

Applied Machine Learning

Dimensionality reduction

Reihaneh Rabbany



COMP 551 (winter 2021)

Learning objectives

What is dimensionality reduction?

What is it good for?

Linear dimensionality reduction:

- Principal Component Analysis
- Relation to Singular Value Decomposition

Motivation

Scenario: we are given high dimensional data and asked to make sense of it!

Real-world data is high-dimensional

- we can't visualize beyond 3D
- features may not have any semantics (value of the pixel vs happy/sad)
- processing and storage is costly
- many features may not vary much in our dataset (e.g., background pixels in face images)



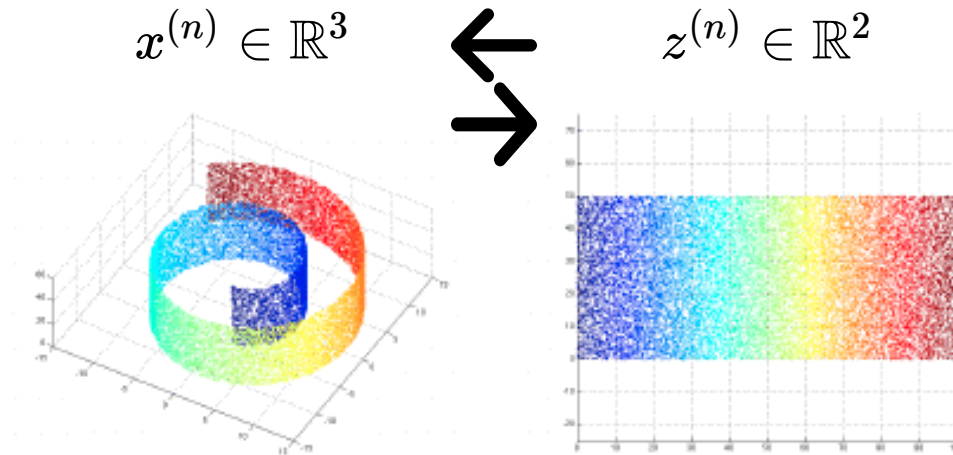
Dimensionality reduction: faithfully represent the data in low dimensions

- We can often do this with real-world data (*manifold hypothesis*)
- finding meaningful low-dimensional structures in high-dimensional observations

Dimensionality reduction

Dimensionality reduction: faithfully represent the data in low dimensions

- learn a mapping between (coordinates) at low-dimension and high-dimensional data

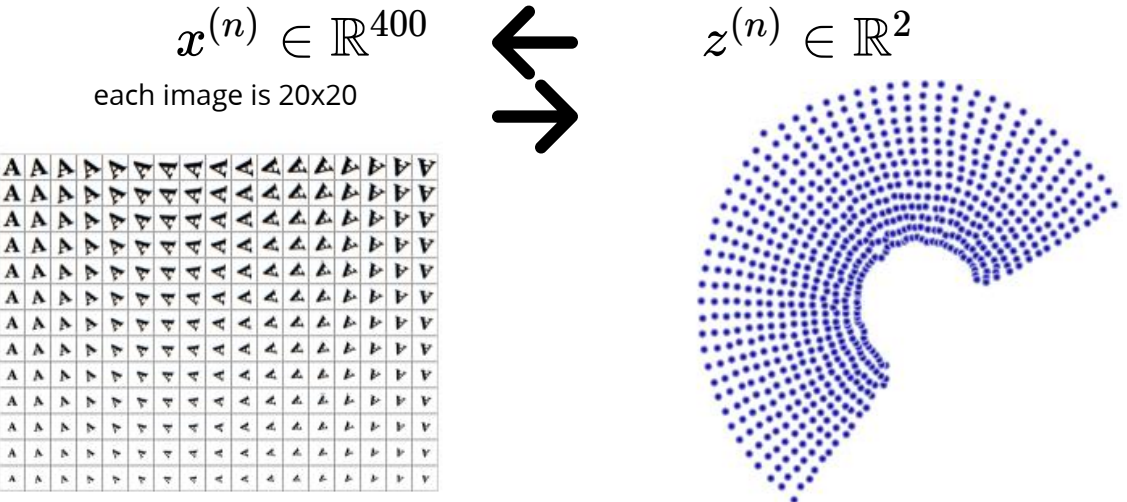


some methods give this mapping in both directions and some only in one direction.

Dimensionality reduction

Dimensionality reduction: faithfully represent the data in low dimensions

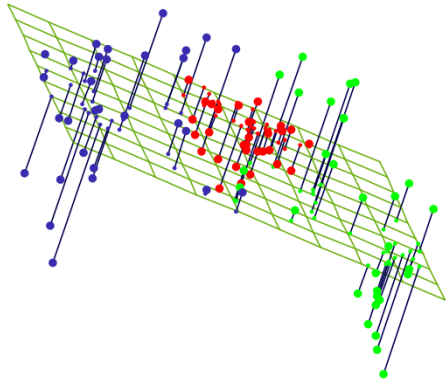
- learn a mapping between (coordinates) at low-dimension and high-dimensional data



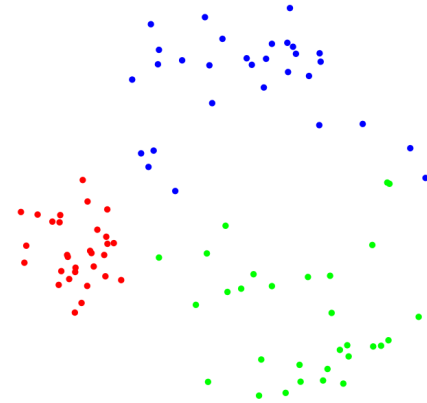
Principal Component Analysis (PCA)

PCA is a **linear** dimensionality reduction method

$$x^{(n)} \in \mathbb{R}^3$$



$$z^{(n)} \in \mathbb{R}^2$$



$$Q \in \mathbb{R}^{3 \times 2}$$
A diagram showing a red arrow pointing right labeled $Q \in \mathbb{R}^{3 \times 2}$ and a blue arrow pointing left labeled Q^T .

