

Introduction to Machine Learning

Lecture 4 - PCA and K-Means

1 The Covariance Matrix

1.1 One-dimensional random variable

Let $X \in \mathbb{R}$ be a one-dimensional random variable and let $\{x_i\}_{i=1}^N$ be N i.i.d realizations from X .

- The first moment of X is given by:

$$\mu_x = \mathbb{E}[X] \approx \frac{1}{N} \sum_{i=1}^N x_i$$

- The variance of X is given by:

$$\sigma_x^2 = \mathbb{E}[(X - \mu_x)^2] \approx \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)^2$$

1.2 High-dimensional random vector

Let $X \in \mathbb{R}^d$ be a high-dimensional random variable and let $\{x_i\}_{i=1}^N$ be N i.i.d realizations from X .

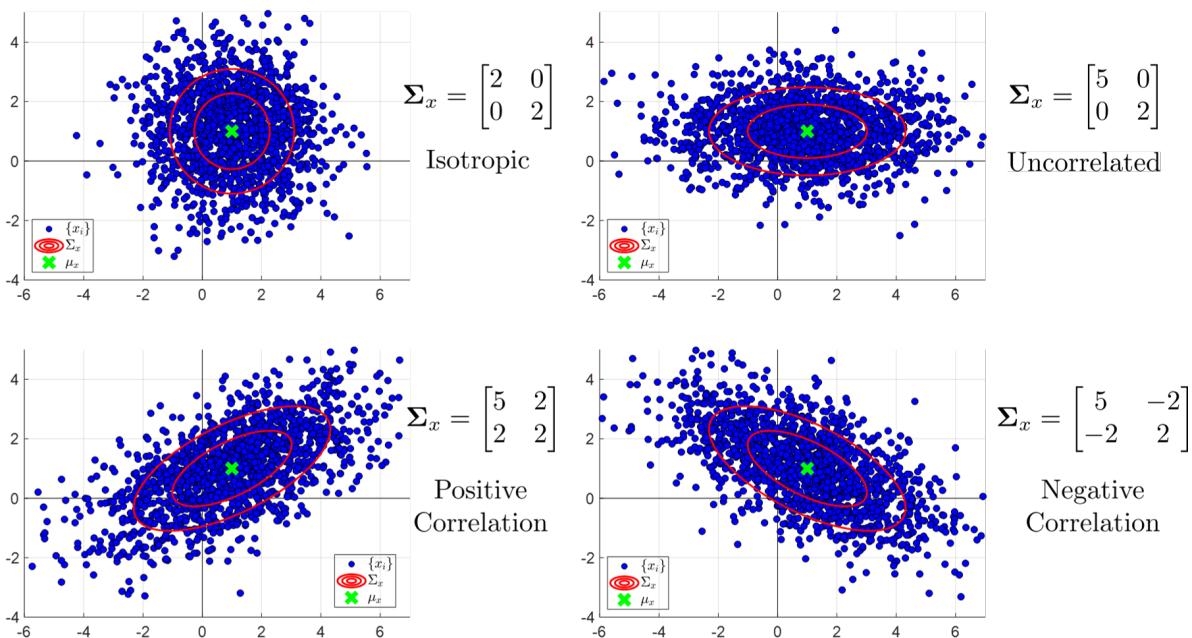
- The first moment of X is given by:

$$\boldsymbol{\mu}_x = \mathbb{E}[X] \approx \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

- The covariance matrix of X is given by:

$$\boldsymbol{\Sigma}_x = \mathbb{E}[(X - \boldsymbol{\mu}_x)(X - \boldsymbol{\mu}_x)^T] \approx \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x)(\mathbf{x}_i - \boldsymbol{\mu}_x)^T$$

The covariance matrix is the generalization of the one-dimensional variance.



2 Principle Component Analysis (PCA)

Let $\{\mathbf{x}_i\}_{i=1}^N$ be a set of vectors such that $\mathbf{x}_i \in \mathbb{R}^D$.

The **Sample (or empirical) Mean** and **Covariance** are given by:

$$\boldsymbol{\mu}_x = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \triangleq \bar{\mathbf{x}}, \quad \boldsymbol{\Sigma}_x = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x)(\mathbf{x}_i - \boldsymbol{\mu}_x)^T$$

- To have an unbiased estimator for $\boldsymbol{\Sigma}_x$, one can use $\frac{1}{N-1}$ instead of $\frac{1}{N}$.

