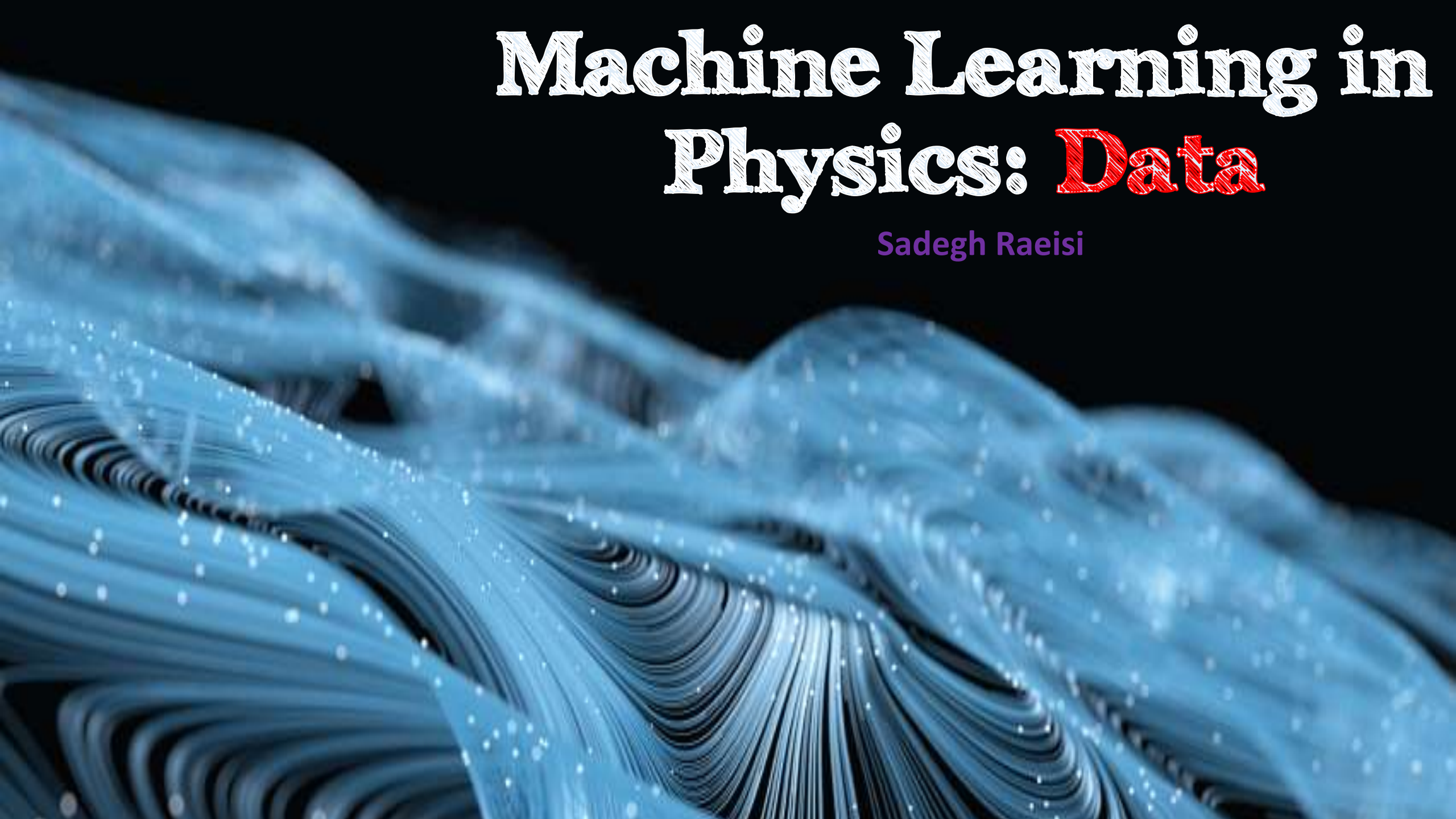
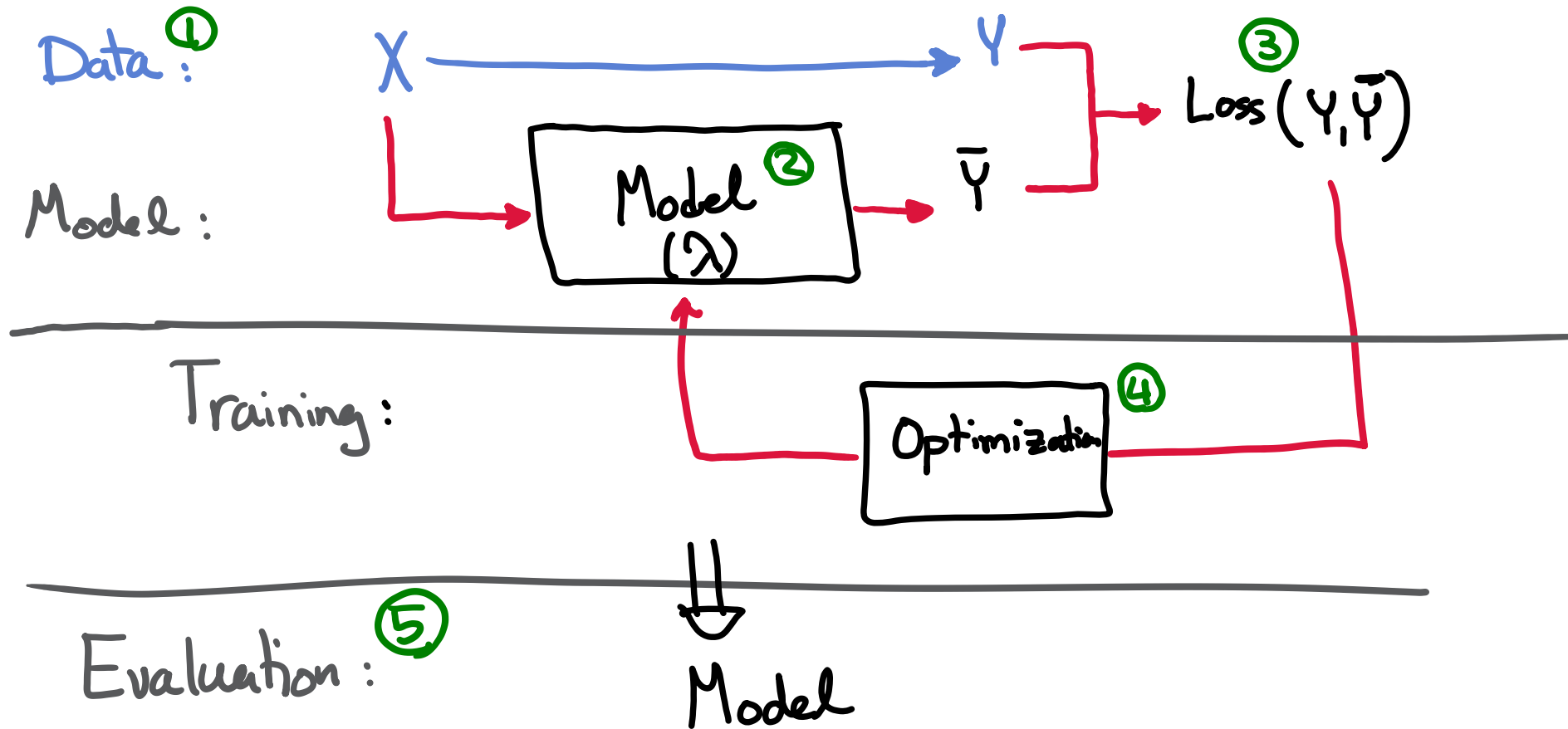


