

CS-GY 6923: Lecture 12

Autoencoders, Principal Component Analysis

NYU Tandon School of Engineering, Prof. Christopher Musco

Optional. Can be completed in place of last assignment, which will contain a mix of written problems and coding problems.

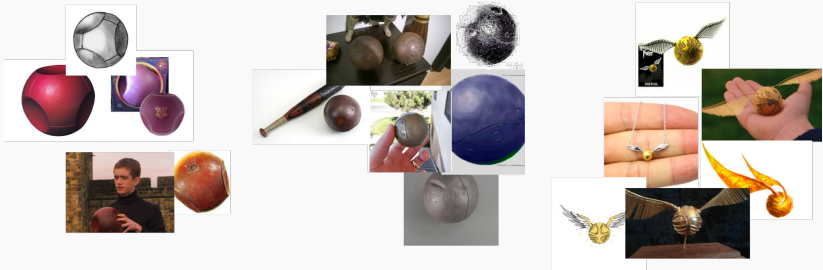
1. Concept Based Project
2. Data set Based Project

Please email me if you want to complete a project and who your team will be (1-3 students allowed). We can also discuss topic ideas.

State-of-the-art supervised learning models like neural networks learn **very good features**.

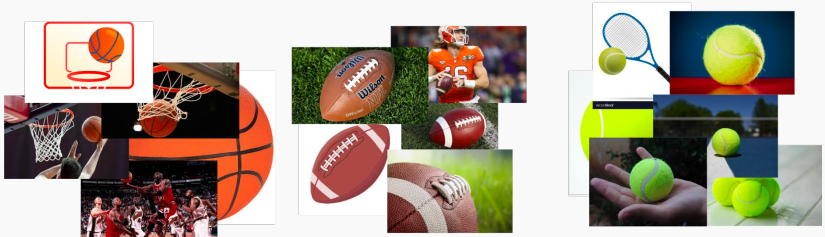
Often these features are useful for tasks other than what the model was trained on. They can be used for other problems where we do not have a lot of training data.

Example: Classify images of different Quidditch balls.



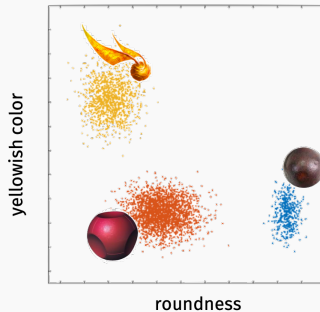
TRANSFER LEARNING

Based on features learned for muggle objects:



FEATURE LEARNING

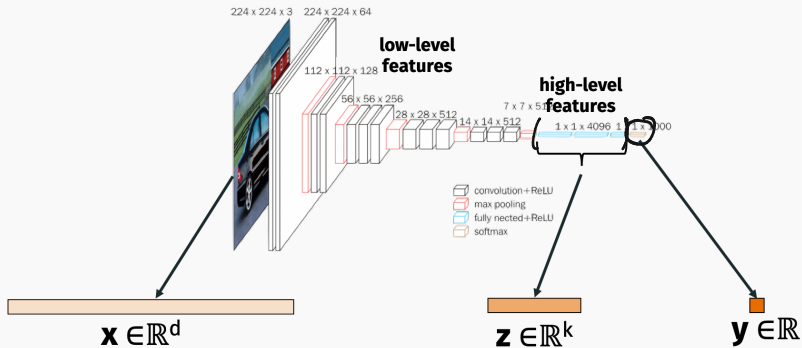
If these features are highly informative (i.e. lead to highly separable data) few training examples are needed to learn.



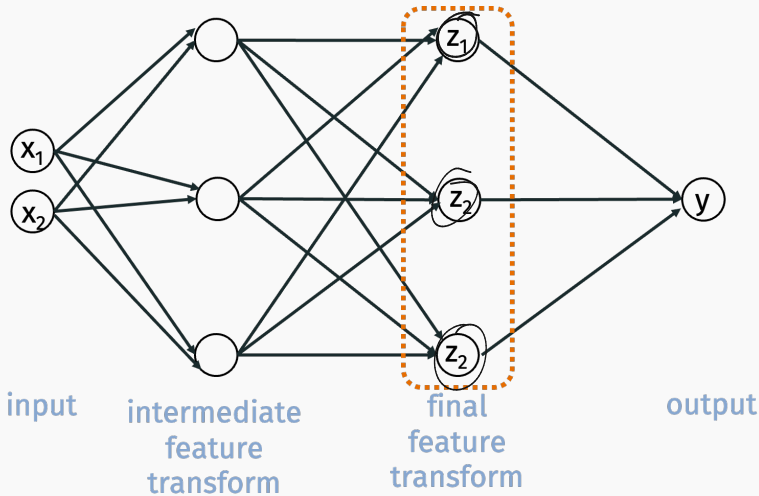
Common technique:

1. Download network trained e.g. on a large image classification dataset (e.g. Imagenet).
2. Extract features \mathbf{z} for any new image \mathbf{x} by running it through the network up until layer before last.
3. Use these features in a simpler machine learning algorithm that requires less data (nearest neighbor, logistic regression, etc.).

