# FYS-2021 Machine Learning

Dimensionality reduction

**Slides: Stine Hansen**
**Guest Lecture: Elisabeth Wetzer**

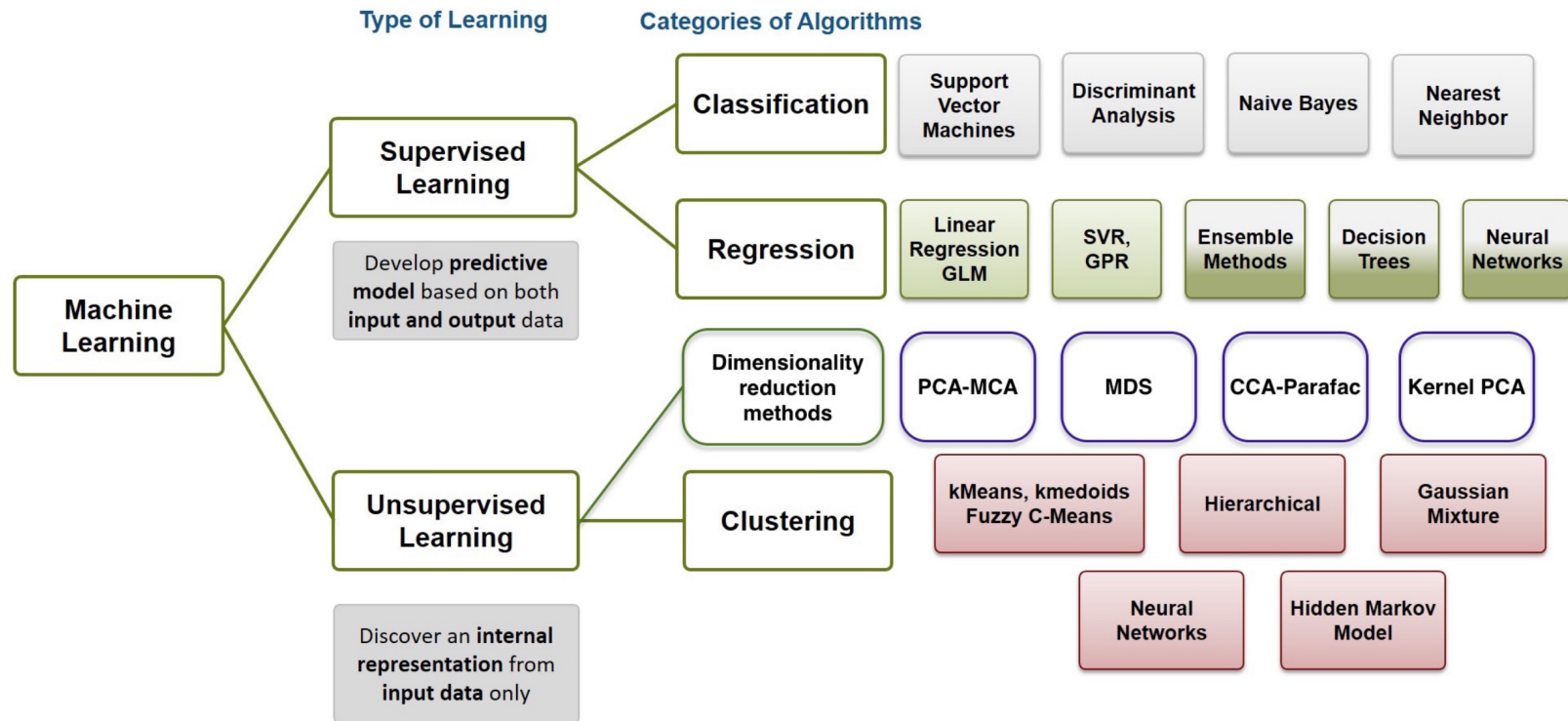UiT The Arctic University of Norway

# Roadmap

1. Dimensionality reduction: what and why
2. Inner product as similarity measure
3. Multidimensional scaling (from samples)
4. Practical example

   Break

5. Multidimensional scaling (from distances)
6. Practical examples
7. Forward selection
8. Practical example

# Machine Learning Taxonomy

**Type of Learning**

**Categories of Algorithms**

**Machine Learning**

**Supervised Learning**

Develop **predictive model** based on both **input and output** data

**Unsupervised Learning**

Discover an **internal representation** from **input data** only

**Classification**

| Support Vector Machines | Discriminant Analysis | Naive Bayes | Nearest Neighbor |

**Regression**

| Linear Regression GLM | SVR, GPR | Ensemble Methods | Decision Trees | Neural Networks |

**Dimensionality reduction methods**

| PCA-MCA | MDS | CCA-Parafac | Kernel PCA |

**Clustering**

| kMeans, kmedoids Fuzzy C-Means | Hierarchical | Gaussian Mixture |

| Neural Networks | Hidden Markov Model |

# Dimensionality reduction

Given a set of $n$ features, the goal in dimensionality reduction is to **reduce the number of features** to $d$ features $(d < n)$, while retaining a meaningful representation of the data.
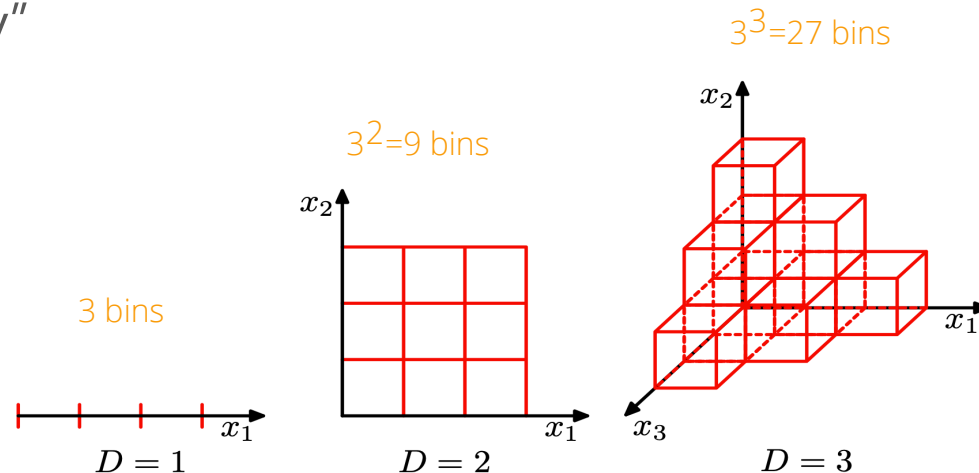
$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \xrightarrow{\text{feature selection}} \begin{bmatrix} X_q \\ X_r \\ \vdots \\ X_d \end{bmatrix} \quad \text{or} \quad \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \xrightarrow{\text{feature extraction}} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_d \end{bmatrix} = f\left( \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \right)$$