where Q has orthonormal columns $Q^T Q = I$

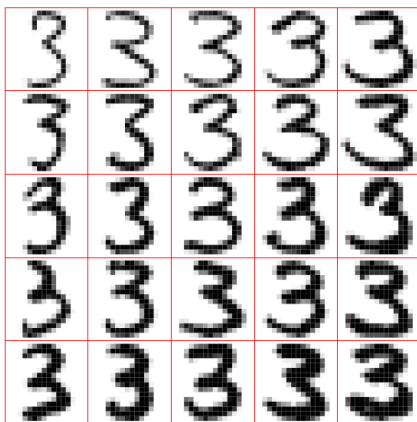
it follows that the pseudo-inverse of Q is $Q^\dagger = (Q^T Q)^{-1} Q^T = Q^T$

PCA: optimization objective

PCA is a **linear** dimensionality reduction method

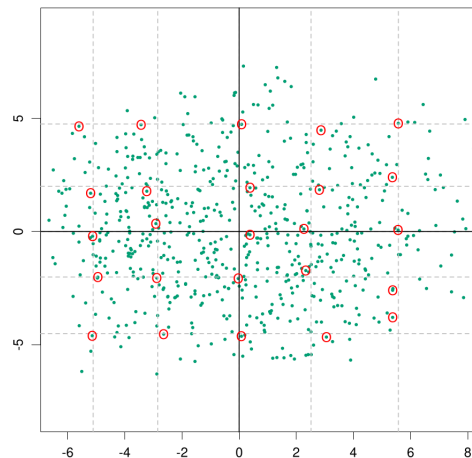
$$x^{(n)} \in \mathbb{R}^{784}$$

each image has $28 \times 28 = 784$ pixels



$$Q \in \mathbb{R}^{784 \times 2}$$
$$Q^T$$

$$z^{(n)} \in \mathbb{R}^2$$



faithfulness is measured by the reconstruction error

$$\min_Q \sum_n \|x^{(n)} - \underbrace{x^{(n)} Q}_{z^{(n)}} Q^T\|_2^2 \quad s.t. \quad Q^T Q = I$$

PCA: optimization objective

PCA is a **linear** dimensionality reduction method

faithfulness is measured by the reconstruction error

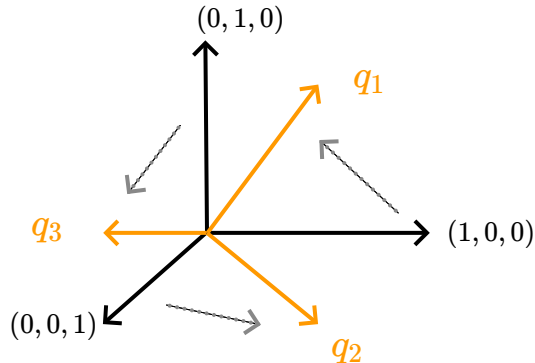
$$\min_Q \sum_n \left\| \underset{z^{(n)}}{x^{(n)}} - x^{(n)\top} Q Q^\top \right\|_2^2 \quad s.t. \quad Q^\top Q = I$$

strategy: find $D \times D$ matrix Q , and only use D' columns

Since Q is orthogonal we can think of it as a change of coordinates

$$Q = \begin{bmatrix} Q_{1,1}, \dots, Q_{1,D} \\ \vdots, \dots, \vdots \\ Q_{D,1}, \dots, Q_{D,D} \end{bmatrix}$$

q_1 q_D



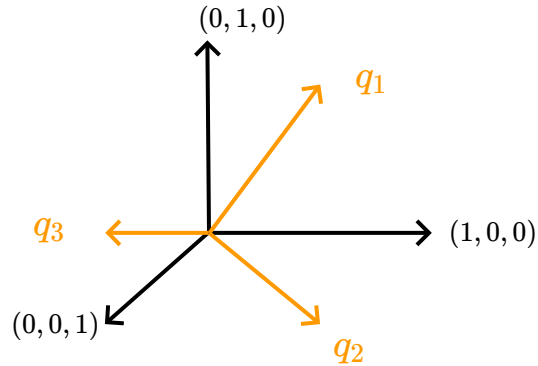
PCA: a change of coordinates

strategy: find $D \times D$ matrix Q , and only use D' columns

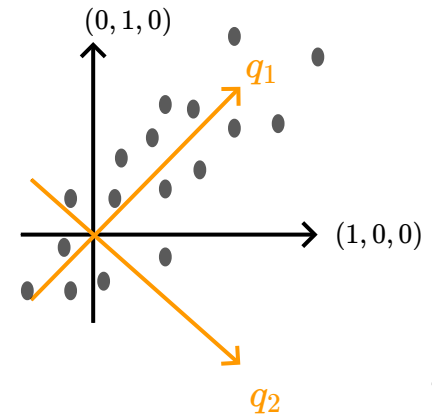
Since Q is orthonormal we can think of it as a change of coordinates

$$Q = \begin{bmatrix} Q_{1,1} & \dots & Q_{1,D} \\ \vdots & \ddots & \vdots \\ Q_{D,1} & \dots & Q_{D,D} \end{bmatrix}$$

q_1 q_D



example $D = 2$



we want to change coordinates such that coordinates $1, 2, \dots, D'$ best explain the data for any given D'

PCA preserves variance

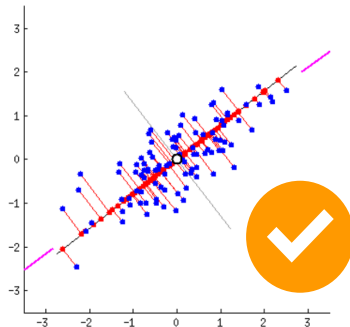
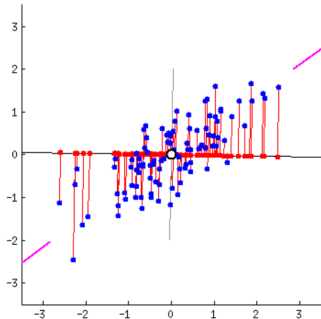
Find a change of coordinate using *orthonormal matrix*

first new coordinate has **maximum variance** (*lowest reconstruction error*)

second coordinate has the next largest variance

...

