Unsupervised learning 18.11.2024

EPFL

Outline

- Hour 1
 - Quiz discussion, review of last lecture
 - Introduction to unsupervised learning:
 - Principal component analysis (PCA)

- Hour 2:
 - PCA continued
 - K-means



Review of data statistics, ex: Quiz 1 grades

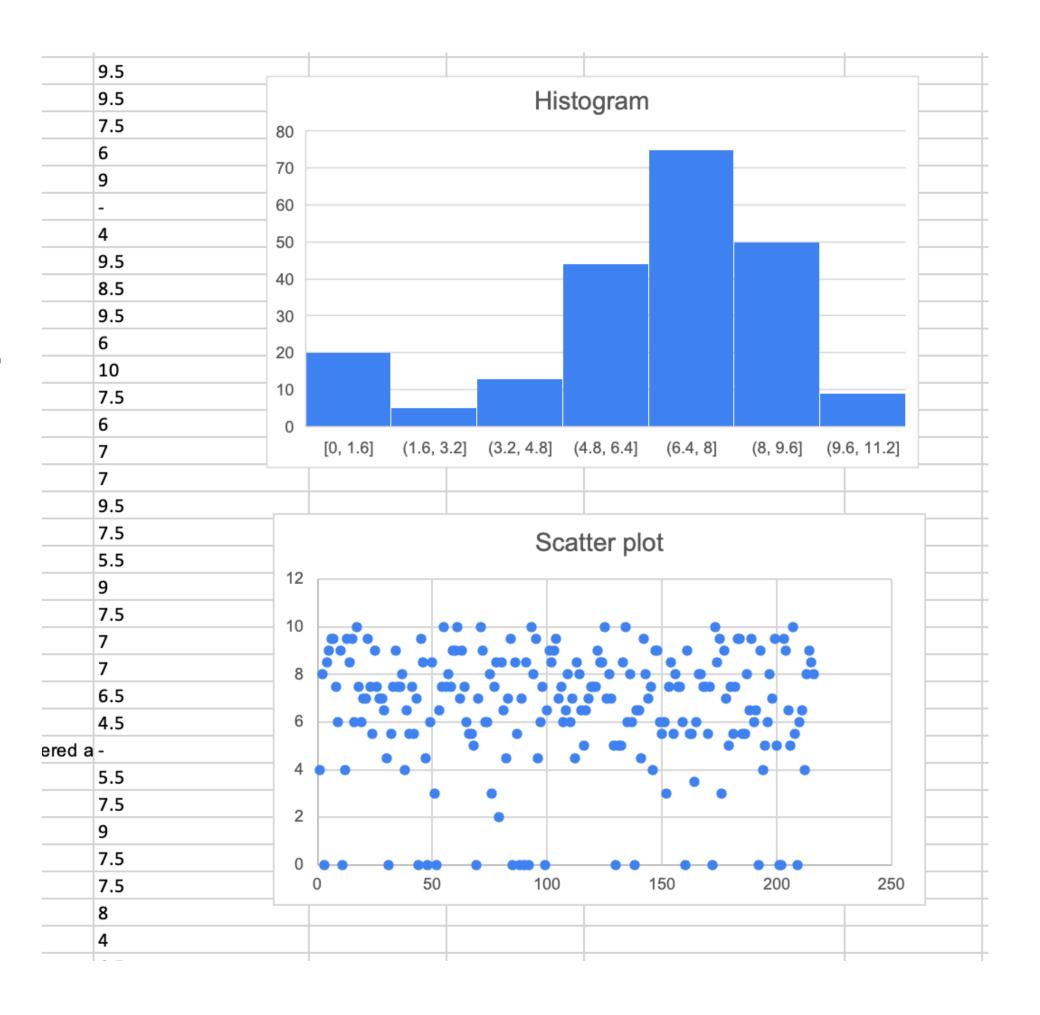
Mean: 7.16, Standard deviation: 1.75, Mode: 7.5

Quartiles:

1st: 6, 25% of grades is below this number

2nd (median): 7.5, 50% of grades is below this number

3rd: 8.5, 75% of grades is below this number



Review last time

Consider MNIST data

Neural networks (NN)

Convolutional neural networks (CNN)

Consider a single layer CNN with 5 filters, each of which is 3x3. How many parameters need to be determined to specify this CNN?