Selects a subset of important features

Extracts a set of new features as a function of the original features
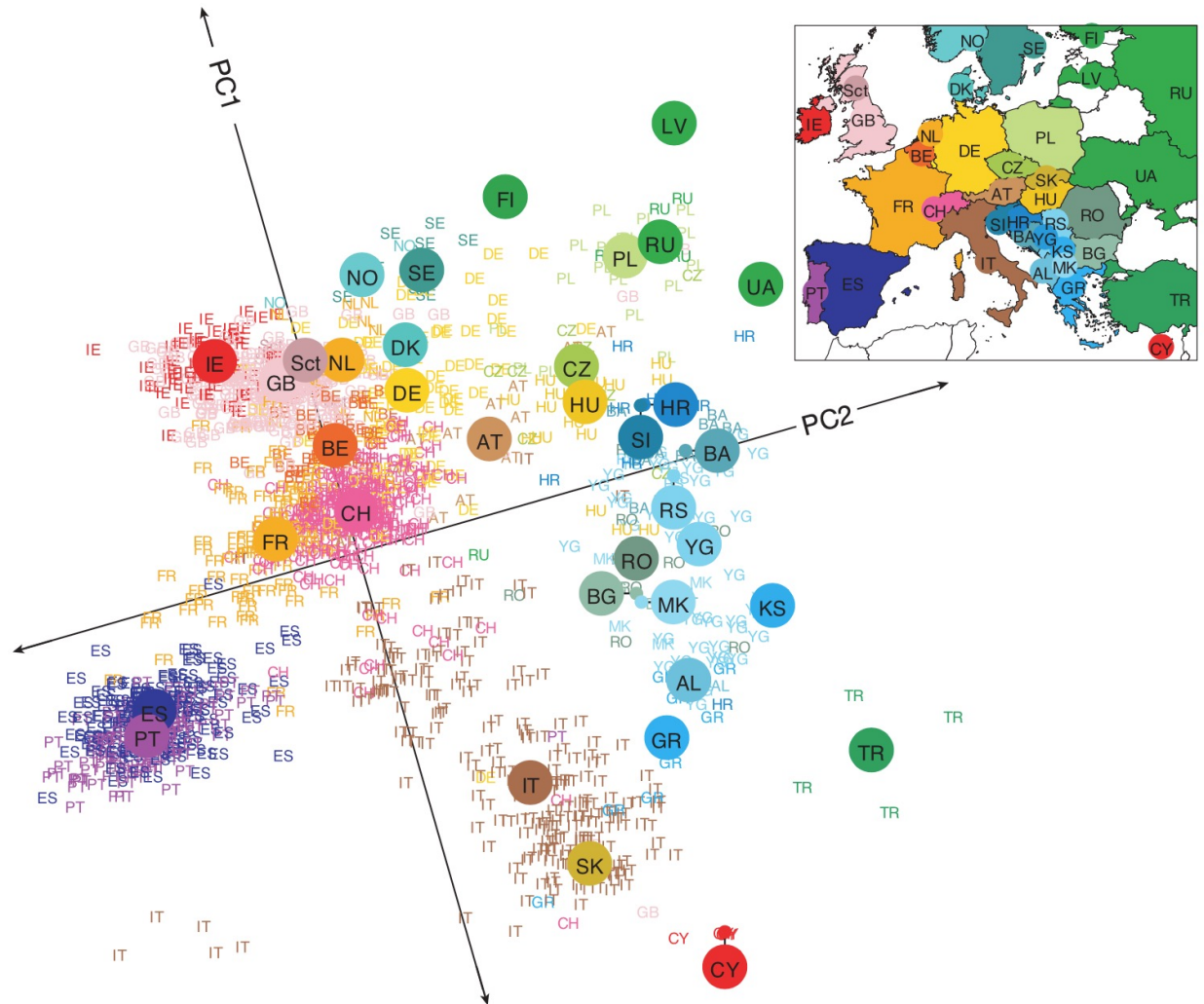
# Why dimensionality reduction?

- Compression:  Reduce computational complexity/memory usage

- Visualization: Enable visualization of high-dimensional data

- "Curse of dimensionality"

$3^3$=27 bins

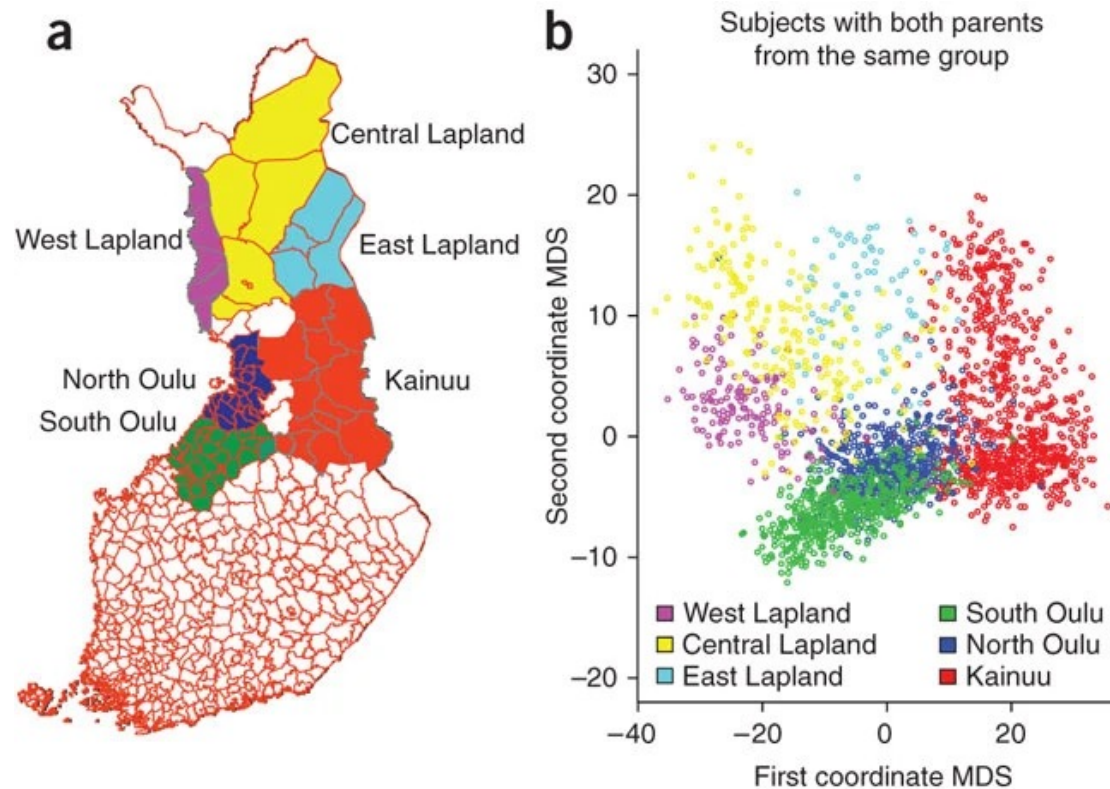$3^2$=9 bins

3 bins

$D = 1$

$D = 2$

$D = 3$

# example: DNA

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_{500\,000} \end{bmatrix} \longrightarrow \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}$$



Example from: Novembre, John, et al. "Genes mirror geography within Europe." Nature 456.7218 (2008): 98-101.