along which one of these directions the data has a higher variance? more spread out?



$$Q = \begin{bmatrix} Q_{1,1}, \dots, Q_{1,D} \\ \vdots, \ddots, \vdots \\ Q_{D,1}, \dots, Q_{D,D} \end{bmatrix}$$

q_1

this direction is the vector q_1

projection is given by $\frac{x^{(n)\top} q_1}{\|q_1\|_2} = x^{(n)\top} q_1$

projection of the whole dataset is $X q_1 = z_1$

$$z_1^\top = [z_1^{(1)}, z_1^{(2)}, \dots, z_1^{(N)}]$$

PCA preserves variance

Find a change of coordinate using *orthonormal matrix*

first new coordinate has maximum variance

projection of the whole dataset is $z_1 = X q_1$

$$\text{Var}(z_1) = \frac{1}{N} \sum_n (z_1^{(n)} - 0)^2$$

assuming features have zero mean, maximize the variance of the projection: $\frac{1}{N} z_1^\top z_1$

$$\max_{q_1} \frac{1}{N} z_1^\top z_1 = \max_{q_1} \frac{1}{N} q_1^\top X^\top X q_1 = \max_{q_1} q_1^\top \Sigma q_1$$

dxd **covariance matrix**

$$\Sigma = \frac{1}{N} X^\top X = \frac{1}{N} \sum_n (x^{(n)} - 0)(x^{(n)} - 0)^\top$$

because the mean is zero

$\Sigma_{i,j}$ is the sample covariance of feature i and j

$$\Sigma_{i,j} = \text{Cov}[X_{:,i}, X_{:,j}] = \frac{1}{N} \sum_n x_i^{(n)} x_j^{(n)}$$

Covariance matrix

variance of a random variable $\text{Var}(x) = \mathbb{E}[(x - \mathbb{E}[x])^2] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$

covariance of two random variable $\text{Cov}(x, y) = \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$

for $x \in \mathbb{R}^D$ we have **covariance matrix**

$$\text{Cov}(x_1, x_1) = \text{Var}(x_1) \quad \text{Cov}(x_1, x_D)$$
$$\Sigma = \begin{bmatrix} \Sigma_{1,1} & \cdots & \Sigma_{1,D} \\ \vdots & \ddots & \vdots \\ \Sigma_{D,1} & \cdots & \Sigma_{D,D} \end{bmatrix} = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^\top] = \mathbb{E}[xx^\top] - \mathbb{E}[x]\mathbb{E}[x]^\top$$

$D \times 1 \quad 1 \times D \quad D \times D \quad D \times D$

given a dataset $\mathcal{D} = \{x^{(1)}, \dots, x^{(N)}\}$ **sample covariance matrix**

$$\hat{\Sigma}^{\text{MLE}} = \mathbb{E}_{\mathcal{D}}[(x - \mathbb{E}_{\mathcal{D}}[x])(x - \mathbb{E}_{\mathcal{D}}[x])^\top]$$

the empirical estimate

$$x - \left(\frac{1}{N} \sum_{x \in \mathcal{D}} x\right)$$

Correlation and dependence

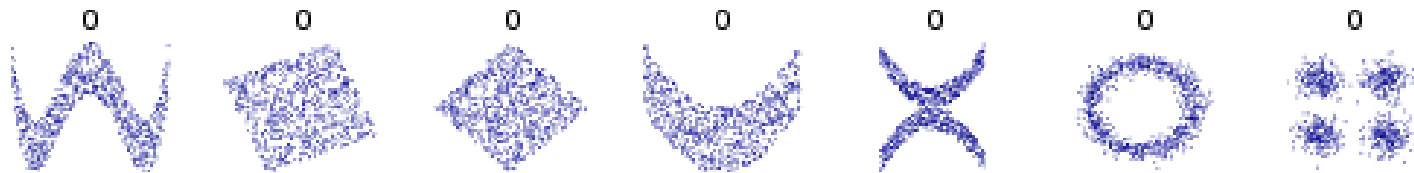
correlation is **normalized covariance**

$$\text{Corr}(x_i, x_j) = \frac{\text{Cov}(x_i, x_j)}{\sqrt{\text{Var}(x_i)\text{Var}(x_j)}} \in [-1, +1]$$

two variables that are **independent are uncorrelated** as well

$$p(x_i, x_j) = p(x_i)p(x_j) \rightarrow \mathbb{E}[x_i x_j] = \mathbb{E}[x_i]\mathbb{E}[x_j] \rightarrow \text{Cov}(x_i, x_j) = 0$$

the inverse is generally not true (zero correlation doesn't mean independence)



in each example above correlation between two coordinates is zero, but they are not independent

Decomposing the covariance matrix

covariance matrix is symmetric positive semi definite

- symmetric
 - $\Sigma_{d,d'} = \text{Cov}(x_d, x_{d'}) = \text{Cov}(x_{d'}, x_d) = \Sigma_{d',d}$
- positive semi definite
 - for any $y \in \mathbb{R}^D$ we have $y^\top \Sigma y = (y^\top \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^\top])y = \text{Var}(y^\top x) \geq 0$

any symmetric positive semi-definite matrix can be decomposed as

$$\Sigma = Q \Lambda Q^\top$$

|
diagonal $D \times D$

orthogonal $Q Q^\top = Q^\top Q = I$ (rotation and reflection)

Spectral Decomposition

PCA with Eigenvalue decomposition

find a change of coordinate using *an orthogonal matrix*

first new coordinate has maximum variance

$$\max_{q_1} q_1^T \Sigma q_1 \quad s.t. \quad \|q_1\| = 1$$

covariance matrix is **symmetric** and **positive semi-definite**

$$(X^T X)^T = X^T X \quad a^T \Sigma a = \frac{1}{N} a^T X^T X a = \frac{1}{N} \|Xa\|_2^2 \geq 0 \quad \forall a$$

any symmetric matrix has the following decomposition

$$\Sigma = \underset{\substack{| \\ |}}{Q} \underset{\substack{| \\ |}}{\Lambda} \underset{\substack{| \\ |}}{Q}^T \quad (\text{as we see shortly using } Q \text{ here is not a co-incidence})$$

