

## Cognitive Algorithms Lecture 5

# Unsupervised Learning

Klaus-Robert Müller, **Ali Hashemi, Lorenz Vaitl,**  
Augustin Krause, Joanina Oltersdorff, Ken Schreiber

Berlin Institute of Technology  
Dept. Machine Learning

Recap  
oo

Unsupervised Learning  
ooooooo

PCA  
oooooooooooo

PCA Applications  
oooooooooooo  
ooo

NMF  
oooooooooooo

Clustering  
oooooooooooo

Summary  
oo

# Today's agenda

Recap

Unsupervised Learning

PCA

PCA Applications

NMF

Clustering

Summary

## Recap: Kernels

What is a kernel function?

$$\rightarrow k(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)^T$$

## Why are kernels useful?

- Implicitly work in high-dimensional space
  - Thanks to representer theorem, work on space spanned by data

How can you show that a function  $k(\cdot, \cdot)$  is a kernel function?

- Find  $\phi(\cdot)$  s.t.  $k(x_i, x_i) = \phi(x_i)^T \phi(x_i)$  or show that  $k(\cdot, \cdot)$  is PSD

## Recap: Covariance Matrices

Given  $n$  data points  $\mathbf{x}_i \in \mathbb{R}^d$  in a data matrix

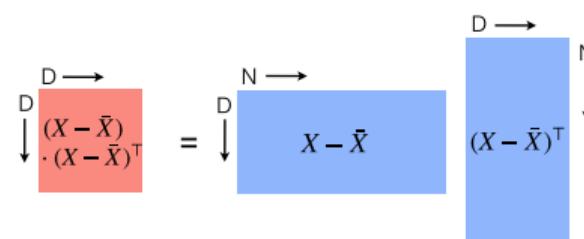
$$X \in \mathbb{R}^{d \times n}$$

the empirical estimate of the **covariance matrix** is defined as

$$\bar{\Sigma} = \frac{1}{n} (X - X)(X - \bar{X})^\top$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

$$\bar{X} = (\bar{x}, \bar{x}, \dots, \bar{x}) \in \mathbb{R}^{d \times n}$$



# Unsupervised learning

## Supervised algorithms

**Classification** and **regression**

Labels for training

But labels not always given e.g.

Mixtures of different speakers in an audio

recording

Complex artefacts in experimental recordings

$y$	
(Binary)	$\in \{+1, -1\}$
Classification	
Regression	$\in \mathbb{R}$
Unsupervised	$\in \emptyset$ <i>don't exist</i>

## Unsupervised algorithms

No labels for training

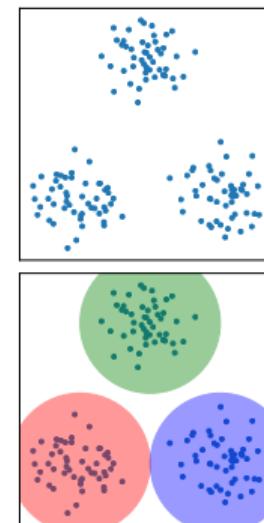
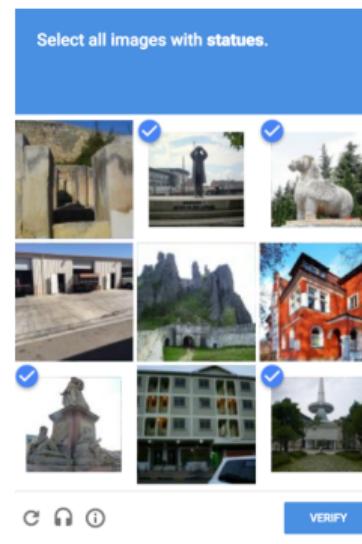
# Unsupervised learning

## Why unsupervised learning

- No need for human made labels
- Algorithm finds structure autonomously

## What can it be used for?

- Clustering
- Dimensionality Reduction
- ⋮

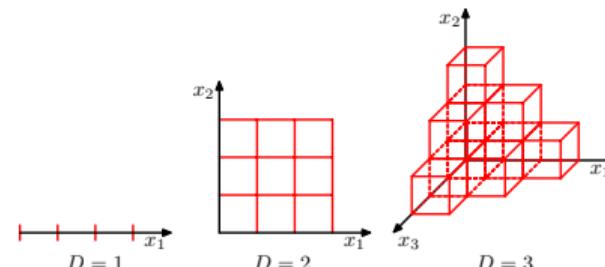


# Dimensionality Reduction

In many applications, we have

- high-dimensional data
  - reason to believe they lie close to a lower dimensional subspace
- Fewer parameters needed to account for the data properties  
*hidden causes or latent variables*

## Curse of Dimensionality



[Bishop, 2007]

解释数据属性、隐藏原因或潜在变量所需的参数更少

# Why Dimensionality Reduction

- **Visualization:**

Insights into high-dimensional structures in the data

- **Better Generalization:**

Fewer dimensions → smaller chance of overfitting

- **Speeding up** learning algorithms:

Most algorithms scale badly with increasing data dimensionality

- **Data compression:**

Less storage requirements

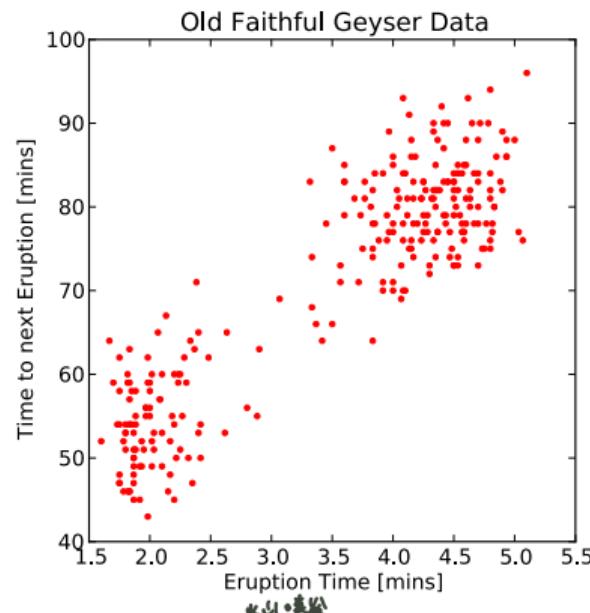


[Xiao et al., 2017]

## Example: Old Faithful Geyser Dataset

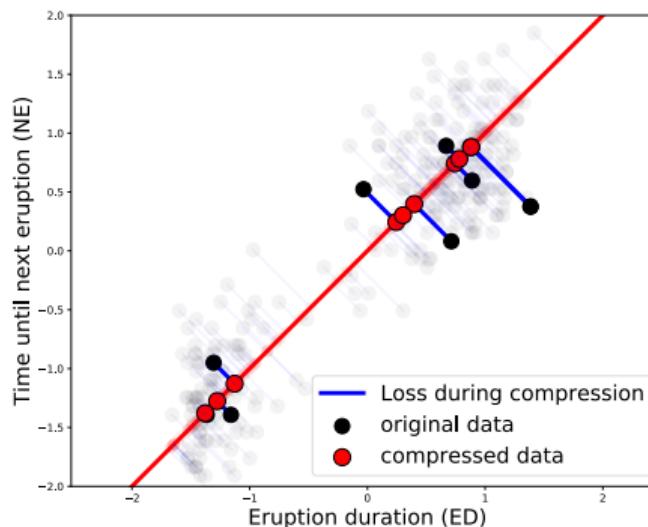


Old Faithful Geyser



## Informal example dimensionality reduction

- High correlation between  $NE$  and  $ED$
- Let's try to project data on  $1D$  subspace
- Relatively good representation