# Machine Learning in Physics: **Data**

Sadegh Raeisi



# Supervised: Ingredients



# Outline

Notation

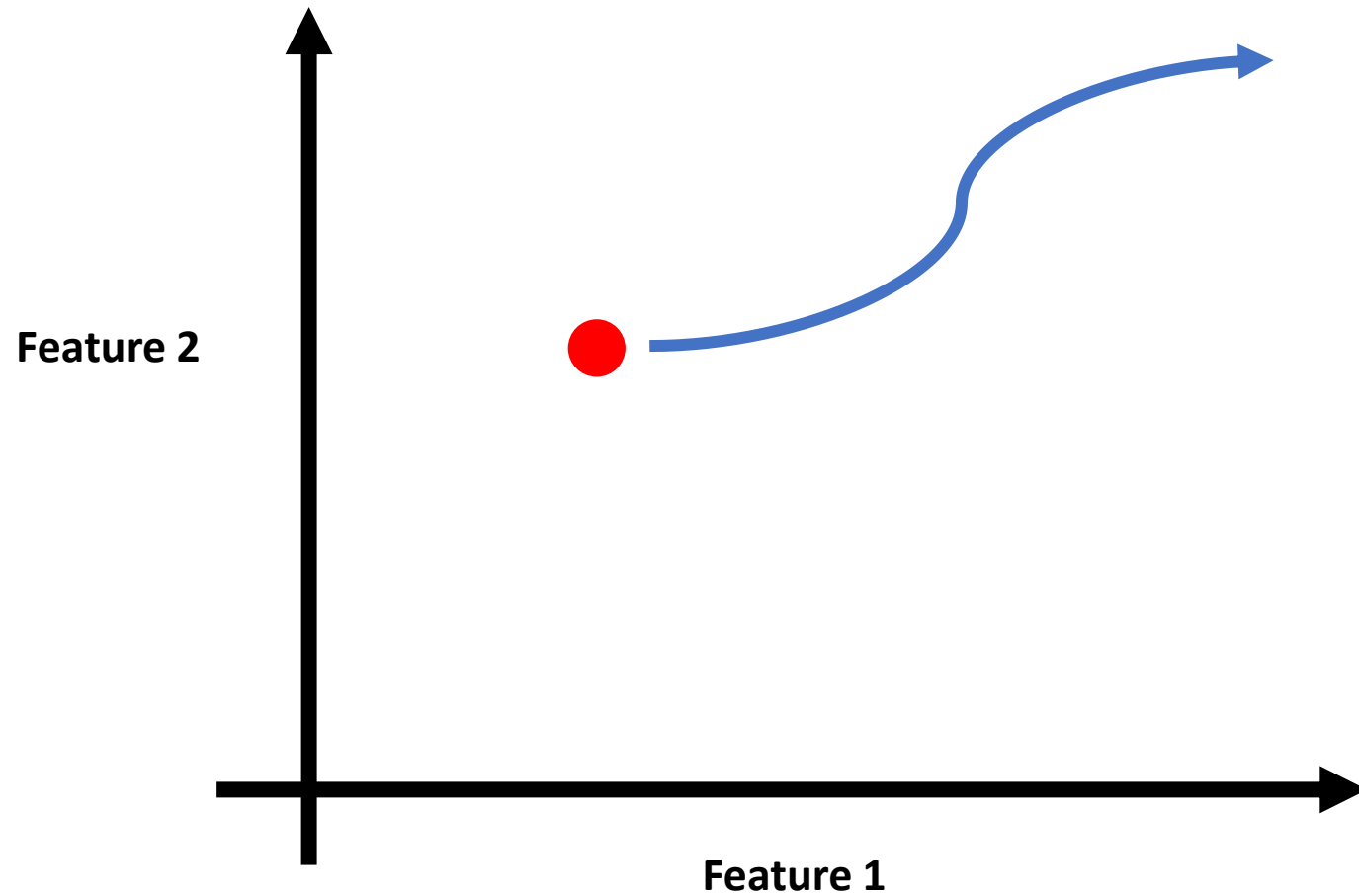
Types of Data

Encoding

Transformations and preprocessing

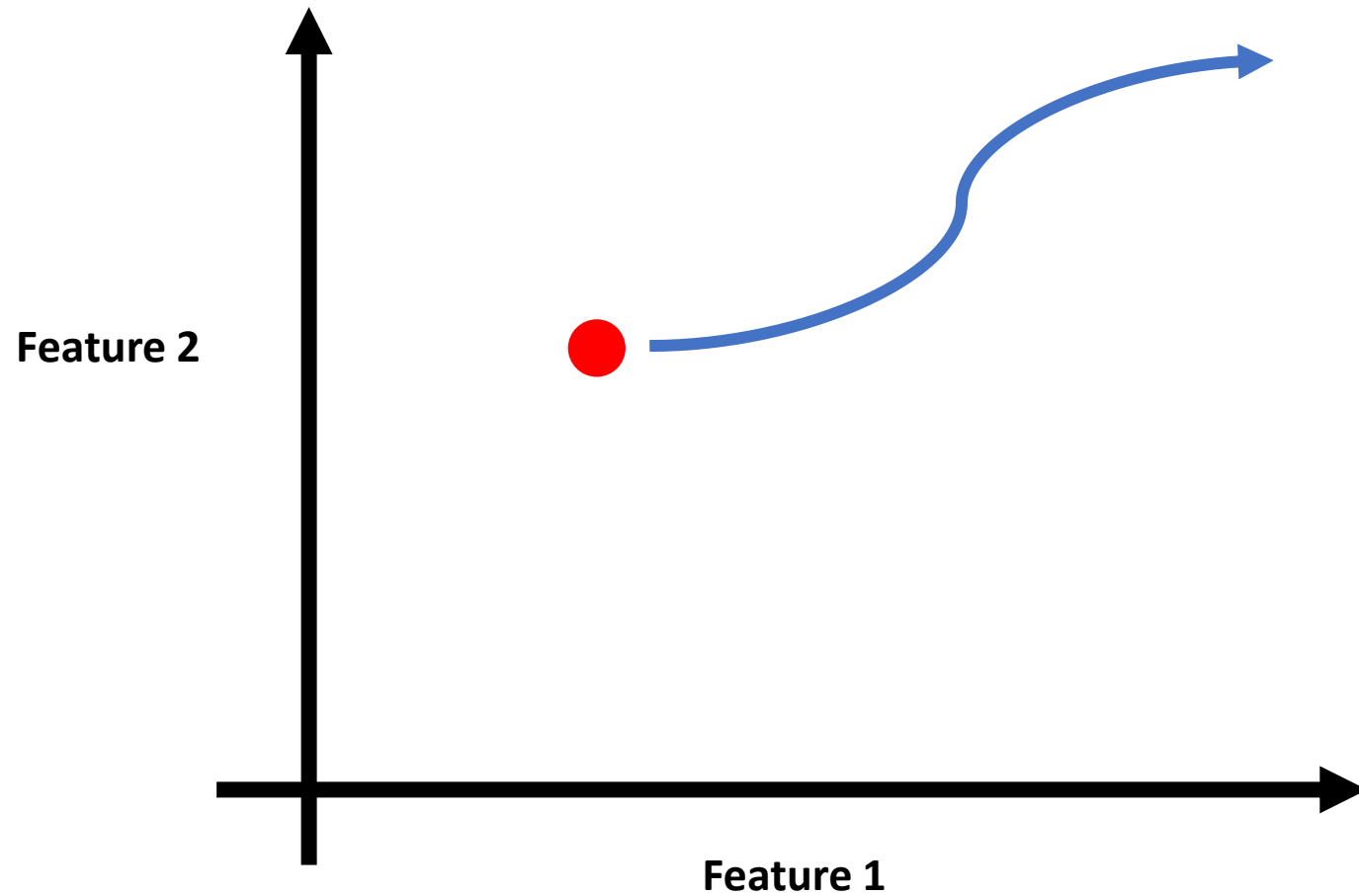
# Notation

# One sample



**Representation of the sample:  
(Feature 1, Feature 2)**

# One sample



**Representation of the sample:  
(Feature 1, Feature 2)**

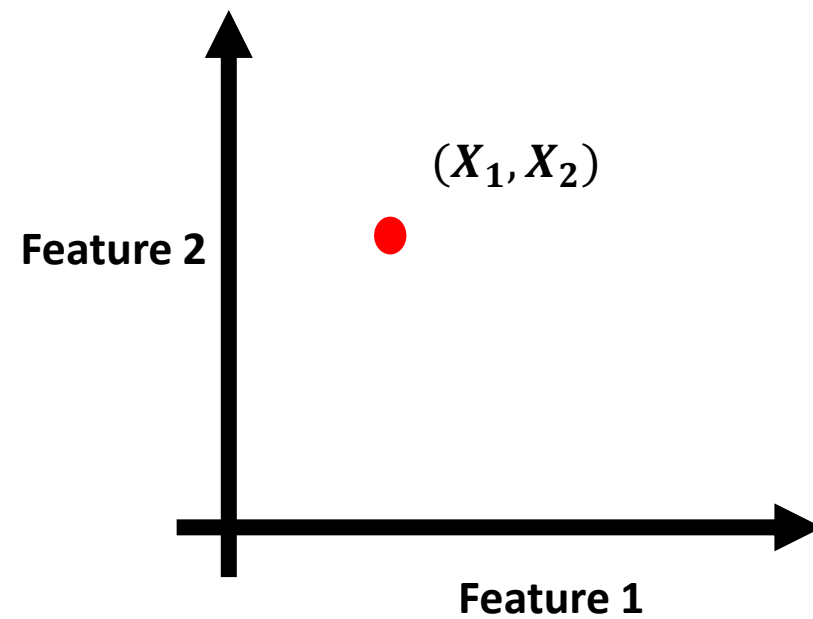


**Feature 1  $\Rightarrow X_1$   
Feature 2  $\Rightarrow X_2$**

One sample



$(X_1, X_2)$

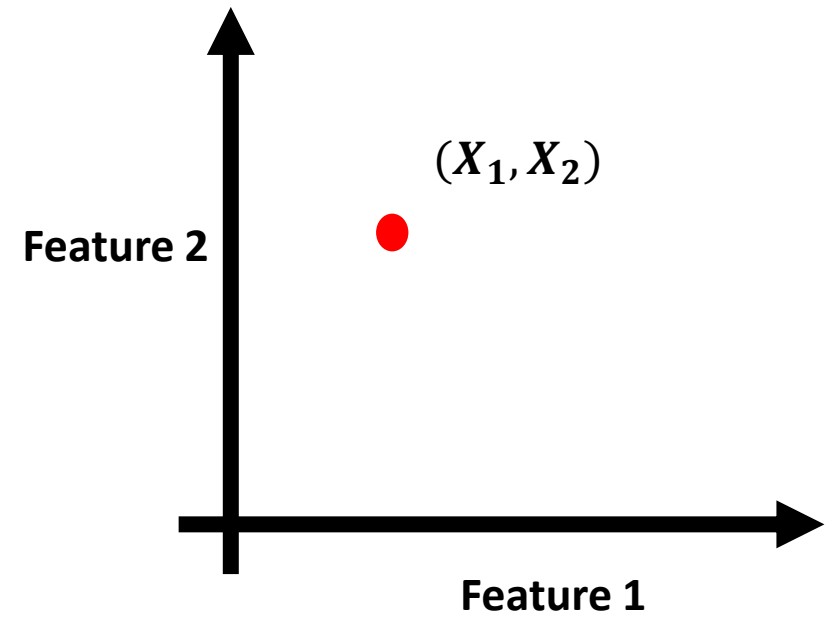


One sample



$$\vec{X} = (X_1, X_2, \dots, X_{nf})$$

*nf*: Dimension of the space





# Collection of samples

Sample 1:  $\vec{X}^1 = (X_1^1, X_2^1, \dots X_{nf}^1)$

Sample 2:  $\vec{X}^2 = (X_1^2, X_2^2, \dots X_{nf}^2)$

Sample ns:  $\vec{X}^{ns} = (X_1^{ns}, X_2^{ns}, \dots X_{nf}^{ns})$

$$X = \begin{pmatrix} X_1^1, X_2^1, \dots X_{nf}^1 \\ X_1^2, X_2^2, \dots X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots X_{nf}^{ns} \end{pmatrix}$$

# Data : Features

$$\mathbf{X} = \begin{pmatrix} X_1^1, X_2^1, \dots, X_{nf}^1 \\ X_1^2, X_2^2, \dots, X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots, X_{nf}^{ns} \end{pmatrix}$$

$nf$ : Dimension of the space

$ns$  : Number of Samples

Data : Labels

$$\mathbf{X} = \begin{pmatrix} X_1^1, X_2^1, \dots, X_{nf}^1 \\ X_1^2, X_2^2, \dots, X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots, X_{nf}^{ns} \end{pmatrix} \begin{matrix} \longrightarrow \\ \longrightarrow \\ \longrightarrow \\ \longrightarrow \end{matrix} \begin{pmatrix} y^1 \\ y^2 \\ \dots \\ y^{ns} \end{pmatrix}$$

$nf$ : Dimension of the space  
 $ns$ : Number of Samples

Data : Labels

$$\mathbf{X} = \begin{pmatrix} X_1^1, X_2^1, \dots, X_{nf}^1 \\ X_1^2, X_2^2, \dots, X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots, X_{nf}^{ns} \end{pmatrix}$$