$Q Q^T = Q^T Q = I$ dxd
each column is an eigenvector

orthogonal matrix
diagonal and sorted ($\lambda_1 > \lambda_2 > \lambda_3 > \dots$)
corresponding eigenvalues are on the diagonal
positive semi-definiteness means these are non-negative

PCA: Principal Component Analysis

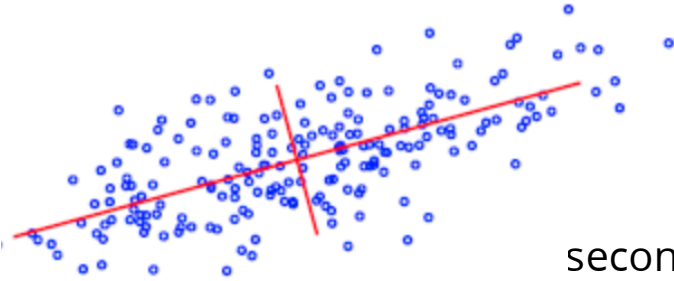
find a change of coordinate using *an orthogonal matrix*

first new coordinate has maximum variance

$$q_1^* = \arg \max_{q_1} q_1^\top \Sigma q_1 \quad s.t. \quad \|q_1\| = 1$$

$\max_{q_1} q_1^\top Q \Lambda Q^\top q_1 = \lambda_1$ using eigenvalue decomposition

maximizing direction is the eigenvector with the largest eigenvalue (first column of Q)



$q_1 = Q_{:,1}$ first principal direction

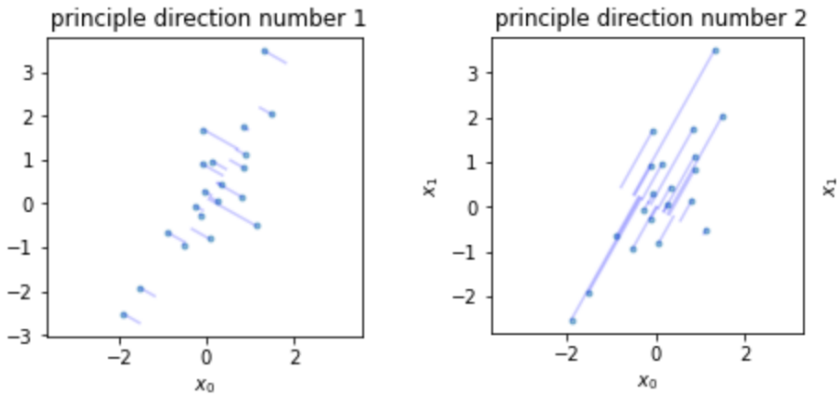
second eigenvector gives the $q_2 = Q_{:,2}$ second principal direction

⋮

so for PCA we need to find the eigenvectors of the covariance matrix

Reducing dimensionality

projection into the principal direction q_i is given by Xq_i



think of the projection XQ as a change of coordinates

we can use the first D' coordinates $Z = XQ_{:,D'}$

to reduce the dimensionality while capturing a lot of the variance in the data

we can project back into original coordinates using $\tilde{X} = ZQ_{:,D'}^T$
reconstruction

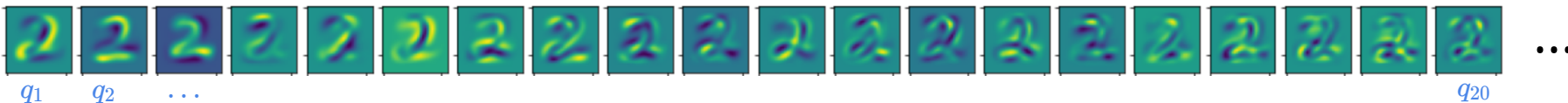
Example: digits dataset

let's only work with digit 2! $x^{(n)} \in \mathbb{R}^{784}$



center the data and form the covariance matrix $\Sigma_{784 \times 784}$

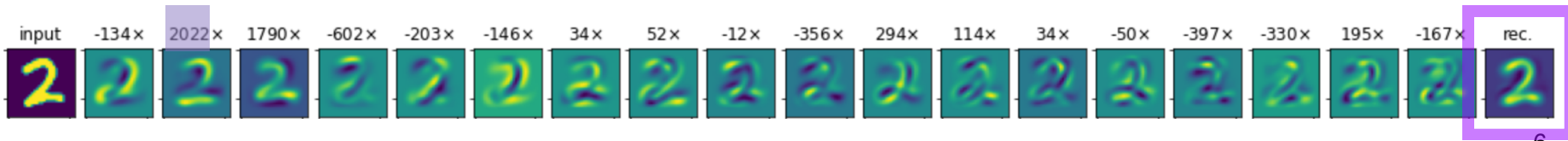
find the **eigenvectors** of the covariance matrix, the principal directions



use the first 20 directions to reduce dimensionality from 784 to 20!

PC coefficient $x^T q_i$ (the new coordinates)