# example: DNA

# Dimensionality reduction methods

**Feature extraction**

- Principal component analysis (PCA)
- Linear discriminant analysis (LDA)
- Multidimensional scaling (MDS)
- t-distributed stochastic neighborhood embedding (t-SNE)
- Uniform manifold approximation and projection (UMAP)
- ...

**Feature selection**

- Forward selection
- Backward elimination
- Variance threshold
- L1 regularization
- ...

# Dimensionality reduction methods

## Feature extraction

- Principal component analysis (PCA)
- Linear discriminant analysis (LDA)
- Multidimensional scaling (MDS)
- t-distributed stochastic neighborhood embedding (t-SNE)
- Uniform manifold approximation and projection (UMAP)
- ...

Will give abstract features which may be non-linear combination of the original features
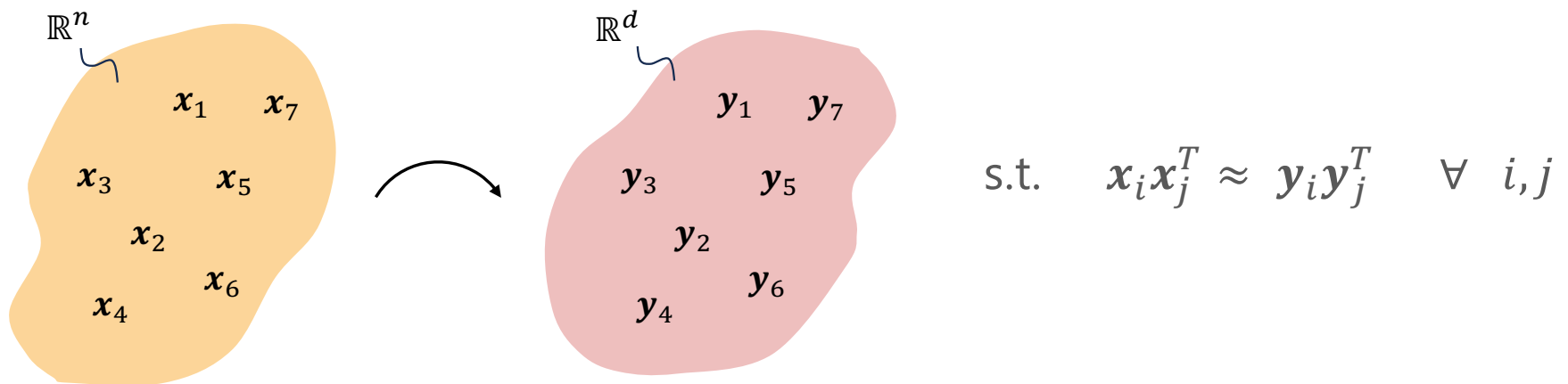
## Feature selection

- Forward selection
- Backward elimination
- Variance threshold
- L1 regularization
- ...

Selects a suitable subset of existing features

# MDS: Multi-Dimensional Scaling

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \longrightarrow \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_d \end{bmatrix} = f\left( \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \right)$$
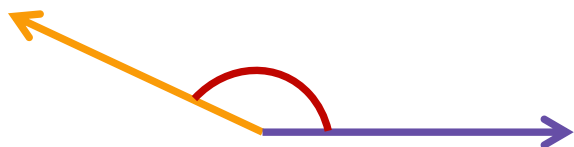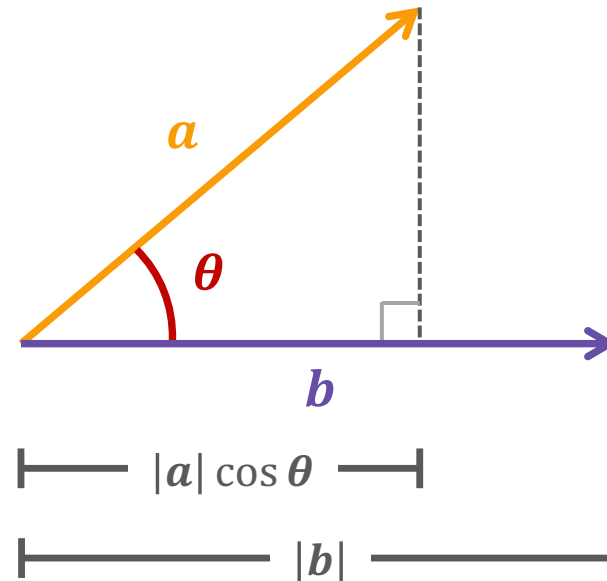
is an unsupervised feature extraction method that aims to **preserve inner products**.

