

Linear Unsupervised Learning

Instructor: Hui Guan

Slides adapted from: https://github.com/jermwatt/machine_learning_refined

Outline

- Background: Fixed Spanning Sets, Orthonormality, and Projections
- The Linear Autoencoder and Principal Component Analysis
- K-means Clustering

8.2 Fixed Spanning Sets, Orthonormality, and Projections

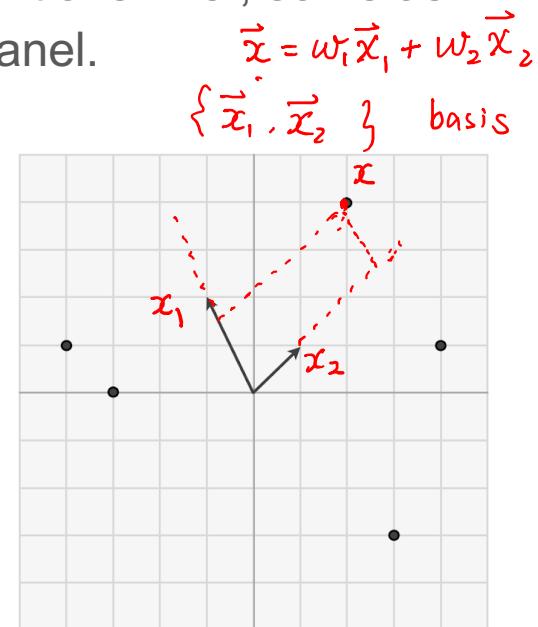
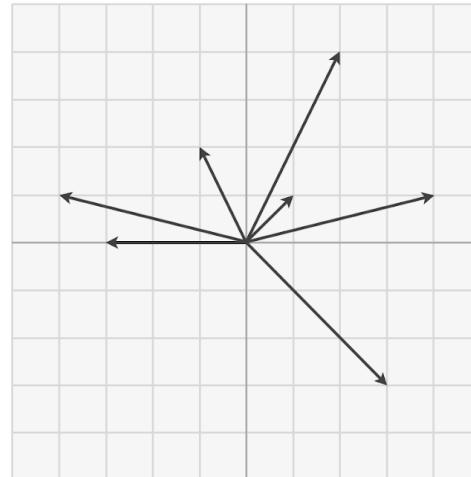
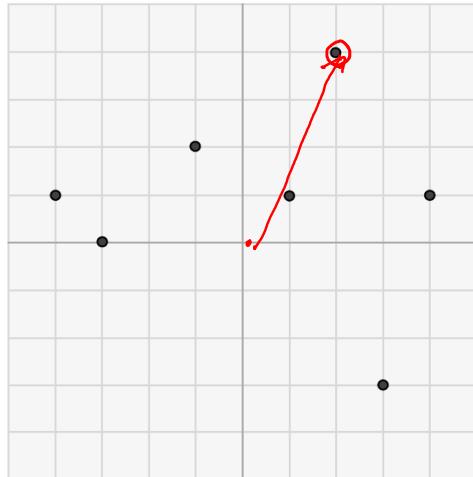
- The fundamental principles of unsupervised learning all lie in some very rudimentary ideas from linear algebra combined with the cost-function based approach to learning we have seen previously with supervised methods.
- This is because unsupervised learning is about effectively representing datasets with no output, i.e., ***only input data***.
- Our goal with such data is simply to effectively represent it, and this is a task that has roots in the notion of ***a spanning set of vectors*** from basic linear algebra.

- In this Section we review the fundamental concepts from linear algebra surrounding the notion of a spanning set of vectors.
- This includes: **linear independence**, learning proper weights for a fixed basis set representation, **orthogonal bases** and representation, and projections.
- Reviewing these topics up front will allow us to much more effectively tackle and intuit unsupervised learning techniques since they can be thought of as direct extensions.

Perfectly representing data via a fixed spanning set

- Remember when thinking about data points in a **multi-dimensional vector space**, we can think of them simultaneously as 'dots' (as shown in the left panel) or as 'arrows' (as in the middle panel below).
- When thinking about *spanning sets of vectors* it is often helpful to visualize points in the same space using both of these conventions - i.e., some as 'dots' and others as 'arrows' as shown in the right panel.

2D



- Those vectors drawn as arrows are particular points - often called a basis or spanning set of vectors - over which we aim to efficiently represent every other point in the space.
- Those other vectors shown as dots are those points we wish to represent over our basis / spanning set.

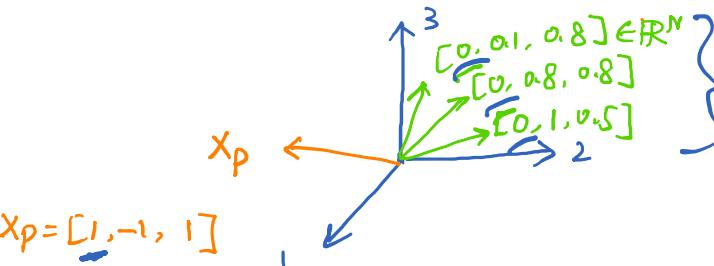
$$\tilde{x}_i = x_i - \frac{1}{P} \sum_{i=1}^P x_i$$

normalization Standardization $\rightarrow \frac{x - \mu}{\sigma} \sim G(0, 1)$

- Let's suppose our set of input points - our *dataset* - is written as $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_P\}$ and lives in N dimensions (that is, each point is N dimensional).
- Moreover throughout we will assume that our dataset has been *mean-centered* - a simple and completely reversible operation that involves subtracting off the mean of the dataset along each input dimension - so that it straddles the origin.
- This sort of normalization is almost always done in practice, and is computationally inexpensive to execute.

- In order for our basis / spanning set to be capable of perfectly representing all P of our points it too must live in the same N dimensional space.
- For a candidate basis / spanning set $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N$ to be capable of perfectly representing such generic N dimensional (input) data means...
- ...that for each data point a set of weights exists so that, in particular linear combination, our basis set can match each data point as

$$\sum_{n=1}^N \underbrace{\mathbf{c}_n}_{\text{vector } C_n \in \mathbb{R}^N} \underbrace{w_{n,p}}_{\text{coefficient}} = \mathbf{x}_p \in \mathbb{R}^N \quad p = 1 \dots P.$$



- Technically speaking in order for this possibility to exist our spanning set must be *linearly independent*
- that is, the spanning vectors do not 'overlap', they point in completely different directions in the space.

$$\vec{x}_p \in \mathbb{R}^N : \underbrace{\vec{c}_1, \dots, \vec{c}_N}_{\geq N} \quad \vec{x}_p = \sum_{n=1}^N c_n \cdot w_{n,p}$$

& at least N vectors are ^{linearly} independent

[given any set of N vectors, we can find $\{v_{jnp}\}_{n=1..N}$ that can

represent \vec{x}_p .] **False**

As the simplest example imagine our spanning set was the set of N *standard basis vectors*. The n^{th} element of a standard basis takes the form of vectors that consist entirely of zeros, except for a 1 in its n^{th} slot

$$(n^{th} \text{ element of the standard basis}) \quad \mathbf{c}_n = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \textcircled{1} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{--- } n^{th}$$

$$\tilde{x}_p = \sum_{n=1}^N c_n \cdot x_{n,p}$$

$$x_{n,p} = \tilde{x}_p [n-1]$$

- To represent a data point \mathbf{x}_p over the standard basis is a trivial affair, and one can easily check that the perfect weights must be defined as

$$w_{n,p} = x_{n,p}$$

- i.e., each weight is simply equal to the value of the data point we aim to represent.
- For most any other spanning set, however these weights must be solved for numerically, which we frame in terms of a cost function minimization after the examples below.

- If indeed the spanning set of vectors is linearly independent, then having properly tuned the weights of point \mathbf{x}_p , denoted as the N length vector

$$\boxed{\mathbf{w}_p} = \begin{bmatrix} w_{1,p} \\ w_{2,p} \\ \vdots \\ w_{N,p} \end{bmatrix}$$

$\vec{\mathbf{x}}_p = \sum_{n=1}^N c_n \cdot \underline{\mathbf{w}_{n,p}}$

vector

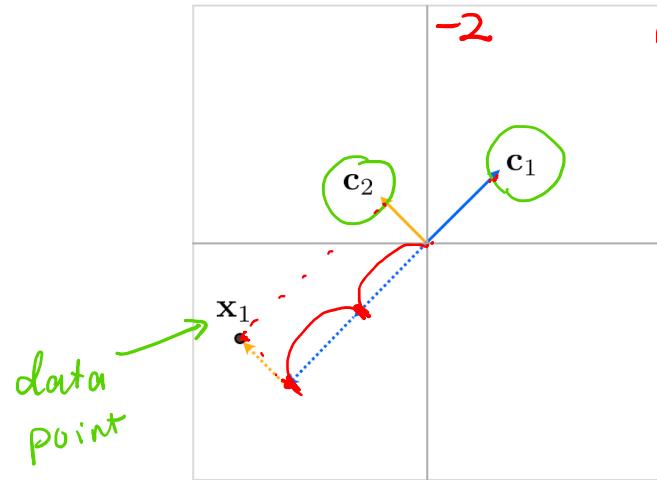
- This weight set provides the **new representation** of \mathbf{x}_p with respect to the spanning set.