using 20 numbers we can represent each image with a good accuracy



example: digits dataset

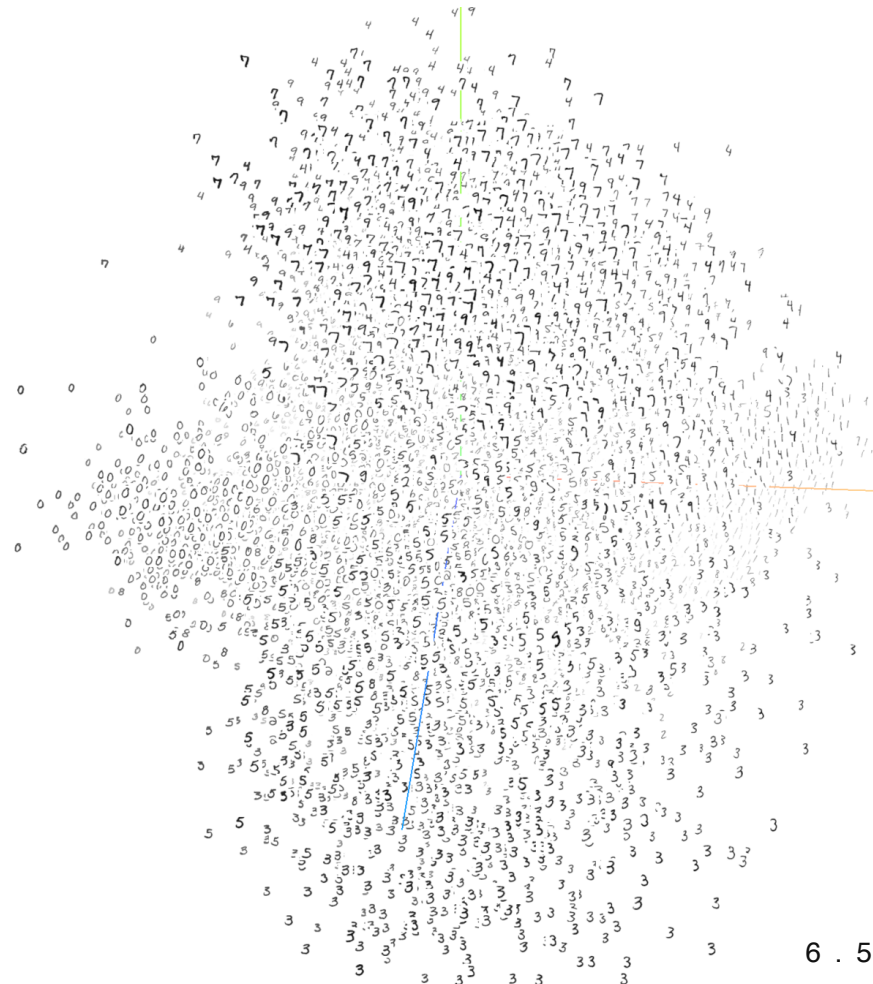
3D embedding of MNIST digits

(<https://projector.tensorflow.org/>)

$$x^{(n)} \in \mathbb{R}^{784}$$

the embedding 3D coordinates are

$$X_{q_1}, X_{q_2}, X_{q_3}$$

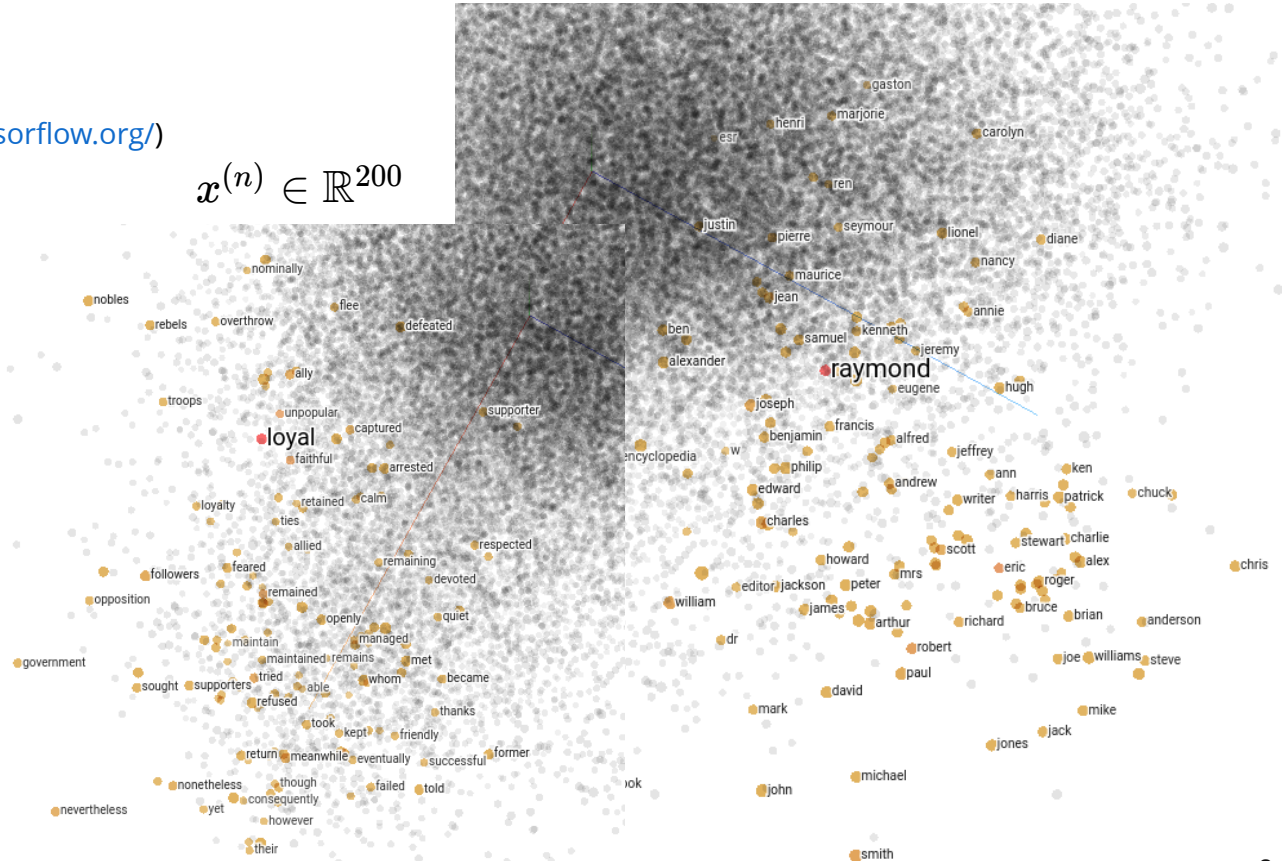


example: text dataset

3D embedding of Word2Vec embeddings (<https://projector.tensorflow.org/>)

$$x^{(n)} \in \mathbb{R}^{200}$$

it is common to use dimensionality reduction to visualize and inspect results of other representation learning methods



example: face dataset

eigenfaces for face recognition
 read more [here](#)

$$x^{(n)} \in \mathbb{R}^{64 \times 64}$$

$$x^{(n)}$$

$$z^{(n)} = x^{(n)\top} Q_{:, :250}$$



Figure #9: n_components=250

q_1, q_2, \dots, q_{15}



Figure #6: Bunch of ghost shaped images. Look at them in the eyes.

