

MATHEMATICS FOR MACHINE LEARNING

Nagoya University, Fall 2023

Lecture 12

Principal Component Analysis (PCA)

This week Tutorial: Friday 19th Dec. 6th period

<https://www.henrikbachmann.com/mml2023.html>

Semester project

Objective: Choose a project topic related to Nagoya or Japan more broadly that can be addressed using machine learning algorithms. Your task is to develop a machine learning model to solve a specific problem or provide insights into an aspect of life, business, environment, culture, etc., in Nagoya/Japan.

Group Size: 1-3 members

Code: Preferably a Google Colab notebook. Exceptions are possible; please provide full documentation for any different technology or package used. If you plan not to submit a Google Colab, please contact us in advance.

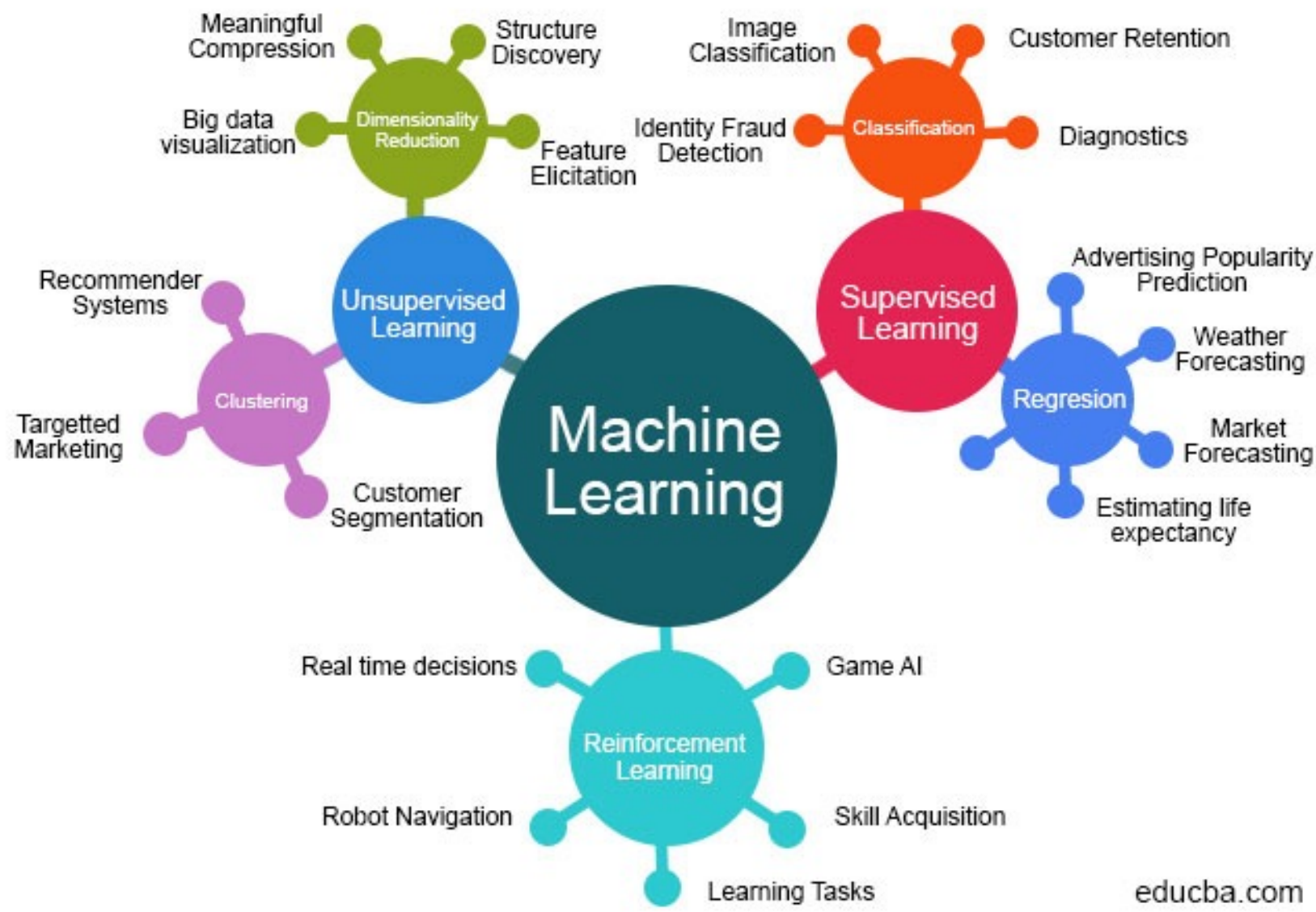
Documentation: **5-10 slides** as if you were going to present the project.