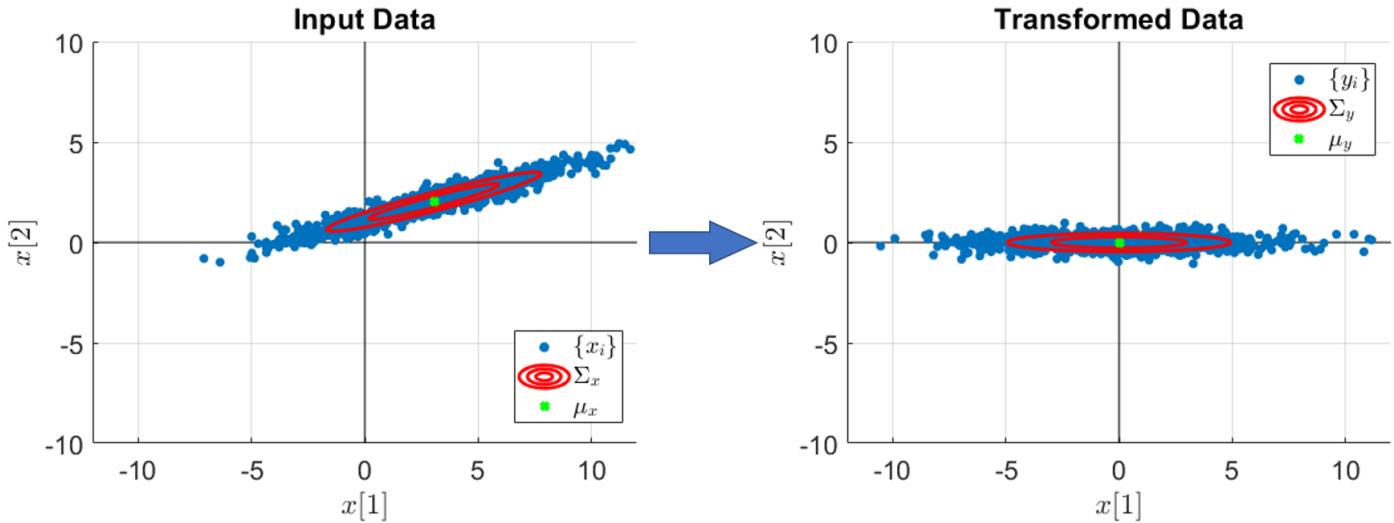
- $\boldsymbol{\Sigma}_x$ is a Positive Semi-Definite (PSD),
that is, $\boldsymbol{\Sigma}_x = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$ where :

 - \mathbf{U} is unitary, namely, $\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I}$

 - $\boldsymbol{\Lambda}$ is a diagonal matrix with non-negative elements, namely $\boldsymbol{\Lambda} = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_D \end{bmatrix}$, and $\lambda_i \geq 0$ for all i .

2.1 Change of basis (axis rotation)

Consider the set $\{\mathbf{x}_i\}_{i=1}^N$, where $N = 1000$, and $\mathbf{x}_i = \begin{bmatrix} \mathbf{x}_i[1] \\ \mathbf{x}_i[2] \end{bmatrix} \in \mathbb{R}^2$ (left figure).



A common preprocess approach is to apply (right figure):

$$\mathbf{y}_i = \phi(\mathbf{x}_i)$$

such that:

- $\{\mathbf{y}_i\}_i$ is centered around zero, namely, $\boldsymbol{\mu}_y = \mathbf{0}$.
- $\boldsymbol{\Sigma}_y$ is a diagonal matrix.
- $\|\mathbf{x}_i - \mathbf{x}_j\|_2 = \|\mathbf{y}_i - \mathbf{y}_j\|_2$ for all i, j .

The following Lemma suggests such map ϕ .

Lemma 1. The following transformation:

$$\mathbf{y}_i = \phi(\mathbf{x}_i) = \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x)$$

where the columns of \mathbf{U} are the eigenvectors of Σ_x satisfies (1), (2) and (3).

Proof.

- 1. By definition:

$$\boldsymbol{\mu}_y = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i = \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) = \mathbf{U}^T \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x) = \mathbf{U}^T \left(\underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i}_{\boldsymbol{\mu}_x} - \boldsymbol{\mu}_x \right) = 0$$

- 2. Consider the eigen-value decomposition: $\Sigma_x = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T$ (\mathbf{U} is unitary) Hence,

$$\begin{aligned} \Sigma_y &= \frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \boldsymbol{\mu}_y) (\mathbf{y}_i - \boldsymbol{\mu}_y)^T = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i \mathbf{y}_i^T \\ &= \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) (\mathbf{x}_i - \boldsymbol{\mu}_x)^T \mathbf{U} = \mathbf{U}^T \underbrace{\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x) (\mathbf{x}_i - \boldsymbol{\mu}_x)^T}_{=\Sigma_x} \mathbf{U} \\ &= \mathbf{U}^T \Sigma_x \mathbf{U} = \mathbf{U}^T \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \mathbf{U} = \boldsymbol{\Lambda} \end{aligned}$$

- 3. Since \mathbf{U}^T is unitary, we have $\|\mathbf{U}^T \mathbf{v}\|_2 = \|\mathbf{v}\|_2$, thus:

$$\|\mathbf{y}_i - \mathbf{y}_j\|_2 = \left\| \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) - \mathbf{U}^T (\mathbf{x}_j - \boldsymbol{\mu}_x) \right\|_2 = \left\| \mathbf{U}^T (\mathbf{x}_i - \mathbf{x}_j) \right\|_2 = \|\mathbf{x}_i - \mathbf{x}_j\|_2$$

□

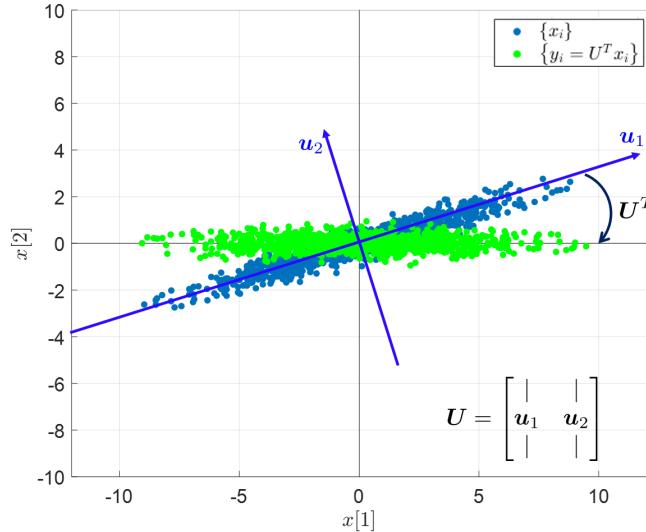
- It is common to arrange the order of $\Sigma_y = \boldsymbol{\Lambda}$ from the largest to the lowest eigenvalue; and respectively, the order of the eigen-vectors \mathbf{U} . Namely, if:

$$\boldsymbol{\Sigma}_y = \text{diag}([\lambda_1, \lambda_2, \dots, \lambda_D])$$

Then:

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$$

- Note that after the mean subtraction $\boldsymbol{\mu}_x$, the matrix \mathbf{U}^T only rotates the data to its main principal components.