mean face used
 for centring the
 data



Figure #5: mean face

there is another way to do PCA

without using the covariance matrix

Singular Value Decomposition (SVD)

any $N \times D$ real matrix has the following decomposition

$$\underset{N \times D}{\mathbf{X}} = \underset{N \times N}{\mathbf{U}} \underset{N \times D}{\mathbf{S}} \underset{D \times D}{\mathbf{V}}^T$$

orthogonal

$$\begin{bmatrix} | & & | \\ u_1 & \dots & u_N \\ | & & | \end{bmatrix}$$

$u_i^\top u_j = 0 \forall i \neq j$
 $\{u_i\}$ left singular vectors

rectangular diagonal

$$\begin{bmatrix} s_1 & & \\ & s_2 & \\ & & \ddots \end{bmatrix}$$

$s_i \geq 0$
 singular values

orthogonal

$$\begin{bmatrix} | & \dots & | \\ v_1 & \dots & v_N \\ | & \dots & | \end{bmatrix}^T$$

$v_i^\top v_j = 0 \forall i \neq j$
 right singular vectors

compressed SVD

assuming $N > D$ we can ignore

- the last $(N-D)$ columns of \mathbf{U} why?
- last $(N-D)$ rows of \mathbf{S}

similarly if $D > N$ we can compress \mathbf{V}, \mathbf{S}

$$\underset{N \times D}{\mathbf{X}} = \underset{N \times D}{\mathbf{U}} \underset{D \times D}{\mathbf{S}} \underset{D \times D}{\mathbf{V}}^T$$

Singular value & eigenvalue decomposition

recall that for PCA we used the eigenvalue decomposition of $\Sigma = \frac{1}{N} X^\top X$

how does it relate to SVD?

$$X^\top X = (USV^\top)^\top (USV^\top) = VS^\top U^\top USV^\top = VS^2V^\top$$

compare to $\frac{1}{N} X^\top X = Q\Lambda Q^\top$

$$(X^\top X)^{-1} = VS^{-2}V^\top$$



eigenvectors of Σ are **right singular vectors** of X $Q = V$

for PCA we could use SVD

- this is the standard computation which works directly with data matrix instead of the covariance matrix

Picking the number of PCs

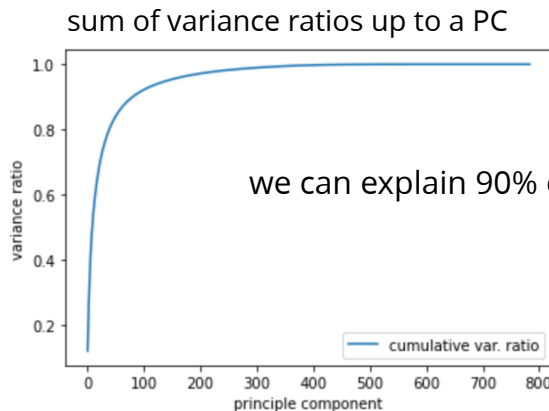
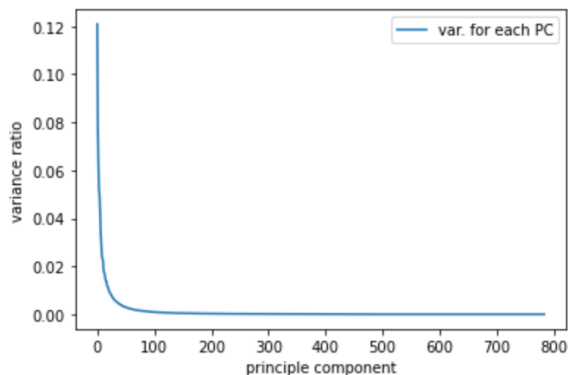
number of PCs in PCA is a hyper-parameter, how should we choose this?

each new principle direction explains some variance in the data $a_d = \frac{1}{N} \sum_n z_d^{(n)2}$

such that we have $a_1 \geq a_2 \geq \dots \geq a_D$ (by definition of PCA)

we can divide by total variance to get a ratio $r_i = \frac{a_i}{\sum_d a_d}$

example for our digits example we get



we can explain 90% of variance in the data using 100 PCs

first few principal directions explain most of the variance in the data!

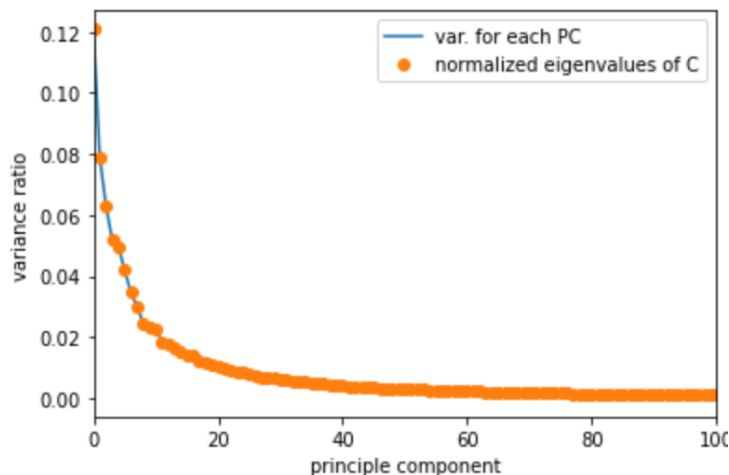
Picking the number of PCs

recall that for picking the principal direction we maximized the variance of the PC

$$\max_{\substack{q \\ \|q\|=1}} \frac{1}{N} q X^\top X q^\top = \max_{\substack{q \\ \|q\|=1}} q \Sigma q^\top = \max_{\substack{q_1 \\ \|q\|=1}} q^\top Q \Lambda Q^\top q = \lambda_1$$

so the variance ratios are also given by $r_i = \frac{\lambda_i}{\sum_d \lambda_d}$