$\mathbb{R}^n$

$x_1 \quad x_7$

$x_3 \qquad x_5$

$x_2$

$x_6$

$x_4$

$\mathbb{R}^d$

$y_1 \quad y_7$

$y_3 \qquad y_5$

$y_2$

$y_6$

$y_4$

s.t. $\quad x_i x_j^T \approx y_i y_j^T \quad \forall \; i, j$

# Inner product: geometric interpretation

measure of the *similarity* between vectors

$$ab^T = \sum_{i=1}^{d} a_i b_i$$

$$= |a||b| \cos(\theta)$$



$|a| \cos \theta$

$|b|$
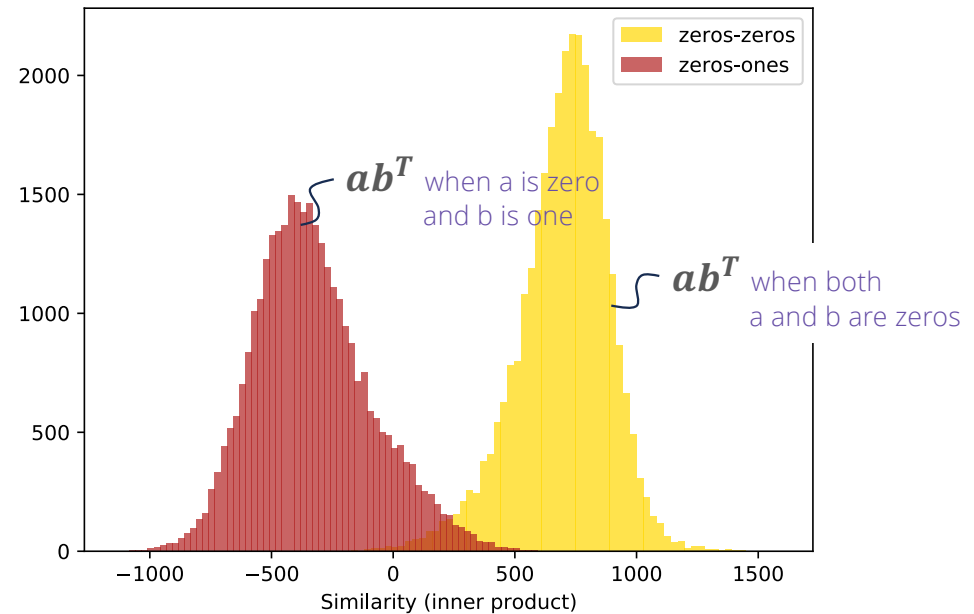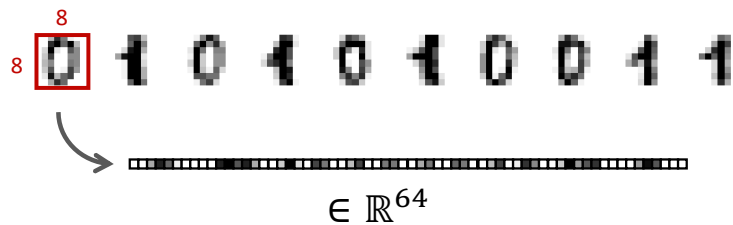
$$ab^T = 3 \cdot 3 \cos \frac{5\pi}{6} \approx -7.8$$

$$ab^T = 3 \cdot 3 \cos \frac{\pi}{12} \approx 8.7$$

$$ab^T = 3 \cdot 4 \cos \frac{\pi}{12} \approx 11.6$$

# example: inner product as similarity measure

vectorize images (size 8x8) of zeros and ones from the digits dataset[1] and compute inner products
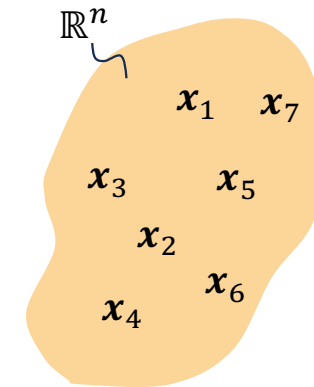


$$\in \mathbb{R}^{64}$$

$ab^T$ when a is zero and b is one

$ab^T$ when both a and b are zeros

zeros-zeros
zeros-ones

Similarity (inner product)

[1] https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load_digits.html

# The Gram matrix

Let

$$X = \begin{bmatrix} x^1 & \rightarrow \\ x^2 & \rightarrow \\ & \vdots \\ x^N & \rightarrow \end{bmatrix} \in \mathbb{R}^{N \times n}$$

<span style="color:red">the rows in **X** are the samples (feature vectors of length n)</span>

$\mathbb{R}^n$

$x_1 \quad x_7$

$x_3 \qquad x_5$

$x_2$

$x_6$

$x_4$

Then

$$B = XX^T = \begin{bmatrix} x^1(x^1)^T & \cdots & x^1(x^N)^T \\ \vdots & \ddots & \vdots \\ x^N(x^1)^T & \cdots & x^N(x^N)^T \end{bmatrix} \in \mathbb{R}^{N \times N}$$

<span style="color:red">element ij in **B** is the inner product between $x_i$ and $x_j$</span>
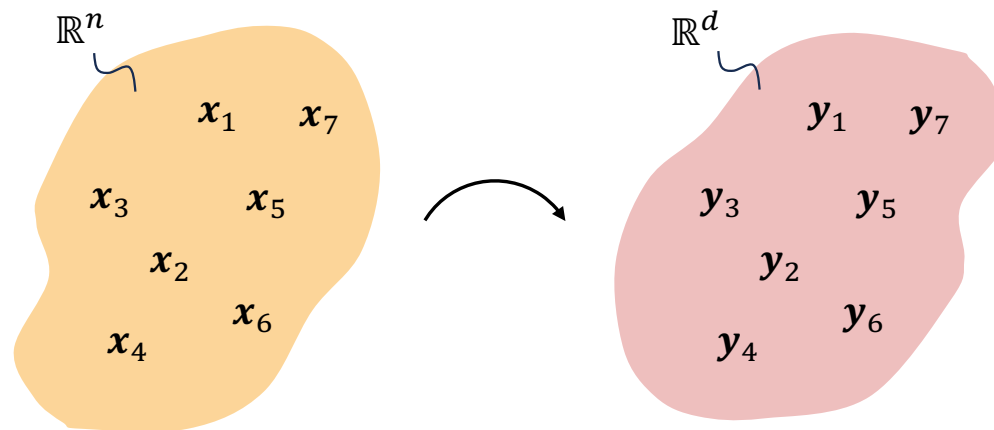
is called the Gram matrix/inner product matrix and contains all possible inner products between the samples.

# MDS: preserving inner products (similarities)

we want to find $\boldsymbol{y}^1, \boldsymbol{y}^2, \ldots, \boldsymbol{y}^N \in \mathbb{R}^d$ such that $\boldsymbol{YY^T} \approx \boldsymbol{XX^T}$

Gram matrix of
Low dim. samples

Gram matrix of
high dim. samples

$\mathbb{R}^n$

$\boldsymbol{x}_1$    $\boldsymbol{x}_7$

$\boldsymbol{x}_3$    $\boldsymbol{x}_5$

$\boldsymbol{x}_2$

$\boldsymbol{x}_6$

$\boldsymbol{x}_4$

$\mathbb{R}^d$

$\boldsymbol{y}_1$    $\boldsymbol{y}_7$

$\boldsymbol{y}_3$    $\boldsymbol{y}_5$

$\boldsymbol{y}_2$

$\boldsymbol{y}_6$

$\boldsymbol{y}_4$

# Properties of the Gram matrix

The Gram matrix $\mathbf{B} = \boldsymbol{X}\boldsymbol{X}^T$ is:

1. Real    all elements are real numbers
2. Symmetric: $\boldsymbol{B}^T = \boldsymbol{B}$    remember: $(A \cdot B)^T = B^T A^T$

Which means that:

1. The eigenvalues of $\mathbf{B}$ are real
2. The eigenvectors of $\mathbf{B}$ are orthonormal

eigenvector of $\mathbf{B}$

$$\mathbf{B}\boldsymbol{e} = \lambda\boldsymbol{e}$$

eigenvalue of $\mathbf{B}$

# Eigen-decomposition of the Gram matrix

Can rewrite the Gram matrix as:

$$\underset{N\times N}{B} = \underset{N\times N}{E^T}\ \underset{N\times N}{\Lambda}\ \underset{N\times N}{E} = \underset{N\times N}{E^T}\ \underset{N\times N}{\Lambda^{1/2}}\ \underset{N\times N}{\Lambda^{1/2}}\ \underset{N\times N}{E}$$

$$\begin{bmatrix} e^1 & e^2 & & e^N \\ \downarrow & \downarrow & \cdots & \downarrow \end{bmatrix} \begin{bmatrix} \lambda^1 & & & \\ & \lambda^2 & & \\ & & \ddots & \\ & & & \lambda^N \end{bmatrix} \begin{bmatrix} e^1 & \rightarrow \\ e^2 & \rightarrow \\ & \vdots \\ e^N & \rightarrow \end{bmatrix} \qquad \begin{bmatrix} \sqrt{\lambda^1} & & & \\ & \sqrt{\lambda^2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda^N} \end{bmatrix} \begin{bmatrix} \sqrt{\lambda^1} & & & \\ & \sqrt{\lambda^2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda^N} \end{bmatrix}$$

# MDS: Putting it together

Given our original samples $x^1, x^2, \ldots, x^N \in \mathbb{R}^n$ MDS aims to find a lower dimensional representation of our data $y^1, y^2, \ldots, y^N \in \mathbb{R}^d$ ($d < n$) while preserving the inner products (similarities) between samples, $YY^T \approx XX^T$.