$$\mathbf{Y} = \begin{pmatrix} y^1 \\ y^2 \\ \dots \\ y^{ns} \end{pmatrix}$$

$nf$ : Dimension of the space  
 $ns$ : Number of Samples

# Types of Data

Quantitative Data

Qualitative Data

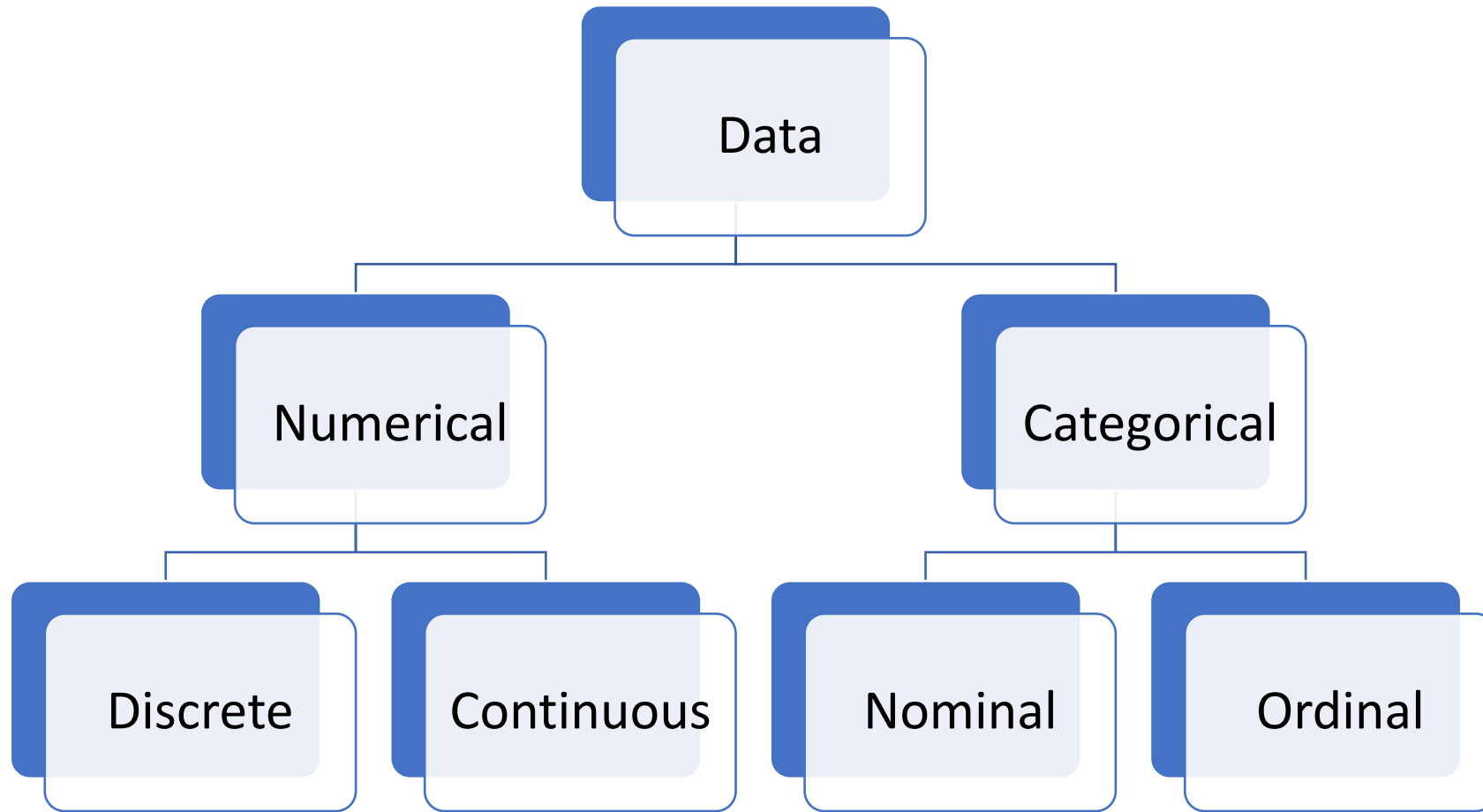
Experimental Data

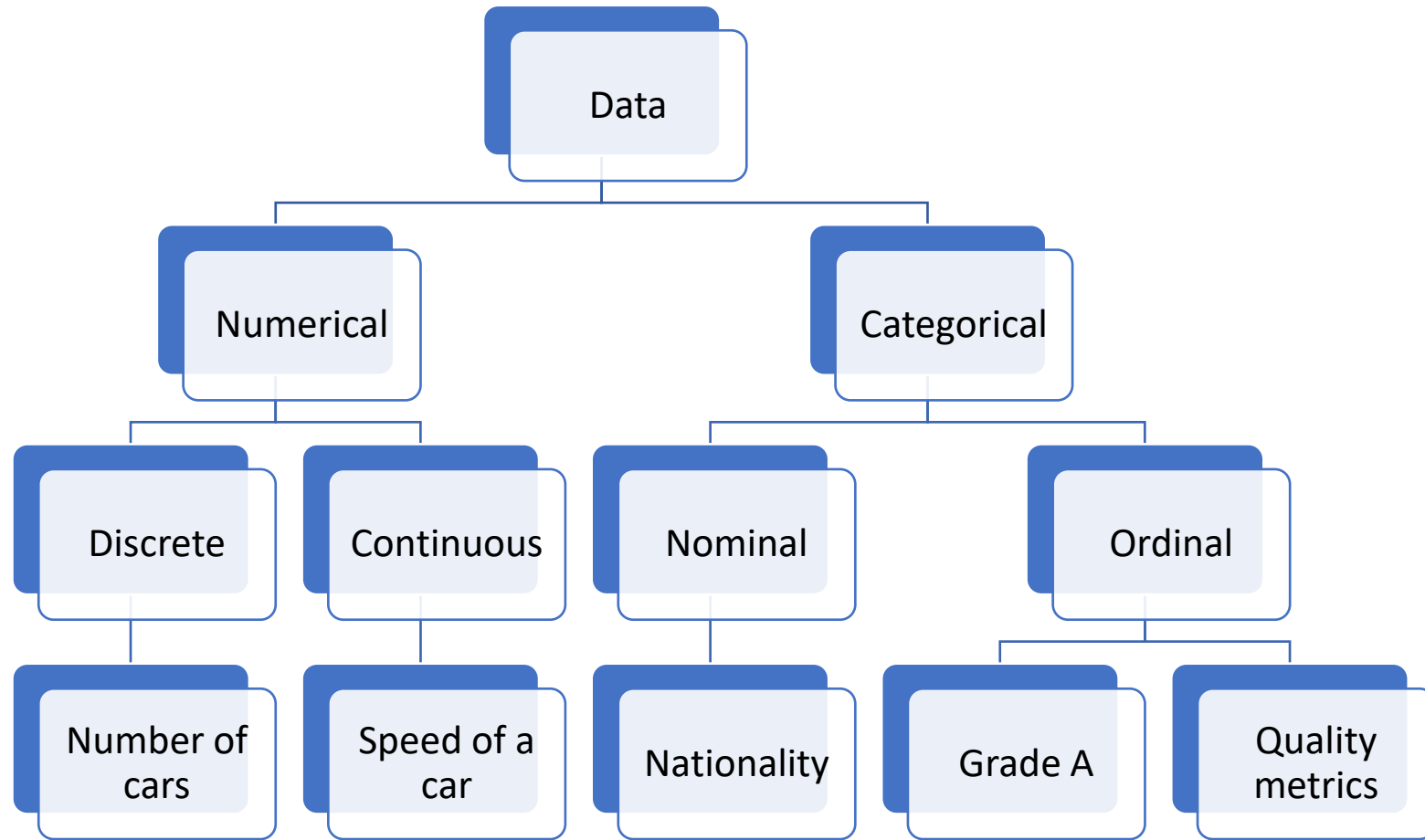
Observational Data

Survey Data

Archival Data

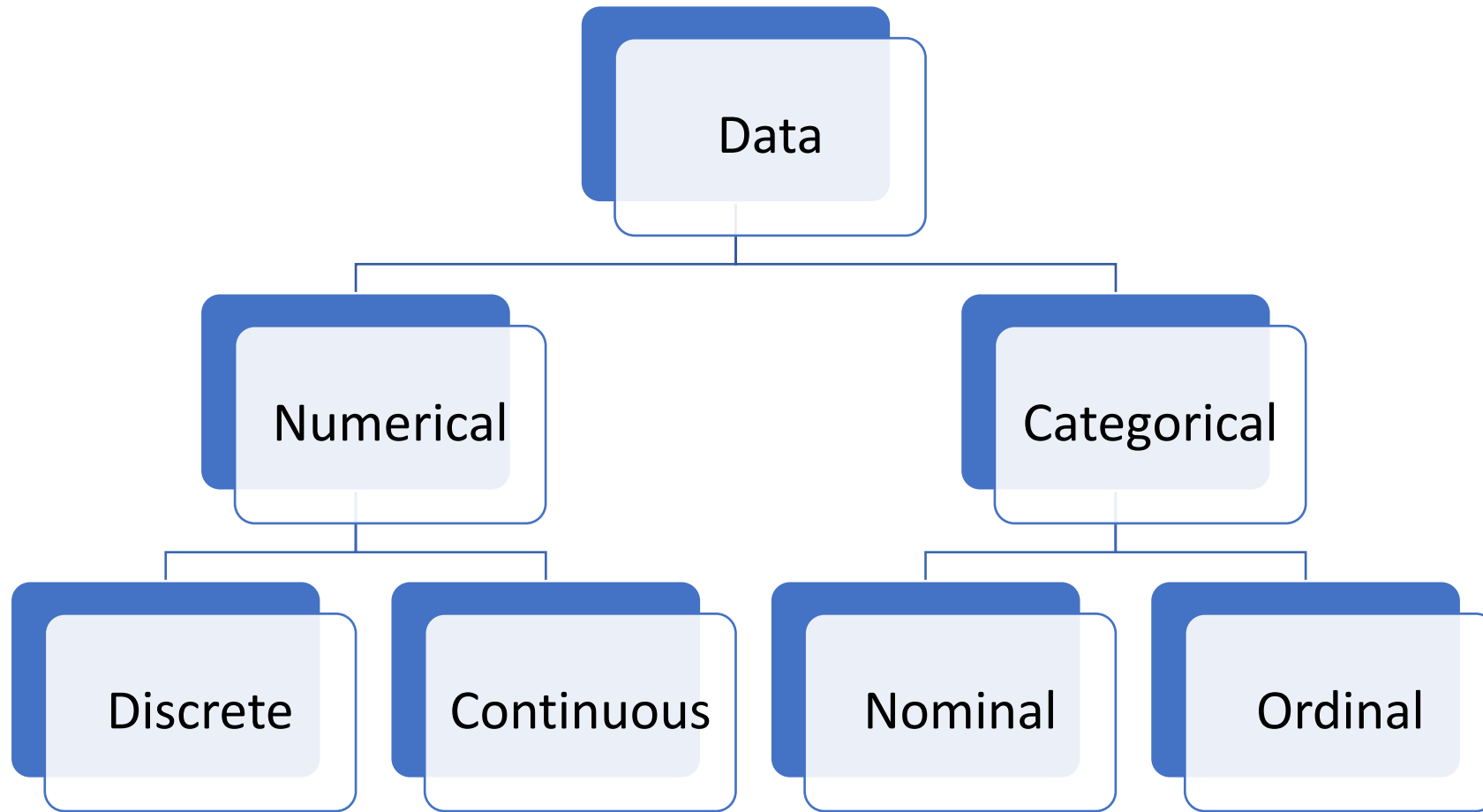
There are more than one categorization ...







# Encoding of Data



## Exercise

**How would you encode  
categorical types of data?**

# Processing the Data

# Missing data

# Missing data

$$\mathbf{X} = \begin{pmatrix} X_1^1, X_2^1, \dots, X_{nf}^1 \\ X_1^2, \textcolor{red}{NA}, \dots, X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots, X_{nf}^{ns} \end{pmatrix}$$

# Missing data

$$\mathbf{X} = \begin{pmatrix} X_1^1, X_2^1, \dots, X_{nf}^1 \\ X_1^2, \textcolor{red}{NA}, \dots, X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots, X_{nf}^{ns} \end{pmatrix}$$

**What can we do about the missing data?**

# Missing data

$$X = \begin{pmatrix} X_1^1, X_2^1, \dots, X_{nf}^1 \\ X_1^2, \text{NA}, \dots, X_{nf}^2 \\ \dots \\ X_1^{ns}, X_2^{ns}, \dots, X_{nf}^{ns} \end{pmatrix}$$

## 1. Get rid of the sample

- What's the disadvantage?