so we can also use eigenvalues to pick the number of PCs

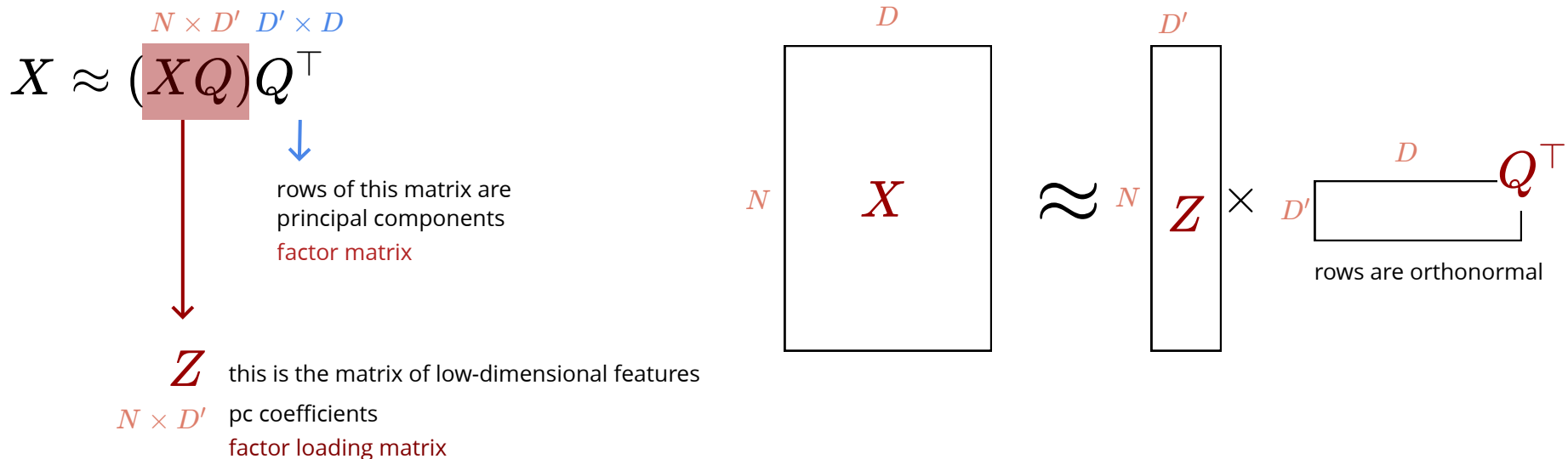


digits **example**:

two estimates of variance ratios do match

Matrix factorization

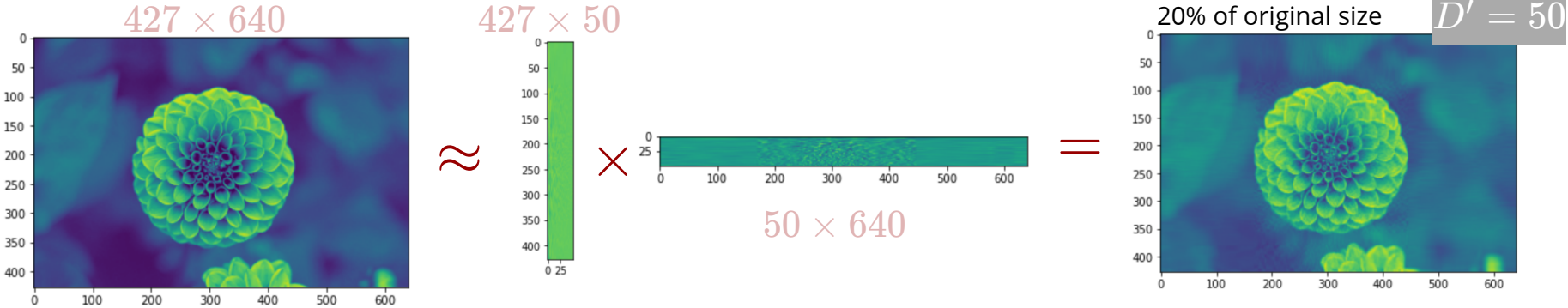
PCA and SVD perform **matrix factorization**



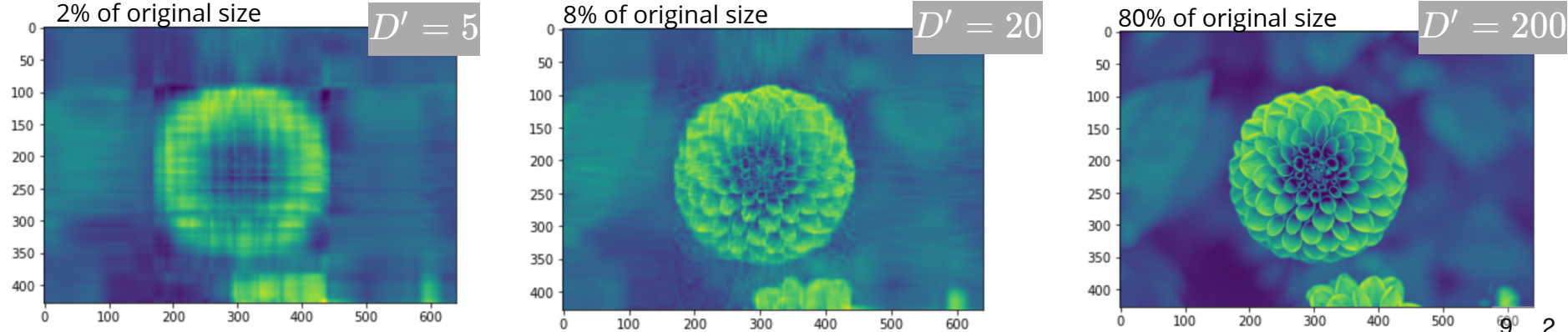
this gives a **row-rank approximation** to our original matrix X

- we can use this to compress the matrix
- we can find give a "smooth" reconstruction of X (remove noise or fill missing values)