$$\vec{\mathbf{x}}_p \underset{\substack{\uparrow \\ \text{original}}}{=} \underbrace{c_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}}_{\substack{\uparrow \\ \text{basis}}} \Rightarrow \boxed{\mathbf{w}_p} : \vec{\mathbf{x}}_p = \{c_n\} \begin{bmatrix} 1.0 \\ 2.0 \\ 3.0 \\ \vdots \end{bmatrix} \cdot \mathbf{w}_p$$

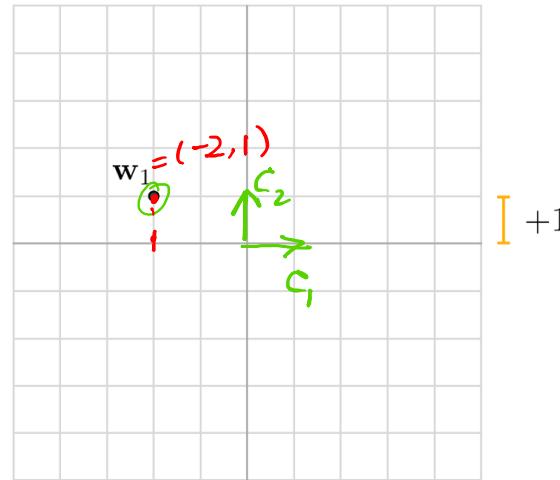
basis

- For example, in the trivial case where our spanning set consists of the standard basis our representation or feature vector $\mathbf{w}_p = \mathbf{x}_p$ is the data point itself!
- Otherwise $\underline{\mathbf{w}_p}$ - typically referred to as the new **encoding** of $\underline{\mathbf{x}_p}$ in the **transformed feature space** whose coordinate axes are defined by the spanning set - is a vector that differs from original data point.

$$\vec{x}_1 = \underline{w^{(1)} c_1} + \underline{w^{(2)} c_2}$$



$$\vec{w}_1 = \begin{bmatrix} w^{(1)} \\ w^{(2)} \end{bmatrix}$$



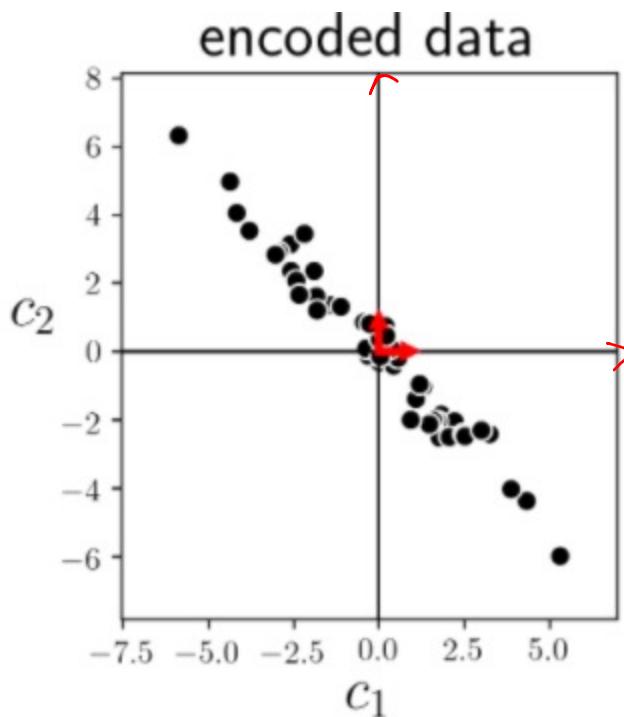
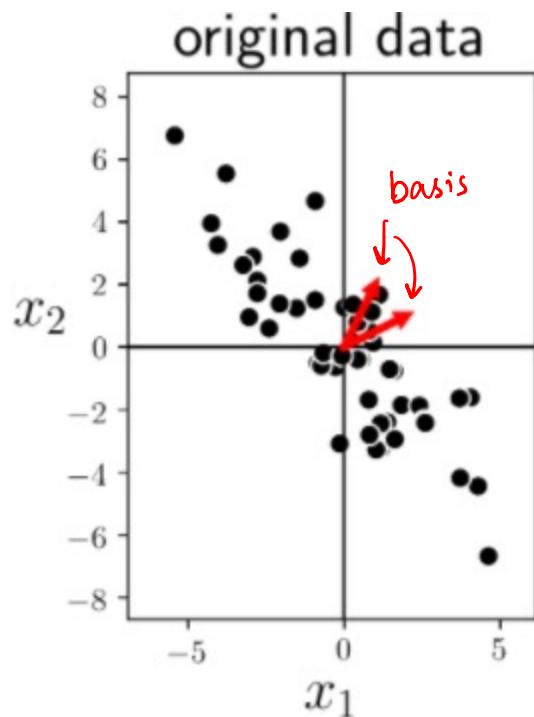
(left) A 2-d point, \mathbf{x}_1 , shown in the space spanned by vectors \mathbf{c}_1 and \mathbf{c}_2 .
 (right) Representation, also known as the *encoding*, of \mathbf{x}_1 in the *transformed feature space*.

Perfect representation with complete spanning set in $N=2$ dimensions

- Given two spanning vectors which are *linearly independent*.

$$\mathbf{c}_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix} \quad \mathbf{c}_2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}.$$

- That is they do not overlap completely and point in different directions in the space, we can perfectly represent any point in the space using some linear combination of them as $w_{1,p}\mathbf{c}_1 + w_{2,p}\mathbf{c}_2 = \mathbf{x}_p$
- Here the weights $w_{1,p}$ and $w_{2,p}$ are unique to each point \mathbf{x}_p .



- One way we can solve for the proper weights is by setting up / minimizing an appropriate cost function.
- It is natural (as in our development of e.g., linear regression) to square the difference as

$$g(w_{p,n}) = \left(\underbrace{\sum_{n=1}^N c_n w_{p,n}}_{\text{Known}} - \underline{x_p} \right)^2$$

- We then tune the corresponding weights to make this squared error as small as possible.

- So for each data point \mathbf{x}_p we have a cost function

$$g(w_{1,p}, \dots, w_{N,p}) = \left(\sum_{n=1}^N \mathbf{c}_n \underline{w_{n,p}} - \mathbf{x}_p \right)^2 \quad p = 1 \dots P$$

- When minimized appropriately provides us with a proper set of parameters for representing \mathbf{x}_p over the spanning set.

- It is perhaps easiest to see this by writing the above cost function more compactly as

$$h(\mathbf{w}_p) = \|\mathbf{C}\mathbf{w}_p - \mathbf{x}_p\|_2^2.$$

- Here the $N \times N$ matrix \mathbf{C} is formed by stacking the spanning set vectors column-wise, and the $N \times 1$ vectors \mathbf{w}_p as

$$\mathbf{C} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{c}_1 & \mathbf{c}_2 & \cdots & \mathbf{c}_N \\ | & | & \cdots & | \end{bmatrix} \quad \mathbf{w}_p = \begin{bmatrix} w_{p,1} \\ w_{p,2} \\ \vdots \\ w_{p,N} \end{bmatrix}$$

- Formally then our optimization problem to recover all of the proper weight vectors can be expressed as

$$\underset{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_P}{\text{minimize}} \quad \frac{1}{P} \sum_{p=1}^P \|\mathbf{C}\mathbf{w}_p - \mathbf{x}_p\|_2^2.$$

Perfect representation using a fixed orthonormal spanning set

A very special kind of basis / spanning set one often sees when representing data points is one whose elements are perpendicular and unit length.

Such a spanning set $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N$ is one that spans an entire space, but whose elements are perpendicular have length one.

Algebraically this is written as

$$\left\{ \begin{array}{l} \mathbf{c}_n^T \mathbf{c}_m = 0 \iff \\ \|\mathbf{c}_n\|_2^2 = 1 \text{ for } n = 1 \dots N. \end{array} \right.$$

$$\mathbf{C}^T = \begin{bmatrix} \mathbf{c}_1 \\ \mathbf{c}_2 \\ \vdots \\ \mathbf{c}_N \end{bmatrix}$$

$$\mathbf{C} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_N]$$

$$\mathbf{C}^T \mathbf{C} = \begin{bmatrix} \mathbf{c}_1 \cdot \mathbf{c}_1 = 1 & & \\ & \mathbf{c}_1 \cdot \mathbf{c}_2 = 0 & \\ & & \mathbf{c}_2 \cdot \mathbf{c}_2 = 1 \end{bmatrix} = \mathbf{I}$$

$\mathbf{c}_1 \cdot \mathbf{c}_2 = 0 \quad \mathbf{c}_1 \cdot \mathbf{c}_3 = 0$

- Note that the same constraints expressed in terms of the concatenated basis matrix \mathbf{C} can be written compactly as

$$\boxed{\mathbf{C}^T \mathbf{C} = \mathbf{I}_{N \times N}}.$$

- Often such a spanning sets are likewise referred to as ***orthonormal*** or - sometimes - ***orthogonal***
- Technically speaking the latter term should only be applied when a set of vectors consists of perpendicular elements that are not necessarily unit-length, it is unfortunately often used to refer to unit-length perpendicular vectors as well.