Your slides should cover (for example):

- Problem Statement
 - Data Collection
 - Data Exploration and Visualization
 - Model Building and Evaluation
 - Conclusion
- **The slides are the documentation and should explain what you did.**
 - **Do not just include examples and diagrams.**
 - **It should be possible to understand your project by just reading the slides.**

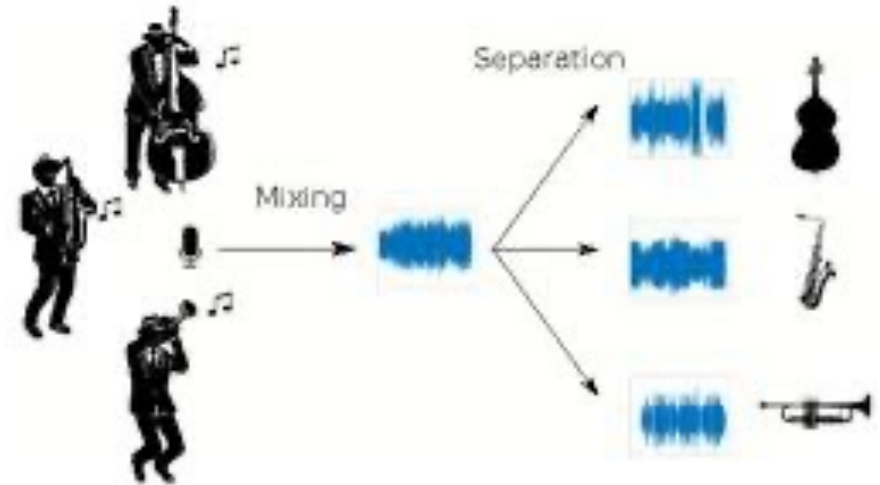
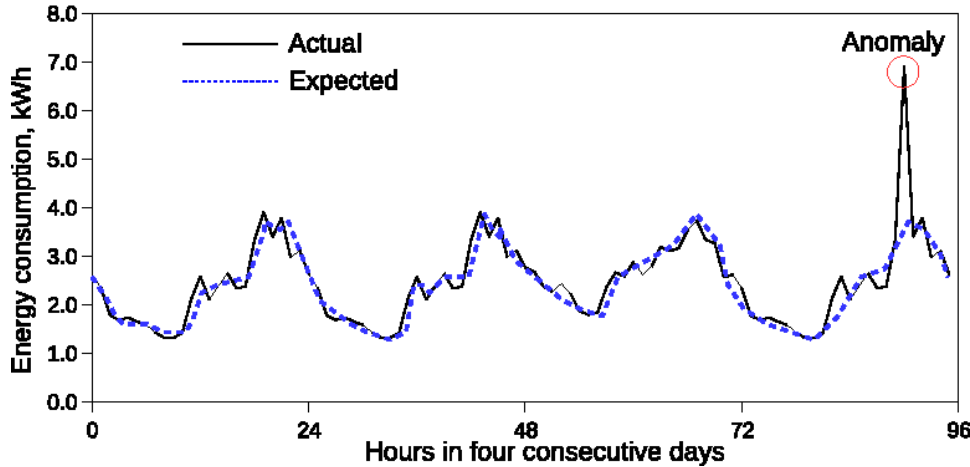
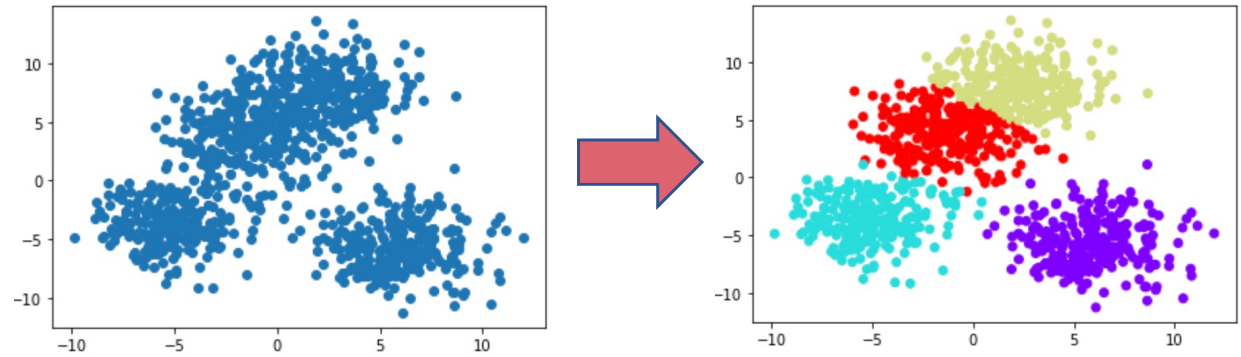
Overview

Machine Learning Algorithms



Unsupervised learning

- **Clustering (Lecture 11)**
- Anomaly detection
- Signal separation
- ...