# The mathematical model for linear dimensionality reduction

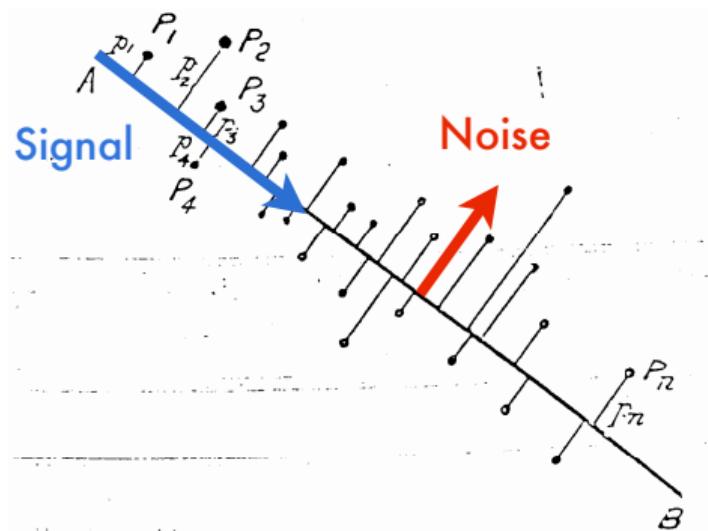
We have

- high-dimensional data  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$
- reason to believe they lie close to a lower dimensional subspace  
→  $m < d$  parameters needed to account for the data properties  
*hidden causes* or *latent variables*  $H = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n] \in \mathbb{R}^{m \times n}$

**Goal:** Find  $m < d$  hidden causes  $H \in \mathbb{R}^{m \times n}$ , that explain the observed data via a mixing  $W \in \mathbb{R}^{d \times m}$ :

$$X \approx WH$$

# Principal Component Analysis (PCA)



Adapted from Pearson [1901]

has the following two equivalent objectives:  
find the direction  $w$  that minimizes the

noise and maximizes the signal

or equivalently

find the line  $w$  that maximizes the variance  
within the data set

## Maximizing variance

We obtained some data  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$

PCA finds a direction  $\mathbf{w} \in \mathbb{R}^d$  such that the estimated variance of the projected data  $\mathbf{w}^\top X$  is maximal

$$\begin{aligned}\overline{\text{Var}}(\mathbf{w}^\top X) &= \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - \frac{1}{n} \sum_{j=1}^n \mathbf{w}^\top \mathbf{x}_j)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - \mathbf{w}^\top \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j)^2 = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}))^2 \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{w}^\top (\mathbf{x}_i - \bar{\mathbf{x}}) \cdot (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{w} \\ &= \mathbf{w}^\top \underbrace{\left( \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) \cdot (\mathbf{x}_i - \bar{\mathbf{x}})^\top \right)}_{\text{Empirical covariance matrix } \bar{\Sigma}} \mathbf{w}\end{aligned}$$

## Maximizing variance

We obtained some data  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$

PCA finds a direction  $\mathbf{w} \in \mathbb{R}^d$  such that the variance of the projected data  $\mathbf{w}^\top X$  is maximal  
Let's assume centered data for easier notation <sup>1</sup>

$$\begin{aligned}\overline{\text{Var}}(\mathbf{w}^\top X) &= \mathbf{w}^\top \underbrace{\left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{x}_i^\top \right)}_{\text{Empirical Covariance matrix } \bar{\Sigma}} \mathbf{w} \\ &\propto \mathbf{w}^\top \underbrace{\mathbf{X}\mathbf{X}^\top}_{\text{Scatter matrix } S} \mathbf{w}\end{aligned}$$

And: we need to constrain  $\mathbf{w}$

---

<sup>1</sup>i.e. we assume  $\bar{\mathbf{x}} = 0$

## Maximizing variance

PCA finds a direction  $\mathbf{w} \in \mathbb{R}^d$  such that the variance of the projected data  $\mathbf{w}^T \mathbf{X}$  is maximal

$$\operatorname{argmax}_{\mathbf{w}} \frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}}$$

This objective function is independent of the scaling of  $\mathbf{w}$ .

Note the similarity to the objective of Linear Discriminant Analysis!

→ Different covariance matrices, different problem, but: same maths solve it

# Maximizing variance

$$\operatorname{argmax}_{\mathbf{w}} \frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}} \quad (1)$$

Set the derivative w.r.t  $\mathbf{w}$  to zero:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}} &= \frac{(\mathbf{w}^T \mathbf{w}) 2 S \mathbf{w} - (\mathbf{w}^T S \mathbf{w}) 2 \mathbf{w}}{(\mathbf{w}^T \mathbf{w})^2} = 0 \rightarrow \text{Eq. 1 can be reduced to the standard eigenvalue problem.} \\ \Rightarrow \underbrace{(\mathbf{w}^T S \mathbf{w})}_{\text{scalar}} \mathbf{w} &= \underbrace{(\mathbf{w}^T \mathbf{w})}_{\text{scalar}} S \mathbf{w} \qquad S \mathbf{w} = \lambda \mathbf{w} \\ \Rightarrow \underbrace{\frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}}}_{\equiv \lambda} \mathbf{w} &= S \mathbf{w} \end{aligned}$$

## Maximizing variance

Setting  $S\mathbf{w} = \lambda\mathbf{w}$  in Eq. 1, we see that the variance in direction  $\mathbf{w}$  is given by:

$$\operatorname{argmax}_{\mathbf{w}} \frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}} = \frac{\mathbf{w}^T \lambda \mathbf{w}}{\mathbf{w}^T \mathbf{w}} = \lambda$$

The variance of the projected data in an eigendirection  $\mathbf{w}$  is given by the corresponding eigenvalue!

The direction of maximal variance in the data is equal to the eigenvector having the largest eigenvalue.

## PCA vs. LDA

Which is which  
and why?

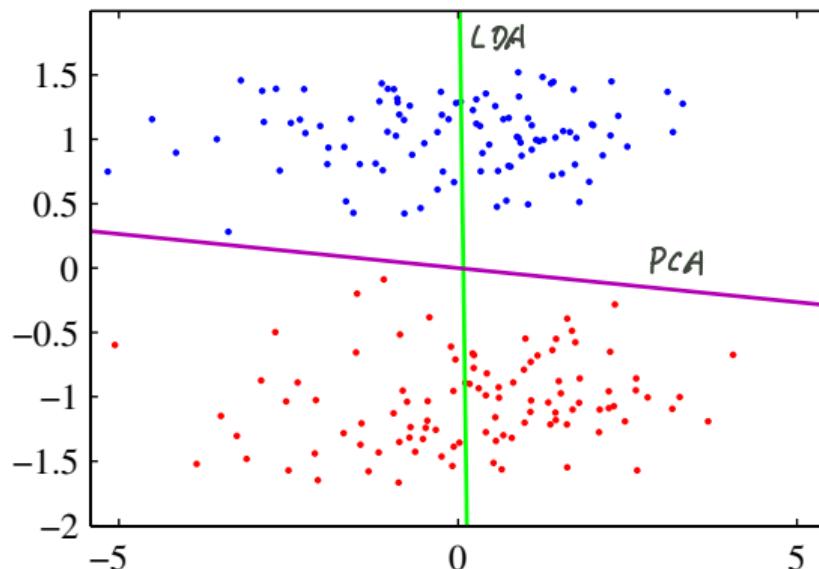


Figure: Directions found by PCA (magenta) and LDA(green) Bishop [2007]

## Finding more Principal Components

Incremental PCA then finds additional principal components, by looking at directions **orthogonal** to previous ones, that maximize variance. It can be shown that:

The  $k$  first PCA basis vectors are the eigenvectors corresponding to the largest  $k$  eigenvalues

$$SW = W\Lambda$$

where  $W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k]$  contains the eigenvectors sorted according to their eigenvalues and  $\Lambda$  is a diagonal matrix containing the corresponding eigenvalues.