Imperfectly representing data using a fixed spanning set

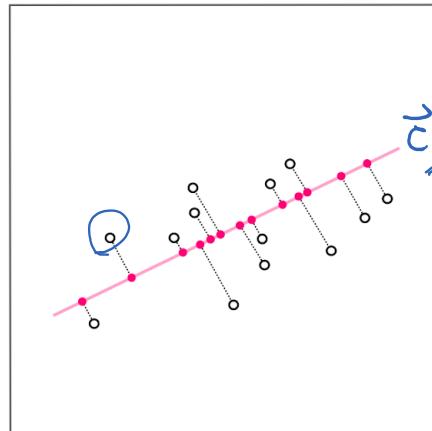
- Above we reviewed two important facts that we require for a spanning set / basis in order for it to be capable of perfectly representing points in a generic N dimensional space:
- The set of vectors are ***linearly independent***, that is they point in different directions in the space
- **There are at least N spanning vectors**

$$x_p \in \mathbb{R}^N \quad \overset{\text{basis}}{\vec{C}} = [\vec{c}_1 \dots \vec{c}_k]$$

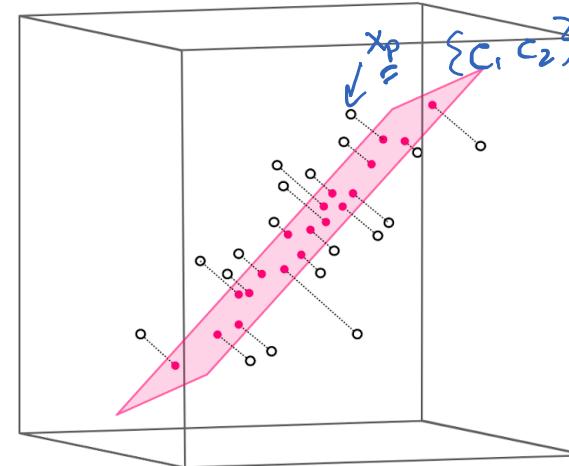
- What happens when we relax the second condition and suppose we more generally have $K \leq N$ spanning vectors / basis elements?
- Unsurprisingly we can no longer perfectly represent a generic set of P points in the space.
- In N dimensions this set of K vectors can - at best - span a **K dimensional subspace**.

- For example as illustrated in the figure below, if $N=3$ the best any set of $K=2$ spanning vectors can do is span a hyperplane.
- A set consisting of a single $K = 1$ spanning vector can only span a line.
- In both instances we clearly will not be able to perfectly represent all possible points in a space, since our reach in each case is restricted to a lower dimensional ***subspace*** of the full space.

$N=2$

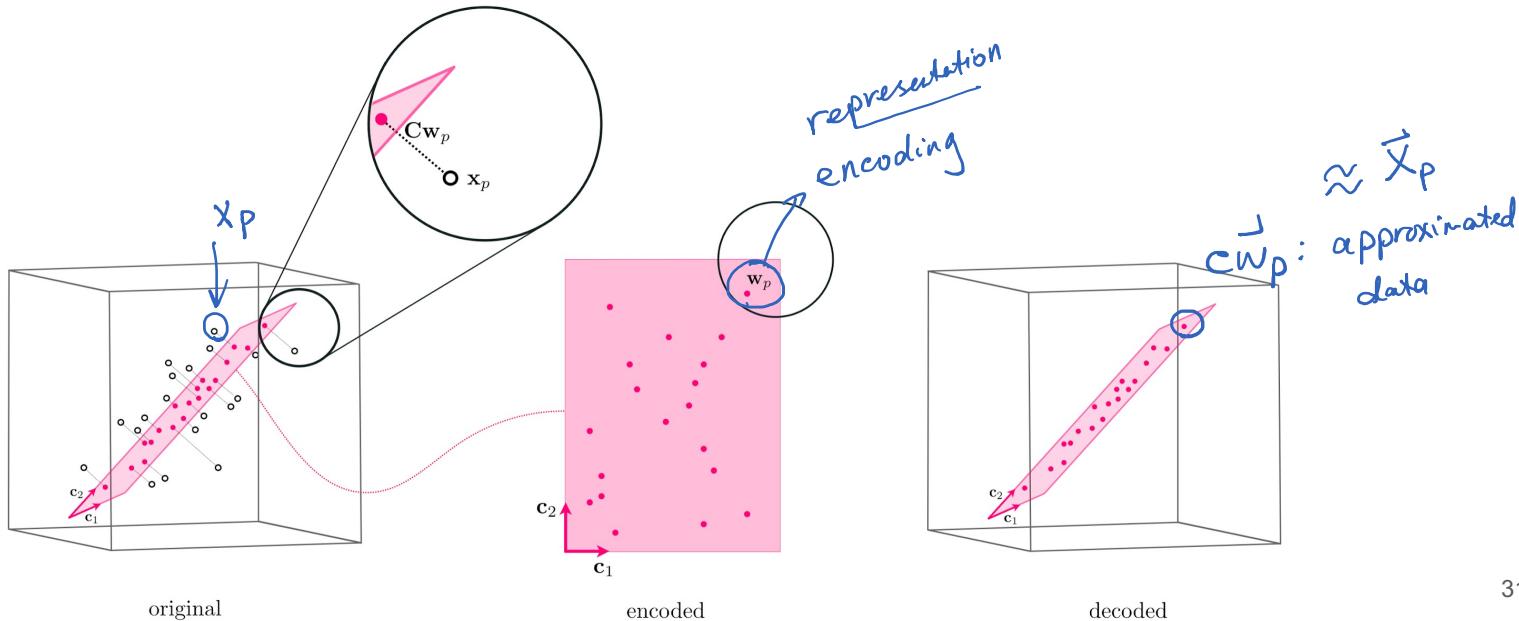


$N=3$



$$\vec{x}_p \approx C \vec{w}_p$$

- The representation over the subspace spanned by C - is Cw_p
- This is called the projection
- The weight vector w_p is the encoding of a point x_p in the (sub)space spanned by K spanning vectors $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_K$



$$C^T C = I_{N \times N}$$

Imperfect representation using a fixed orthonormal spanning set

$$C = [\vec{c}_1, \vec{c}_2, \dots, \vec{c}_N] \Rightarrow \vec{x}_p = C \cdot \vec{w}_p + \text{C orthonormal} \Rightarrow C^T \vec{x}_p = \underline{C^T C} \cdot \vec{w}_p$$

↓
encoding

basis

$$C^T \vec{x}_p = \underline{I} \cdot \vec{w}_p$$

- If our spanning set / basis of K elements is **orthonormal**, each weight vector encoding \mathbf{w}_p is:

$$\mathbf{w}_p = \mathbf{C}^T \mathbf{x}_p \quad p = 1 \dots P$$

- since we will have

$$\text{encoding } \mathbf{C}^T \mathbf{C} = \mathbf{I}_{K \times K}.$$

"decoding" $\mathbf{w}_p \rightarrow \mathbf{x}_p$: $\vec{x}_p = C \cdot \vec{w}_p$

if $C \in \mathbb{R}^{N \times N}$, orthonormal

$$\min_{C} g(C) = \| C^T x_p - x_p \|_2$$

PCA:

$C \in \mathbb{R}^{K \times N}$

Orthogonal

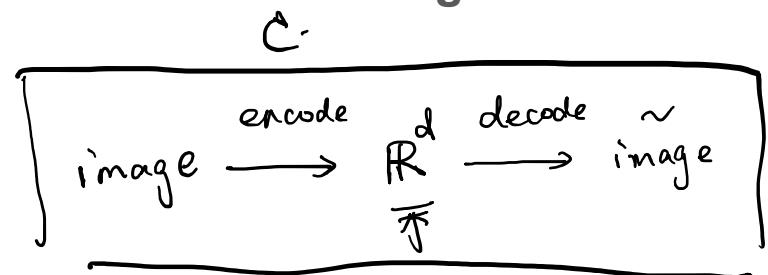
$\star C = [C_1, C_2, \dots, C_K]$

- Our representation is given simply in terms of the incomplete spanning set and datapoints as

$$C C^T x_p \approx x_p$$

$$p = 1 \dots P.$$

- The above formula is called **autoencoder** formula.
- In this set of lectures, our aim is to *learn the best basis for a given set of data* the assumption of orthogonality.



8.3 The Linear Autoencoder and Principal Component Analysis

- The most fundamental unsupervised learning method is known as **Principal Component Analysis** or **PCA** for short.

- Learn the proper weights w to best represent input data over a given fixed spanning set
- Learn a proper spanning set as well.

c

$$\left\{ \begin{array}{l} g(c, w) = \|c \cdot c^T x - x\|_2^2 \\ \text{s.t. } w = c^T x \\ c^T c = I \end{array} \right.$$

Learning proper spanning sets

- Given a fixed basis or spanning set of $\boxed{K} \leq N$ vectors $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_K$ - we can represent a set of P mean-centered points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_P$ as

$$\sum_{n=1}^K \mathbf{c}_n w_{n,p} \approx \mathbf{x}_p \quad p = 1 \dots P$$

- Here we tune the weights $w_{p,n}$ by minimizing a Least Squares cost function

$$g(\mathbf{w}_1, \dots, \mathbf{w}_P) = \frac{1}{P} \sum_{p=1}^P \|\mathbf{C}\mathbf{w}_p - \mathbf{x}_p\|_2^2.$$

- In general when we used a full set of *linearly independent* $K = N$ fixed basis vectors, we can learn corresponding weights that drive this cost function to zero, and give us strict equality in the desired equation.