https://en.wikipedia.org/wiki/Signal_separation

https://en.wikipedia.org/wiki/Cocktail_party_effect

Lecture 11: k-means clustering

k-means algorithm

1. Initialize the means $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ with some starting value.

(a) **Forgy method:** Choose randomly k different numbers $\{s_1, \dots, s_k\} \subset \{1, \dots, n\}$ and set

$$\mu_i = p_{s_i}$$

for $i = 1, \dots, k$

(b) **Random partition:** Choose $c : P \rightarrow \{1, \dots, k\}$ randomly and set for $i = 1, \dots, k$

$$\mu_i = \frac{1}{|C_i|} \sum_{p \in C_i} p.$$

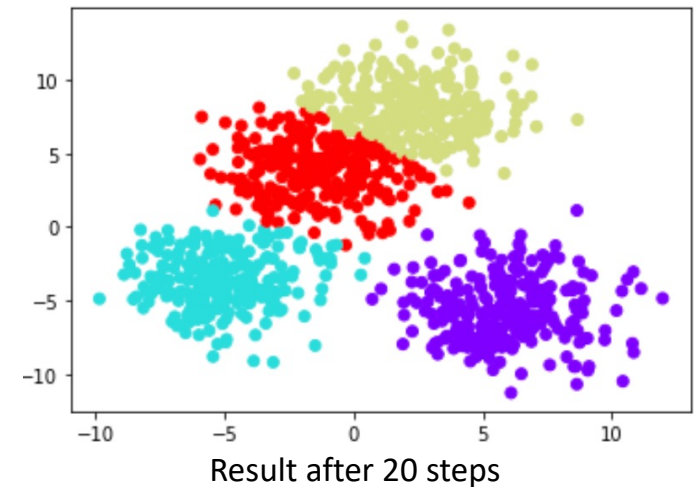
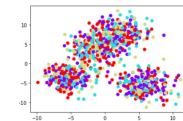
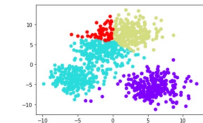
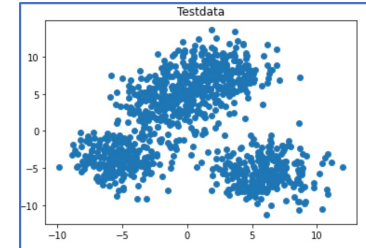
2. Define the clusters C_i for $i = 1, \dots, k$ by

$$C_i = \{p \in P \mid \|p - \mu_i\| \leq \|p - \mu_j\| \text{ for } j = 1, \dots, k\}$$

3. Recalculate the means μ_i for $i = 1, \dots, k$ by

$$\mu_i = \frac{1}{|C_i|} \sum_{p \in C_i} p.$$

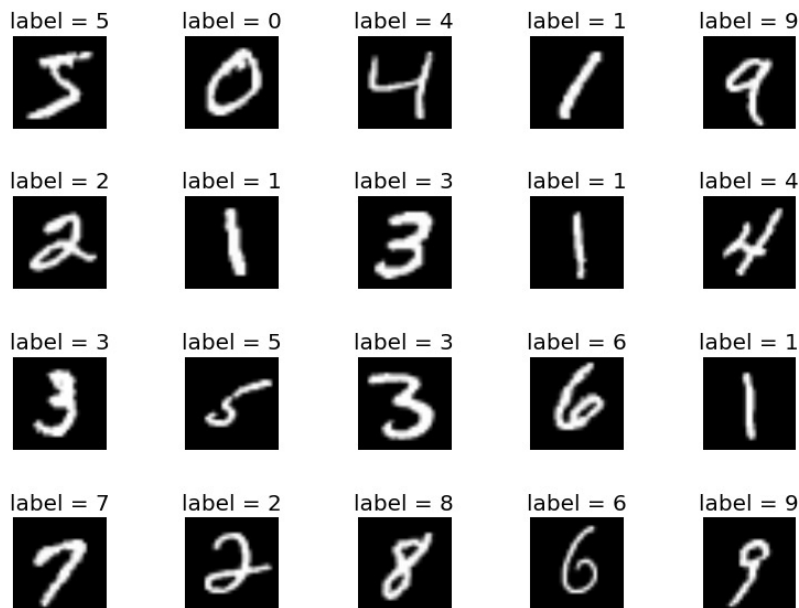
4. Repeat with step 2.



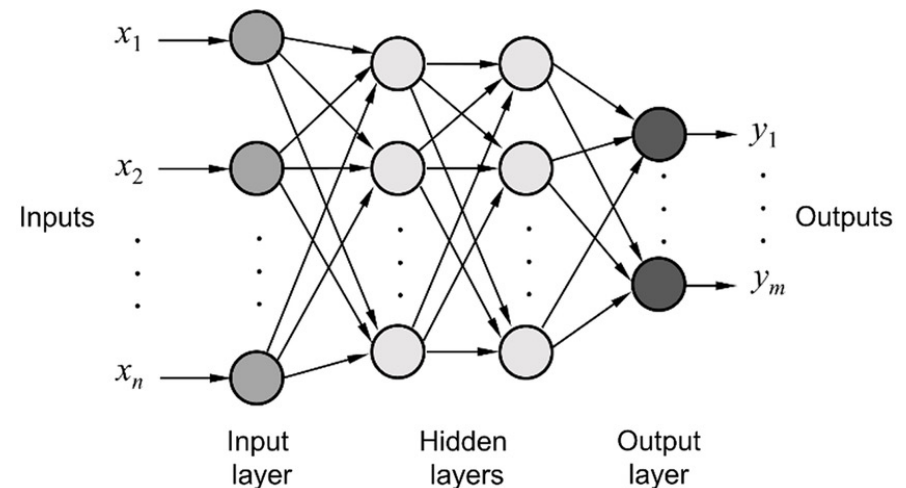
Unsupervised learning: Dimensionality reduction