- Note that \mathbf{x}_i can be recovered from \mathbf{y}_i by:

$$\mathbf{x}_i = \mathbf{U} \mathbf{y}_i + \boldsymbol{\mu}_x$$

2.2 Linear (affine) dimensionality reduction (using PCA)

- Let D be the dimension of the input data, which also known as the **features vector**: $\mathbf{x}_i \in \mathbb{R}^D$.
- D might be a vary large number.
 - For example, D can be the number of pixels in a High Definition (HD) image, $D = 1920 \times 1080$.
- To reduce the length of the features vector we can apply PCA but consider only the first $d < D$ components of \mathbf{y}_i .

2.2.1 Example 1

- Let $\{\mathbf{x}_i\}_i$ be a $2D$ data set, namely $\mathbf{x}_i \in \mathbb{R}^2$.
For simplicity we set $\mu_x = \mathbf{0}$ and let $\Sigma_x = \mathbf{U}\Lambda\mathbf{U}^T$.
- The full PCA (namely, only a rotation) is given by:

$$\mathbf{y}_i = \mathbf{U}^T \mathbf{x}_i = \begin{bmatrix} -\mathbf{u}_1^T \\ -\mathbf{u}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{x}_i \\ | \\ | \end{bmatrix} \in \mathbb{R}^2$$

and we can perfectly reconstruct \mathbf{x}_i by:

$$\mathbf{x}_i = \mathbf{U} \mathbf{y}_i = \begin{bmatrix} \mathbf{u}_1 \\ | \\ \mathbf{u}_2 \\ | \end{bmatrix} \begin{bmatrix} \mathbf{y}_i \\ | \\ | \end{bmatrix} \in \mathbb{R}^2$$

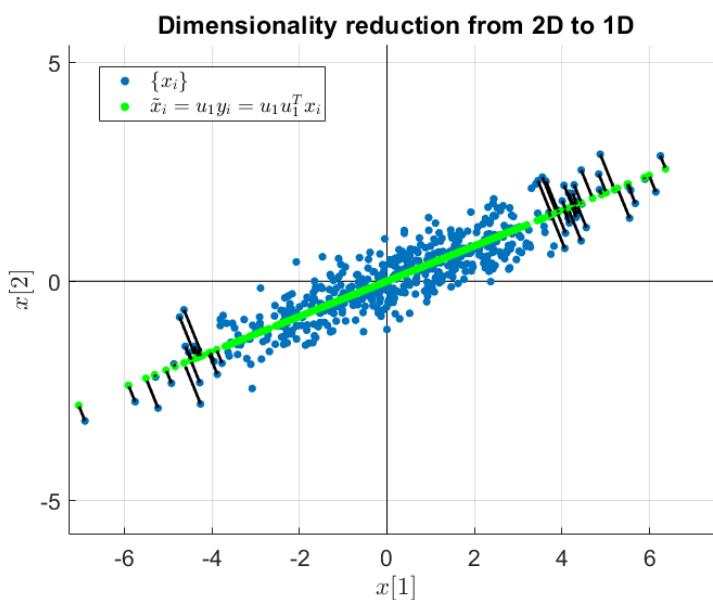
- One can reduce the dimension of $\mathbf{y}_i \in \mathbb{R}^2$ from 2 to 1 by consider only the first eigenvector:

$$y_i = \mathbf{u}_1^T \mathbf{x}_i = [-\mathbf{u}_1^T] \begin{bmatrix} \mathbf{x}_i \\ | \\ | \end{bmatrix} \in \mathbb{R}$$

where \mathbf{u}_1 is the first leading eigen-vector of Σ_x (corresponding to the largest eigen-vlaue λ_1)

- The imperfect reconstruction is given by:

$$\tilde{\mathbf{x}}_i = \mathbf{u}_1 y_i = \begin{bmatrix} \mathbf{u}_1 \\ | \end{bmatrix} y_i = \begin{bmatrix} \mathbf{u}_1 \\ | \end{bmatrix} [-\mathbf{u}_1^T] \begin{bmatrix} \mathbf{x}_i \\ | \\ | \end{bmatrix} \in \mathbb{R}^2$$



Reconstruction error Let us compute the error between an original \mathbf{x}_i and its reconstructed version $\tilde{\mathbf{x}}_i$

$$\epsilon_i^2 \triangleq \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2^2 = ?$$

$$\begin{aligned}\epsilon_i^2 &= \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\|_2^2 \\ &= \left\| \mathbf{U}^T \mathbf{x}_i - \mathbf{U}^T \tilde{\mathbf{x}}_i \right\|_2^2 \\ &= \left\| \begin{bmatrix} \mathbf{u}_1^T \mathbf{x}_i \\ \mathbf{u}_2^T \mathbf{x}_i \end{bmatrix} - \mathbf{U}^T \mathbf{u}_1 \mathbf{u}_1^T \mathbf{x}_i \right\|_2^2 \\ &= \left\| \begin{bmatrix} \mathbf{u}_1^T \mathbf{x}_i \\ \mathbf{u}_2^T \mathbf{x}_i \end{bmatrix} - \begin{bmatrix} 1 \\ 0 \end{bmatrix} \mathbf{u}_1^T \mathbf{x}_i \right\|_2^2 \\ &= \left\| \begin{bmatrix} 0 \\ \mathbf{u}_2^T \mathbf{x}_i \end{bmatrix} \right\|_2^2 \\ &= (\mathbf{u}_2^T \mathbf{x}_i)^2\end{aligned}$$