- Here we will *learn* a proper basis (along with the weights).
- This slight thematic twist - where we *learn* the basis along with the proper weights - is called *Principal Component Analysis*.

$$g(\underline{\mathbf{w}_1, \dots, \mathbf{w}_P}, \mathbf{C}) = \frac{1}{P} \sum_{p=1}^P \|\mathbf{C}\underline{\mathbf{w}_p} - \underline{\mathbf{x}_p}\|_2^2.$$

- This cost function can be properly minimized using any number of standard approaches like e.g., gradient descent or coordinate / block-coordinate descent.

- if our K spanning vectors concatenated column-wise to form the spanning matrix \mathbf{C} are **orthonormal**, then the encoding of each \mathbf{x}_p may be written simply as $\boxed{\mathbf{w}_p = \mathbf{C}^T \mathbf{x}_p}$
- If we plug in this simple solution for \mathbf{w}_p into the pth summand of the Least Squares cost above, we get a cost that is a function of \mathbf{C}

$$g(\mathbf{C}) = \frac{1}{P} \sum_{p=1}^P \left\| \underbrace{\mathbf{C} \mathbf{C}^T \mathbf{x}_p}_{\text{decode}} - \mathbf{x}_p \right\|_2^2.$$

\mathbf{w}_p : encoding
 decode

It is often referred to as the **linear Autoencoder**

Principal Component Analysis

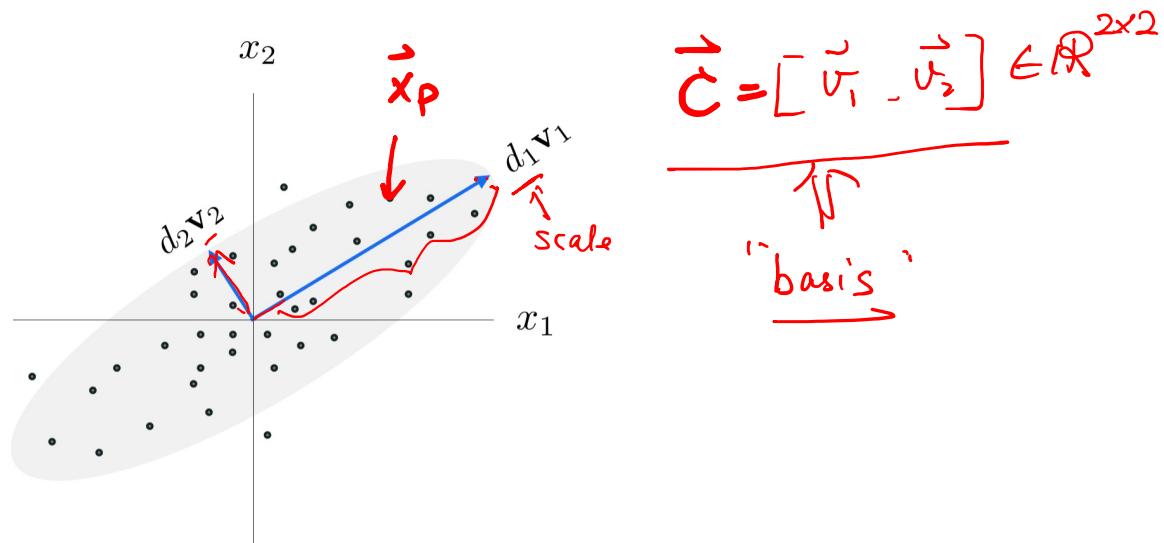
$$\rightarrow C = [\vec{c}_1 \ \vec{c}_2 \ \dots \ \vec{c}_k]$$

\vec{c}_k : principal component

- The linear Autoencoder cost may have many minimizers, of which the set of principal components is a particularly important one.
- The spanning set of principal components always provide a consistent skeleton for a dataset, with its members pointing in the dataset's largest directions of orthogonal variance.
- Employing this particular solution to the linear Autoencoder is often referred to as Principal Component Analysis, or PCA for short, in practice.

- A scaled version of the first principal component of this dataset points in the direction in which the dataset is most spread out, also called its largest direction of variance.
- A scaled version of the second principal component points in the next most important direction in which the dataset is spread out that is orthogonal to the first.

$N=2$



- This special orthonormal minimizer of the linear Autoencoder is given by the eigenvectors of the so-called covariance matrix of the data.
- Denoting by \mathbf{X} the $N \times P$ data matrix consisting of our P mean-centered input points stacked column-wise

"mean-centered"

$\mathbf{X} \in \mathbb{R}^{N \times P}$

$\mathbf{x}' \in \mathbb{R}^{N \times P'} \quad (P' > P)$

$\mathbf{C} \rightarrow$

$\mathbf{C} \uparrow$

$\mathbf{C} = \begin{bmatrix} \mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k \end{bmatrix}$ "eigenvector"

$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_P \end{bmatrix}$ 1st dimension

$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_P \end{bmatrix}$ $\in \mathbb{R}^{N \times P}$ dimension

$\Rightarrow \begin{cases} \mathbb{R}^{N \times N} & \text{①} \\ \mathbb{R}^{P \times P} & \text{②} \end{cases}$ #samples

$\mathbf{V} = \begin{bmatrix} \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_k \end{bmatrix}$ eigenvalue

$\frac{1}{P} \mathbf{X} \mathbf{X}^T = \mathbf{C} \cdot \mathbf{V} \mathbf{C}^T$ basis

"Symmetric"

44

- The orthogonal basis provided by this special solution (called the *Principal Components* of a dataset) can be computed (as a minimum of the Autoencoder cost function) as the *eigenvectors* of the corresponding correlation matrix of this data

$$\text{covariance matrix of } \mathbf{X} := \frac{1}{P} \mathbf{X} \mathbf{X}^T$$

- Denoting the eigenvector/value decomposition of the covariance matrix $\frac{1}{P} \mathbf{X} \mathbf{X}^T$ is given as

$$\frac{1}{P} \mathbf{X} \mathbf{X}^T = \mathbf{V} \mathbf{C} \mathbf{C}^T \mathbf{V}^T$$

eigenvector

- then above the orthonormal basis we recover is given precisely by the eigenvectors above, i.e., $\mathbf{C} = \mathbf{V}$
- Again, these are referred to in the jargon of machine learning as the *principal components* of the data.
- Moreover, the variance in each (principal component) direction is given precisely by the corresponding **eigenvalue** in \mathbf{D} .

Python Implementation of PCA

- First we center the data



```
def center(X):
    ...
    A function for normalizing each feaure dimension of an input array, mean-centering
    and division by its standard deviation
    ...
    X_means = np.mean(X, axis=1)[:,np.newaxis]
    X_normalized = X - X_means

    return X_normalized
```

- Next, we compute the principal components of the mean-centered data.

```
def compute_pcs(X, lam):
    """
    A function for computing the principal components of an input data matrix. Both
    principal components and variance parameters (eigenvectors and eigenvalues of XX^T)
    are returned
    """
    # create the correlation matrix
    P = float(X.shape[1])
    Cov = 1/P*np.dot(X,X.T) + lam*np.eye(X.shape[0])λ · I numerical stability
    .
    # use numpy function to compute eigenvalues / vectors of correlation matrix
    D,V = np.linalg.eigh(Cov)eigen value C
    return D,V
```

- *Put it together*

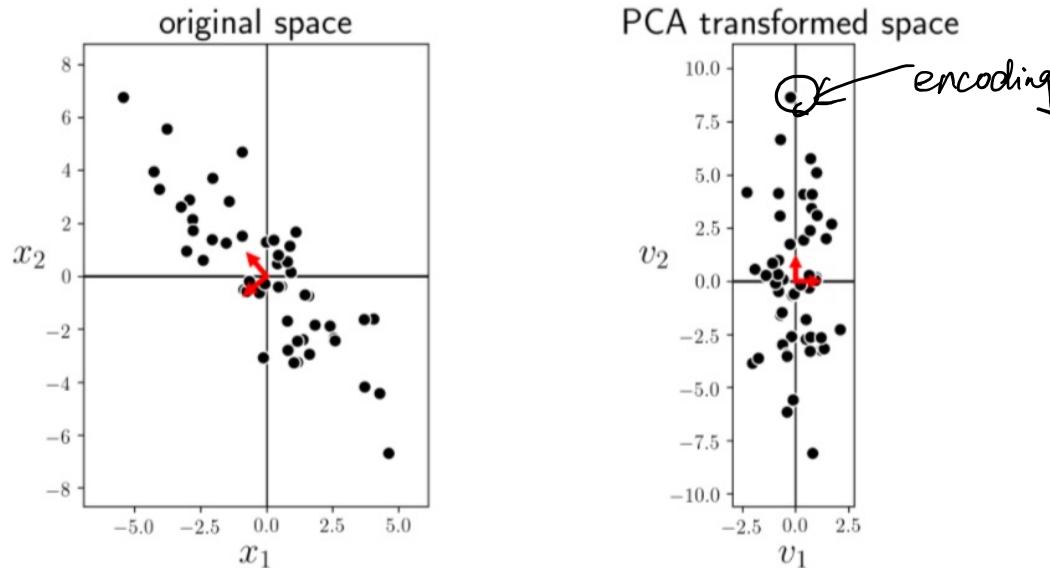
```
def pca_transform_data(X,**kwargs):
    """
    A function for producing the full PCA transformation on an input dataset X.
    """

    # user-determined number of principal components to keep, and regularizer penalty parameter
    num_components = X.shape[0]
    if 'num_components' in kwargs:
        num_components = kwargs['num_components']
    lam = 10**(-7)
    if 'lam' in kwargs:
        lam = kwargs['lam']

    # compute principal components
    D,V = compute_pcs(X, lam)
    V = V[:, -num_components:]
    D = D[-num_components:]

    # compute transformed data for PC space: V^T X
    W = np.dot(V.T,X)
    return W,V
```