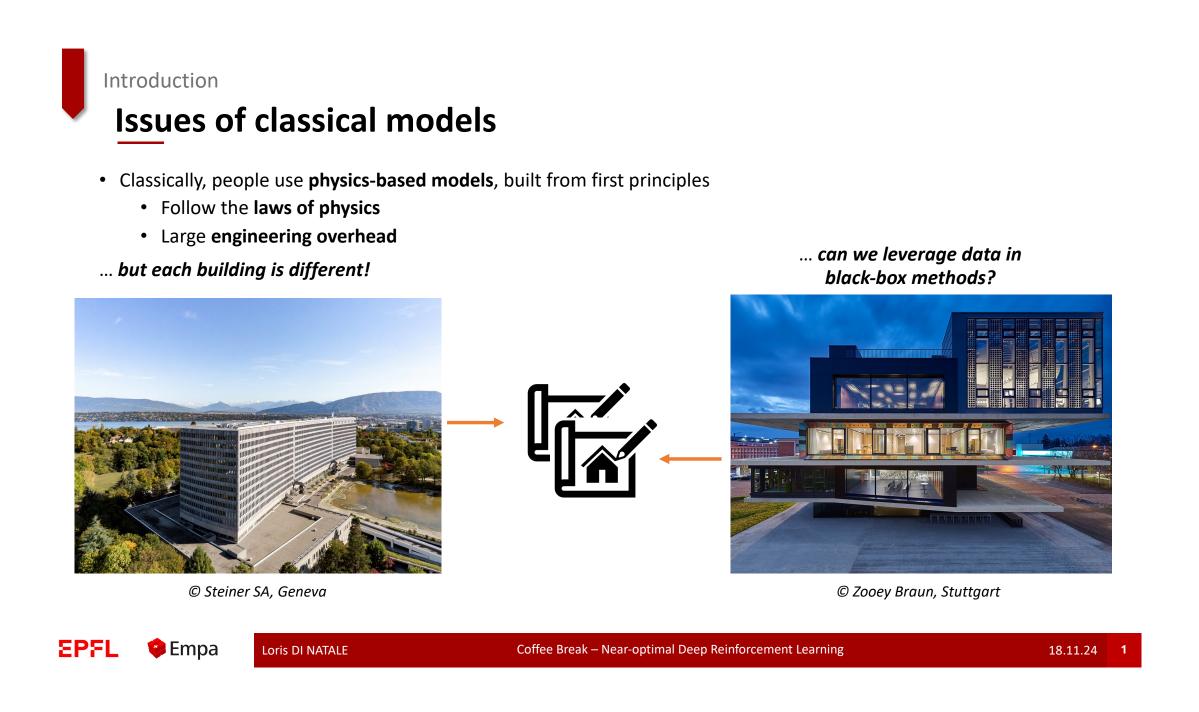
$$5 \times (3 \times 3 + 1) = 50$$
 parameters.