Recall theory of eigenvectors:

Since  $S$  is symmetric, there exist  $d$  orthogonal eigenvectors!

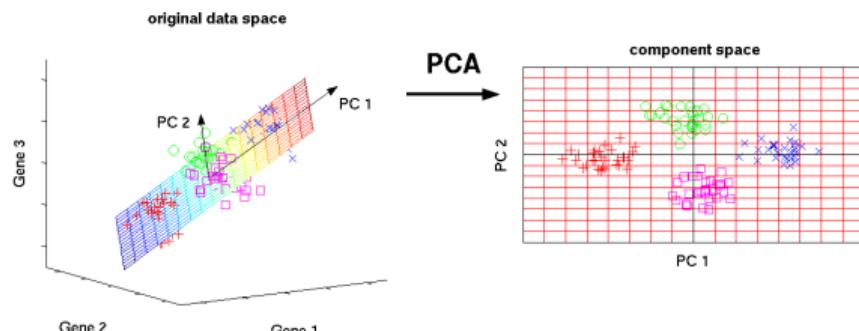
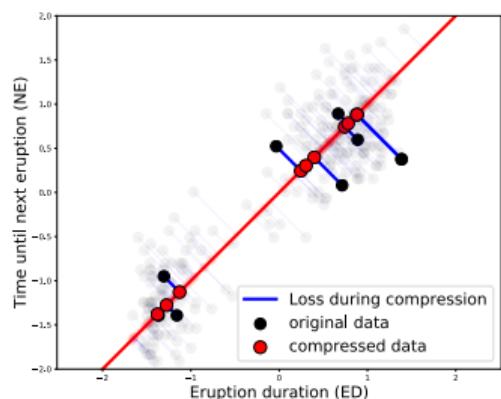
$\rightarrow \mathbf{w}_i \perp \mathbf{w}_j$

# Encoding (from real data to latent representation)

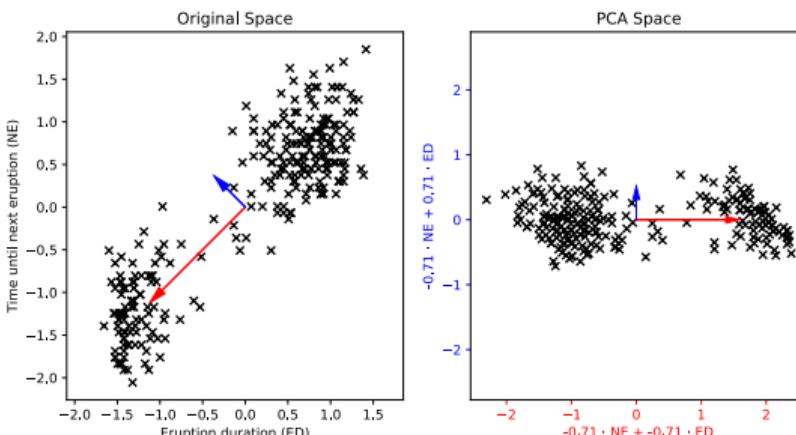
Now that we have

$W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$ , we project each data point  $\mathbf{x}$  onto  $W$

$$h_i = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x}_i \\ \vdots \\ \mathbf{w}_k^T \mathbf{x}_i \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \\ \vdots \\ \mathbf{w}_k^T \end{bmatrix} \mathbf{x}_i = W^T \cdot \mathbf{x}_i$$



# Principal Component Analysis



- PCA aligns maximum variance directions with standard basis
  - Variance along each dimension is **uncorrelated**
  - PCs  $w_j$  and latent representation  $h_{ij}$  do not change with increasing  $k$

# Summary: Principal Component Analysis

1. Estimate the covariance matrix  $S$  of the data  $X \in \mathbb{R}^{d \times n}$
2. Compute the eigenvectors of  $S$
3.  $W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$  where  $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$  are the eigenvectors corresponding to the  $k$  largest eigenvalues
4. Project the data onto  $W$ :  $H = W^\top \cdot X$   
\* If needed<sup>2</sup>: reconstruct data by  $X \approx \tilde{X} = WH$

---

<sup>2</sup>This holds for all Linear Matrix Factorization methods

# Summary: Principal Component Analysis

---

## Algorithm 1: Principal Component Analysis

---

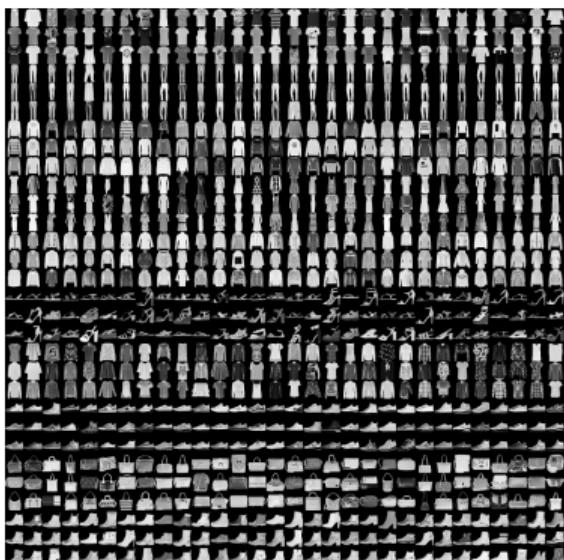
Require: data  $x_1, \dots, x_n \in \mathbb{R}^d$ , number of principal components  $k$

- 1: # Compute Sample Covariance Matrix
  - 2:  $C = 1/n (X - 1/n \sum_i x_i)(X - 1/n \sum_i x_i)^\top$
  - 3: # Compute eigenvectors corresponding to the  $k$  largest eigenvalues
  - 4:  $W = \text{eig}(C)$
  - 5: # Project data onto  $W$
  - 6:  $H = W^\top X$
  - 7: return  $W, H$
-

# PCA on Fashion MNIST

Fashion-MNIST is a dataset of Zalando's article images consisting of 70.000 examples. Each example is a  $28 \times 28$  grayscale image, associated with a label from 10 classes.

