{y}]^\top$$



# Mean and (Co)Variance

Mean and (co)variance are often useful to describe properties of data distributions (expected values and spread).

## Summary

$$\mathbb{E}_x[\mathbf{x}] = \int \mathbf{x}p(\mathbf{x})d\mathbf{x} =: \boldsymbol{\mu} \quad \left( \frac{1}{N} \sum_{n=1}^N x_n \right)$$

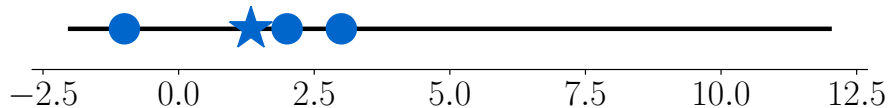
$$\mathbb{V}_x[\mathbf{x}] = \mathbb{E}_x[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] = \mathbb{E}_x[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top =: \boldsymbol{\Sigma}$$

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{x,y}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}_x[\mathbf{x}]\mathbb{E}_y[\mathbf{y}]^\top$$

Compute the mean and (co)variance of the following datasets

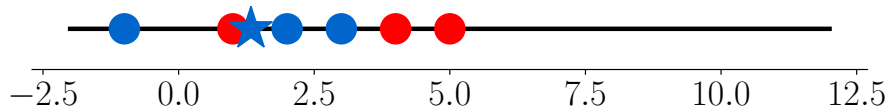
$$\mathcal{D}_1 := \{-2, -1, 2\} \quad \mathcal{D}_2 := \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 5 \\ 4 \end{bmatrix} \right\}$$

## Translation: Effect on the Mean



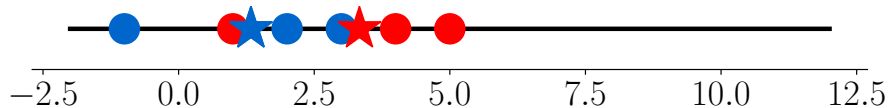
- ▶ What happens to the mean of a dataset if we shift/translate it by 2?

## Translation: Effect on the Mean



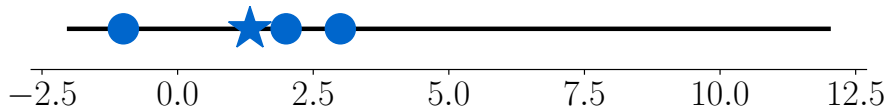
- ▶ What happens to the mean of a dataset if we shift/translate it by 2?

## Translation: Effect on the Mean



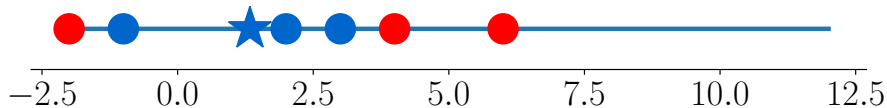
- ▶ What happens to the mean of a dataset if we shift/translate it by 2?

## Scaling: Effect on the Mean



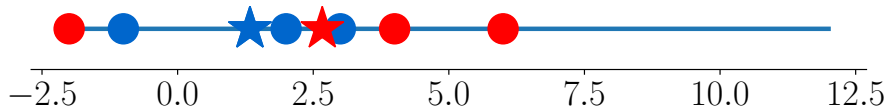
- ▶ What happens to the mean of a dataset if we scale it by 2?

## Scaling: Effect on the Mean



- ▶ What happens to the mean of a dataset if we scale it by 2?

## Scaling: Effect on the Mean



- ▶ What happens to the mean of a dataset if we scale it by 2?

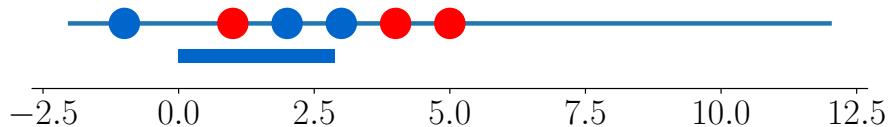
## Translation: Effect on the Variance



- ▶ What happens to the variance of a dataset if we shift it by 2?