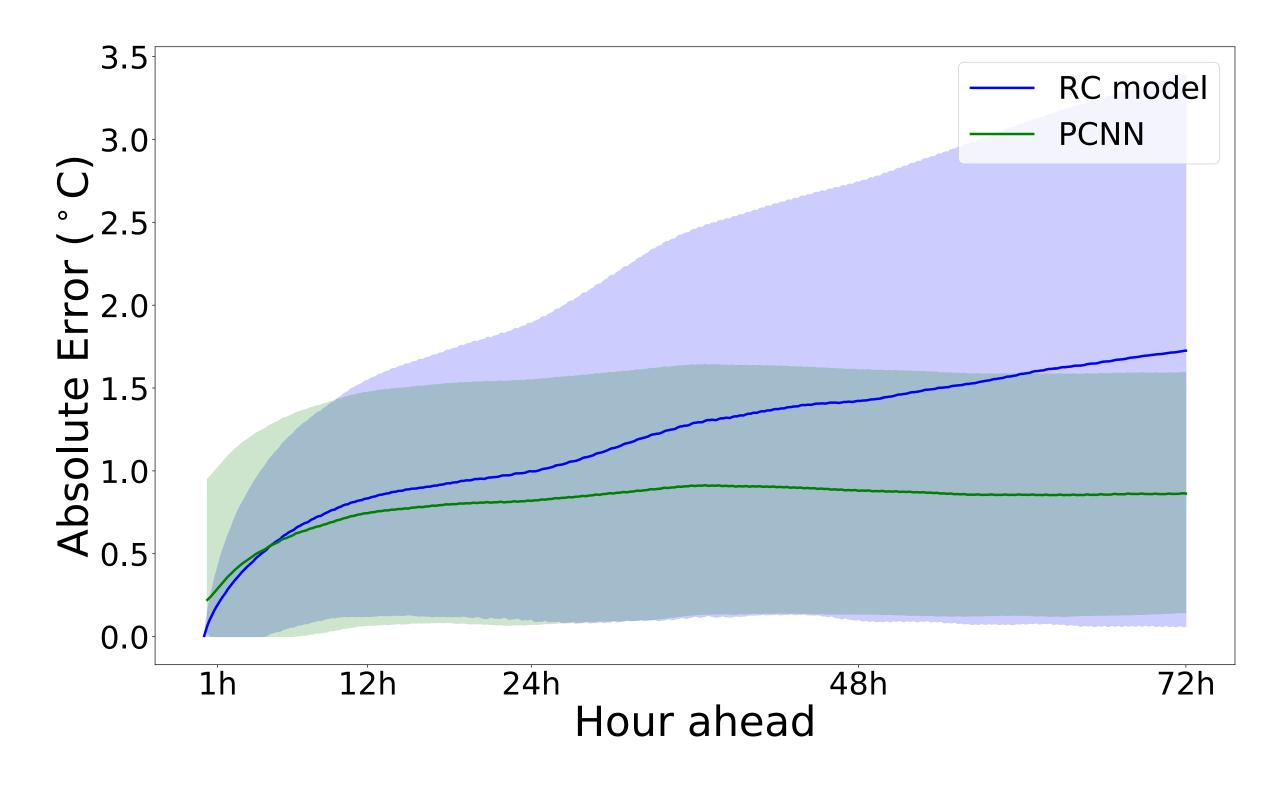


Example applications of NN in IGM

Prof. Colin Jones: Predictive Control Lab

Goal: develop a dynamical model of the temperature evolution in a building for control (Control objective: minimize energy consumption while ensuring comfort of occupants)







Introduction Unsupervised learning

Unsupervised learning is a type of machine learning that looks for previously undetected patterns in a **data set with no pre-existing labels** and with a minimum of human supervision.

in the next techniques, we don't use labels anymore Note: the objective is vague but we will consider 2 concrete instances



Dimensionality reduction

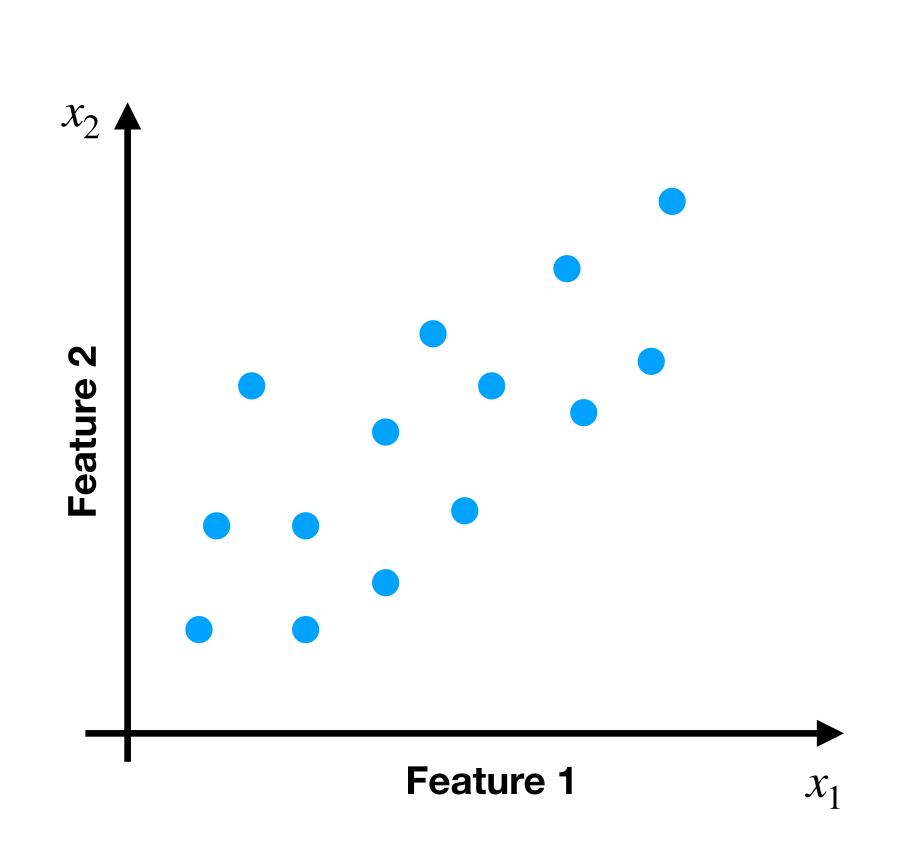
Through

principal component analysis



Motivation Intuition

We have samples described by a series of features



We want to find a smaller set of new features that explain our sample because:

Less features is easier to visualize

Some of the current feature can be redundant

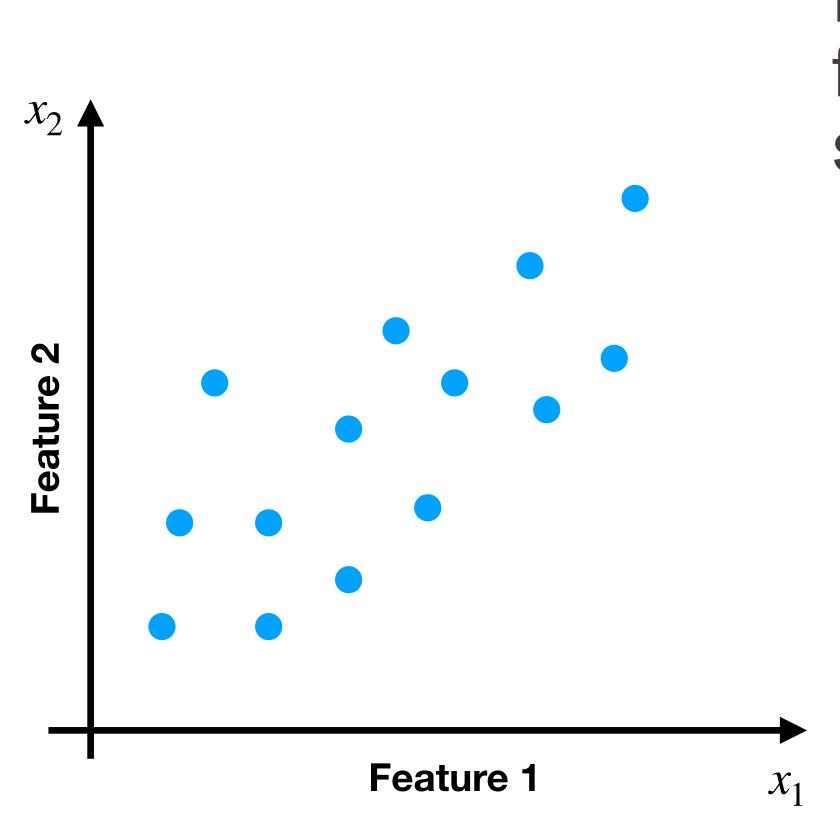
Some of the current features are not very useful to describe our samples



Principal component analysis (PCA)

Approach to dimensionality reduction

How to find this smaller set of new features?



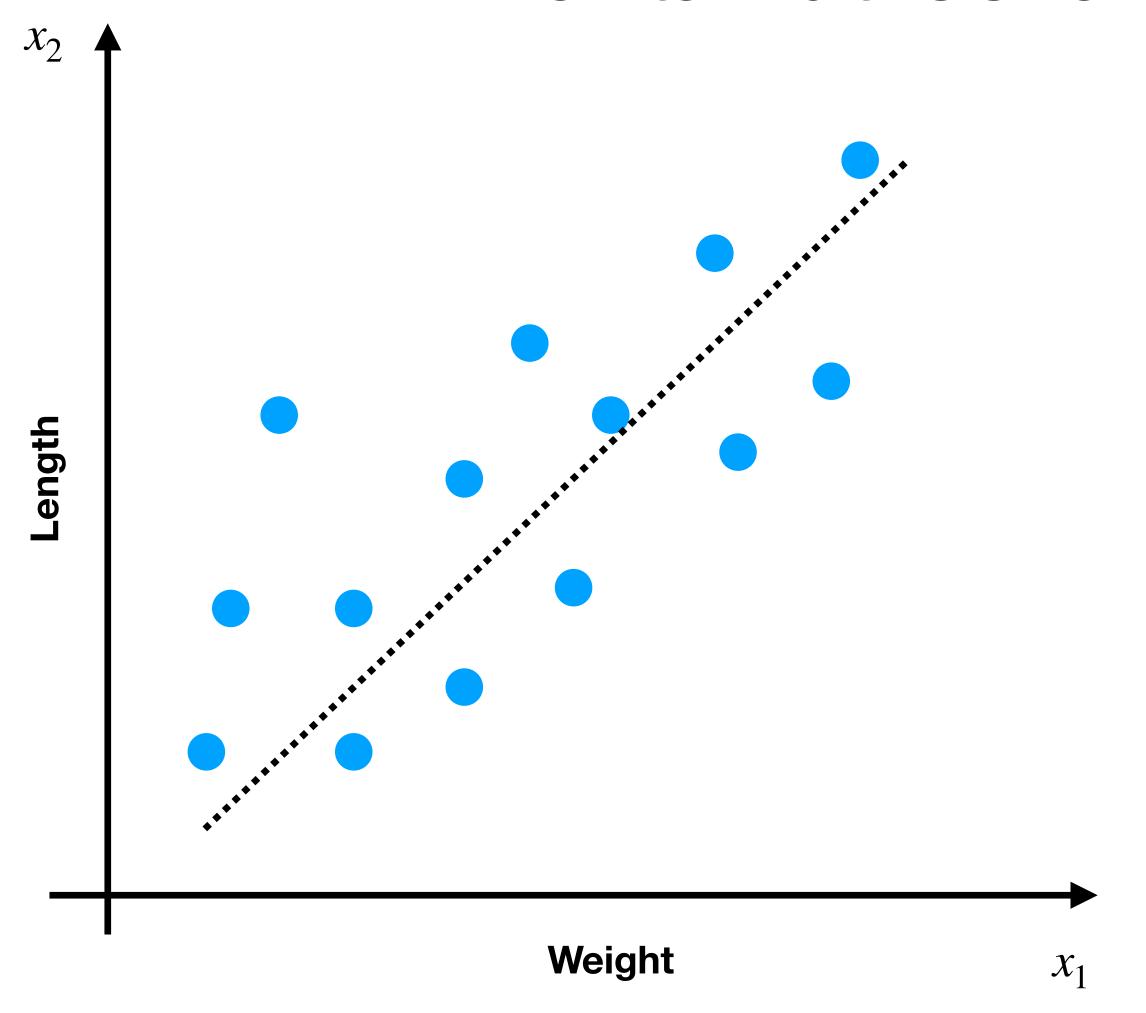
PCA: Find the best linear combination of features to create new features that explain our samples better