# PCA on Fashion MNIST



# PCA on Fashion MNIST: Shirts

Figure: Eigenvectors ( $W$ )

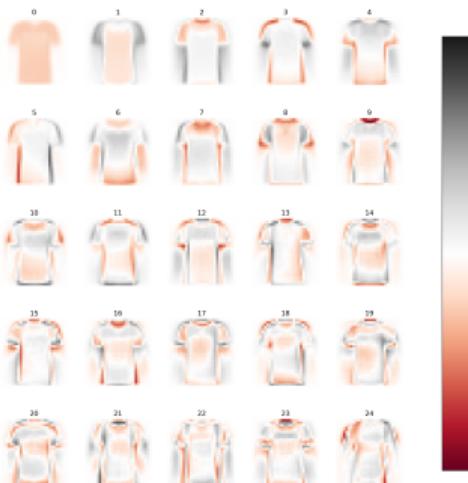
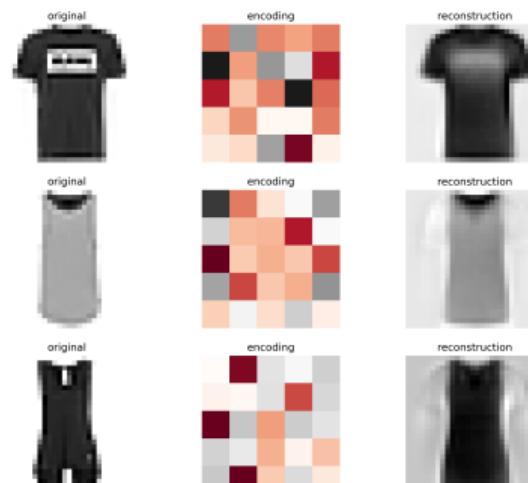


Figure: 3 examples ( $X, H, WH$ )



Here we compress images to 25 dimensions instead of 784 ( $28 \times 28$ )  
→ more than factor 31

# PCA For High-Dimensional Data

Input: centered data  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$  with  $n \ll d$

- Covariance matrix  $XX^\top$  will be very large ( $d$ -by- $d$ )
- Too few samples for a robust covariance matrix estimate

We know the direction of maximal variance  $\mathbf{w}$  must lie in the span of the data (for  $\lambda \neq 0$ ):

$$\lambda \mathbf{w} = XX^\top \mathbf{w} \rightarrow \mathbf{w} = X\alpha$$

where  $\alpha = 1/\lambda \cdot X^\top \mathbf{w}$  is a weighting of each data point

## PCA For High-Dimensional Data

We can plug  $w = X\alpha$  in and obtain

$$\alpha = \frac{1}{\lambda} X^T w$$

$$\begin{aligned} XX^T w &= \lambda w \\ XX^T X\alpha &= \lambda X\alpha \\ \underbrace{X^T X}_{K} \underbrace{X^T X}_{K} \alpha &= \lambda \underbrace{X^T X}_{K} \alpha \\ KK\alpha &= \lambda K\alpha \end{aligned}$$

which can be solved by [Schölkopf et al., 1998]

$$K\alpha = \lambda\alpha.$$

Solving PCA via  $X^T X$  instead of  $XX^T$  is called **linear kernel PCA**

Note: if we want to use other kernels, we have to take care that data is centered in features space → more complicated

## When should we use a kernel for PCA ?

If there are more dimensions than samples ( $n \ll d$ )

→ Compute PCA on linear kernel matrix  $X^\top X \in \mathbb{R}^{n \times n}$

If there are more samples than dimensions ( $d \ll n$ )

→ Compute PCA on covariance matrix  $XX^\top \in \mathbb{R}^{d \times d}$

# Wrap-up: Linear Kernel PCA

---

## Algorithm 2: Linear Kernel PCA

---

**Require:** data  $x_1, \dots, x_n \in \mathbb{R}^d$ ,  $n \ll d$ , number of principal components  $k$

- 1: # Compute Linear Kernel
  - 2:  $K = (X - 1/n \sum_i x_i)^\top (X - 1/n \sum_i x_i)$
  - 3: # Compute eigenvectors corresponding to the  $k$  largest eigenvalues
  - 4:  $\alpha = \text{eig}(K)$
  - 5:  $W = X\alpha$
  - 6: # Project data onto  $W$
  - 7:  $H = W^\top X$
  - 8: **return**  $W, H$
-

# Trends in Text Data

Let's look at some more applications of PCA!

**Extra' Host Maria Menounos' YouTube Channel Gives TV Fans an Online Forum**

**Items (1)**  
6 minutes ago  
[View more popular stories on BetaNews this past week: December 29 — January 4](#)

co-host of syndicated entertainment newsmagazine "Extra," Maria Menounos talks a lot about TV, but she's even charter online. The multiplatform host behind an uncanny reality show has joined forces with YouTube company and a bunch of others, as well as spokespersons like Parsons, to quickly built YouTube network Ferme into a platform... Read more

**live from the Engadget CES Stage: Pebble CEO Eric Migicovsky**

**Items (1)**  
7 minutes ago  
[Steve Larson: What doctors can learn from each other - Steve Larson \(2011\)](#)

clickbuster success story Pebble was the darling of last year's CES, helping to usher in a year in which wearables were all the rage. The smartwatch maker's CEO Eric Migicovsky will be joining us to discuss what the company has up its proverbial sleeve.

**Finally There Is An "Alien" Game That Is Actually Like The Movies**

**Items (1)**  
1 hour ago  
[TV: Malcolm Gladwell: The unheard story of David and Goliath - Malcolm Gladwell \(2013\)](#)

announced today, Alien: Isolation is out this year. The announcement trailer is slow, quiet and terrifying—everything we love about the series.

**Venrock VC leaves to launch China clean energy platform**

**Items (1)**  
1 hour ago  
[An Acer user's a small and affordable Android tablet -- ironia AI and B1](#)

two venture firm Venrock's energy investors, Matt Theodore and Matthew Nordam, have left the firm. Theodore is an accomplished investors who was one of my first interviews after I launched the investment site. We've published many pieces from Nordam on the state of cleantech venture.

**Apple's App Store revenue hit record \$10 billion in 2013**

**Items (1)**  
2 minutes ago on [Seen: whistleblower 2014 tech predictions](#)

Apple announced today that sales from its iOS App Store reached \$10 billion last year, which is up 50 percent from 2012. That's made 2013 Apple's most successful year ever since its launch of the App Store in July of 2008, and it positions the company as the undisputed leader in mobile... Read More.