Matrix factorization



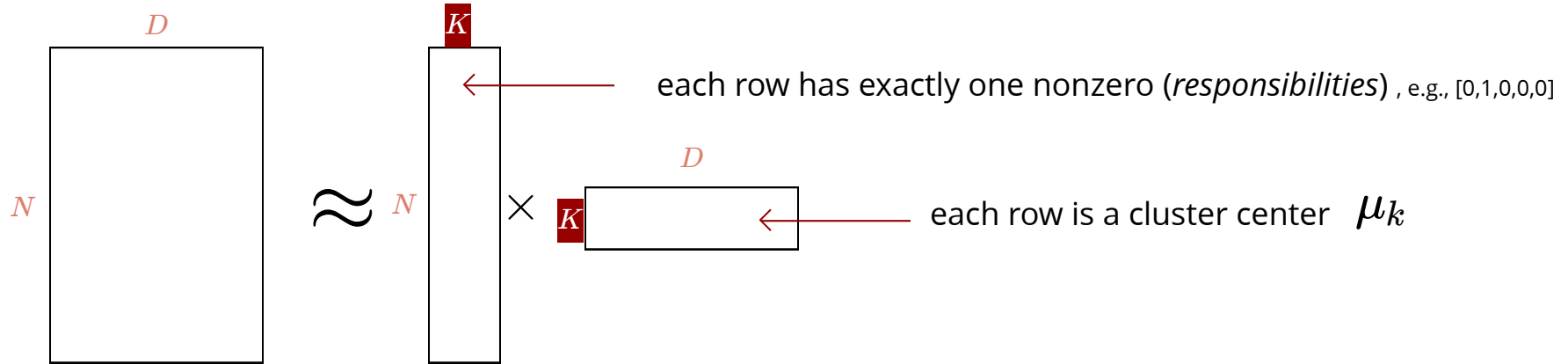
changing the rank D' gives different amount of compression



Matrix factorization

relationship to **K-means**

K-means also can be seen as matrix factorization



matrix product simply equates each row of X with one row of the factor matrix

- instead of principal components \rightarrow cluster centers
- factor loading matrix \rightarrow one nonzero per row of Z (each node belongs to one cluster)

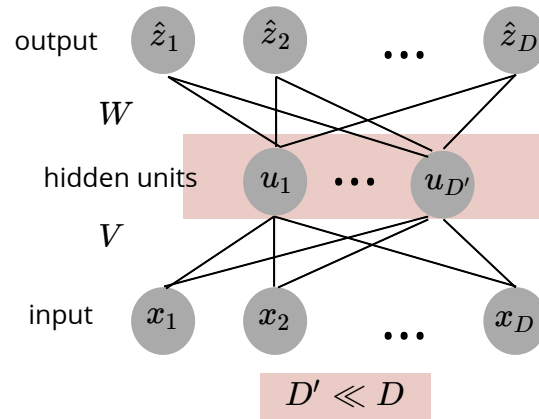
Autoencoders

a feed-forward neural net which predicts its input

- can be trained with reconstruction loss
 - e.g. mean squared error: $\sum_n \|\mathbf{x}^{(n)} - \mathbf{z}^{(n)}\|_2^2$

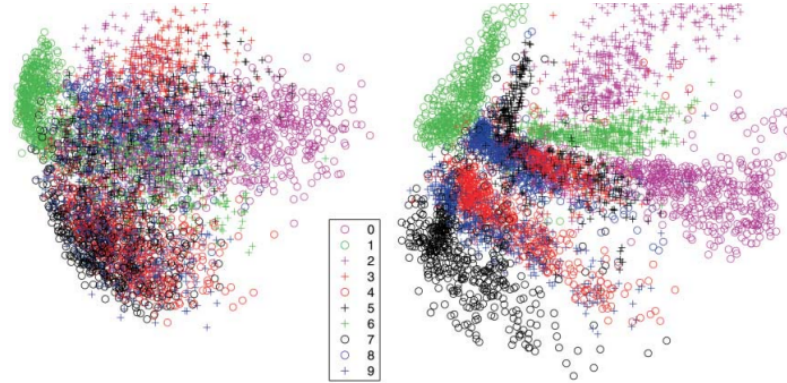
dimensionality reduction with **a bottleneck layer**
much smaller than input

- optimal weights for linear autoencoder are the principal components
- we can do nonlinear dimensionality reduction if activations are not all linear
 - projecting the data on a non-linear manifold
 - deep autoencoders are very powerful



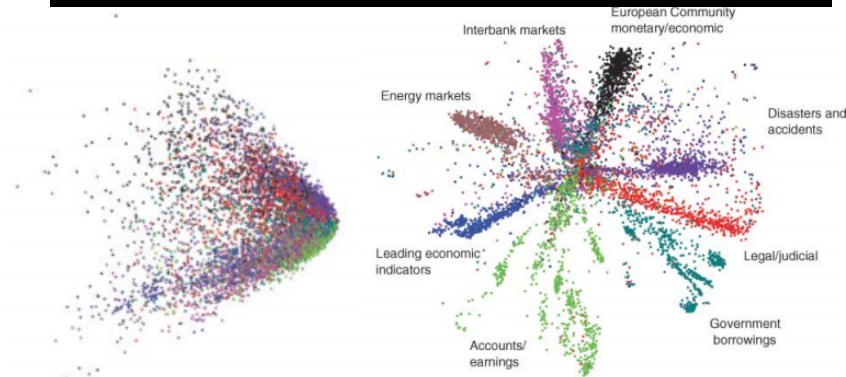
Autoencoders: **example**

MNIST digits



PCA v.s. Autoencoder

newswire stories



read the paper [here](#)

Summary

Dimensionality reduction helps us:

- visualize our data
- compress it
- simplify the computational need of further analysis (clustering, supervised learning etc.)
- also can be used for anomaly detection (not discussed)

PCA is a linear dimensionality reduction method

- projects the data to a linear space (spanned by D' principal directions)
 - directions are eigenvectors of the covariance matrix
 - the projection has maximum variance (minimum reconstruction error)
 - eigenvalues tell us about the contribution of each new principal direction
- PCA using Singular Value Decomposition
- Model selection for PCA
- PCA as matrix factorization and its relationship to k-means
- practical note: don't forget to subtract the mean!