- In the left panel we show the mean-centered data along with its two principal components (pointing in the two orthogonal directions of greatest variance in the dataset) shown as red arrows.
- In the right panel we show the encoded version of the data in a space where the principal components are in line with the coordinate axes.



$$\vec{x}_p \in \mathbb{R}^N \quad P: \text{sampler} \Rightarrow X \in \mathbb{R}^{N \times p}$$

Usage: ① dimension reduction $x_p \in \mathbb{R}^N \Rightarrow w_p \in \mathbb{R}^k$
 $k < N$

Recap

basis/a spanning set of vectors: $\{\vec{c}_n\}_{n=1 \dots N} \in \mathbb{R}^N$ $\Rightarrow C \in \mathbb{R}^{N \times N}$

② $k=2, 3$ ③ fill in missing values
 linearly independent } \Rightarrow span a space
 N dimension

encoding: $\vec{x}_p = C \cdot \vec{w}_p$ \vec{w}_p : representation
 feature space $\mathbb{R}^{N \times N}$ encoding transformed feature space \mathbb{R}^N

Orthonormal: $\vec{c}_i \cdot \vec{c}_j = 0$ if $i \neq j$, $i, j \in \{1 \dots N\} \Rightarrow C^T C = I_{N \times N}$

decoding (C is orthonormal): $\vec{w}_p = C^T \vec{x}_p$ $\vec{x}_p = C \cdot \vec{w}_p \Rightarrow C^T \vec{x}_p = C^T C \vec{w}_p = \vec{w}_p$
 C^T C^T $I_{N \times N}$

Autoencoder: $C(C^T \vec{x}_p) = \vec{x}_p$
 $| C \in \mathbb{R}^{N \times k}, k \leq N \Rightarrow C(C^T \vec{x}_p) \approx \vec{x}_p$

PCA: $C^*, \{\vec{w}_p^*\} = \arg \min_p \| C \vec{w}_p - \vec{x}_p \|_2^2$. S.t. $C^T C = I_{k \times k}$

Numerical Solution
 $\frac{1}{p} X X^T = V D V^T$
 $C^* = V$ eigenvector
 D : eigenvalues

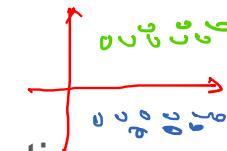
True/False Statement

Given the above context, give me 10 True/False statements.

Report whether the statement is True or False in the following format: statement - [True/False].

Explain the reason why a statement is correct or wrong .

- ✗ PCA, as an unsupervised learning technique, can be used to reduce the dimension of a dataset by reducing the number of inputs. #P N: Feature dimension
- ✓ Spanning vectors that match the data dimension and are linearly independent can perfectly represent all data points.
- ✗ An orthonormal basis consists of vectors that have varying lengths and are orthogonal to each other.
- ✗ The principal components point in random directions in a dataset.
- ✗ Using PCA for dimension reduction always results in better classification results
- ✓ Principal components are derived from the eigendecomposition of a dataset's covariance matrix.



$C^* \in \mathbb{R}^{N \times K}$: principal components = eigen vectors

$$\frac{1}{p} X X^T = V D C^T$$

basis = { }_{1...k}

$$D = \begin{bmatrix} 2 & & \\ & 1.2 & \\ & & 0.8 \end{bmatrix}$$

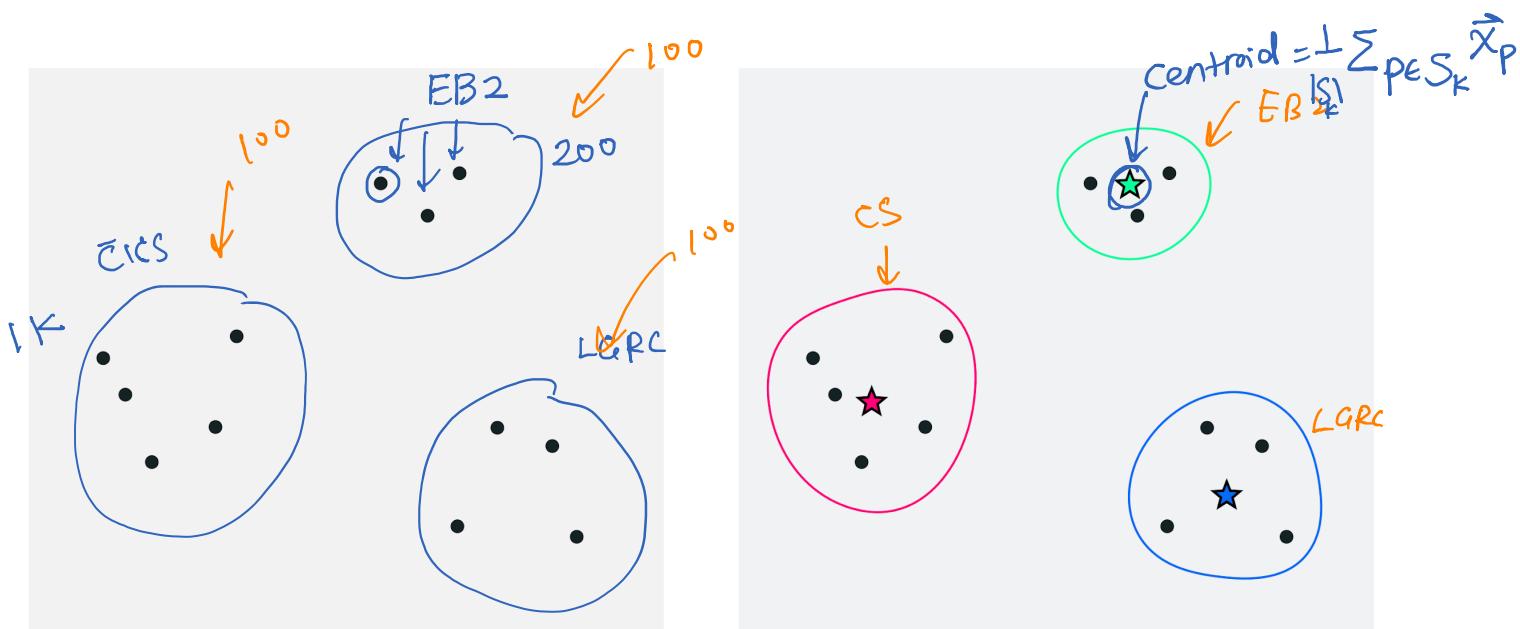
close to zero

- PCA, as an unsupervised learning technique, can be used to reduce the dimension of a dataset by reducing the number of inputs. – [FALSE]
 - Explanation: Should be the number of features instead of inputs.
- Spanning vectors that match the data dimension and are linearly independent can perfectly represent all data points. – [TRUE]
- An orthonormal basis consists of vectors that have varying lengths and are orthogonal to each other. – [FALSE]
 - **Correct Statement:** An orthonormal basis consists of vectors that have unit length and are orthogonal to each other.
- The principal components point in random directions in a dataset. - [False]
 - **Correct Statement:** The principal components point in the dataset's largest directions of orthogonal variance.

- Using PCA for dimension reduction always results in better classification results. - [False]
 - Correct Statement: PCA can potentially destroy inter-class separation, so it doesn't always guarantee better classification results.
- Principal components are derived from the eigendecomposition of a dataset's covariance matrix. - [True]

8.5 K-means Clustering

- The focus of this Section - the *K-means algorithm* - is an elementary example of another set of unsupervised learning methods called *clustering algorithms*.
- These algorithms are designed to (properly) reduce the number of points in a dataset, which we refer to as the *data dimension* of a dataset...
- ...and in doing so help us understand the structure of our data.



(left) A 2-dimensional toy dataset with $P=10$ data points.

(right) The data shown naturally clustered into $K=3$ clusters. Each **cluster center** - also called a **centroid** - is marked by a star symbol colored to match its cluster boundary.

$$\mathcal{S}_k^* = \{ p \mid \text{if } \mathbf{x}_p \text{ belongs to the } \underline{k}^{\text{th}} \text{ cluster} \} .$$

- Algebraically, each centroid is represented as:

$$\underline{\mathbf{c}_k^*} = \frac{1}{|\mathcal{S}_k|} \sum_{p \in \mathcal{S}_k} \underline{\mathbf{x}_p}.$$

- This formula confirms the intuition that each centroid represents a chunk of the data - **the average of those points belonging to each cluster**.
- Next we can state mathematically an obvious and implicit fact about the clustering scenario above: that each point belongs to the cluster whose centroid it is closest to.