**best pictures of the day - live**

**Items (1)**  
10 days ago on [Tech: Chris Downey: Design with the blind in mind - Chris Downey \(2013\)](#)

Guardian's photo team brings you a daily round up from the world of photography. Joana

**Democracy needs whistleblowers. That's why I broke into the FBI in 1971**

**Items (1)**  
2 minutes ago on [Tech: Chris Downey: Design with the blind in mind - Chris Downey \(2013\)](#)

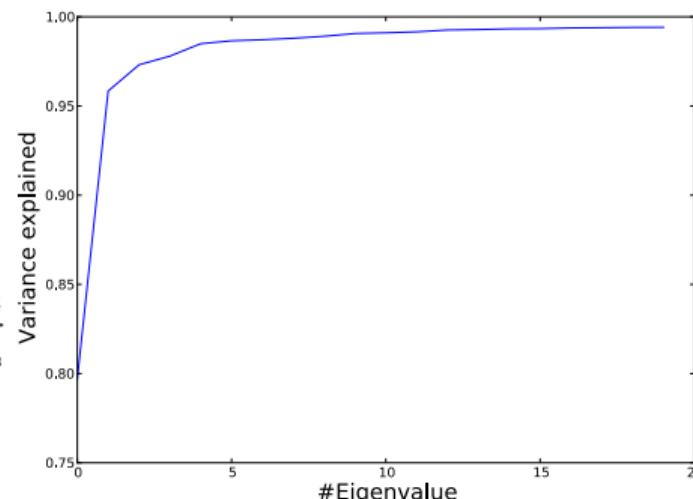
Steubenville, Ohio, had laws to reveal something that was more dangerous. We wanted to warn the public about that. So I got a job at the FBI and I became a informant and I broke into an FBI office in Medina, Pennsylvania, in March 1971 and removed about 1,000 documents from the filing cabinets. We had a hunch that there would be incriminating evidence in those files. And there was. And we found out that they had been doing every single thing that went on under his world would be recorded. But we could not be sure, and until we did, we were on tenterhooks. A shroud went up around the group of eight of us. One of us had to go to the FBI office and get the records. And we had to do it in secret because if we had intruded the bureau's agents to set up interviews of anti-war activists as it would enhance the war on terrorism. So we had to do it in secret. And we had to do it in secret because every agent there was really mailed. That's why the first place of evidence was across there is an FBI office in Medina, Pennsylvania.

That's why I broke into the FBI in 1971. I'm a former partner with Edward Snowden has done in releasing National Security Agency documents that show the NSA's blanket surveillance of Americans. I think Snowden's a legitimate whistleblower, and I guess

# Trends in Text Data

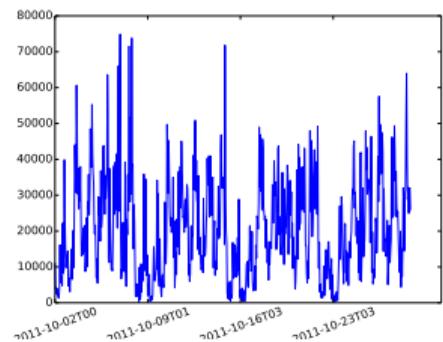
We are looking at Bag-Of-Words data from news web pages  
wunderfacts.com, in October 2011

- We store the data in a matrix  
 $X \in \mathbb{R}^{d \times t}$   
 $X_{ij} = 10$   
The entry  $X_{dt} = 10$  tells us: word  $d$  i  
was counted 10 times in time bin  $t$  j
- Let's apply PCA and look at the  
explained variance ( $ev_i = \frac{\sum_{j=1}^i \lambda_j}{\sum_{l=1}^d \lambda_l}$ )  
 $EV_m = \frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^d \lambda_i}$
- We only need 15 principal directions  
to explain  $>99\%$  of the data

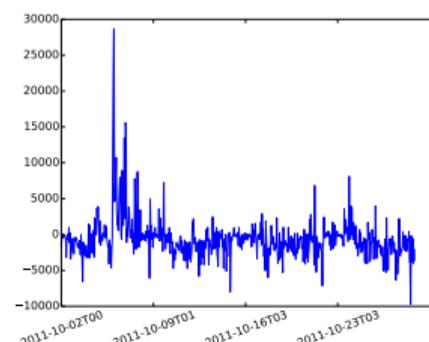


# Trends in Text Data

## First Principal Component



## Second Principal Component



Main Variance due to weekly/daily publishing activity

Spike on day of Steve Job's death

→ We can use PCA as a tool for analysing big unlabelled data

# Non-Negative Matrix Factorization (NMF)

- For some data PCA is not intuitive
- Example: Non-negative data
  - Principal directions will have negative entries
  - This can be hard to interpret
- Many data sets are strictly non-negative
  - Text data
  - Image data
  - Probabilistic data
- NMF is straightforward to implement
- Applicable to Recommender Systems (more info here)

# Non-Negative Matrix Factorization

Given non-negative data  $X \in \mathbb{R}_+^{d \times n}$  we want to find  $W \in \mathbb{R}_+^{d \times m}$ ,  $H \in \mathbb{R}_+^{m \times n}$  such that

$$\operatorname{argmin}_{W,H} \|X - WH\|_{\text{Fro}}^2 = \operatorname{argmin}_{W,H} \sum_{d=1}^d \sum_{n=1}^n (X_{dn} - (WH)_{dn})^2$$

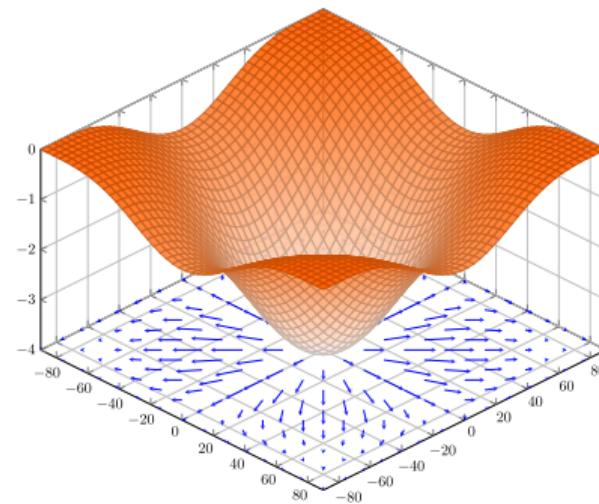
# Non-Negative Matrix Factorization

$$\underset{W,H}{\operatorname{argmin}} \quad ||X - WH||_{\text{Fro}}^2 \quad \text{s.t. } W \geq 0 \text{ and } H \geq 0$$

Note that the constraints make the problem NP-hard and Ill-posed [Gillis, 2014].

## Note

Whenever the problem is too hard just follow the gradients!

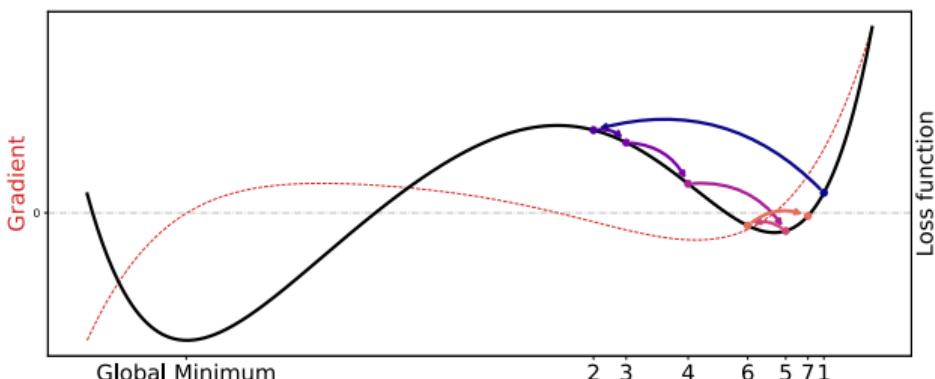


from Wikipedia

# Gradient Descent

Gradient descent finds **an** optimum of Loss  $\mathcal{L}(\alpha)$  by iterating

$$\alpha \leftarrow \alpha - \eta \frac{\partial \mathcal{L}(\alpha)}{\partial \alpha}$$



- **Not guaranteed to find global optimum**
- Adaptive  $\eta$  can yield improvements
- Stochastic Gradient Descent is basis for State of the Art ML