Using the Pythagorean theorem:

$$\begin{aligned}\|\mathbf{x}_i\|^2 &= (\mathbf{u}_1^T \mathbf{x}_i)^2 + (\mathbf{u}_2^T \mathbf{x}_i)^2 \\ \Rightarrow \epsilon_i^2 &= (\mathbf{u}_2^T \mathbf{x}_i)^2 = \|\mathbf{x}_i\|^2 - (\mathbf{u}_1^T \mathbf{x}_i)^2\end{aligned}$$

The overall MSE is obtained by averaging over all points in $\{\mathbf{x}_i\}$:

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N \epsilon_i^2 = \frac{1}{N} \sum_{i=1}^N \left(\|\mathbf{x}_i\|^2 - (\mathbf{u}_1^T \mathbf{x}_i)^2 \right)$$

Exercise: Prove that the first leading eigen-vector \mathbf{u}_1 is indeed the minimizer of the MSE. Formally, show that

$$\mathbf{u}_1 = \arg \min_{\|\mathbf{u}\|=1} \text{MSE}(\mathbf{u})$$

Solution: Using Lagrangian multipliers ($\|\mathbf{u}\| = \|\mathbf{u}\|^2 = \mathbf{u}^T \mathbf{u} = 1$), the Lagrangian is:

$$\begin{aligned}\mathcal{L}(\mathbf{u}, \lambda) &= \text{MSE}(\mathbf{u}) + \lambda (\mathbf{u}^T \mathbf{u} - 1) \\ &= \frac{1}{N} \sum_{i=1}^N \left(\|\mathbf{x}_i\|^2 - (\mathbf{u}^T \mathbf{x}_i)^2 \right) + \lambda (\mathbf{u}^T \mathbf{u} - 1) \\ \Rightarrow \nabla_{\mathbf{u}} \mathcal{L}(\mathbf{u}, \lambda) &= 0 \\ -\frac{1}{N} \sum_{i=1}^N 2\mathbf{x}_i \mathbf{x}_i^T \mathbf{u} + 2\lambda \mathbf{u} &= 0 \\ \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T}_{\Sigma_x} \mathbf{u} &= \lambda \mathbf{u} \\ \Sigma_x \mathbf{u} &= \lambda \mathbf{u}\end{aligned}$$

$$\Rightarrow [\mathbf{u} \text{ is an eigen-vector of } \Sigma_x]$$

Now:

$$\begin{aligned}
 \text{MSE}(\mathbf{u}) &= \frac{1}{N} \sum_{i=1}^N \left(\|\mathbf{x}_i\|_2^2 - \|\mathbf{x}_i^T \mathbf{u}\|_2^2 \right) \\
 &= C - \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i^T \mathbf{u}\|_2^2 \\
 &= C - \frac{1}{N} \sum_{i=1}^N \mathbf{u}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} \\
 &= C - \mathbf{u}^T \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T}_{=\Sigma_x} \mathbf{u} \\
 &= C - \mathbf{u}^T \Sigma_x \mathbf{u} \\
 &= C - \lambda \underbrace{\mathbf{u}^T \mathbf{u}}_{=1} \\
 &= C - \lambda
 \end{aligned}$$

Thus, to minimize the MSE, we should choose \mathbf{u} which corresponds to the largest eigen-value.

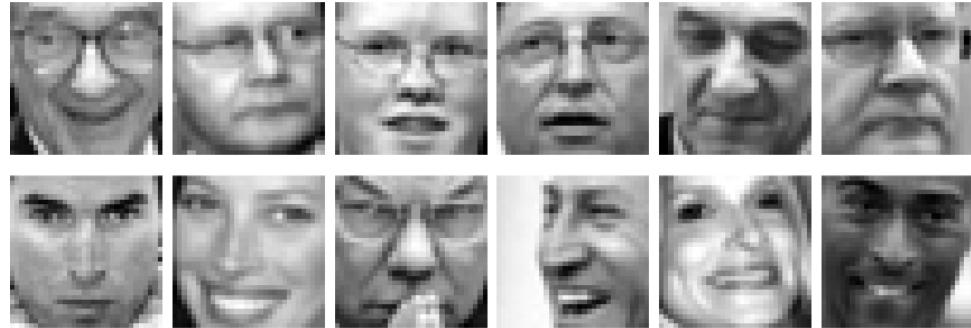
■

Note: C is in fact the trace of Σ_x (which is the sum of the eigen-values):

$$C = \text{tr} \{ \Sigma_x \} = \sum_{i=1}^D \lambda_i (\Sigma_x)$$

2.2.2 Example 2

- Consider the set $\{\mathbf{x}_i\}_{i=1}^N$ of $N = 5,000$ images of faces.
- Each image has 32×32 pixels.
- Each \mathbf{x}_i is a column stack of all pixels, that is, $\mathbf{x}_i \in \mathbb{R}^{1,024}$.