## 2. Get rid of the feature

- When does it make sense to do this?

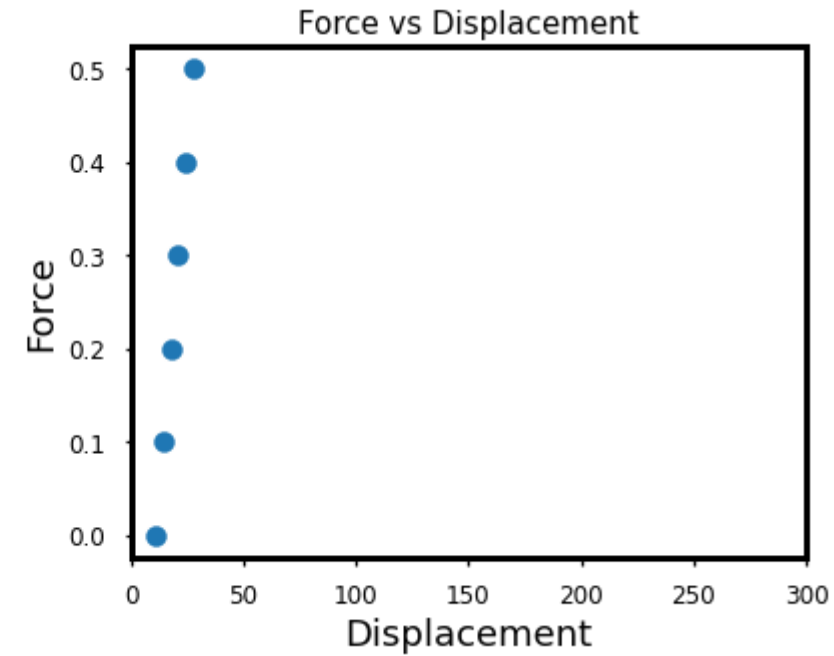
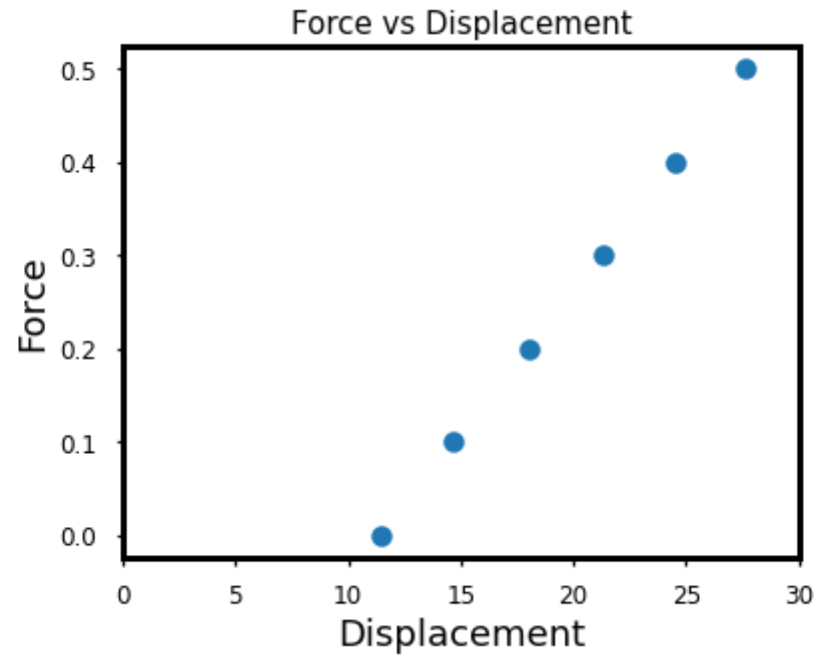
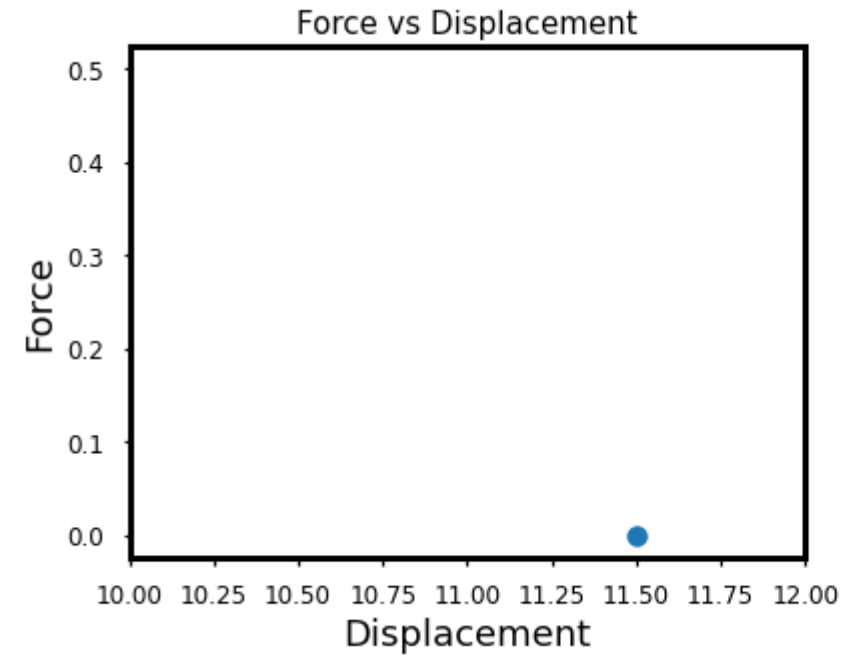
## 3. Assign a value to it?

- How?

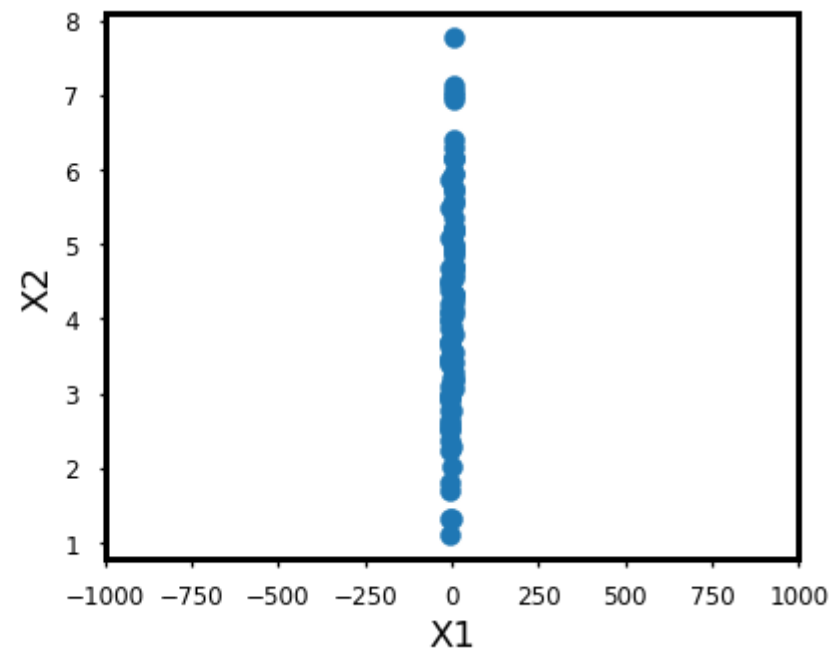
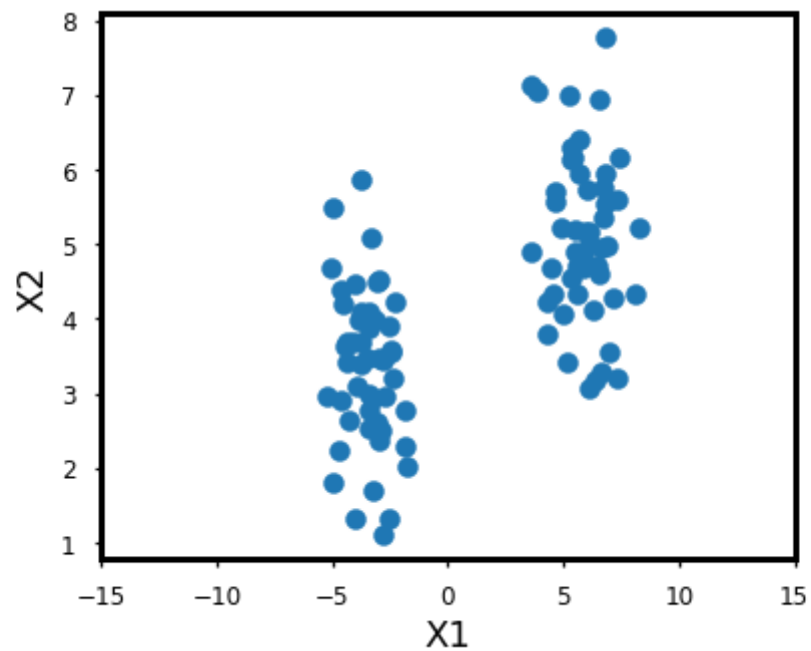
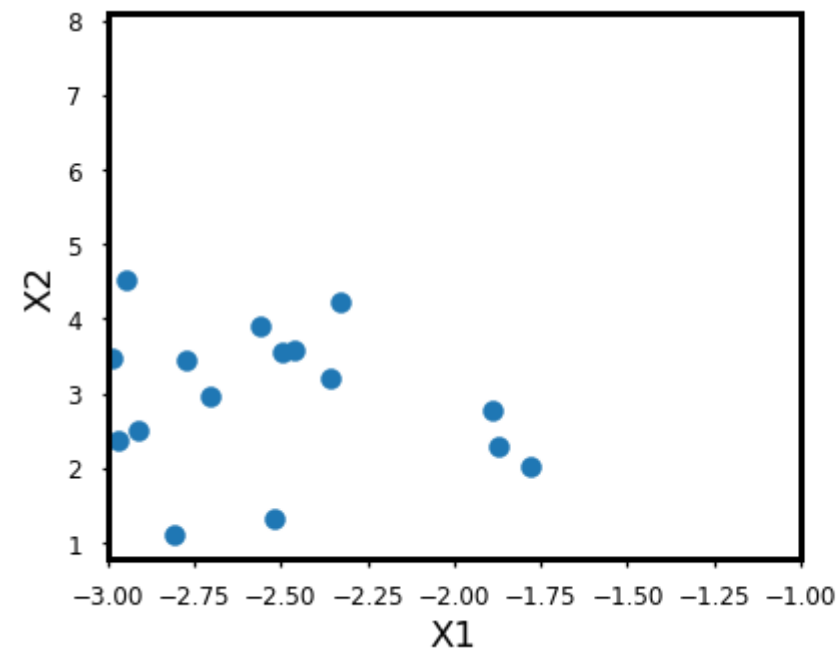