$\text{Tr}$  is the trace of the matrix

$$\text{Tr} \begin{pmatrix} A_{11} & & & \\ \vdots & A_{1d} & & \\ A_{d1} & & & A_{dd} \end{pmatrix} = A_{11} + A_{22} + \dots + A_{dd}$$

## Using

$$\frac{\partial}{\partial X} \|X\|_{\text{Fro}}^2 = \frac{\partial}{\partial X} \text{Tr}(XX^T)$$

$$\text{Tr}(A) = \sum_i A_{ii}$$

$$\frac{\partial \text{Tr}(AXX^TA^T)}{\partial X} = 2A^TAX$$

$$\frac{\partial \text{Tr}(XAX^T)}{\partial X} = XA^T + XA$$

$$\frac{\partial \text{Tr}(AXB)}{\partial X} = \frac{\partial \text{Tr}(B^TX^TA^T)}{\partial X} = A^TB^T$$

$$\frac{\partial \|X - WH\|_{\text{Fro}}^2}{\partial H} = 2(W^TWH - W^TX)$$

$$\frac{\partial \|X - WH\|_{\text{Fro}}^2}{\partial W} = 2(WHH^T - XH^T)$$

# Non-Negative Matrix Factorization

Gradient descent finds an optimal solution by iterating

$$H \leftarrow H - \eta (W^T W H - W^T X)$$

$$W \leftarrow W - \eta (W H H^T - X H^T)$$

$$\eta_{ij}^H := \frac{H_{ij}}{(W^T W H)_{ij}}$$

- By choosing  $\eta$  wisely one can transform the additive updates into multiplicative ones [Lee and Seung, 1999, 2000]:

$$H = H \odot W^T X \oslash W^T W H \quad H_{ij} \leftarrow H_{ij} \frac{(W^T X)_{ij}}{(W^T W H)_{ij}}$$

$$W = W \odot X H^T \oslash W H H^T$$

$$W_{ij} \leftarrow W_{ij} \frac{(X H^T)_{ij}}{(W H H^T)_{ij}}$$

where

- is *element-wise* multiplication (in python '\*')
- is *element-wise* division (in python '/')

# NMF Algorithm

---

## Algorithm 3: Non-negative Matrix Factorization

---

Require: data  $X = [x_1, \dots, x_n] \in \mathbb{R}_+^{d \times n}$ , number of factors  $k$

- 1: # Initialize  $W \in \mathbb{R}_+^{d \times k}$ ,  $H \in \mathbb{R}_+^{k \times n}$  randomly
  - 2: # Add a small constant  $\epsilon = 10^{-19}$  to  $X$  to avoid zero-divisions
  - 3: **for**  $it \leq \text{Iterations}$  **do**
  - 4:      $H = H \odot W^\top X \oslash W^\top WH$
  - 5:      $W = W \odot XH^\top \oslash WHH^\top$
  - 6: **end for**
  - 7: **return**  $W, H$
-

## NMF on Fashion MNIST: Shirts

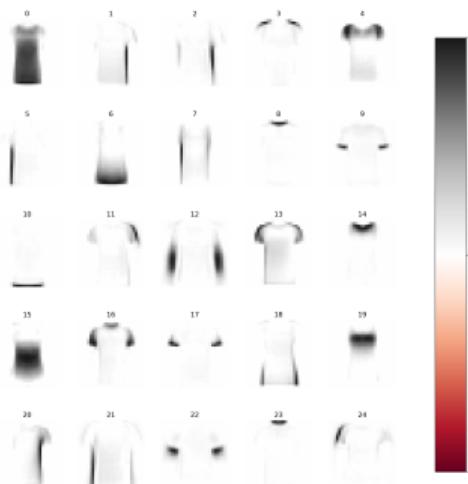


Figure: 25 main components ( $W$ )

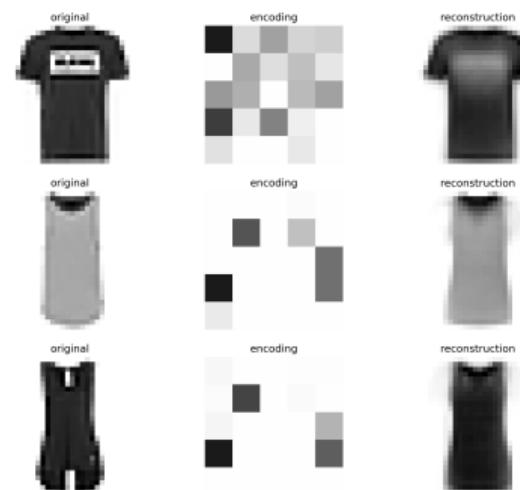


Figure: 3 examples ( $X, H, WH$ )

# PCA vs. NMF - Fashion Mnist

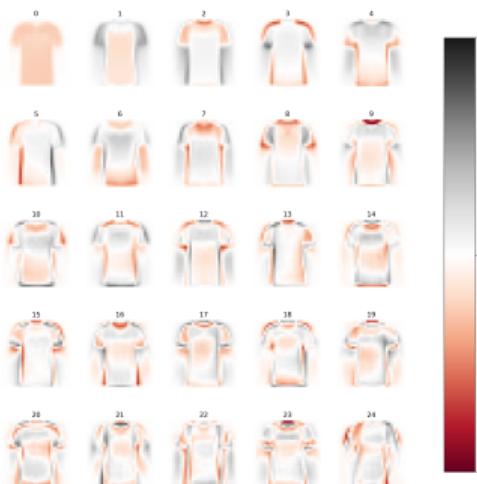


Figure:  $W_{PCA}$

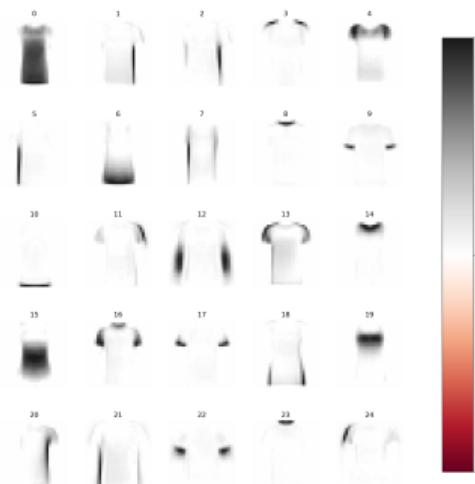
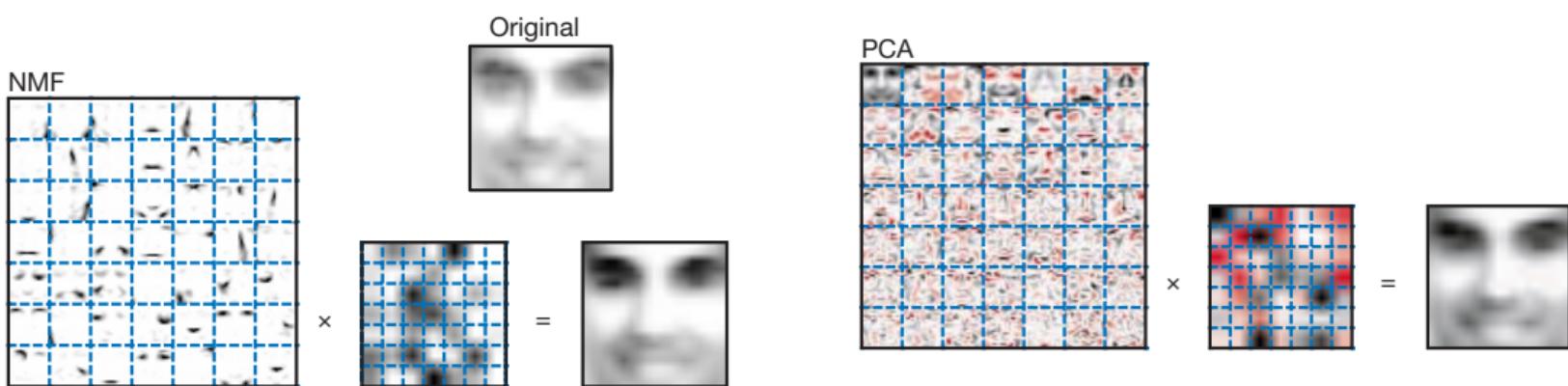