Example: Digit recognition (MNIST Dataset)

Pictures of size $28 \times 28 = 784$ pixels



Neural network with input layer size 784

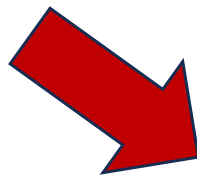
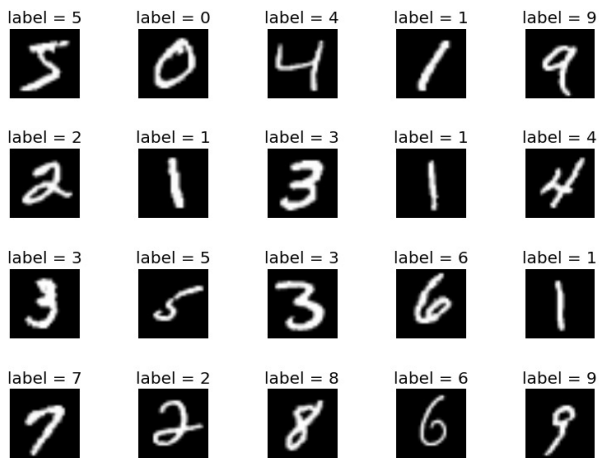


Question: How can we reduce the size of the input layer without losing a lot of information?

Unsupervised learning: Dimensionality reduction

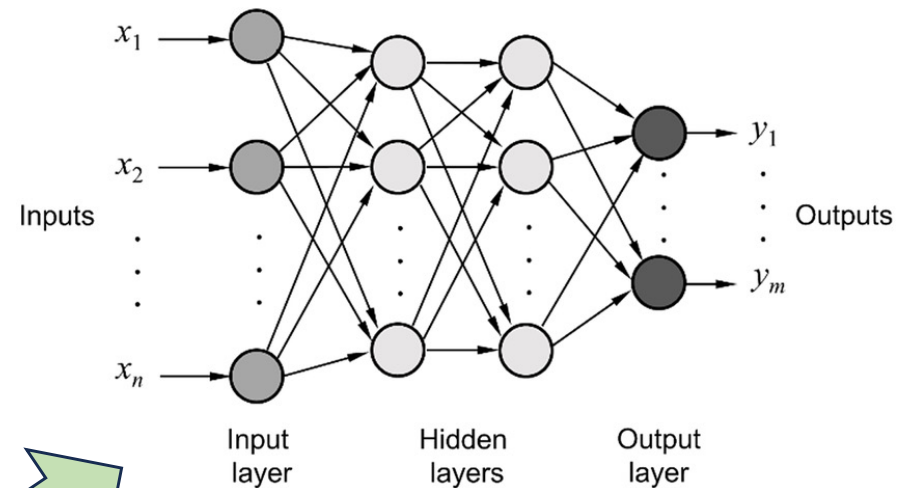
Example: Digit recognition (MNIST Dataset)