Scale of the data

# Scale of the data for regression



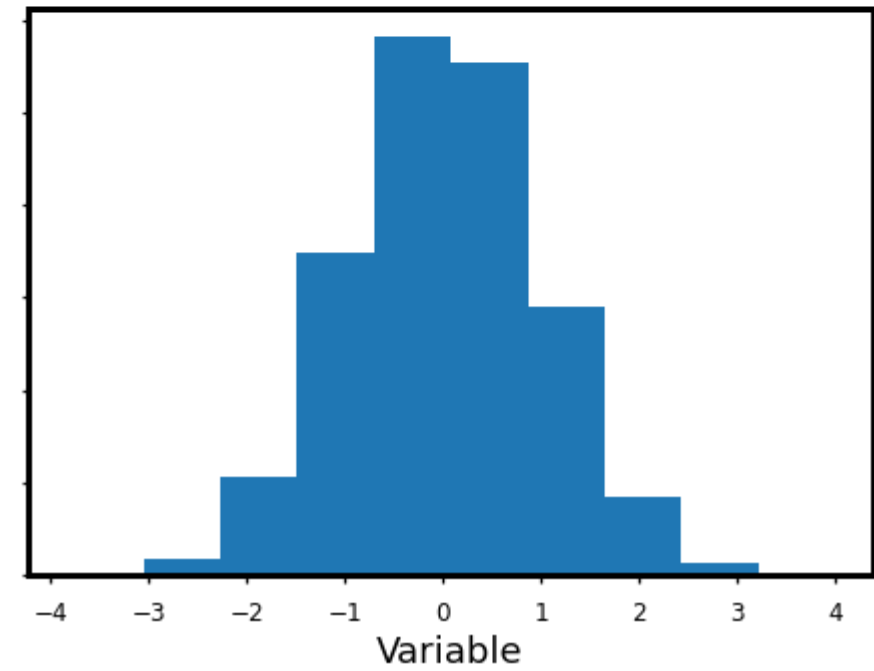
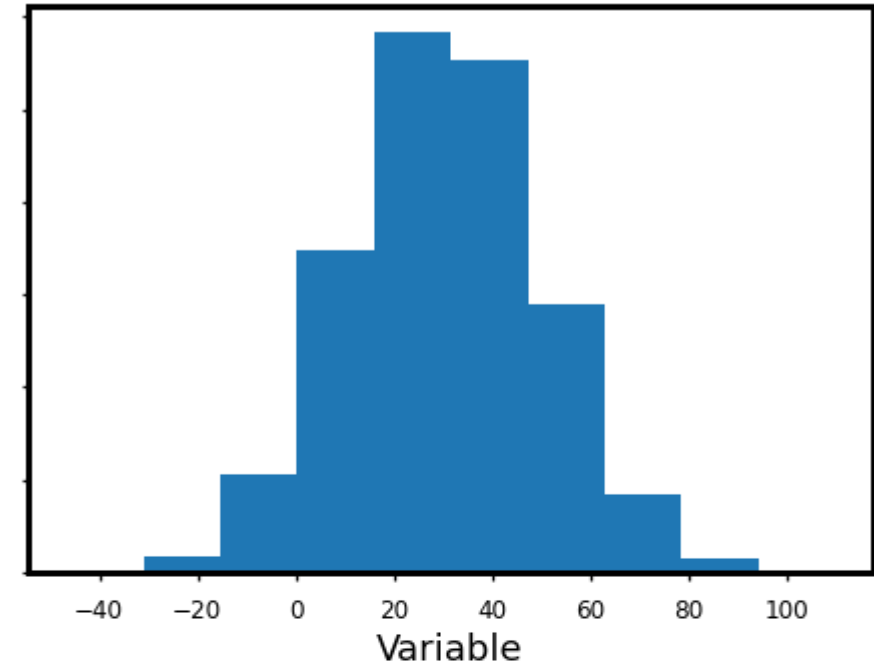
# Scale of the data: Clustering



# Re-scaling the data

1. Physical Scale

2. Scale on which it is changing



# Questions

- What other ways are there to scale the data? What are they good for?
- Does scaling the data affect optimization/training models?

# Data Reduction

# Why

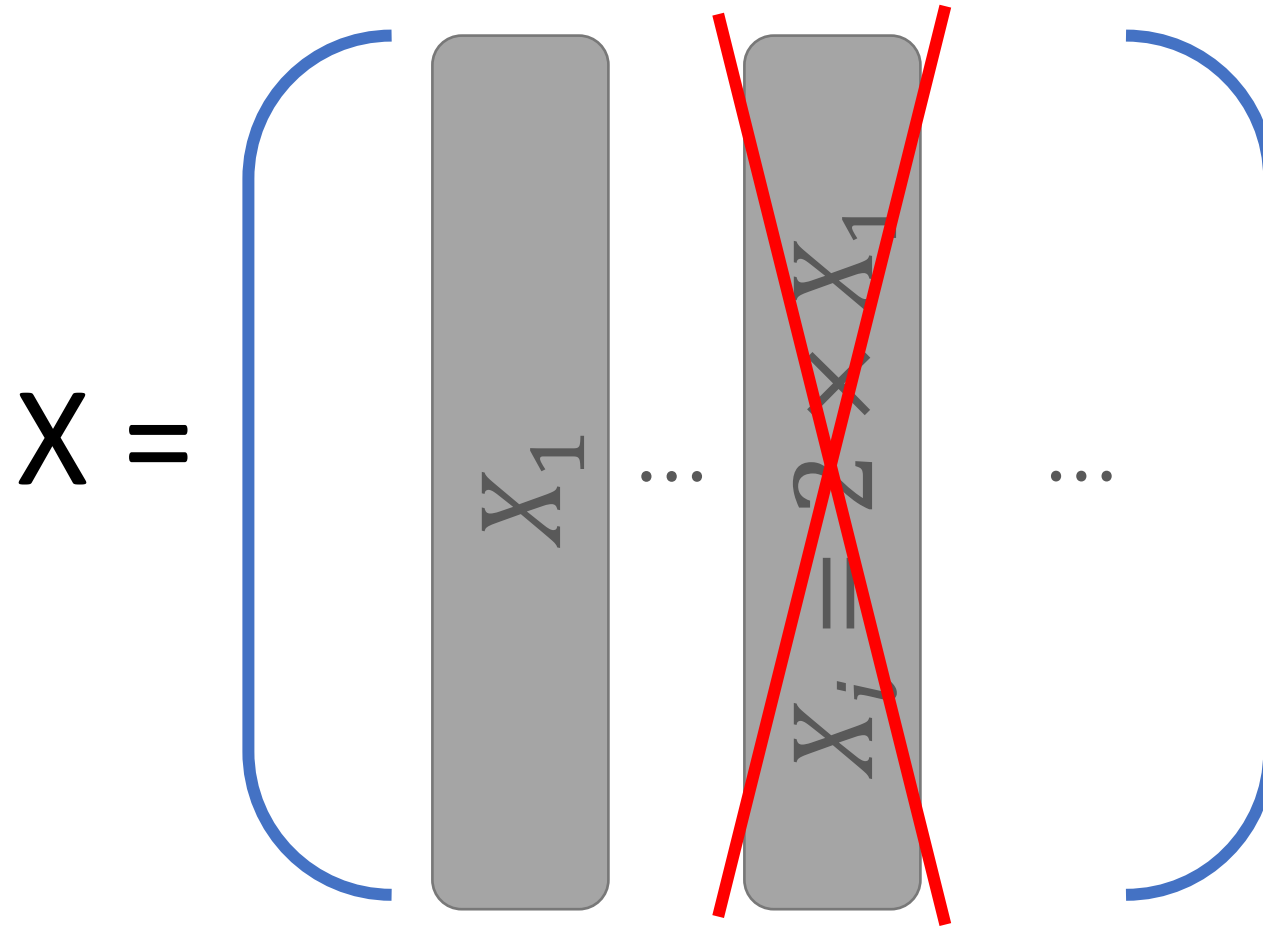
- Visualization
- Computational efficiency
- Curse of dimensionality

# Techniques

- Feature selection
- Feature transformation
  - Linear transformation
  - Manifold learning



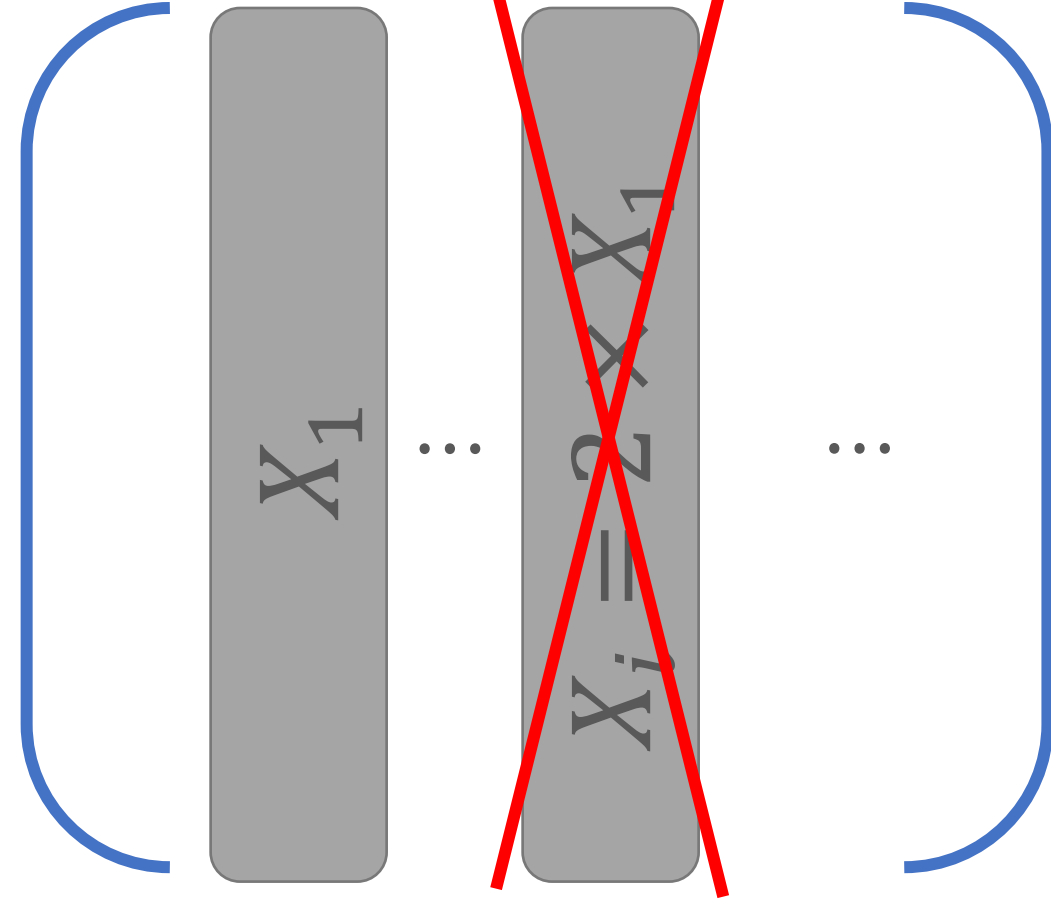
# Feature Selection



When?