- To express this algebraically for a given point \mathbf{x}_p is simply say that the point must belong to the cluster where the distance to the centroid $\|\mathbf{x}_p - \mathbf{c}_k\|_2$ is minimal.
- In other words, the point \mathbf{x}_p belongs to or is *assigned* to cluster k^* if

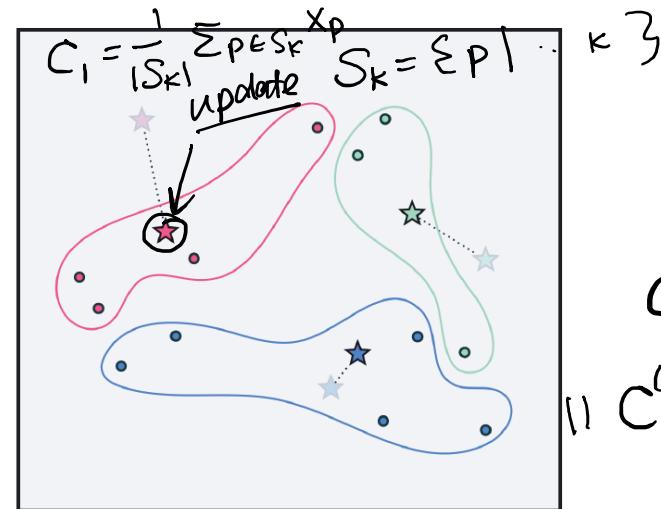
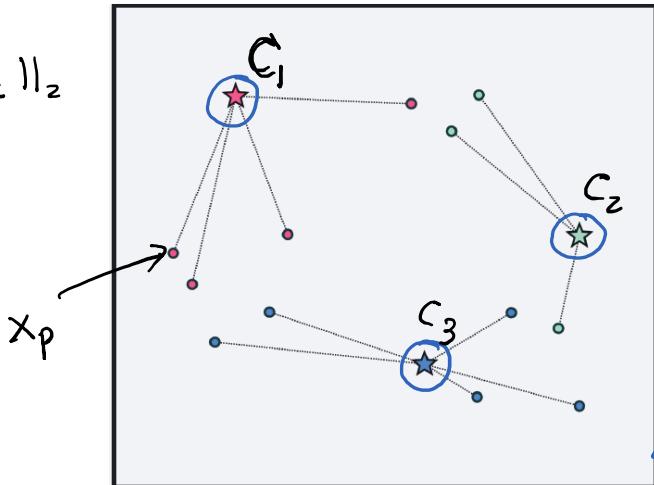
$$a_p = \operatorname{argmin}_{k=1,\dots,K} \|\mathbf{x}_p - \mathbf{c}_k\|_2$$

- In the jargon of machine learning these are called *cluster assignments*.

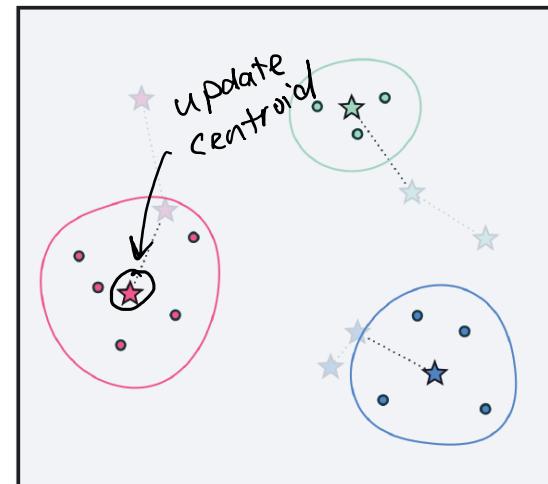
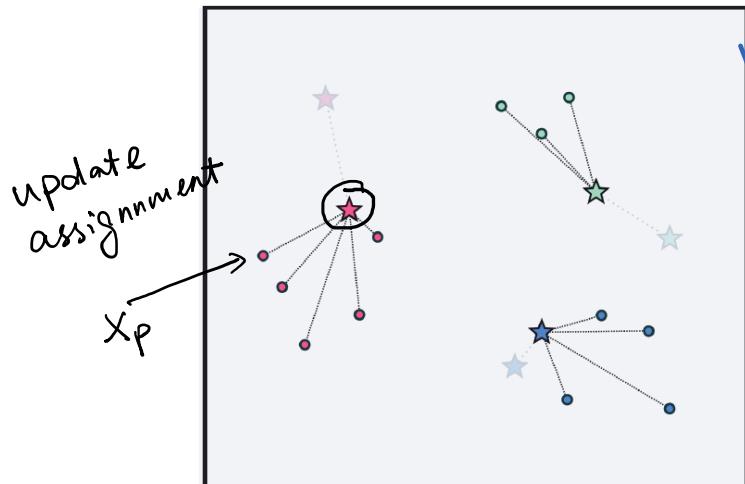
- We do not want to have to visually identify clusters in a dataset ourselves...
- ...and in any being constrained three dimensions if the dimension of of a dataset is greater than $N=3$ there is no way for us humans to do this anyway.
- Instead we want an algorithm that will do this for us automatically.

- Thankfully we can do this rather easily using the framework detailed above for mathematically describing clusters...
- ...the resulting algorithm being called the ***K-means clustering algorithm***.
- As with many of the algorithms we have seen K-means is an ***iterative method***, meaning that we will refine the ideal location for our cluster centroids / cluster assignments over a number of update steps.

$$\min_k \|x_p - c_k\|_2$$



① $\|c^{(t+1)} - c^{(t)}\| < \text{abs}$



② $t < 100$

The K-means algorithm

1: **input:** dataset $\mathbf{x}_1, \dots, \mathbf{x}_P$, initializations for centroids $\mathbf{c}_1, \dots, \mathbf{c}_K$, and maximum number of iterations J

2: **for** $j = 1, \dots, J$

3: **# Update cluster assignments**

4: **for** $p = 1, \dots, P$

5: $a_p = \operatorname{argmin}_{k=1, \dots, K} \|\mathbf{c}_k - \mathbf{x}_p\|_2$

6: **end for**

7: **# Update centroid locations**

8: **for** $k = 1, \dots, K$

9: denote S_k the index set of points \mathbf{x}_p currently assigned to the k^{th} cluster

10: update \mathbf{c}_k via
$$\mathbf{c}_k = \frac{1}{|S_k|} \sum_{p \in S_k} \mathbf{x}_p$$

11: **end for**

12: **end for**

13: **# Update cluster assignments using final centroids**

14: **for** $p = 1, \dots, P$

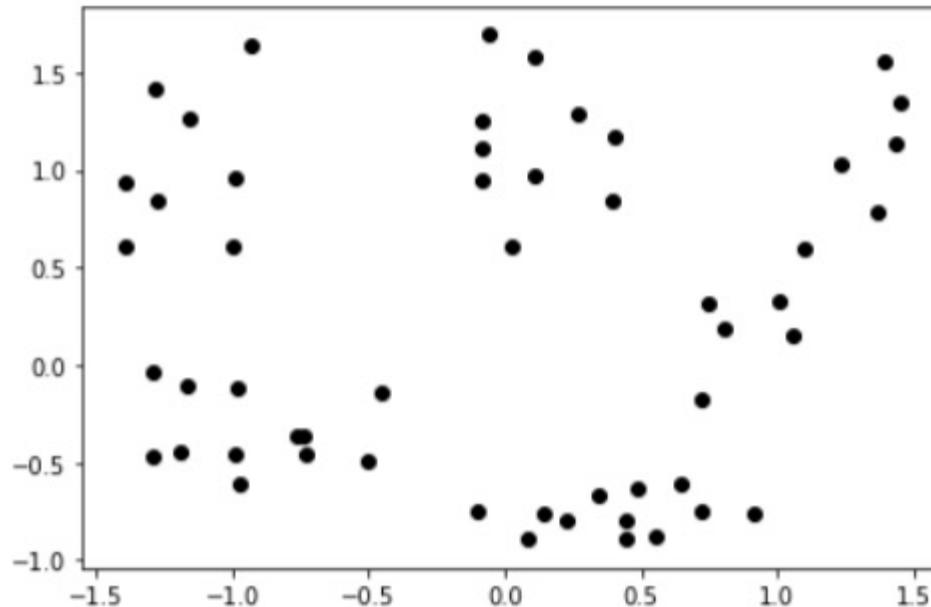
15: $a_p = \operatorname{argmin}_{k=1, \dots, K} \|\mathbf{c}_k - \mathbf{x}_p\|_2$

16: **end for**

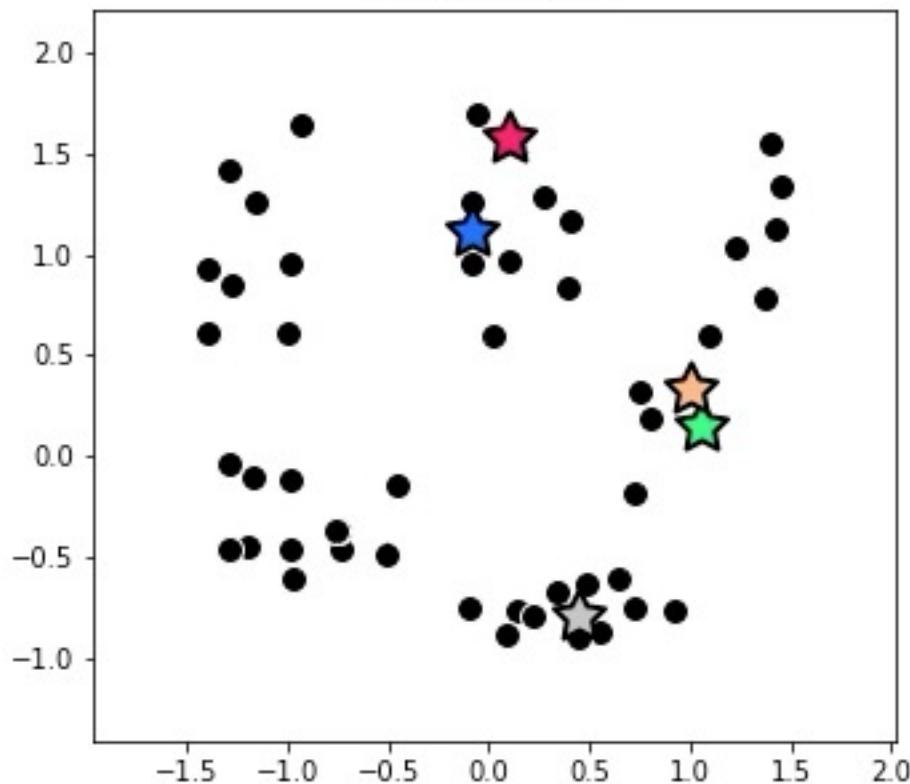
17: **output:** optimal centroids and assignments