PCA

Projection of points onto a lower dimensional subspace

How to find this smaller set of new features?



$$w_1 x_1 + w_2 x_2$$

The new feature

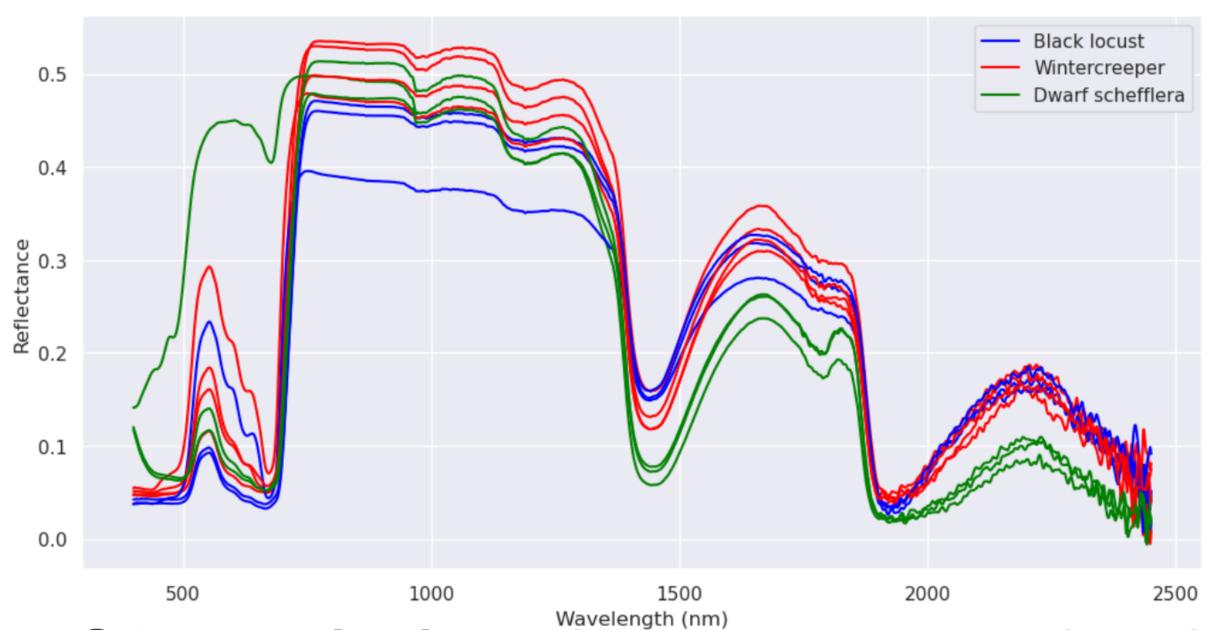
A mix of length and weight that describe our samples better