Have that

$$(A \cdot B)^T = B^T A^T$$

$$XX^T = \mathbf{B} = \left(E^T \Lambda^{1/2}\right)\left(\Lambda^{1/2} E\right) = \left(E^T \Lambda^{1/2}\right)\left(E^T \Lambda^{1/2}\right)^T$$

If we now let $\mathbf{Y} = E^T \Lambda^{1/2}$, then $XX^T = YY^T$!    Y as expression of eigenvectors and eigenvalues of $XX^T$

$$\mathbf{Y} = \begin{bmatrix} e^1 & e^2 & & e^N \\ \downarrow & \downarrow & \cdots & \downarrow \end{bmatrix} \begin{bmatrix} \sqrt{\lambda^1} & & & \\ & \sqrt{\lambda^2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda^N} \end{bmatrix} = \begin{bmatrix} \sqrt{\lambda^1} e^1 & \sqrt{\lambda^2} e^2 & & \sqrt{\lambda^N} e^N \\ \downarrow & \downarrow & \cdots & \downarrow \end{bmatrix} \in \mathbb{R}^{N \times N}$$

$N \times N$       $N \times N$

remember: we wanted to reduce the dimensionality to d (d<n). How?

# MDS: Putting it together

To obtain $d$ dimensional representations for $\boldsymbol{y}^1, \boldsymbol{y}^2, \dots, \boldsymbol{y}^N$ such that $\boldsymbol{Y}\boldsymbol{Y}^T$ is *as close as possible* to $\boldsymbol{X}\boldsymbol{X}^T$, we can choose to set $\quad \mathbf{Y} = \boldsymbol{E}_d^T \boldsymbol{\Lambda}_d^{1/2}$, where $\boldsymbol{E}_d^T$ is the matrix of eigenvectors corresponding to the $d$ largest eigenvalues and $\boldsymbol{\Lambda}_d^{1/2}$ is the diagonal matrix with the $d$ largest eigenvalues.

$$\mathbf{Y} = \begin{bmatrix} \boldsymbol{e}^1 & \boldsymbol{e}^2 & & \boldsymbol{e}^d \\ \downarrow & \downarrow & \cdots & \downarrow \end{bmatrix} \begin{bmatrix} \sqrt{\lambda^1} & & & \\ & \sqrt{\lambda^2} & & \\ & & \ddots & \\ & & & \sqrt{\lambda^d} \end{bmatrix} = \begin{bmatrix} \sqrt{\lambda^1}\boldsymbol{e}^1 & \sqrt{\lambda^2}\boldsymbol{e}^2 & & \sqrt{\lambda^d}\boldsymbol{e}^d \\ \downarrow & \downarrow & \cdots & \downarrow \end{bmatrix} \in \mathbb{R}^{N \times d}$$

$N \times d \qquad\qquad d \times d$

# MDS: algorithm

Given our original samples $x^1, x^2, \dots, x^N \in \mathbb{R}^n$ MDS aims to find a lower dimensional representation of our data $y^1, y^2, \dots, y^N \in \mathbb{R}^d$ $(d < n)$ while preserving the inner products (similarities) between samples, $YY^T \approx XX^T$.





**Note**

$[\lambda^1, \lambda^2, \dots, \lambda^N] \qquad E^T$
eigenvalues, eigenvectors = np.linalg.eig(B)

column eigenvectors[:,i] is the eigenvector corresponding to the eigenvalue eigenvalues[i]
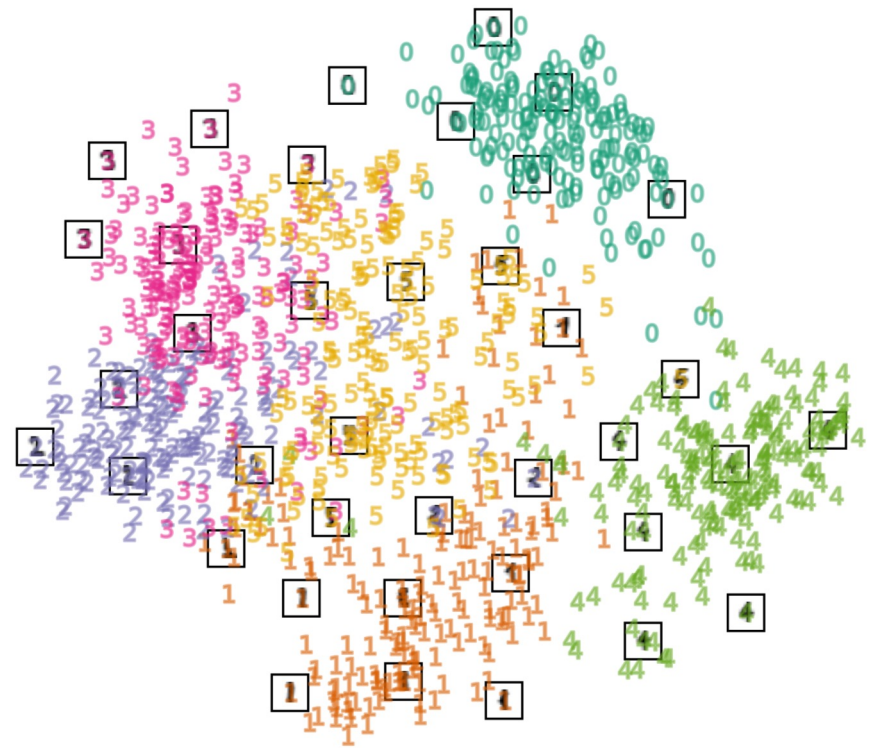
1. Compute the Gram matrix $\mathbf{B} = XX^T$