We can apply PCA by:

$$\mathbf{y}_i = \mathbf{U}_1^T (\mathbf{x}_i - \boldsymbol{\mu}_x) \in \mathbb{R}^d$$

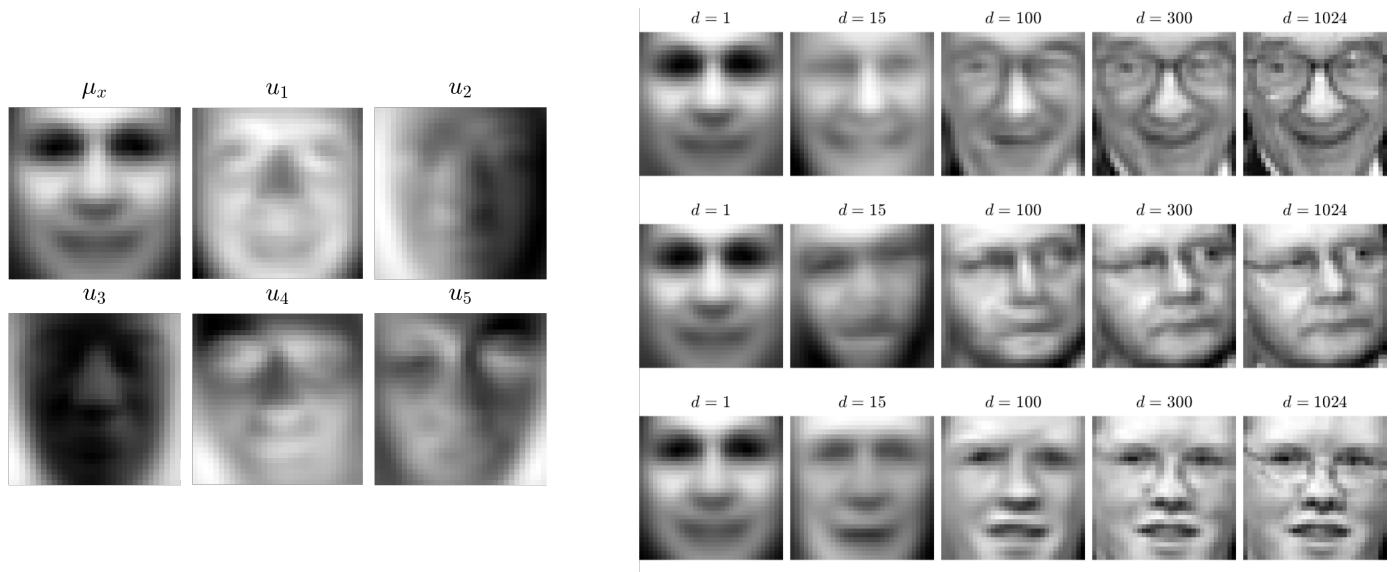
where \mathbf{U}_1 are the first d eigen-vectors of $\Sigma_x = \mathbf{U}\Lambda\mathbf{U}^T$:

$$\mathbf{U} = [\mathbf{U}_1 \quad \mathbf{U}_2], \quad \mathbf{U}_1 \in \mathbb{R}^{D \times d}$$

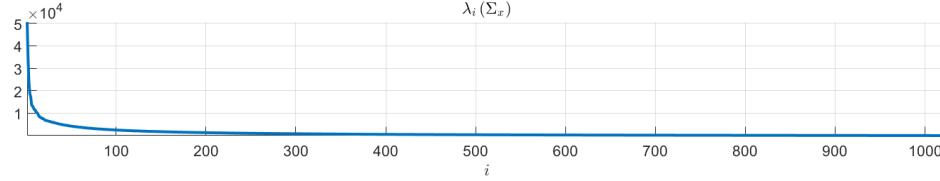
The reconstruction is given by:

$$\tilde{\mathbf{x}}_i = \mathbf{U}_1 \mathbf{y}_i + \boldsymbol{\mu}_x$$

We present the reconstruction $\tilde{\mathbf{x}}_i$ from the low dimensionality representation \mathbf{y}_i for several values of d .



The eigenvalues $\{\lambda_i(\Sigma_x)\}_{i=1}^D$ of Σ_x can tell us about the reconstruction error:



Notice that most of the data variance is captured by the first 200 eigen-vectors.

2.3 Comments

1. Σ_x is a Positive Semi Definite (PSD), namely:

$$\mathbf{v}^T \Sigma_x \mathbf{v} \geq 0, \quad \forall \mathbf{v}$$

2. Any PSD matrix has non-negative eigenvalues and orthogonal eigen-vectors:

$$\Sigma_x = \mathbf{U} \Lambda \mathbf{U}^T, \quad \mathbf{U} \mathbf{U}^T = \mathbf{U}^T \mathbf{U} = \mathbf{I}$$

3. Let the columns of $\mathbf{X} \in \mathbb{R}^{D \times N}$ contain the set $\{\mathbf{x}_i\}_i^D$:

$$\mathbf{X} \triangleq \begin{bmatrix} | & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_N \\ | & & | \end{bmatrix} \in \mathbb{R}^{D \times N}$$

assume zero mean ($\mu_x = \mathbf{0}$).