PCA - example application in python

- **Dataset:** Leaves' optical reflectance measured at each nm in the range of 450-2500 nm Wavelength, see more details about the experiment <u>here</u>
- Question: what is the size of feature vector for each leaf? $x \in \mathbb{R}$ (24 Sa -4 Sc = 2000)

950

Goal: can we reduce the dimensionality of feature vector and still capture the distinguishing features of each leaf



• Other applications: dynamical system model order reduction, audio compression for recommendation algorithms, text processing for news recommendation

EPFL

Finding distance of a point to a subspace

Subspace
$$\Theta_{i} \in \mathbb{R}^{d}$$
, $\alpha_{i} \in \mathbb{R}$, $i = 1, 2, ..., \Gamma$
 $\alpha_{i} = 1, 2, ..., \Gamma$

Distance of a point to a subspace $\alpha_{i} = 1, 2, ..., \Gamma$
 $\alpha_{i} = 1, 2, ..., \Gamma$

Distance of a point to a subspace $\alpha_{i} = 1, 2, ..., \Gamma$
 $\alpha_{i} = 1, 2, .$

dist
$$(x, S) = \min_{x \in S} \|x - S\|_2 = \min_{x \in S} \|x - Ga\|, \text{ where } \alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_r \end{bmatrix},$$

Let us consider min
$$||x-0\alpha||^2$$

$$|| \times - \Theta \alpha ||^{2} = || \times ||^{2} - 2 \langle x, \Theta \alpha \rangle + || \Theta \alpha ||^{2}$$

$$= || \times ||^{2} - 2 a^{T} \Theta^{T} x + a^{T} \Theta^{T} \Theta \alpha$$



Projection of a point onto a subspace

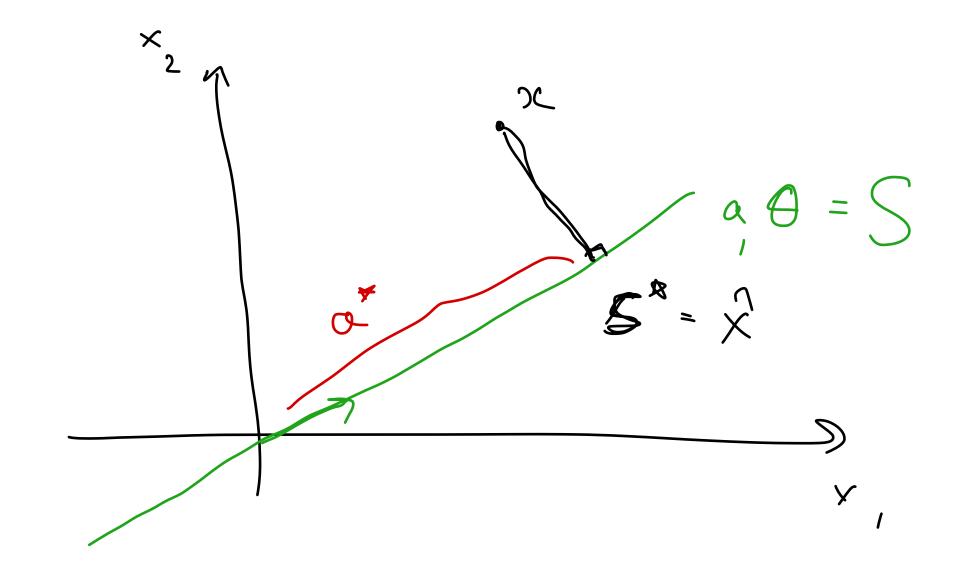
Distance of a point onto a subspace $\left(\text{clistance}\right)^2: \|x\|^2 - 2\alpha^\top \Theta^\top x + \alpha^\top \Theta^\top \Theta \alpha = : \overline{J}(\alpha)$

$$\frac{\partial}{\partial \alpha} J(\alpha) = 0 \implies \alpha^* = (\theta^{T}\theta)^{T}\theta^{T}x$$
Choosing orthonormal vectors $\{\theta_{i}\}_{i=1}^{T}$ are orthonormal $\Rightarrow \theta^{T}\theta = T_{i}x^{T}$

$$\alpha^* = \Theta^T x$$
, $S^* = closest point to x, $\Theta \alpha^* = \Theta \Theta^T x \in \mathbb{R}^{d \times 1}$$

Projection of a point onto a subspace

$$\stackrel{\wedge}{\times}$$
 = $S = \bigcirc \bigcirc \stackrel{\top}{\bigcirc} \stackrel{\top}{\sim}$