2. Compute eigenvectors and eigenvalues of $\mathbf{B}$

3. Sort eigenvalues and eigenvectors

4. Construct a matrix $\Lambda_d$ with the $d$ largest eigenvalues on the diagonal

5. Construct a matrix $E_d^T$ where the columns are the eigenvectors corresponding to the $d$ largest eigenvalues

6. Obtain $y^1, y^2, \dots, y^N$ as the rows of $\mathbf{Y} = E_d^T \Lambda_d^{1/2}$

# example: image dataset



$\in \mathbb{R}^{64}$

# Roadmap

1. Dimensionality reduction: what and why
2. Inner product as similarity measure
3. Multidimensional scaling (from samples)
4. Practical example
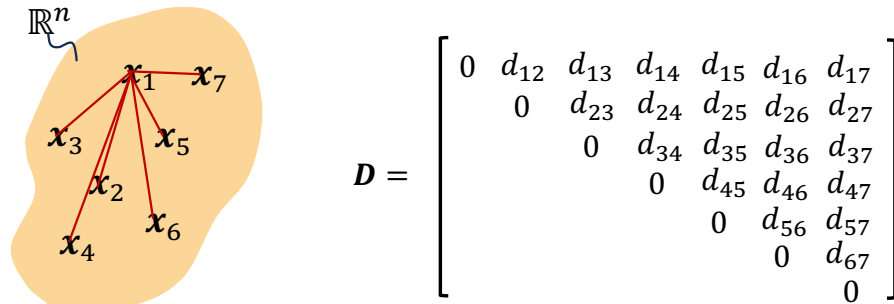
   <span style="color:red">Break</span>

5. Multidimensional scaling (from distances)
6. Practical examples
7. Forward selection
8. Practical example

# Gram matrix from pairwise distances

can **not** compute inner products directly

What if we don't have access to the original samples $x^1, x^2, \ldots, x^N \in \mathbb{R}^n$, but only to the matrix $D \in \mathbb{R}^{N \times N}$ containing all squared pairwise distances $d_{ij} = \left\| x^i - x^j \right\|^2$?

$\mathbb{R}^n$

$x_1$ $x_7$
$x_3$ $x_5$
$x_2$
$x_4$ $x_6$

$$D = \begin{bmatrix} 0 & d_{12} & d_{13} & d_{14} & d_{15} & d_{16} & d_{17} \\ & 0 & d_{23} & d_{24} & d_{25} & d_{26} & d_{27} \\ & & 0 & d_{34} & d_{35} & d_{36} & d_{37} \\ & & & 0 & d_{45} & d_{46} & d_{47} \\ & & & & 0 & d_{56} & d_{57} \\ & & & & & 0 & d_{67} \\ & & & & & & 0 \end{bmatrix}$$

We can compute the Gram matrix $B \in \mathbb{R}^{N \times N}$ consistent with these distances as a function of elements in $D$!

$$B = XX^T = -\frac{1}{2}CDC^T$$

assumption: the data is centered around the origin ("double-centering")

where $C = I - \frac{1}{N}J_N$ , and $I$ is the identity matrix and $J_N$ is an $N x N$ matrix of all ones. Thus use MDS to obtain a set of samples $y^1, y^2, \ldots, y^N \in \mathbb{R}^d$ such that $YY^T \approx XX^T$.

# MDS: algorithm

1. **If** *samples $\boldsymbol{X}$ are available:*
   Compute the Gram matrix $\mathbf{B} = \boldsymbol{X}\boldsymbol{X}^T$

   **else if** *only pairwise distances $\boldsymbol{D}$ are available:*
   Compute the Gram matrix $\mathbf{B} = -\frac{1}{2}\boldsymbol{C}\boldsymbol{D}\boldsymbol{C}^T$

2. Compute eigenvectors and eigenvalues of $\mathbf{B}$

3. Sort eigenvalues and eigenvectors

4. Construct a matrix $\boldsymbol{\Lambda}_d$ with the $d$ largest eigenvalues on the diagonal

5. Construct a matrix $\boldsymbol{E}_d^T$ where the columns are the eigenvectors corresponding to the $d$ largest eigenvalues

6. Obtain $\boldsymbol{y}^1, \boldsymbol{y}^2, \dots, \boldsymbol{y}^N$ as the rows of $\mathbf{Y} = \boldsymbol{E}_d^T \boldsymbol{\Lambda}_d^{1/2}$

$\mathbb{R}^n$

$\boldsymbol{x}_1 \quad \boldsymbol{x}_7$

$\boldsymbol{x}_3 \quad \boldsymbol{x}_5$

$\boldsymbol{x}_2$

$\boldsymbol{x}_6$

$\boldsymbol{x}_4$

$\mathbb{R}^d$

$\boldsymbol{y}_1 \quad \boldsymbol{y}_7$

$\boldsymbol{y}_3 \quad \boldsymbol{y}_5$

$\boldsymbol{y}_2$

$\boldsymbol{y}_6$

$\boldsymbol{y}_4$

# example: city distances

# example: letter recognition

| Letter | C | D | G | H | M | N | Q | W |
|--------|---|---|---|---|---|---|---|---|
| C | — | | | | | | | |
| D | 5 | — | | | | | | |
| G | 12 | 2 | — | | | | | |
| H | 2 | 4 | 3 | — | | | | |
| M | 2 | 3 | 2 | 19 | — | | | |
| N | 2 | 4 | 1 | 18 | 16 | — | | |
| Q | 9 | 20 | 9 | 1 | 2 | 8 | — | |
| W | 1 | 5 | 2 | 5 | 18 | 13 | 4 | — |

Number of times "G" and "C" were confused

...these numbers represent similarities

$$d_{ij} = \begin{cases} c - sim_{ij}, & if\ i \neq j \\ 0, & if\ i = j \end{cases}$$



MDS can work with pairwise distances – they do **not** have to be metrics! Can be (dis-)similiarities!

# Drawback of feature extraction

The new features are some combination of the original features, but what do they represent?

Depending on application, interpretability can be important!

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \xrightarrow{\text{feature extraction}} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_d \end{bmatrix} = f\left(\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}\right)$$

**Feature selection**: Selects a subset of the original features.

the features still have an interpretation!

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \xrightarrow{\text{feature selection}} \begin{bmatrix} X_q \\ X_r \\ \vdots \\ X_d \end{bmatrix}$$

27

# Feature types

1. **Relevant features**: important for the predictive task     weight [kg], blood pressure, age, ...

2. **Irrelevant features**: unimportant for the predictive task     favorite movie

3. **Redundant features**: become irrelevant in the presence of other features

             weight [g]

**Feature selection** algorithms aim to find the **relevant features**!

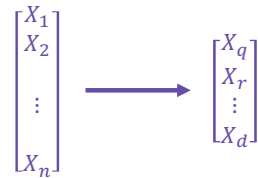# Dimensionality reduction methods

## Feature extraction

- Principal component analysis (PCA)
- Linear discriminant analysis (LDA)
- Multidimensional scaling (MDS)
- t-distributed stochastic neighborhood embedding (t-SNE)
- Uniform manifold approximation and projection (UMAP)
- ...