} optional

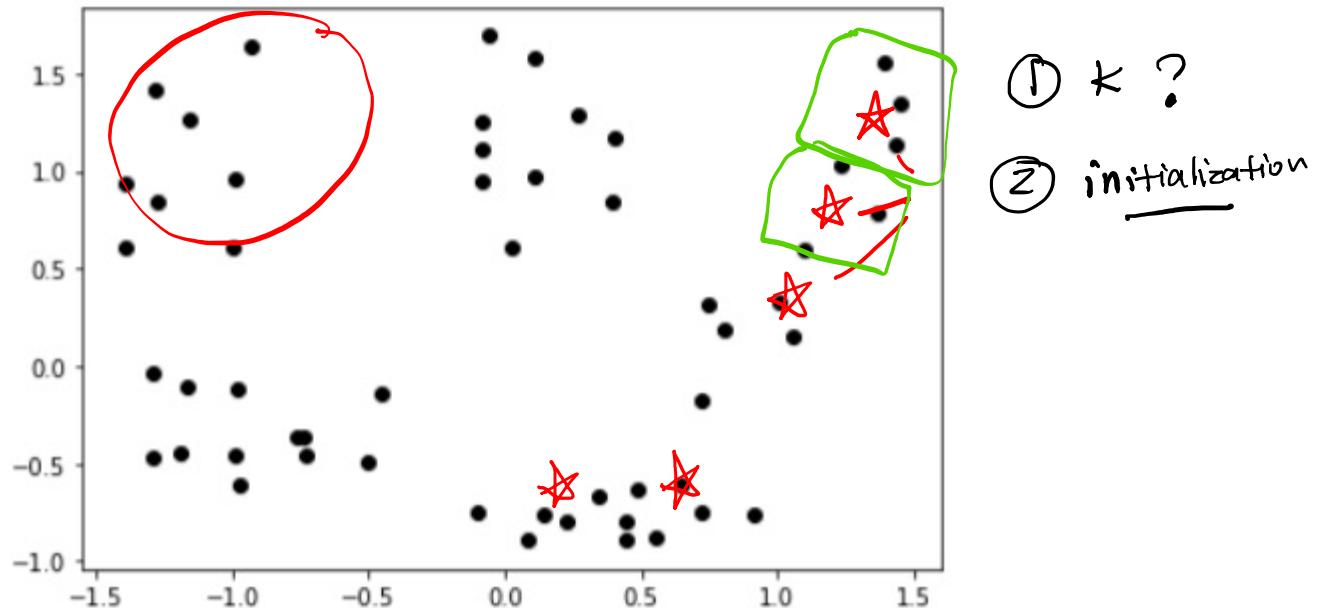
Q: How many clusters does the dataset have?



iteration 1

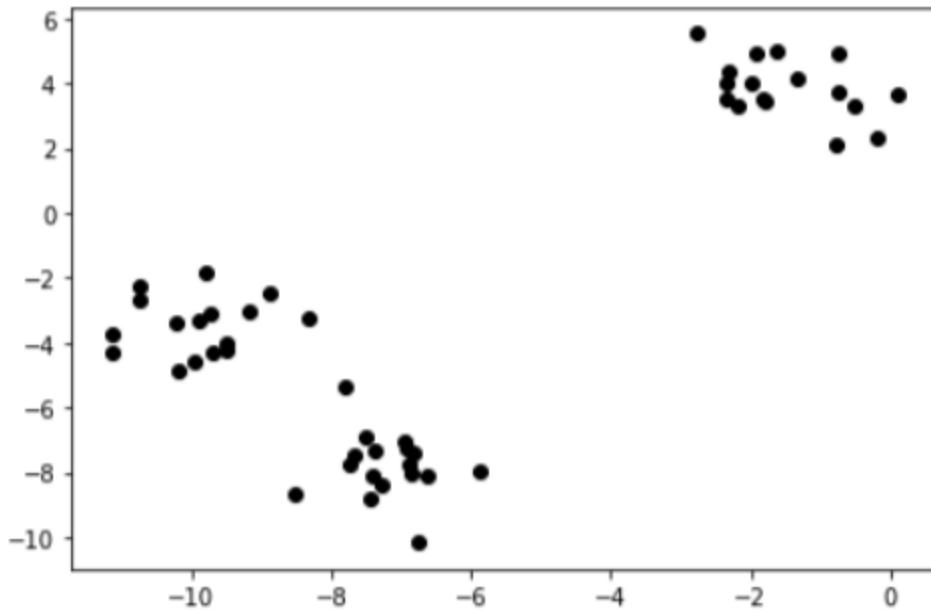


Q: What are the potential problems with k-means clustering?



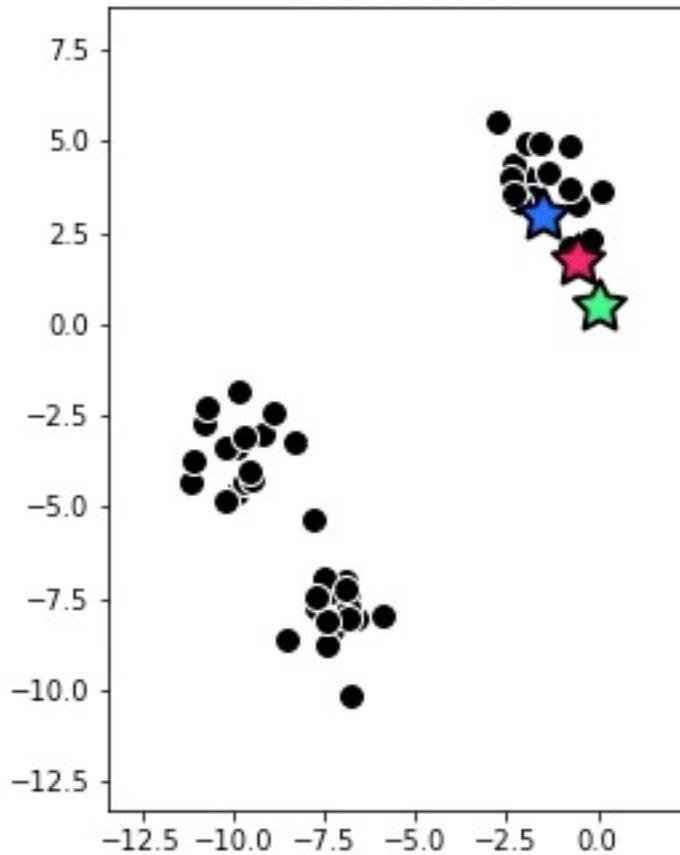
Example: K-means failures due to a poor choice of initialization

- Generally speaking, a good initialization is one where the initial centroids are spread evenly throughout the distribution of the data.
- Conversely, a poor initialization is typically one where all of the initial centroids are bunched up together in a small region of the space.
- This can make it difficult for the centroids to spread out effectively, leading to less than optimal clusterings or even 'empty clusters' (i.e., those to which no points are finally assigned).



- The first and most fundamental issue is that of ***empty clusters***, meaning clusters with no points assigned to them.
- It is indeed possible, for some unfortunate initialization choices, for clusters to end up being empty when running K-means.
- This is not an overly common occurrence, but is still worth noting.