- (a) The eigen-vectors of Σ_x are equal to the left singular vectors of \mathbf{X} .
- (b) The eigen-values of Σ_x are equal to the singular values (squared) of \mathbf{X} , namely:

$$\sigma_i^2(\mathbf{X}) = \lambda_i(\Sigma_x)$$

4. The reconstruction error as a function of d (the number of dimensions) is given by the eigen-values:

$$\text{err}(d) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \tilde{\mathbf{x}}_i\| = \sum_{i=d+1}^D \lambda_i$$

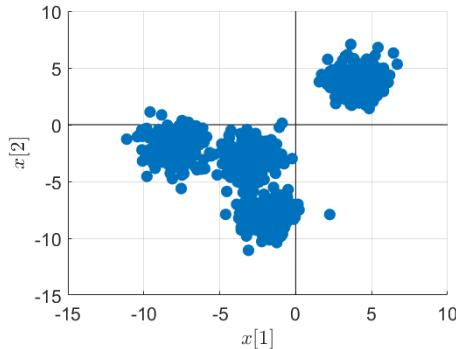
where

$$\tilde{\mathbf{x}}_i = \mathbf{U}_1 \mathbf{U}_1^T \mathbf{x}_i, \quad \mathbf{U}_1 \in \mathbb{R}^{D \times d}$$

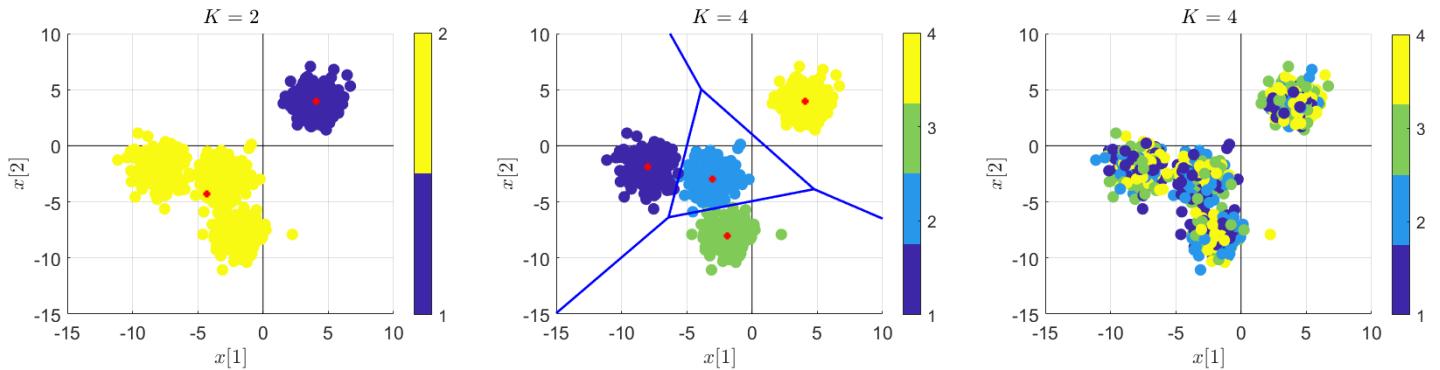
3 K-Means (Clustering)

- Consider the set $\{x_i\}_{i=1}^N$.
- We want to cluster (split) the set $\{x_i\}_i$ into $K > 1$ subsets.

Consider, for example, the following set:



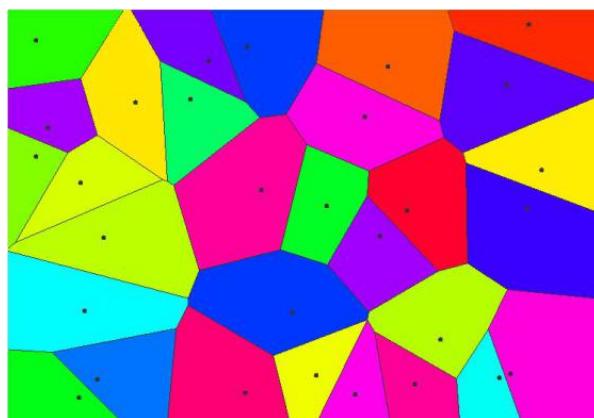
We can split it into $K = 2$ subsets or $K = 4$ subsets like so:



- The red markers are the cluster **centroid**.
- The last clustering (right) makes less sense (but it is a valid clustering into $K = 4$ subsets).
- In order to have a meaningful clustering we need a measure of similarity (or dissimilarity).

3.1 Voronoi

Given a set of centroids $\{\mu_k\}_{k=1}^K$,
the Voronoi diagram plots the regions consisting of all points closer to their centroid than to any other.



3.2 Metric Function

A function $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is called a metric function if it satisfies the following:

1. $d(\mathbf{x}_i, \mathbf{x}_j) = d(\mathbf{x}_j, \mathbf{x}_i)$ (Symmetry).
2. $d(\mathbf{x}_i, \mathbf{x}_j) = 0$ if and only if $\mathbf{x}_i = \mathbf{x}_j$.
3. $d(\mathbf{x}_i, \mathbf{x}_k) \leq d(\mathbf{x}_i, \mathbf{x}_j) + d(\mathbf{x}_j, \mathbf{x}_k)$ (triangle inequality).

3.2.1 Some Examples

1. L_2 (Euclidean):

$$d_E(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|_2 = \sqrt{\sum_{i=1}^D (\mathbf{x}_1[i] - \mathbf{x}_2[i])^2}$$

2. L_1 :

$$d_M(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|_1 = \sum_{i=1}^D |\mathbf{x}_1[i] - \mathbf{x}_2[i]|$$

3. log-Euclidean:

$$d(x_1, x_2) = |\log(x_1) - \log(x_2)|, \quad x_1, x_2 > 0$$

Exercise: Prove that any metric function d , satisfies:

$$d(\mathbf{x}_i, \mathbf{x}_j) \geq 0$$

for all $\mathbf{x}_i \in \mathcal{X}$ and $\mathbf{x}_j \in \mathcal{X}$.