EPFL Find a subspace to minimize average distance of data to it

Sum of distances of all data points to a subspace $S = \{ \Theta a \mid a \in \mathbb{R}^r, \Theta \in \mathbb{R}^r \}$ objective function in PCA is to find S such that

$$J(S) = \frac{1}{N} \sum_{i=1}^{N} d_{i}st(x^{i}, S)$$
 is minimized

PCA chooses a subspace to minimize

$$J(S) = \frac{1}{N} \sum_{i=1}^{N} || x^{i} - \hat{x}^{i}|| = \frac{1}{N} \sum_{i=1}^{N} || x^{i} - oo^{7}x^{i}||$$

for any given
$$S$$
 spanned by $\{0\}_{i=1}^{r}$ $\Theta = [0,(0,2),(0,1)]$



Formulation of PCA objective using the data matrix

Standardize data: for each feature, subtract mean & divide by it stol.

Recall, we wanted to find $d\Theta: \xi$ to min $\frac{1}{N} \sum_{i=1}^{N} || x^i - \hat{x}^i||$

PCA objective:

$$\min_{\Theta} \frac{1}{N} \| \times - \times \Theta O T \|_{F}^{2}, \quad \Theta \in \mathbb{R}^{d \times r}$$

We are trying to find $\Theta = [0, |0, |0]$

EPFL

Eigenvalue decomposition of a symmetric matrix

Let
$$C = x^T x \in \mathbb{R}^{d \times d}$$
 $x \in \mathbb{R}^{d \times d}$

observe C is symmetric: $C = C^T$.

First: any symmetric matrix $C \in \mathbb{R}^{d \times d}$ has an eigenvalue elecamposition as follows:

 $C = V D V^T$. $D = \begin{cases} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \end{cases}$
 $V = [V_1 \mid V_2 \mid \dots \mid V_d] \in \mathbb{R}^{d \times d}$, where $\{\lambda_i\}_{i=1}^d \subset \mathbb{R}$ are the eigenvalues of C and $\{V_i\}_{i=1}^d \subset \mathbb{R}^d$ are the corresponding eigenvectors.

Solution to PCA using eigenvalue decomposition

EPFL

Principal components

Example: projection of data onto 2 principal components

Suppose
$$x_1, x_2$$
 of $C = X^TX$ are much larger than x_3, x_4, \dots, x_d . Then, we can let $x = 2$, and define the subspace S as span of v_1, v_2 , where v_1, v_2 are eigenvectors corresponds to x_1, x_2 , respectively. $S = \{v_1 \mid v_2\} \{a_1\} \{a_2\} \{a_1, a_2 \in \mathbb{R}\}$ compressed features: $A = XO \in \mathbb{R}^N$, $O = \{v_1 \mid v_2\} \{a_2\} \{a_3\} \{a_4\} \{a_4$

PCA

Pseudo-code

- 1. Standardize data (subhact mean, alivide by stal) X EIR
- 2. Eigenvalue decomposition of the data covariance matrix

- 3. To reduce the dimension to r << d , order eigenvalues in decreasing value, and choose the eigenvectors corresponds to r largest eigenvalues.
- 4. Compressed feature matrix

$$A = \times 0 \in \mathbb{Z}$$

 $C = X \times C | A \times c |$

PCA example - distinguishing texts

Defining features

Each data sample is a document with a large not be a large not be

There are d unique words in all the documents $x \in \mathbb{R}^{c}$

Feature j is positive for document i if word j is in document \mathfrak{g}

```
ocj >0 ED word j 15 in clocument i
```

TFIDF feature definition for documents

Herm

Prequing Inverse chocument Prequency

Term frequency of word in document 24

Document frequency of word ...j.

TFIDF for each word i.. in document of

$$x_{j} = f(0_{j}) \times \log \frac{1}{f_{doc}}, \quad i = 1, ..., N$$

$$j = 1, ..., d$$



PCA example

Based on "Principal Component Analysis" lecture of Stanford EE104: https://ee104.stanford.edu/lectures/pca.pdf

Distinguishing text: The Critique of Pure Reason by Immanuel Kant and The Problems of Philosophy by Bertrand Russell

d can be very high dimensional



Dimensionality reduction Other techniques

There are several other techniques for dimensionality reduction:

Linear discriminant analysis (LDA)

Generalized discriminant analysis (GDA)

T-distributed Stochastic Neighbor Embedding (t-SNE)

Autoencoders (neural networks)

Autoencoder Introduction

- An autoencoder is a type of neural network often used for dimensionality reduction.
- Autoencoders are trained in an unsupervised manner, by minimizing the