Figure:  $W_{NMF}$

# PCA vs. NMF - Face Parts

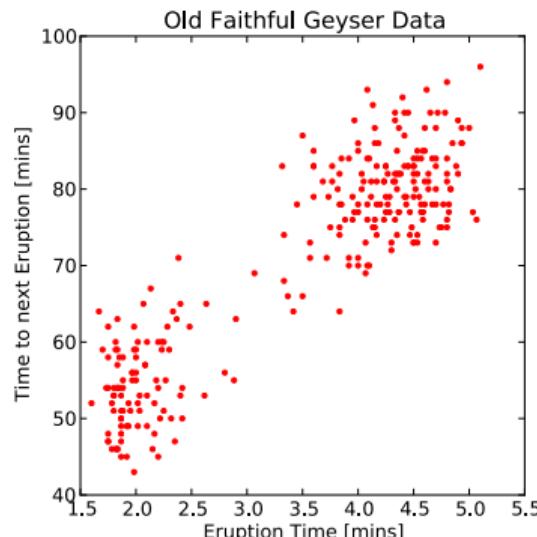


Taken from Lee and Seung [1999]

# Clustering

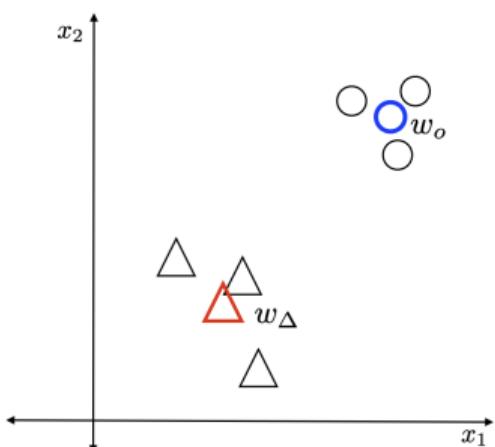
PCA/NMF are applied to reduce dimensionality  
Often the problem setting is different:

You want to categorize eruptions without labels



# Clustering

Remember Lecture 2:  
Psychological Models of Categorization: Prototypes



Prototypes  $\mu_\Delta$  and  $\mu_o$ :

$$\mu_\Delta = \frac{1}{n_\Delta} \sum_n^{n_\Delta} x_{\Delta,n}$$

$$\mu_o = \frac{1}{n_o} \sum_n^{n_o} x_{o,n}$$

New data points  $x$  are assigned to their closest cluster center  $\mu^*$

$$\mu^* = \operatorname{argmin}_i (\|\mu_i - x\|_2)$$

# K-means Clustering

**Objective:** Find cluster centers  $\mu_1, \dots, \mu_k$  such that the sum of distances of data points to their respective cluster centers (WCSS<sup>3</sup>) is minimized

$$L(\{\mu_1, \dots, \mu_k\}, r) = \sum_{i=1}^n \|x_i - \mu_{r_i}\|^2$$

where  $r_i$  : cluster index of data point  $i$

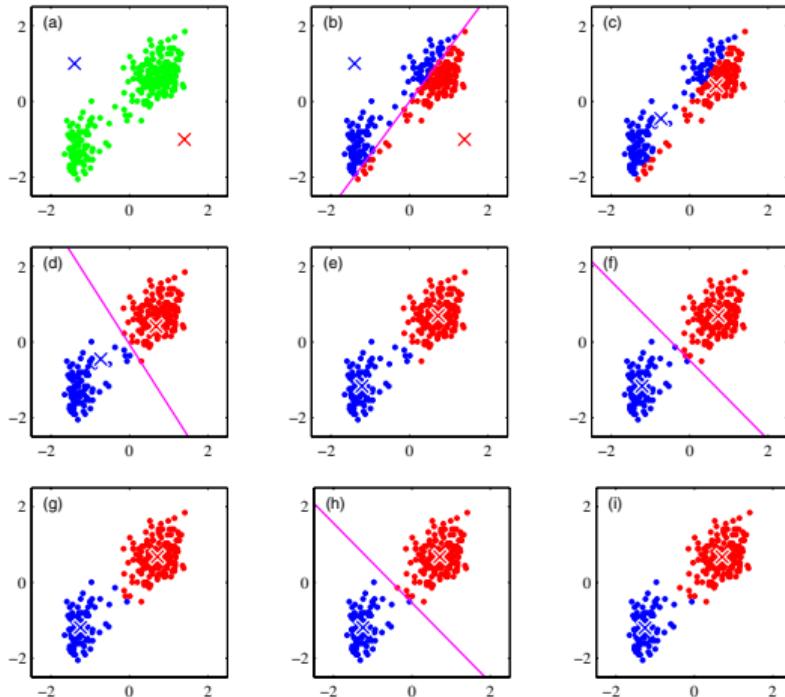
We minimize  $L$  by re-iterating two steps:

1. Assign each data point  $x_i$  to their closest cluster  $\mu_{r_i}$
2. Update  $\mu_r$  to the mean of the members in that cluster  $r$

---

<sup>3</sup>Within Cluster Sum of Squares

# K-means Clustering Step-by-Step



Re-iterate two steps:

1. Assign each  $x_i$  to closest cluster  $\mu_r$ ,
2. Update  $\mu_r$  to mean of members in cluster  $r$

# K-means Clustering Algorithm

---

**Require:** data  $x_1, \dots, x_n \in \mathbb{R}^d$ , number of clusters  $k$ , iterations  $m$ .

```
1: Choose random data points as initial cluster centers  $\mu_1 \leftarrow x_{i_1}, \dots, \mu_k \leftarrow x_{i_k}$  where  $i_j \neq i_l$  for all  $j \neq l$ .
2:  $r \leftarrow \mathbf{0}_n$ 
3:  $r' \leftarrow \mathbf{0}_n$ 
4:  $i \leftarrow 0$ 
5: while  $i < m$  do
6:   for  $j \leftarrow 1$  to  $n$  do
7:     Find nearest cluster center  $r'_j \leftarrow \operatorname{argmin}_{1 \leq l \leq k} \|x_j - \mu_l\|_2$ 
8:   end for
9:   for  $j \leftarrow 1$  to  $k$  do
10:    Compute new cluster center  $\mu_j \leftarrow \frac{1}{|\{l:r'_l=j\}|} \sum_{l:r'_l=j} x_l$ 
11:   end for
12:   if  $r = r'$  then
13:     break
14:   end if
15:    $r \leftarrow r'$ 
16:    $i \leftarrow i + 1$ 
17: end while
18: return cluster centers  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$ , assignment vector  $r \in \mathbb{R}^n$ 
```