Pictures of size $28 \times 28 = 784$ pixels



Pictures of size $8 \times 8 = 64$ pixels

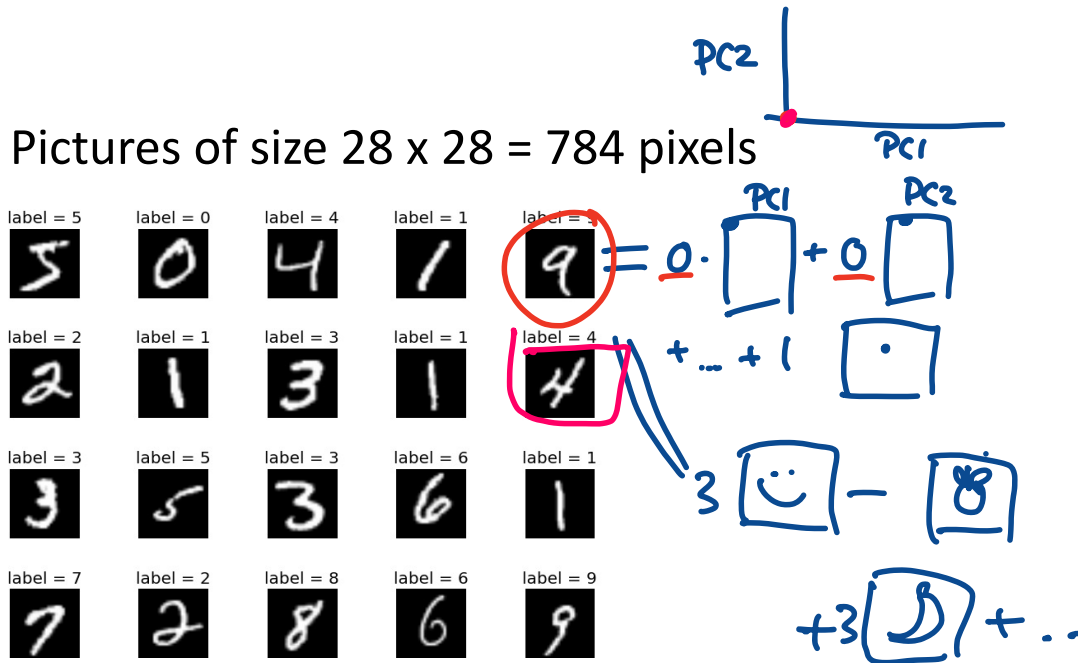
Neural network with input layer size 64



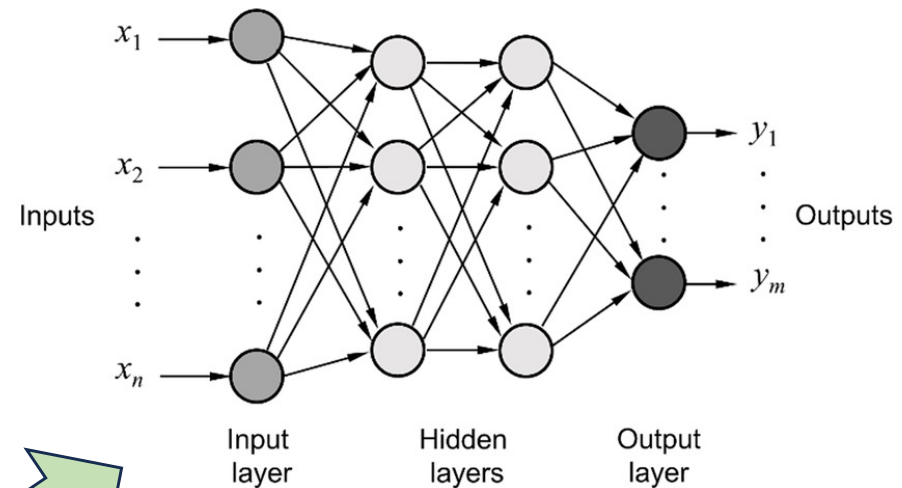
Naive approach: Scale the image down or ignore pixels

Unsupervised learning: Dimensionality reduction

Example: Digit recognition (MNIST Dataset)



Neural network with input layer size X

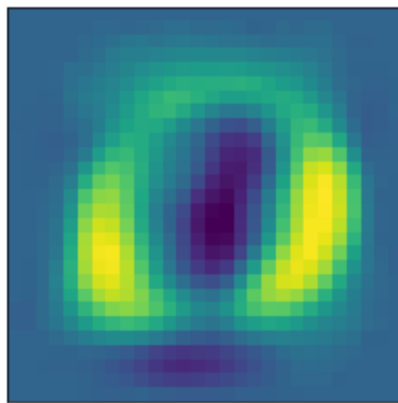


784 Datapoints ordered by “importance”.
Choose the first X of them.

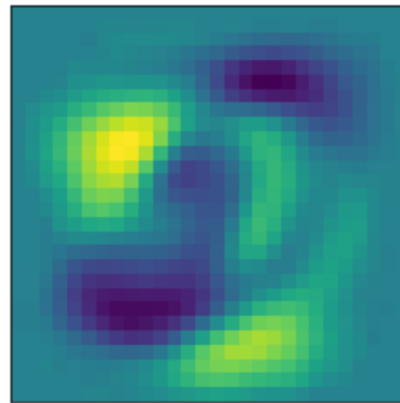
Better approach: Find a new representation of the picture into principal components.

Unsupervised learning: Dimensionality reduction

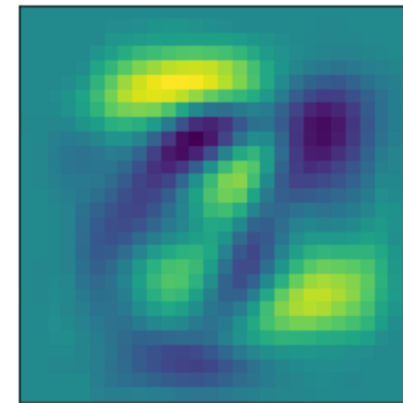
The first three principal components (PC) of the MNIST dataset