# Feature Selection

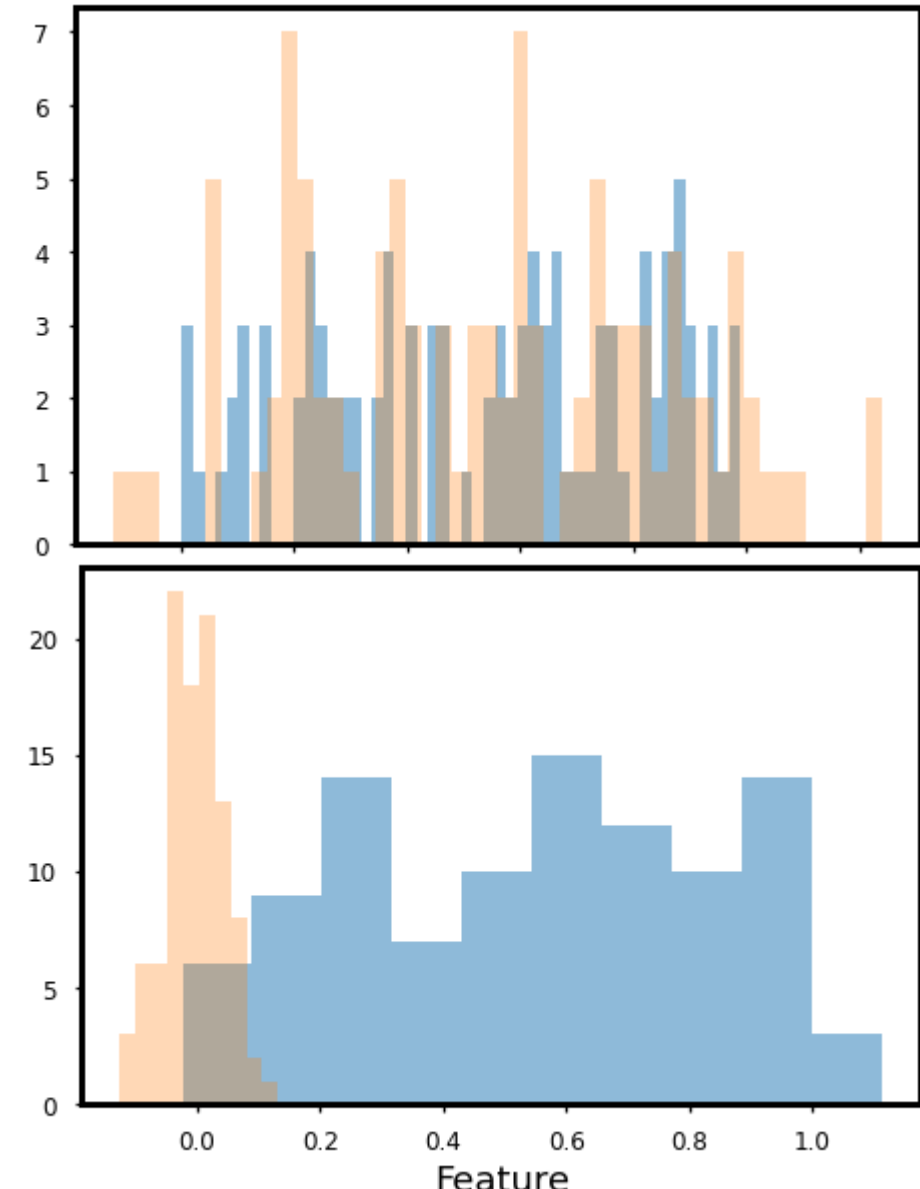
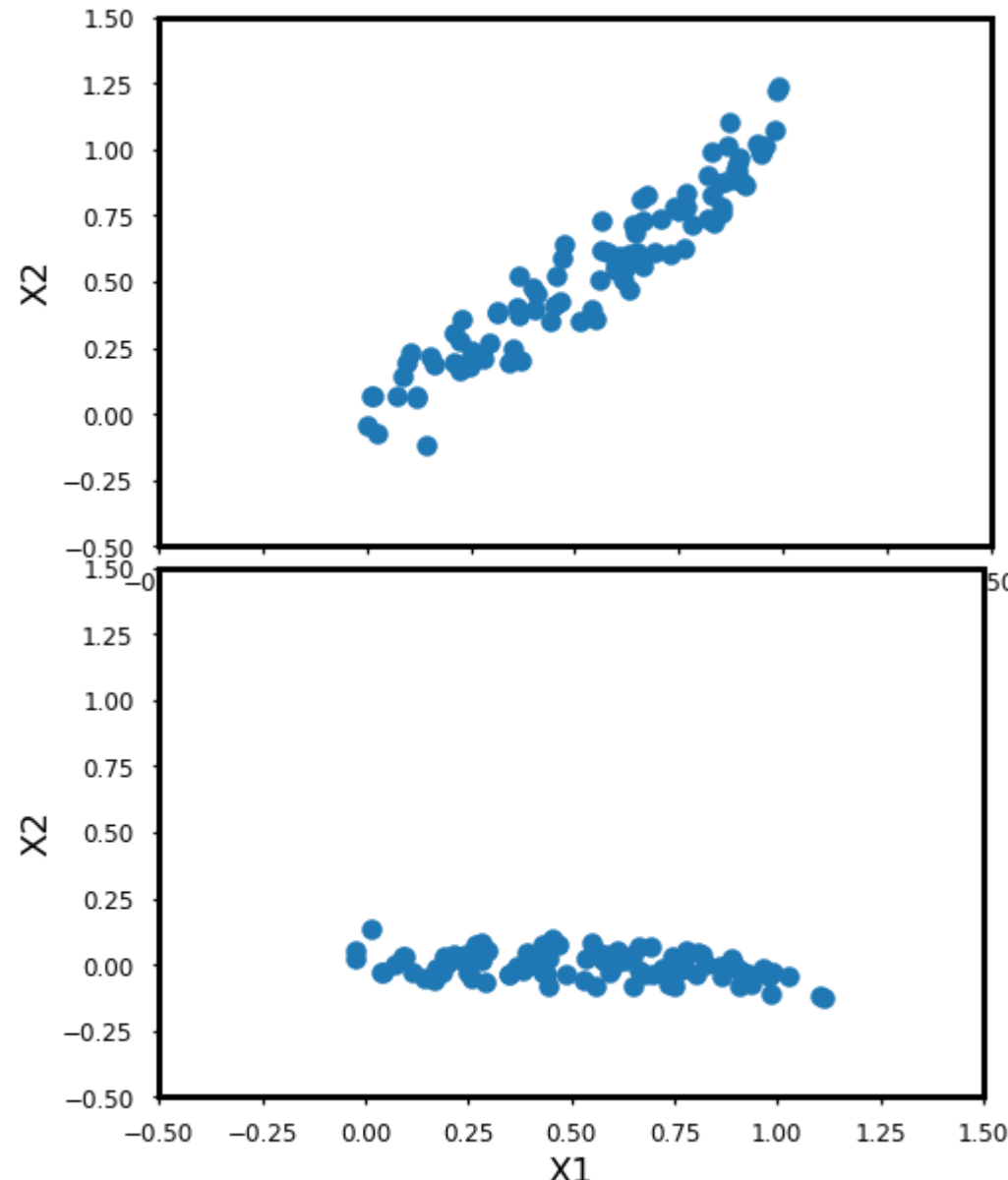
$X =$



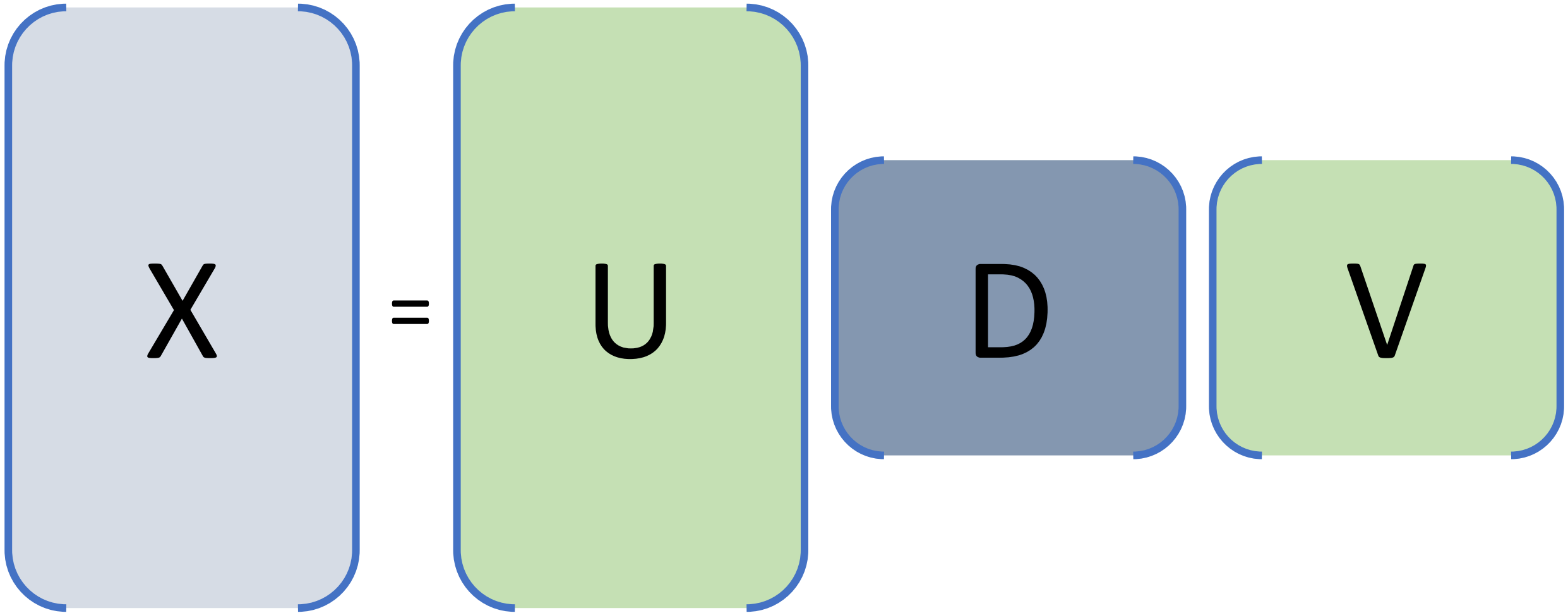
## How to:

- Un-Supervised
  - Variance
  - Importance/significance
- supervised

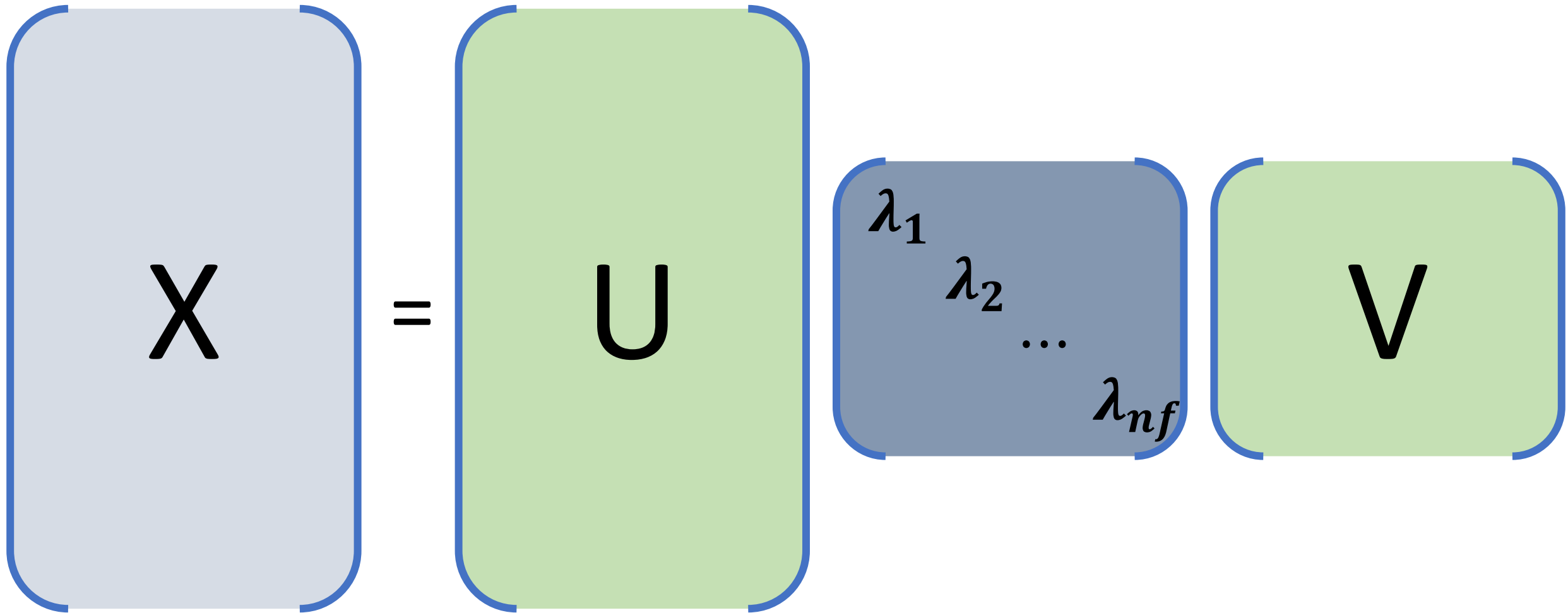
# Feature Transformation



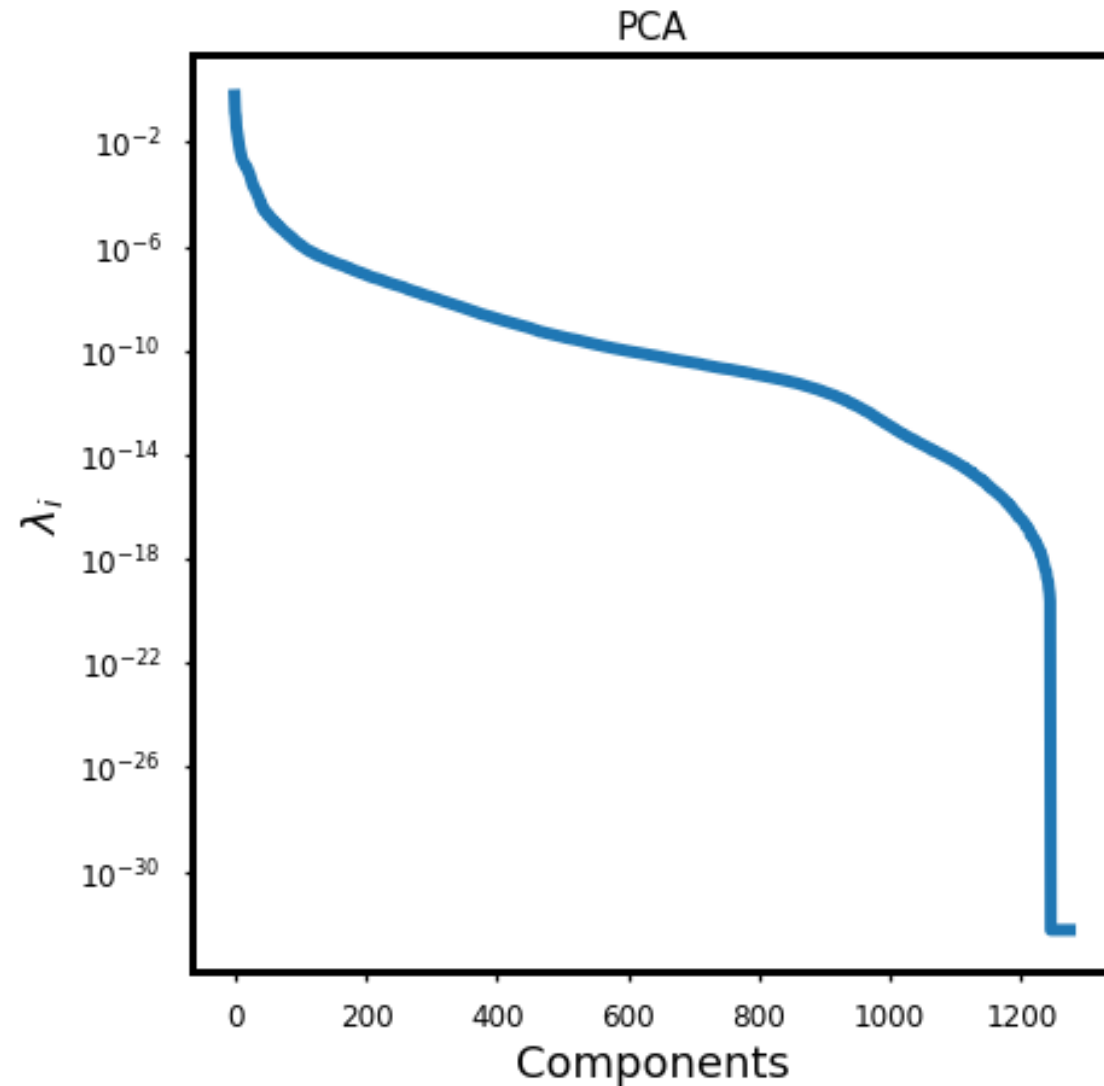
# Feature Transformation: PCA



# Feature Transformation: PCA

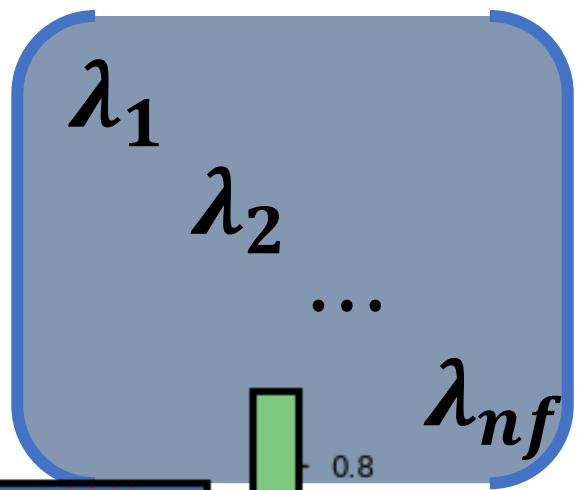
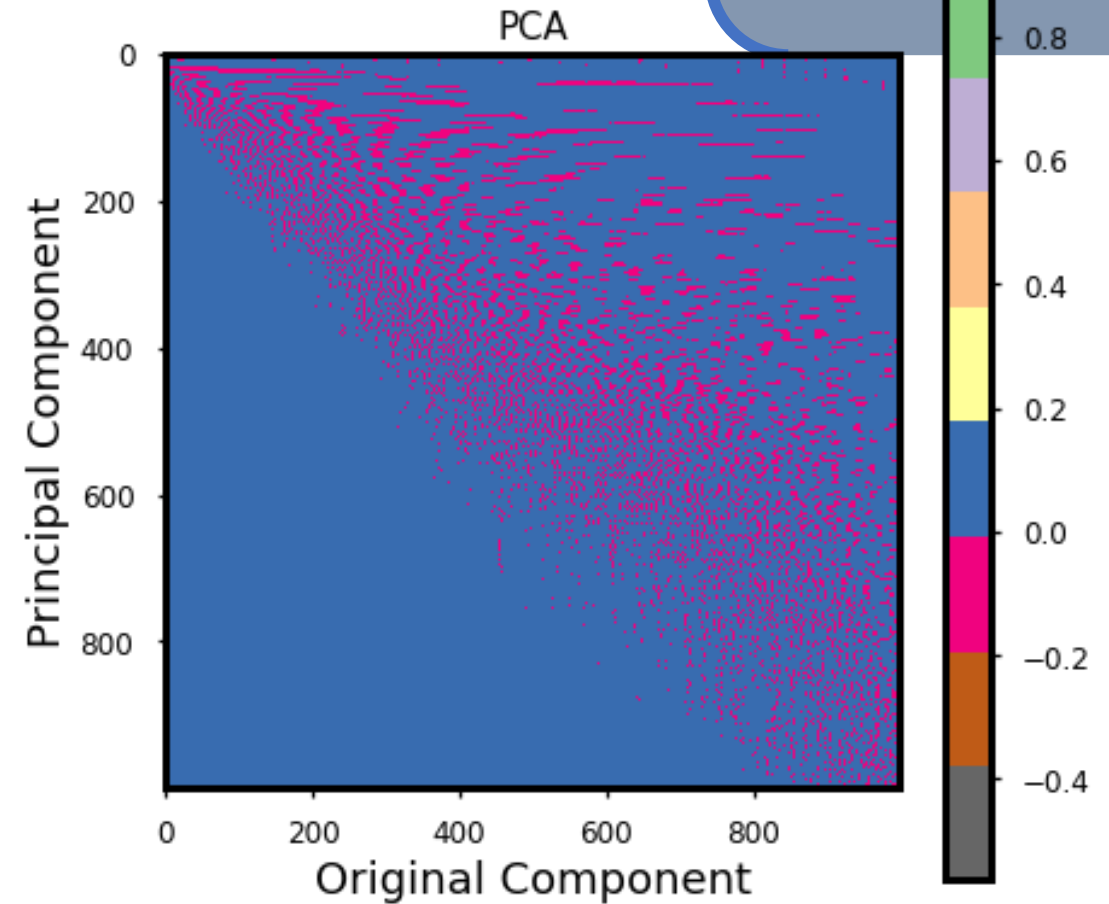
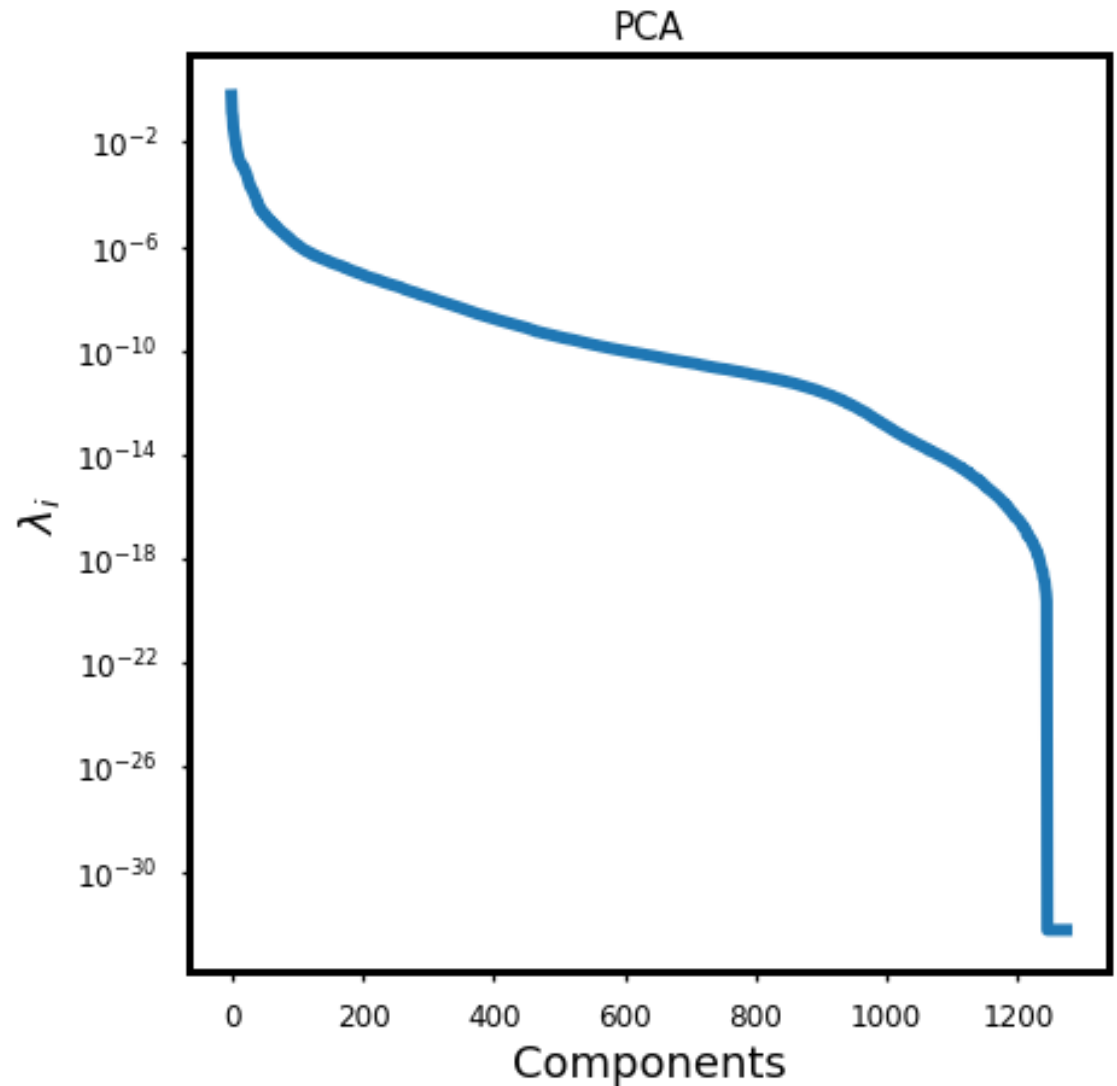