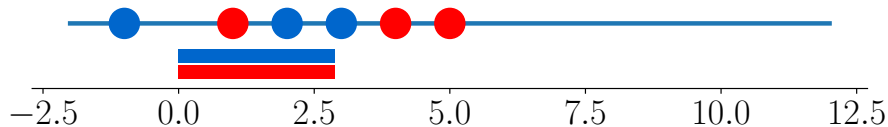


## Translation: Effect on the Variance



- What happens to the variance of a dataset if we shift it by 2?

## Translation: Effect on the Variance



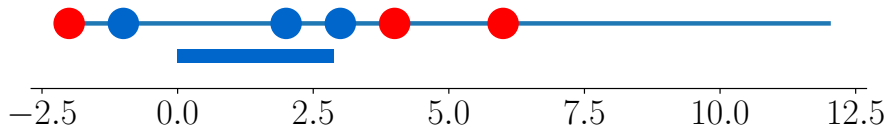
- ▶ What happens to the variance of a dataset if we shift it by 2?

## Scaling: Effect on the Variance



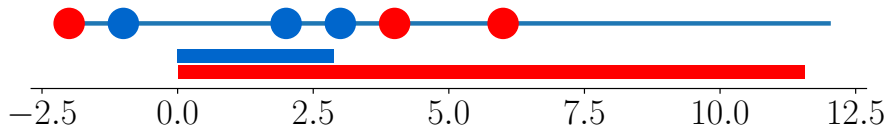
- ▶ What happens to the variance of a dataset if we scale it by 2?

## Scaling: Effect on the Variance



- ▶ What happens to the variance of a dataset if we scale it by 2?

## Scaling: Effect on the Variance



- ▶ What happens to the variance of a dataset if we scale it by 2?

# Linear/Affine Transformations

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}, \quad \text{where } \mathbb{E}_x[\mathbf{x}] = \boldsymbol{\mu}, \mathbb{V}_x[\mathbf{x}] = \boldsymbol{\Sigma}$$
$$\mathbb{E}[\mathbf{y}] =$$
$$\mathbb{V}[\mathbf{y}] =$$

# Linear/Affine Transformations

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}, \quad \text{where } \mathbb{E}_x[\mathbf{x}] = \boldsymbol{\mu}, \mathbb{V}_x[\mathbf{x}] = \boldsymbol{\Sigma}$$
$$\mathbb{E}[\mathbf{y}] = \mathbb{E}_x[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}_x[\mathbf{x}] + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$$
$$\mathbb{V}[\mathbf{y}] =$$

# Linear/Affine Transformations

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}, \quad \text{where } \mathbb{E}_x[\mathbf{x}] = \boldsymbol{\mu}, \mathbb{V}_x[\mathbf{x}] = \boldsymbol{\Sigma}$$
$$\mathbb{E}[\mathbf{y}] = \mathbb{E}_x[\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\mathbb{E}_x[\mathbf{x}] + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$$
$$\mathbb{V}[\mathbf{y}] =$$



# Linear/Affine Transformations

$$\begin{aligned} \mathbf{y} &= \mathbf{Ax} + \mathbf{b}, & \text{where } \mathbb{E}_x[\mathbf{x}] &= \boldsymbol{\mu}, \mathbb{V}_x[\mathbf{x}] = \boldsymbol{\Sigma} \\ \mathbb{E}[\mathbf{y}] &= \mathbb{E}_x[\mathbf{Ax} + \mathbf{b}] = \mathbf{A}\mathbb{E}_x[\mathbf{x}] + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \\ \mathbb{V}[\mathbf{y}] &= \mathbb{V}_x[\mathbf{Ax} + \mathbf{b}] = \mathbb{V}_x[\mathbf{Ax}] = \mathbf{A}\mathbb{V}_x[\mathbf{x}]\mathbf{A}^\top = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top \end{aligned}$$

# Sum of Random Variables

$\mathbf{x}, \mathbf{y} \in \mathbb{R}^D$  random variables. Then

$$\mathbb{E}[\mathbf{x} \pm \mathbf{y}] = \mathbb{E}_x[\mathbf{x}] \pm \mathbb{E}_y[\mathbf{y}]$$

$$\mathbb{V}[\mathbf{x} \pm \mathbf{y}] = \mathbb{V}[\mathbf{x}] + \mathbb{V}[\mathbf{y}] \pm \text{Cov}[\mathbf{x}, \mathbf{y}] \pm \text{Cov}[\mathbf{y}, \mathbf{x}]$$