PC1

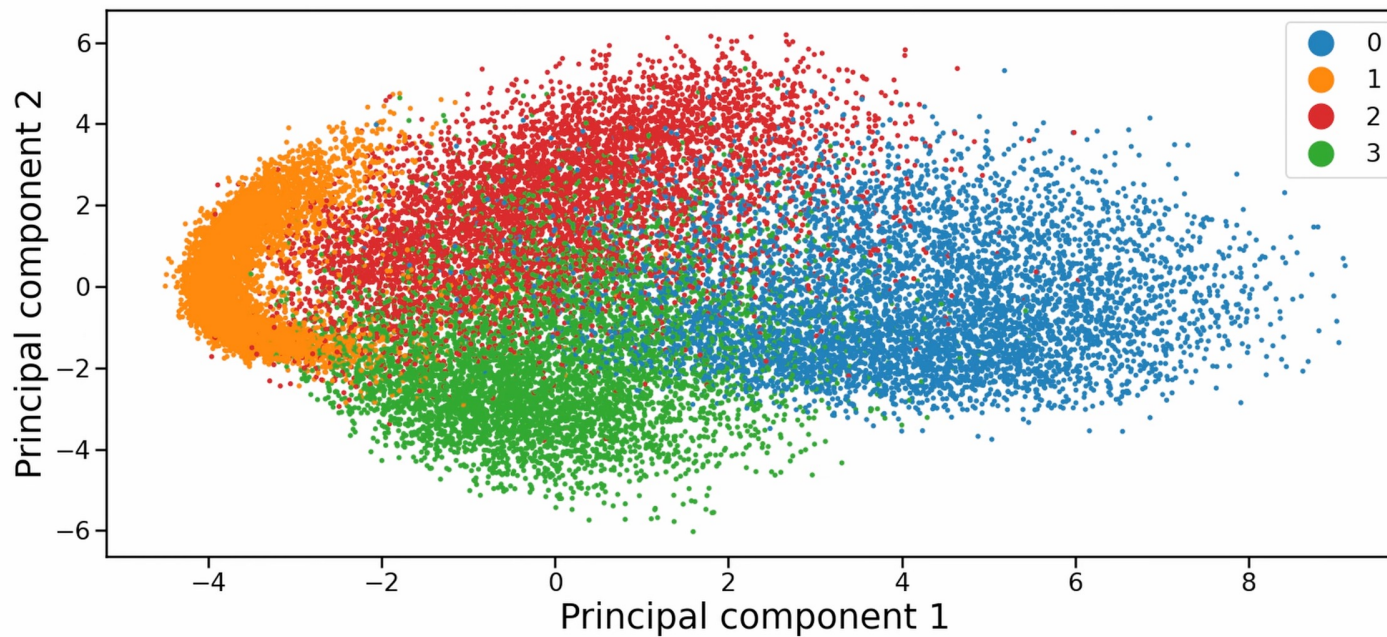


PC2



PC3

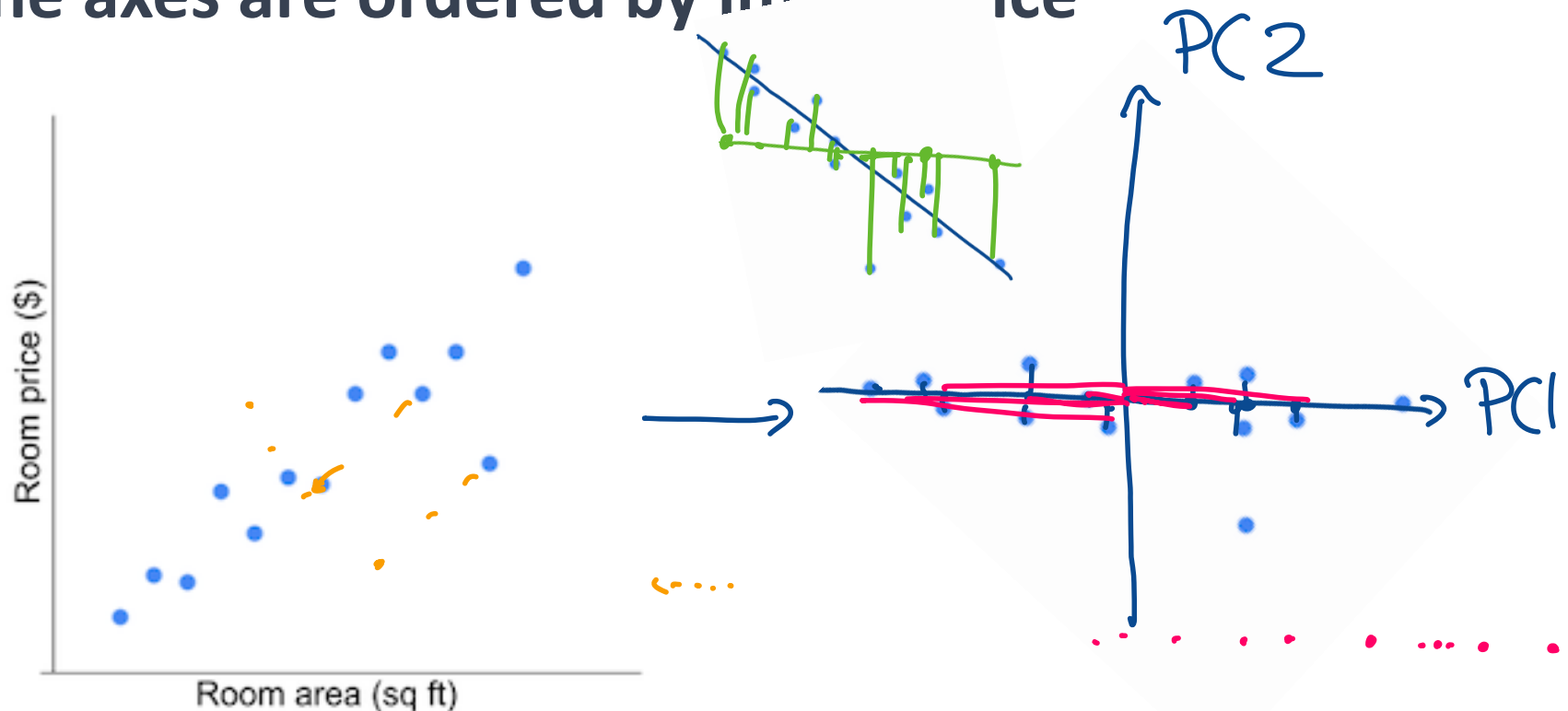
...