## Feature selection

- Forward selection
- Backward elimination
- Variance threshold
- L1 regularization
- ...

# Forward selection

$$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \longrightarrow \begin{bmatrix} X_q \\ X_r \\ \vdots \\ X_d \end{bmatrix}$$

is a supervised feature selection method that aims to find the **best subset of features**.

Best subset: the least number of features that most contribute to accuracy

forward selection is done in the context of a model and a prediction task!

# Forward selection: algorithm

Aim: Given a model $f$ and a labeled dataset $\{x^i, y^i\}_{i=1}^N$, where $x^i \in \mathbb{R}^n$, select the set of features $\mathcal{F}$ that "work best" together.

0.  Shuffle the data and split into training, validation, and test

1.  Init: $\mathcal{F} = \emptyset$   starting with an empty set

2.  Repeat:

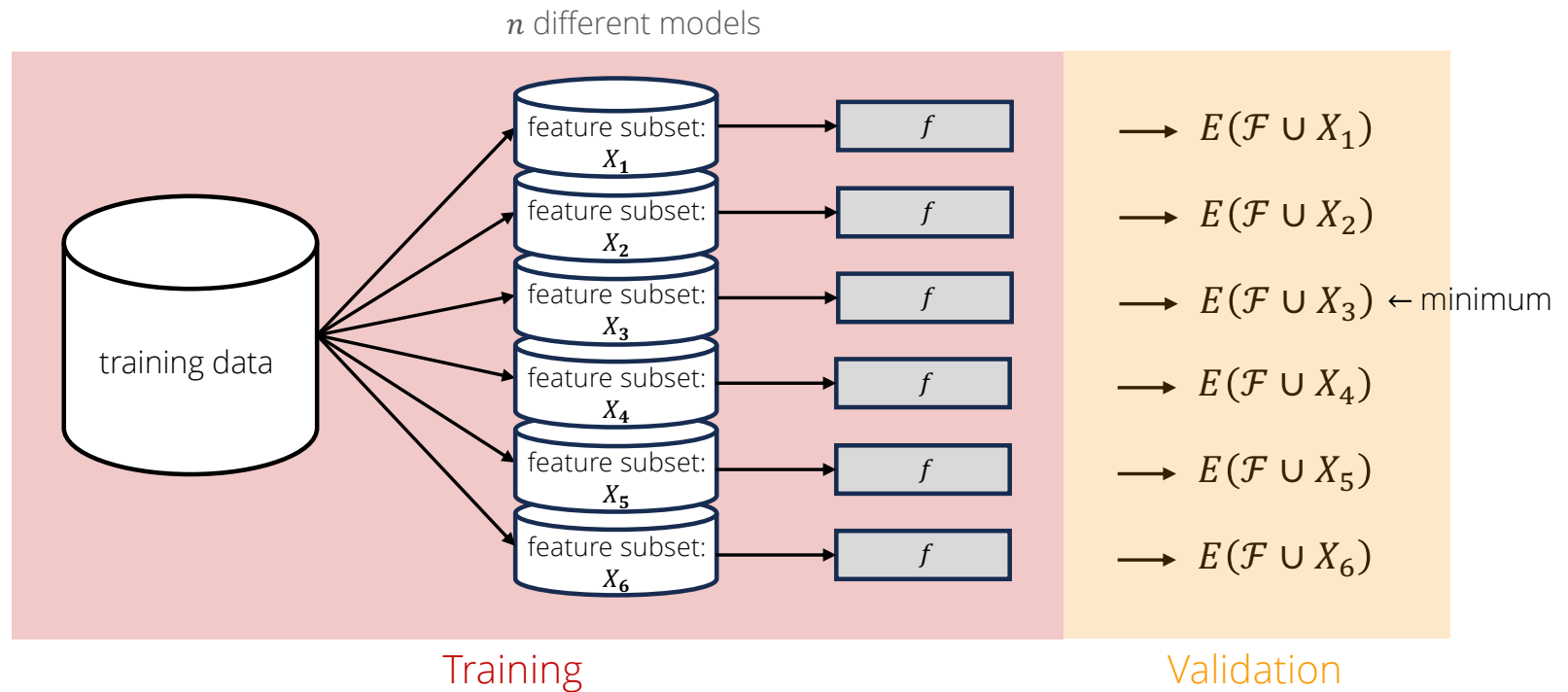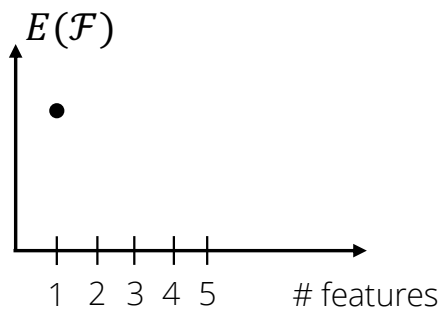    - for each feature $x_i \notin \mathcal{F}$, train $f$ on the set of features $\mathcal{F} \cup x_i$  (training data)

    - find the feature yielding the lowest validation error: $j = \arg\min_i E(\mathcal{F} \cup x_i)$

    - add $x_j$ to $\mathcal{F}$ if $E(\mathcal{F} \cup x_j) < E(\mathcal{F})$

    Until: $E(\mathcal{F} \cup x_j) \geq E(\mathcal{F})$   stop when adding a feature does not decrease the validation error

# Forward selection: example

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix}$$ 
age
hair color
glucose level
BMI
height
weight