# Feature Transformation: PCA

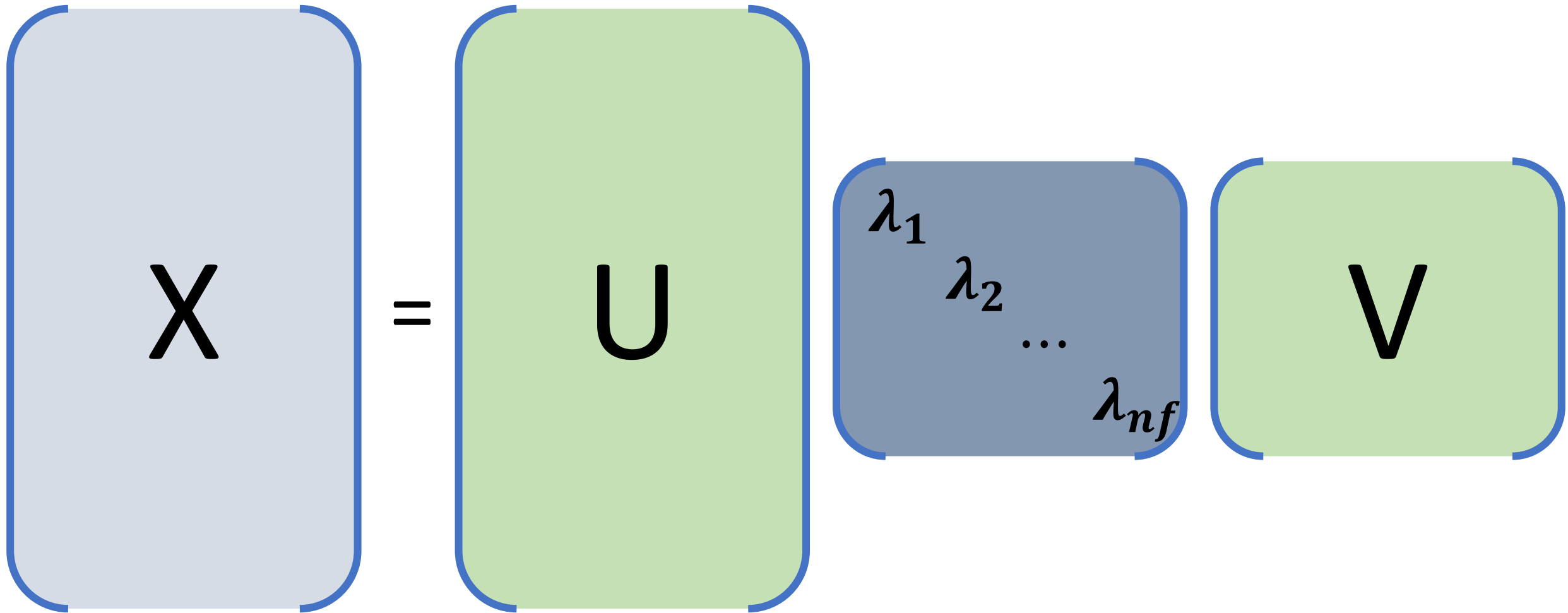


$$\lambda_1 \quad \lambda_2 \quad \dots \quad \lambda_{nf}$$

# Feature Transformation: PCA

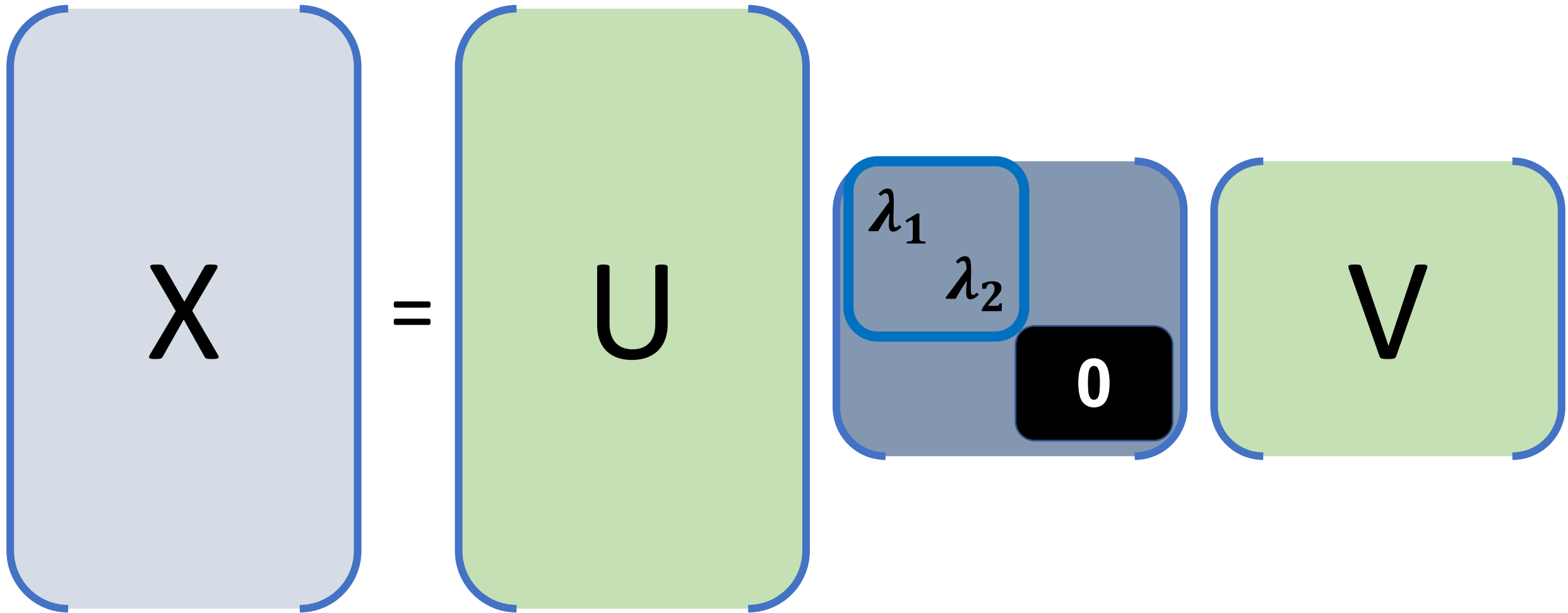


# Feature Transformation: PCA





# Feature Transformation: PCA



# Feature Transformation: PCA



How much  
**information**  
do we loose?

# Non-linear transformations: Manifold learning

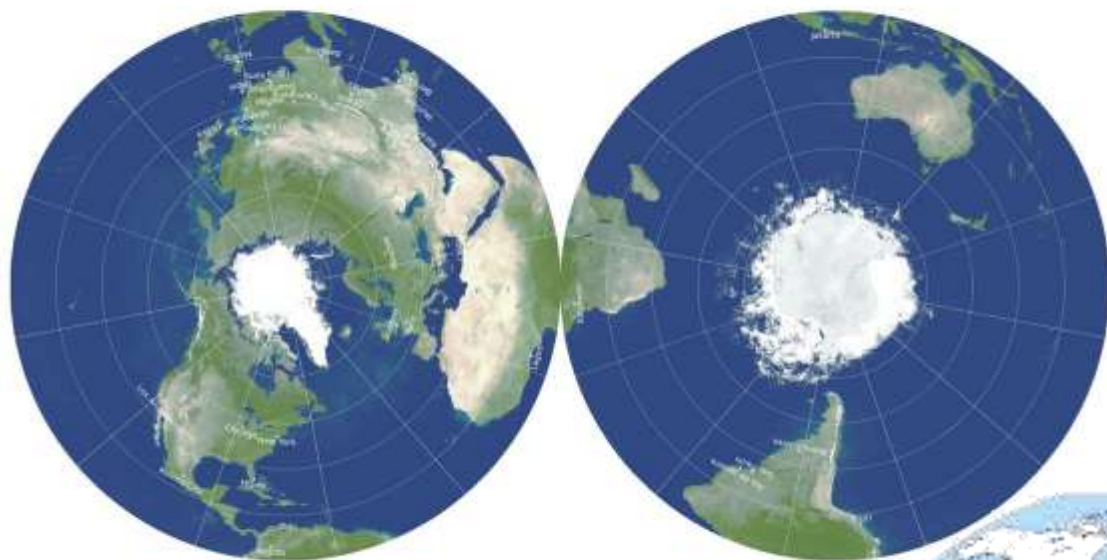


**In three dimension**



**2D reduction**

# Non-linear transformations: Manifold learning



**2D reduction**

# Summary

Notation

Types of Data

Encoding

Transformations and preprocessing

- Missing data
- Scaling the data
- Data reduction

# Supervised: Ingredients