Principal Component Analysis (PCA)

Principal Component Analysis (PCA) is a method that reduces the dimensionality of data by transforming it into principal components, each representing unique variance, while retaining the most significant information.

“Finds a better coordinate system for given data, such that the axes are ordered by importance”



PCA – Variance

Given numbers $X = \{x_1, \dots, x_n\} \subset \mathbb{R}$

Average:
$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

Variance:
$$\text{Var}(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{Var}\left(\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\right) > \text{Var}\left(\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\overset{\cdot}{\text{---}}\right)$$

PCA – Orthogonal projection

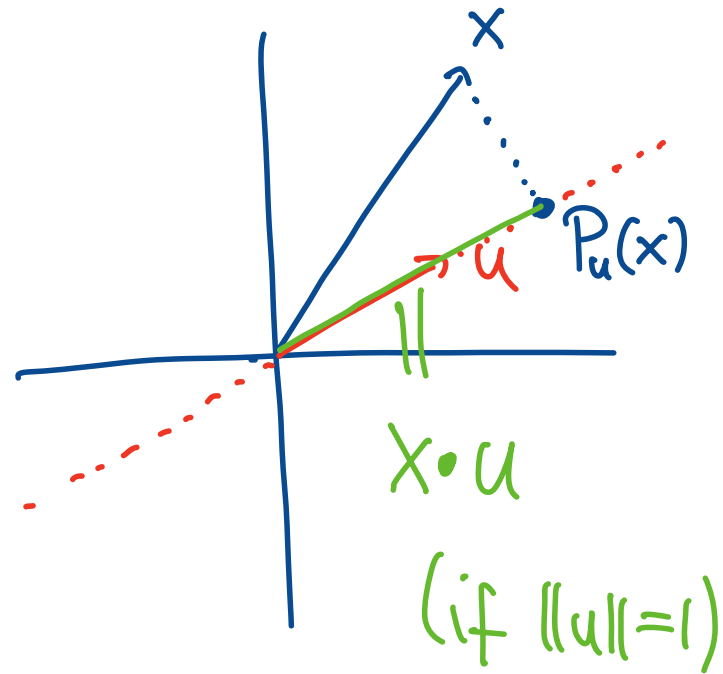
Given $u \in \mathbb{R}^n$ $x \cdot u = x^T u$

orthogonal projection:

$$P_u(x) = \frac{x \cdot u}{u \cdot u} u$$

If $\|u\|=1$ ($u \cdot u = 1$)

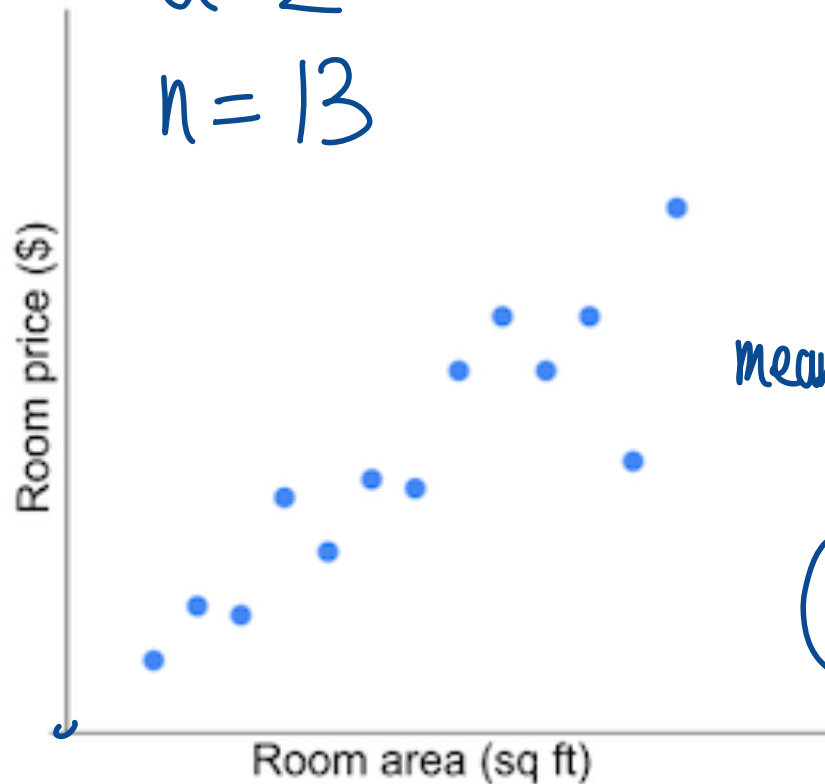
$$P_u(x) = (x \cdot u) u$$



PCA – First principal component

How to find the first “principal component”