---

# Application Example: Image compression

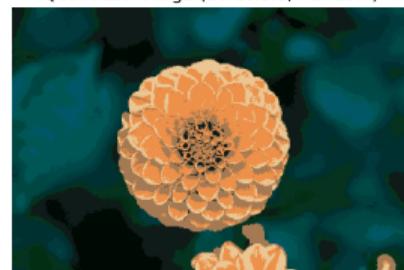
Original image (colors: RGB 3 channels [0, 255])



Quantized image (16 colors, K-Means)



Quantized image (16 colors, Random)

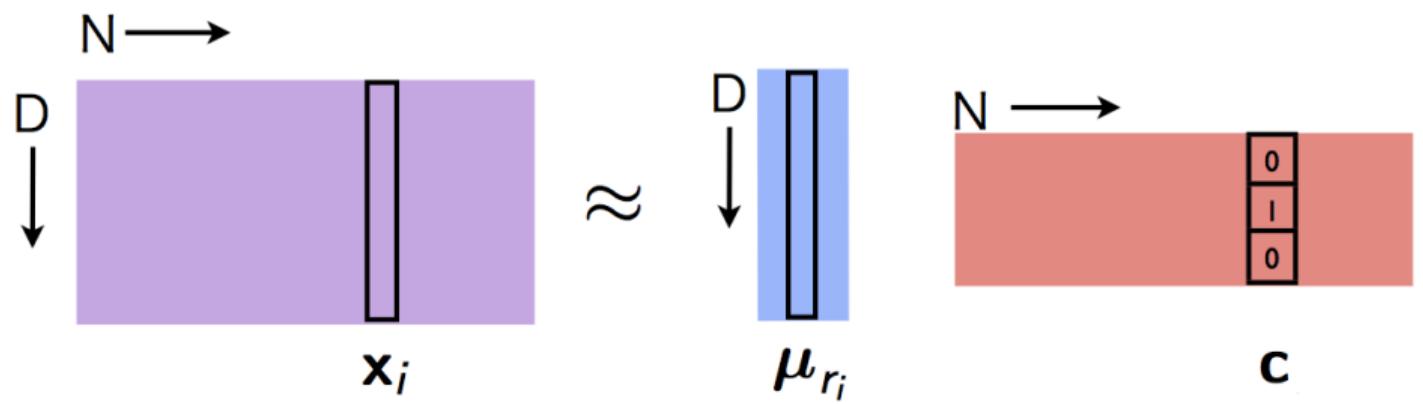


Adapted from sklearn

Encode color with one nibble (4 bit) instead of 3 Byte (24 bit)  
→ only need  $1/6^{th}$  of the bandwidth (+ dictionary of colors)

## Clustering can be seen as Matrix Factorization

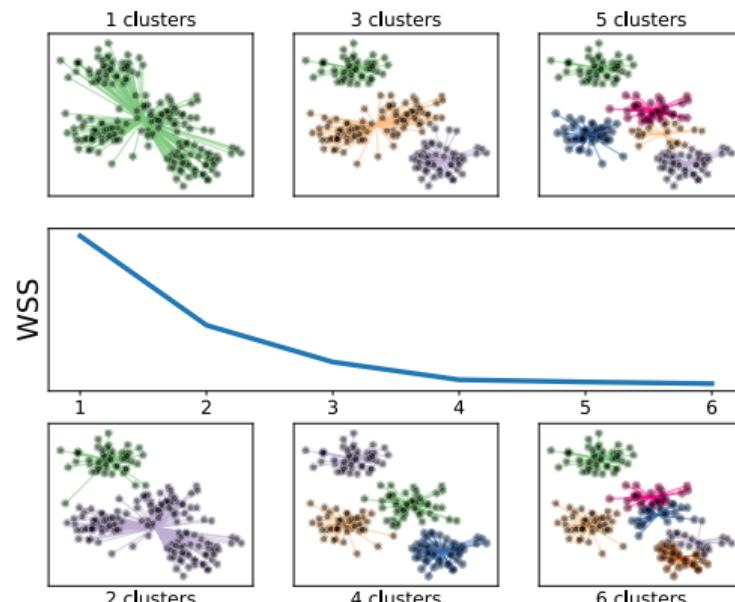
Clustering finds an **optimal partitioning** of a data set <sup>4</sup>



<sup>4</sup>In the previous example, we actually used this by approximating the original colors  $x_i$  with the corresponding  $\mu_{r_i}$

## How to choose $k$

- Number of clusters  $k$  is critical hyper-parameter
- In supervised settings we use model selection (grid search) to optimize hyper-parameters for accuracy on test data
- How can we optimize the number of clusters?



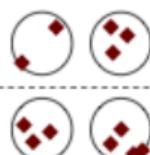
→ One approach: find "elbow"; Lowest  $k$  after which no real change

## Clustering Instability

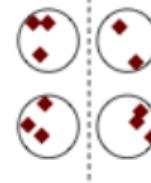
Number of Clusters is a critical parameter

k = 2;

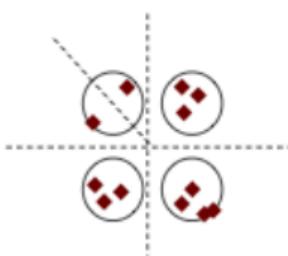
## Sample 1



## Sample 2



$$k = 5:$$



Clusterings are unstable (i.e. converge to different results) if number of clusters is too small or too large

# Summary I

## Matrix Factorization Methods

PCA and NMF belong to this class

Linearly approximate original data  $X$  with  $\approx WH$

## Principal Component Analysis

is a popular dimensionality reduction tool

aligns directions of maximal variance with standard basis

finds orthogonal directions

finds optimal matrix factorization

## Non-negative Matrix Factorization

works for non-negative data (count data, probabilistic data)

does not find orthogonal directions / uncorrelated factors

NMF encoding typically more sparse than PCA encoding

# Summary II

## Gradient Descent

useful for non-convex optimization  
work-horse of Machine Learning

## K-Means Clustering

finds an optimal partitioning of a data set  
K-Means requires

Good initialization  
Knowledge of optimal  $k$

# References

- C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 2007.
- N. Gillis. The why and how of nonnegative matrix factorization. *Regularization, Optimization, Kernels, and Support Vector Machines*, 12(257):257–291, 2014.
- D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–91, 1999. doi: 10.1038/44565.
- D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *In NIPS*, pages 556–562. MIT Press, 2000.
- K. Pearson. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901.
- B. Schölkopf, A. J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(6):1299–1319, 1998.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. 08 2017.

# References

- C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. 2007.
- N. Gillis. The why and how of nonnegative matrix factorization. *Regularization, Optimization, Kernels, and Support Vector Machines*, 12(257):257–291, 2014.
- D. D. Lee and H. S. Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755):788–91, 1999. doi: 10.1038/44565.
- D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *In NIPS*, pages 556–562. MIT Press, 2000.
- K. Pearson. Liii. on lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572, 1901.
- B. Schölkopf, A. J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(6):1299–1319, 1998.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. 08 2017.