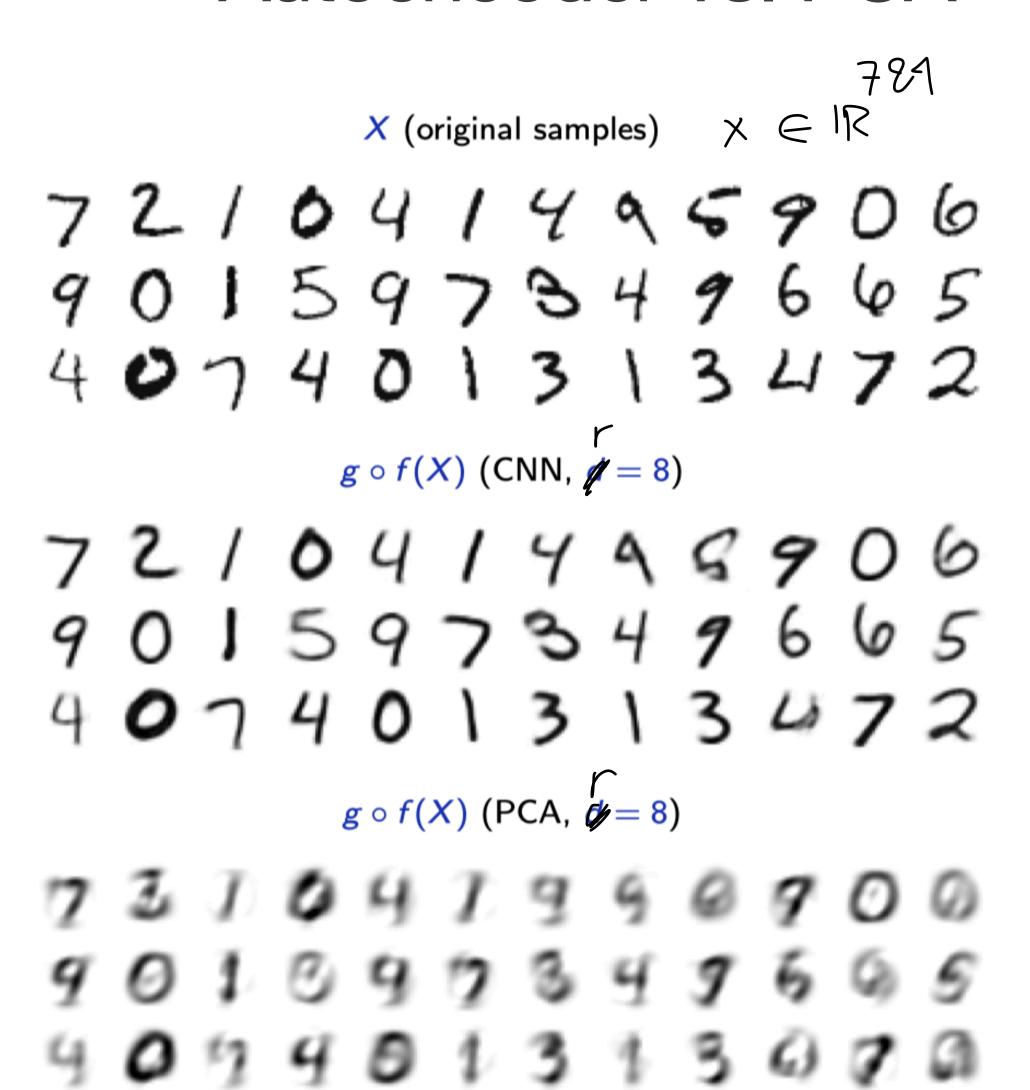
reconstruction error / loss
$$\sum_{i=1}^{N} L(x^{i}, \hat{x}^{i})$$
 $\approx = 9 (x^{i}) \implies$

Example: Squared error



Autoencoder

Autoencoder vs. PCA



Top: Some examples of the original MNIST test samples

Middle: Reconstructed output from an autoencoder with a latent space of 8 dimensions

This auto-encoder uses convolutional layers, and was trained on the MNIST training set

Bottom: Reconstructed output from PCA with 8 reduced dimensions

Image credit: F. Fleuret, Deep Learning (EPFL)



Summary - dimensionality reduction

Used for

- Exploratory data analysis
- Visualizing data
- Help reduce overfitting by reducing feature dimension

PCA: an approach to dimensionality reduction

- Projects data onto a linear subspace
- Useful in case there is approximately linear dependence between different features
- Easy to compute
- Connection to singular value decomposition (see Problem set 2)