$d=2$
 $n=13$



Given $x^{(i)} \in \mathbb{R}^d \quad 1 \leq i \leq n$
 $X = \{x^{(i)}\}$

Step 1: Scale s.t.

$$\text{mean } \frac{1}{n} \sum_{i=1}^n x^{(i)} = 0$$

(standard deviation = 1)

PCA -

How to find the first "principal component" $\|u\|=1$

after scaling:

Want $u \in \mathbb{R}^d$ such that

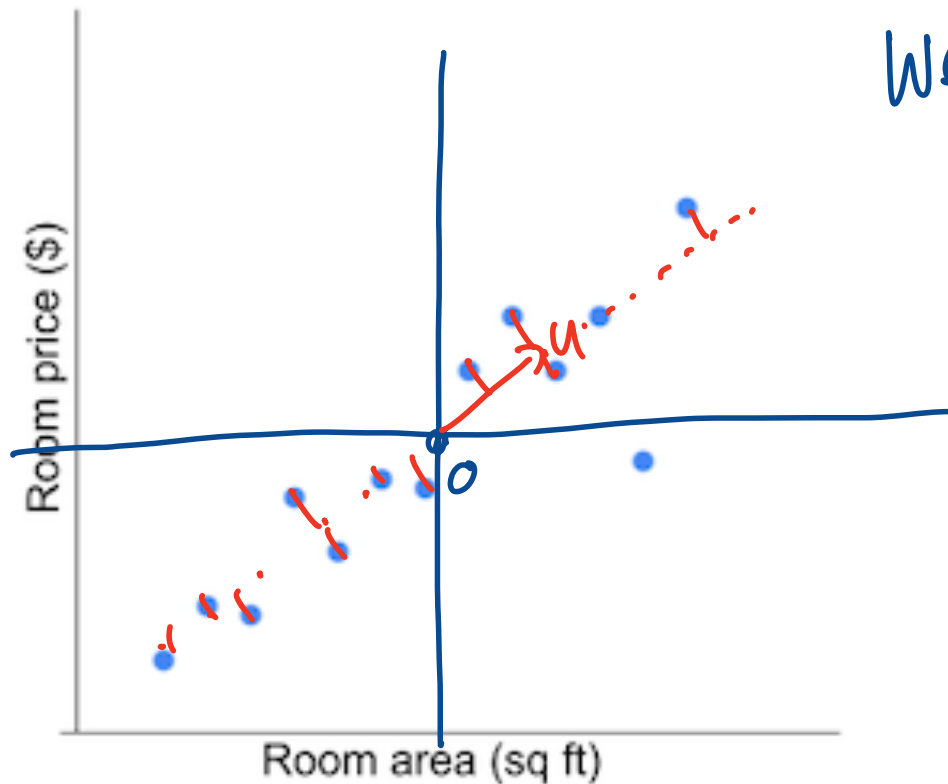
We maximize:

$$\frac{1}{n} \sum_{i=1}^n (x^{(i)} \cdot u)^2$$

$$\mathbb{R}^{d \times d} \quad (x^{(i)T} u)^T x^{(i)T} u$$

$$u^T x^{(i)} x^{(i)T} u$$

$$= u^T \left[\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)T} \right] u$$



$$= u^T \left[\frac{1}{n} \sum_{i=1}^n x^{(i)} x^{(i)T} \right] u \quad X = \begin{matrix} | & & | \\ \hline & d & \\ & x^{(1)} & \\ & x^{(2)} & \\ & \vdots & \\ & x^{(n)} & \\ \hline & & | \end{matrix}$$

Σ : covariance matrix

$$\frac{1}{n} X X^T \quad \text{In example}$$

$$\Sigma = \begin{pmatrix} & \\ & \end{pmatrix}$$

Covariance
of room size & price

$$A \in \mathbb{R}^{n \times n}, x \neq 0$$

$$Ax = \lambda x$$

λ : eigenvalue, x : eigenvector

$$A = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}, \quad A \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 1 \cdot \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$A \begin{pmatrix} 0 \\ 1 \end{pmatrix} = 2 \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

Summary: For symmetric M we want to maximize $u^T M u$.

Theorem: For any sym. $M \in \mathbb{R}^{d \times d}$

the term $\frac{u^T M u}{u^T u}$ (Rayleigh quotient)

is maximal if u is an eigenvector of M for the largest eigenvalue.

Proof: By the spectral theorem

$$M = Q D Q^T \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$$

$$D = \begin{pmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_d \end{pmatrix}$$

Assume $\|u\|=1$

$$u^T M u = u^T Q D Q^T u$$

$$u = Qy = \underbrace{(Q^T u)^T}_y D Q^T u$$

$$Q^T u = y = y^T D y$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_d \end{pmatrix}$$

$$\|y\| = 1$$

$$y_1^2 + y_2^2 + \dots + y_d^2 = 1$$

$$y^T D y = \sum_{i=1}^d \lambda_i y_i^2$$

$\begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{pmatrix}$

maximal for $y = \begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix}$
 (then it is λ_1)

$$u = Qy$$

$Q \begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix} =$ first column of Q
 $=$ eigenvector for λ_1 .