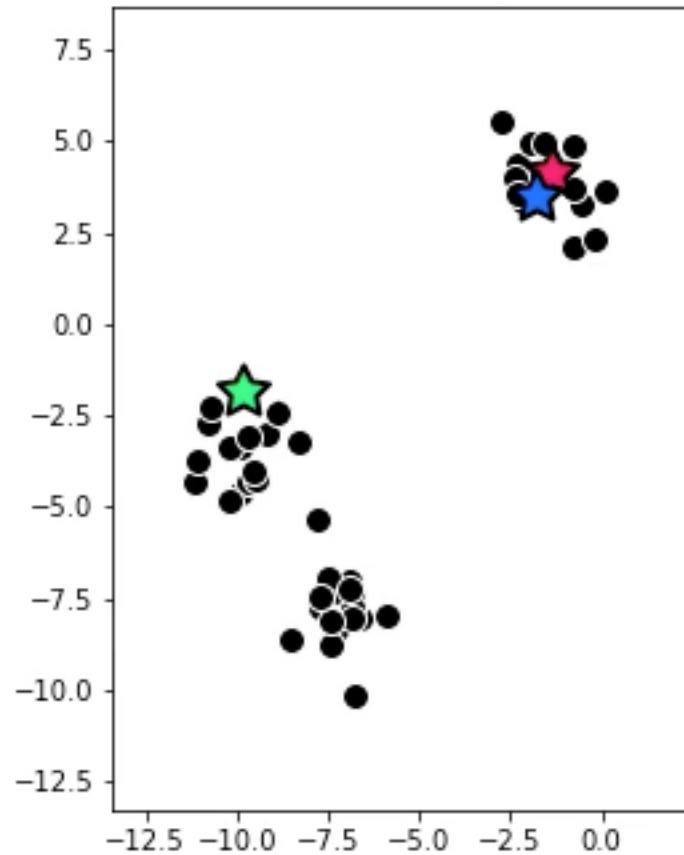
iteration 1



- Notice here how - due to the centroid initializations - the blue cluster very quickly becomes empty (every point is closer to one of the other two centroids) and never recovers any assigned points.
- This problem can be easily detected in practice (simply count up the number of assignments at the end of a run of K-means), and **re-running with a different initialization is typically the anecdote**.

- The second issue to keep note of with K-means is the *sub-optimal clustering*.
- This means that although the proper number of clusters K was chosen the final clustering itself fails to fully capture the cluster-structure of the data.
- The animation below illustrates this failure on the same dataset used above.
- Here we can see that - due to the initialization - the single cluster in the top of the image is cleaved in two, while the two lower clusters are identified a single massive cluster.

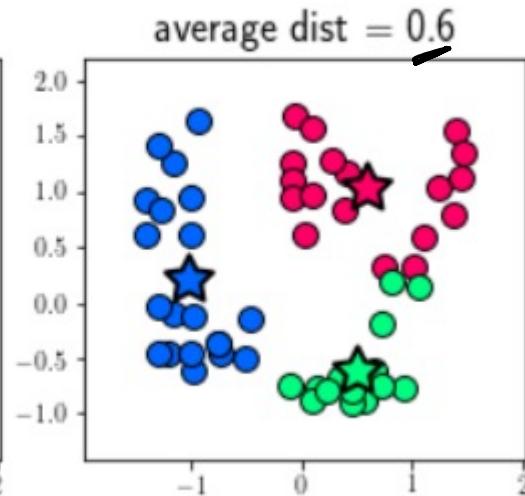
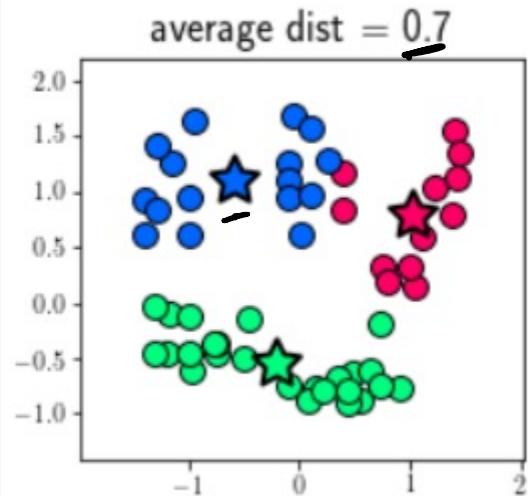
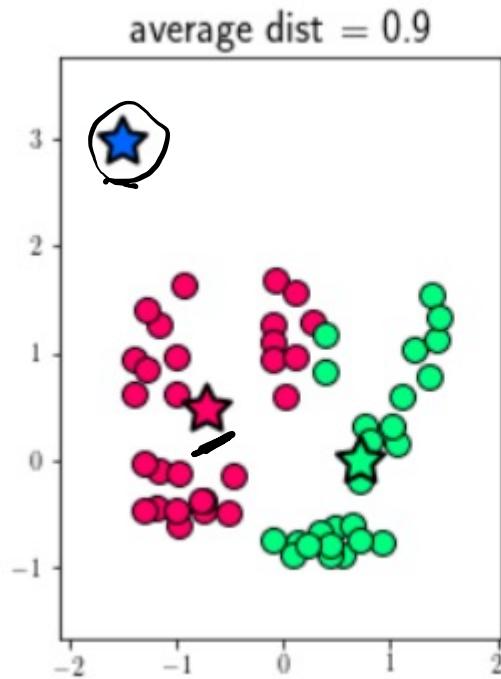
iteration 1



- Once again the best antidote for this issue is to simply re-run the algorithm with a different initialization.
- To determine the best clustering from set of runs we can use the *average distance of each point to its cluster centroid* - called the average intra-cluster distance - as an objective measure to rank the runs.
- Denoting \mathbf{c}_{k_p} the final cluster centroid of the p^{th} point \mathbf{x}_p , then the average distance from each point to its respective centroid can be written as

$$\text{average intra-cluster distance} = \underbrace{\frac{1}{P} \sum_{p=1}^P \| \mathbf{x}_p - \mathbf{c}_{k_p} \|_2 }_{\text{Diagram: A horizontal line segment with arrows at both ends, indicating summation over } p \text{ from 1 to } P.}$$

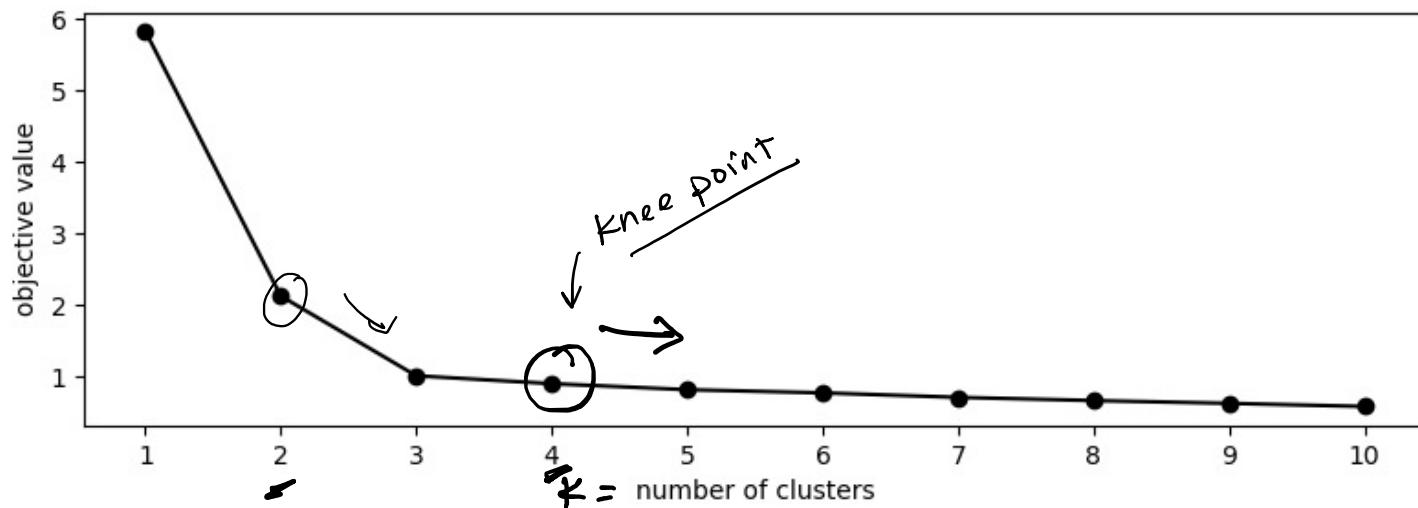
- Computing this for each run of K-means we choose the final clustering that achieves the ***smallest*** such value as the best clustering arrangement.



Example: Choosing the ideal number of clusters K

- To determine the optimal setting of the parameter K - i.e., the number of clusters in which to cluster the data - **we typically must try a range of different values for K**, run the K-means algorithm in each case, and compare the results using the average intra-cluster distance.

Scree plot



- As one should expect, the intra-cluster distance decreases **monotonically** as we increase K.
- Notice, however that the scree plot above has an ***elbow*** at K=3, meaning that increasing the number of clusters from 3 to 4 and onwards reduces the objective value by very little.
- Because of this **we can argue that K=3 is a good choice for the number of clusters for this particular dataset** (as we saw in the prior example, it is indeed the optimal number) since any fewer clusters and the intra-cluster distance is comparatively large, while adding additional clusters does not decrease the total intra-cluster distance too much.

Fun Stuff

Table 8.1 Common matrix factorization problems $CW \approx X$ subject to possible constraints on C and W .

Matrix factorization problem	Constraints on C and W
Principal Component Analysis	C is orthonormal $\min_{C, W} \ CW - X\ $
Recommender systems	No constraint on C or W ; X is only partially known
K-means clustering	Each column of W is a standard basis vector $S_k \Rightarrow W = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$
Sparse dictionary learning	Each column of W is sparse $W = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ $C = \begin{bmatrix} & & & & \end{bmatrix}$ $\in N \times K$ $(K > N)$
Nonnegative matrix factorization	Both C and W are nonnegative $\uparrow \quad \uparrow$

Summary

- Background: Fixed Spanning Sets, Orthonormality, and Projections
- The Linear Autoencoder and Principal Component Analysis
 - Reduce the feature dimension N to the number of principle components
- K-means Clustering
 - Reduce the data dimension P to a set of centroids
 - Sensitive to initialization
 - Sensitive to the choice of K