$\mathcal{F} = \{X_3\}$

$E(\mathcal{F})$

$\bullet$

1 2 3 4 5    # features

$n$ different models

training data

feature subset: $X_1$ → $f$ → $E(\mathcal{F} \cup X_1)$

feature subset: $X_2$ → $f$ → $E(\mathcal{F} \cup X_2)$

feature subset: $X_3$ → $f$ → $E(\mathcal{F} \cup X_3)$ ← minimum

feature subset: $X_4$ → $f$ → $E(\mathcal{F} \cup X_4)$

feature subset: $X_5$ → $f$ → $E(\mathcal{F} \cup X_5)$

feature subset: $X_6$ → $f$ → $E(\mathcal{F} \cup X_6)$

Training

Validation

32

# Forward selection: example

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix} \begin{matrix} \text{age} \\ \text{hair color} \\ \text{glucose level} \\ \text{BMI} \\ \text{height} \\ \text{weight} \end{matrix}$$
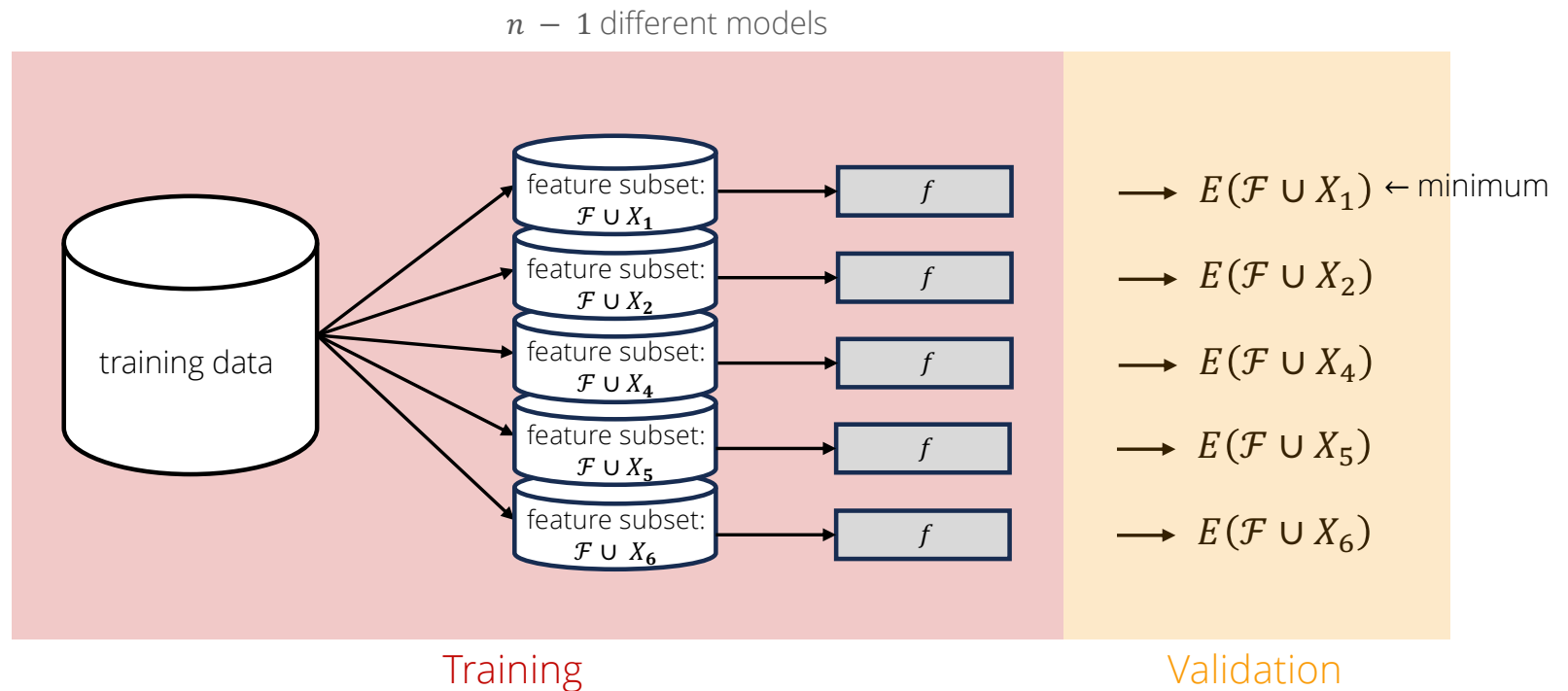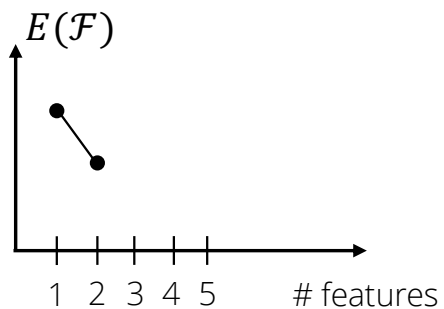
$$\mathcal{F} = \{X_3, X_1\}$$

$E(\mathcal{F})$



# features
1  2  3  4  5

$n - 1$ different models

training data

| feature subset: $\mathcal{F} \cup X_1$ | $\rightarrow$ | $f$ | $\longrightarrow$ | $E(\mathcal{F} \cup X_1)$ ← minimum |
| feature subset: $\mathcal{F} \cup X_2$ | $\rightarrow$ | $f$ | $\longrightarrow$ | $E(\mathcal{F} \cup X_2)$ |
| feature subset: $\mathcal{F} \cup X_4$ | $\rightarrow$ | $f$ | $\longrightarrow$ | $E(\mathcal{F} \cup X_4)$ |
| feature subset: $\mathcal{F} \cup X_5$ | $\rightarrow$ | $f$ | $\longrightarrow$ | $E(\mathcal{F} \cup X_5)$ |
| feature subset: $\mathcal{F} \cup X_6$ | $\rightarrow$ | $f$ | $\longrightarrow$ | $E(\mathcal{F} \cup X_6)$ |

Training

Validation

# Forward selection: example

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix}$$

$X_1$ — age
$X_2$ — hair color
$X_3$ — glucose level
$X_4$ — BMI
$X_5$ — height
$X_6$ — weight

$\mathcal{F} = \{X_3, X_1, X_4\}$

$E(\mathcal{F})$

$n - 2$ different models

training data

feature subset: $\mathcal{F} \cup X_2$ → $f$ → $E(\mathcal{F} \cup X_2)$

feature subset: $\mathcal{F} \cup X_4$ → $f$ → $E(\mathcal{F} \cup X_4)$ ← minimum

feature subset: $\mathcal{F} \cup X_5$ → $f$ → $E(\mathcal{F} \cup X_5)$

feature subset: $\mathcal{F} \cup X_6$ → $f$ → $E(\mathcal{F} \cup X_6)$

Training

Validation

1  2  3  4  5    # features

# Forward selection: example

$$X = \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \end{bmatrix}$$
age
hair color
glucose level
BMI
height
weight

$\mathcal{F} = \{X_3, X_1, X_4\}$

$E(\mathcal{F})$

1 2 3 4 5  # features

$n - 3$ different models

training data

feature subset:
$\mathcal{F} \cup X_2$ → $f$ → $E(\mathcal{F} \cup X_2)$

feature subset:
$\mathcal{F} \cup X_5$ → $f$ → $E(\mathcal{F} \cup X_5)$

feature subset:
$\mathcal{F} \cup X_6$ → $f$ → $E(\mathcal{F} \cup X_6)$ ← minimum

Training

Validation

# Forward selection: remarks

- To go from $n$ to $d$ features, we need to train $n + (n-1) + (n-2) + \ldots + (n-d)$ models, which might become costly

- The selected features depend heavily on the model and the prediction task

- The selected features might depend on the train/validation split

  improve robustness through cross validation!

- Finding the optimal subset is **not** guaranteed (context/dependencies)

# Conclusion

- Visited two methods for <span style="color:red">dimensionality reduction:</span>

  - Unsupervised feature **extraction** method: Multidimensional scaling
  - Supervised feature **selection** method: Forward selection

- MDS from high dimensional data (samples)

- MDS from pairwise distances

- Forward selection (in the context of a model)

- Choice of method depends on the data and the problem setting