Solution:

$$0 = d(\mathbf{x}_i, \mathbf{x}_i) \leq d(\mathbf{x}_i, \mathbf{x}_j) + d(\mathbf{x}_j, \mathbf{x}_i) = 2d(\mathbf{x}_i, \mathbf{x}_j)$$

$$\Rightarrow \boxed{d(\mathbf{x}_i, \mathbf{x}_j) \geq 0}$$

■

3.2.2 Dissimilarity Function

Dissimilarity function $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is usually symmetric and non-negative but does not have to be a metric function. For example, the squared Euclidean:

$$d_E^2(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|^2 = \sum_{i=1}^D (\mathbf{x}_i[i] - \mathbf{x}_j[i])^2$$

does not satisfy the triangle inequality:

$$100 = d_E^2(10, 0) \not\leq d_E^2(10, 5) + d_E^2(5, 0) = 25 + 25 = 50$$

Another example is the cosine dissimilarity:

$$d(\mathbf{x}_i, \mathbf{x}_j) = 1 - \frac{\langle \mathbf{x}_i, \mathbf{x}_j \rangle}{\|\mathbf{x}_i\|_2 \|\mathbf{x}_j\|_2}$$

Formulation Using a dissimilarity function d , we write the clustering problem as follows:

$$\min_{\{\mathcal{C}_k\}, \{\boldsymbol{\mu}_k\}} \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} d(\mathbf{x}_i, \boldsymbol{\mu}_k)$$

where $\mathcal{C}_k \subseteq \{\mathbf{x}_i\}$ are the **disjoint** subsets such that $\bigsqcup_k \mathcal{C}_k = \{\mathbf{x}_i\}$, and $\{\boldsymbol{\mu}_k\}$ are the clusters centroids.

- Choosing the dissimilarity function will affect the obtained clustering.
- This problem is not convex, hence, finding the global minimum is not guaranteed.

3.3 K-Means Using Squared Euclidean Dissimilarity

Using the squared Euclidean measure we have:

$$\{\mathcal{C}_k\}, \{\boldsymbol{\mu}_k\} = \arg \min_{\{\mathcal{C}_k\}, \{\boldsymbol{\mu}_k\}} \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2 \quad (1)$$

3.3.1

For fixed centroids $\{\boldsymbol{\mu}_k\}$.

What are the optimal $\{\mathcal{C}_k\}$?

$$\min_{\{\mathcal{C}_k\}} \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$$

Note that each point \mathbf{x}_i is assigned only to a **single** centroid (one subset).

Thus, to minimize this problem, one should set:

$$\boxed{\mathcal{C}_s = \left\{ \mathbf{x}_i \mid \|\mathbf{x}_i - \boldsymbol{\mu}_s\|_2^2 \leq \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2 \right\}, \quad k \in \{1, 2, \dots, K\}}$$

- In case of equality $\|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2 = \|\mathbf{x}_i - \boldsymbol{\mu}_{k'}\|_2$ we can assign \mathbf{x}_i to the cluster with the smaller index $\min(k, k')$.

In other words, we can write:

$$\boxed{\mathbf{x}_i \in \mathcal{C}_s \iff s = \arg \min_k \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2}$$

3.3.2

For fixed clusters $\{\mathcal{C}_k\}$..

What are the optimal $\{\boldsymbol{\mu}_k\}$?

We can compare the gradient with respect to $\boldsymbol{\mu}_s$ to zero:

$$\begin{aligned} \nabla_{\boldsymbol{\mu}_s} \left(\sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2 \right) &= \mathbf{0} \\ -2 \sum_{\mathbf{x}_i \in \mathcal{C}_s} (\mathbf{x}_i - \boldsymbol{\mu}_s) &= \mathbf{0} \end{aligned}$$

$$\Rightarrow \boxed{\boldsymbol{\mu}_s = \frac{1}{|\mathcal{C}_s|} \sum_{\mathbf{x}_i \in \mathcal{C}_s} \mathbf{x}_i}$$

In words, the optimal centroid $\boldsymbol{\mu}_k$ of the k th cluster \mathcal{C}_k is the mean of the cluster.

3.3.3 Algorithm

We can find a local minimum in Eq. (1) using the following alternating iterative approach:

Algorithm 1 K-Means

1. Set initial centroids $\{\mu_k\}$
2. For $s = 1, 2, \dots, K$, find the clusters \mathcal{C}_s by:

$$\mathcal{C}_s = \left\{ \mathbf{x}_i \mid \|\mathbf{x}_i - \mu_s\|_2^2 \leq \|\mathbf{x}_i - \mu_k\|_2^2, \quad \forall k \right\}$$

3. Update the centroids:

$$\mu_s = \frac{1}{|\mathcal{C}_s|} \sum_{\mathbf{x}_i \in \mathcal{C}_s} \mathbf{x}_i$$

4. Repeat 2 - 3 until convergence.
-

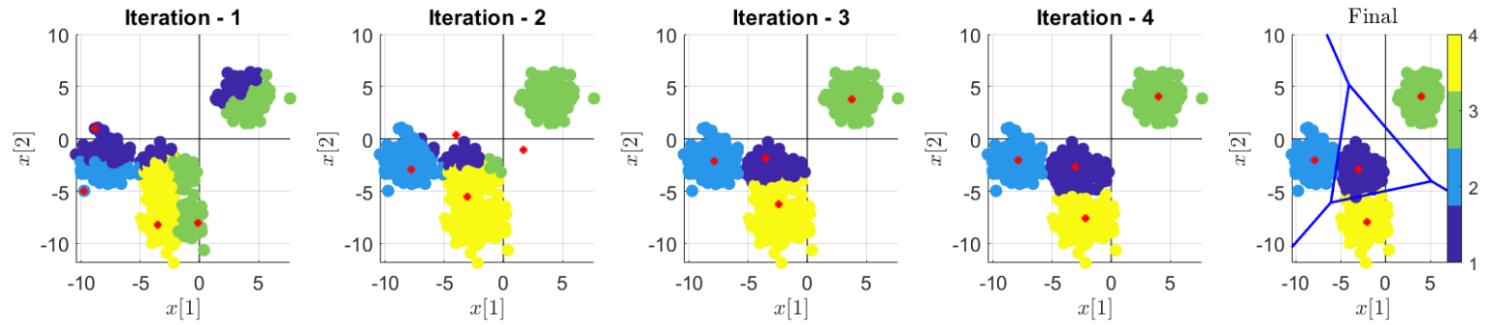
Notes:

1. Choosing the number of clusters K is an important issue.
2. Different initial conditions might lead to different clustering.

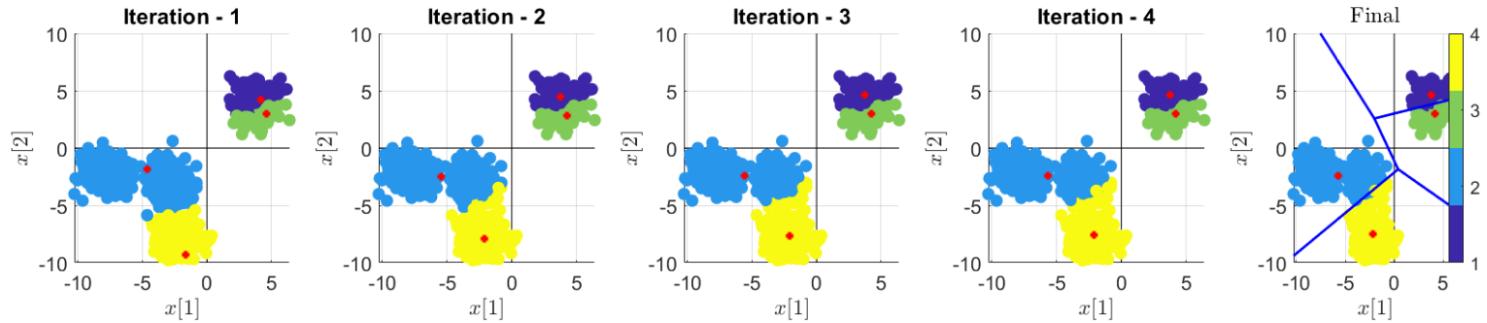
3.3.4 Example 1 - sensitivity to initial conditions

Given the set $\{\mathbf{x}_i\}_{i=1}^N$, we set the initial conditions $\{\mu_k\}$ twice:

- First run:



- Second run:



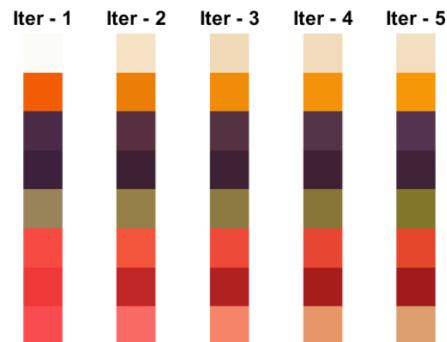
3.3.5 Example 2

Consider the following image $I \in \mathbb{R}^{384 \times 512 \times 3}$:

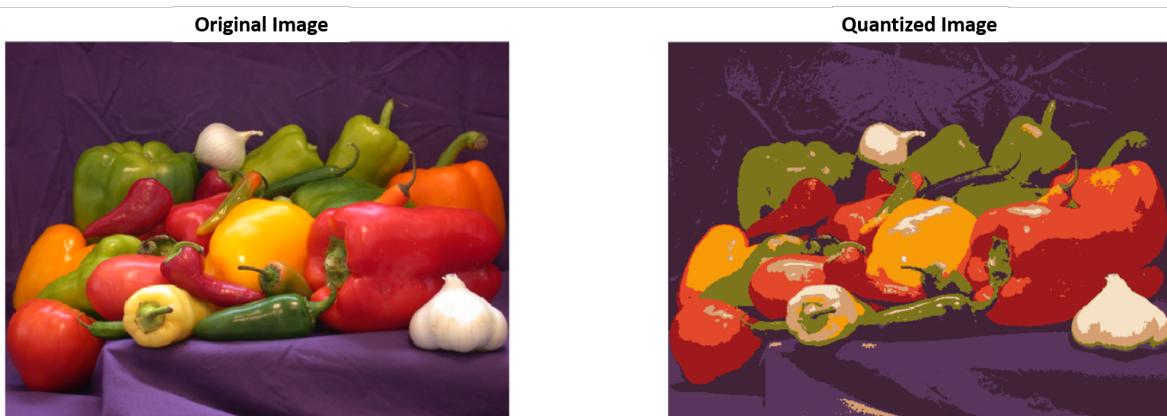


- The color image is a combination of 3 images: a red level image, a blue level image, and a green level image.
- Each color (red, blue and green) is represented by $2^8 = 256$ levels.
- Each pixel $x_i \in \mathbb{R}^3$ contains the values of the three colors (Red, Green and Blue).
- Each pixel x_i can represent $2^8 \cdot 2^8 \cdot 2^8 = 2^{24}$ different colors.
- Thus, every pixel requires 24 bits in the memory.

We can cluster the 2^{24} colors into $2^3 = 8$ centroids using the K-Means algorithm:



We can represent each pixel by its associated centroid:



The clustered (quantized) image requires only 3 bits per pixel (instead of 24).