TRANSFER LEARNING



TRANSFER LEARNING



Transfer learning: Lots of labeled data for one problem makes up for little labeled data for another.

But what if we don't even have labeled data for a sufficiently related problem?

How to extract features in a data-driven way from unlabeled data is one of the central problems in unsupervised learning.

Simple but clever idea: If we have inputs $\underline{x}_1, \dots, \underline{x}_n \in \mathbb{R}^d$ but few or no targets $\underline{y}_1, \dots, \underline{y}_n$, just make the inputs the targets.

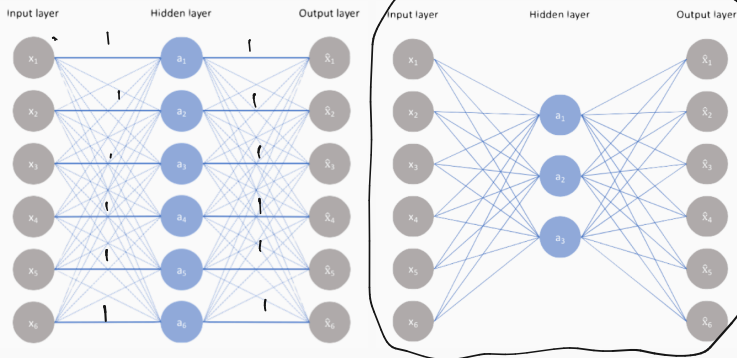
- Let $\underline{f}_\theta: \mathbb{R}^d \rightarrow \mathbb{R}^d$ be our model.
- Let L_θ be a loss function. E.g. squared loss:

$$L_\theta(\underline{x}) = \|\underline{x} - \underline{f}_\theta(\underline{x})\|_2^2.$$
- Train model: $\underline{\theta}^* = \min_{\theta} \sum_{i=1}^n L_\theta(\underline{x}_i).$

If \underline{f}_θ is a model that incorporates feature learning, then these features can be used for supervised tasks.

\underline{f}_θ is called an **autoencoder**. It maps input space to input space (e.g. images to images, french to french, PDE solutions to PDE solutions).

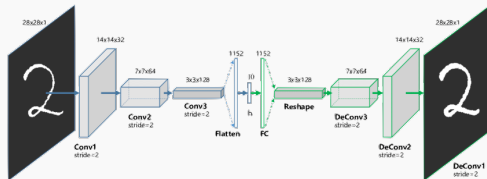
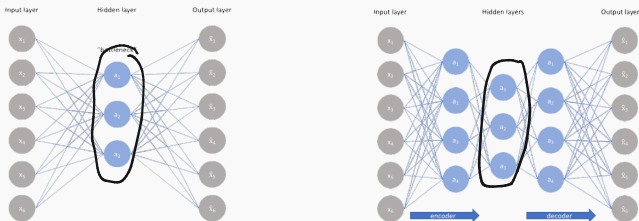
Two examples of autoencoder architectures:



Which would lead to better feature learning?

AUTOENCODER

Important property of autoencoders: no matter what architecture is use, there must always be a **bottleneck** with fewer parameters than the input. The bottleneck ensures information is “distilled” from low-level features to high-level features.



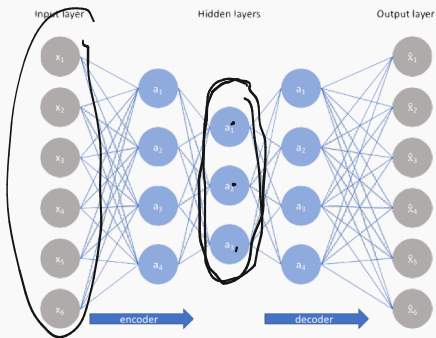
AUTOENCODER

Architecture typically split into two parts:

Encoder: $e : \mathbb{R}^d \rightarrow \mathbb{R}^k$ $\mathbb{R}^6 \rightarrow \mathbb{R}^3$

Decoder: $d : \mathbb{R}^k \rightarrow \mathbb{R}^d$ $\mathbb{R}^3 \rightarrow \mathbb{R}^6$

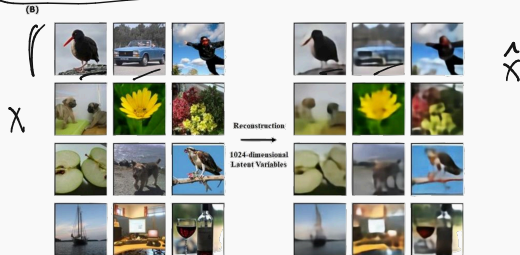
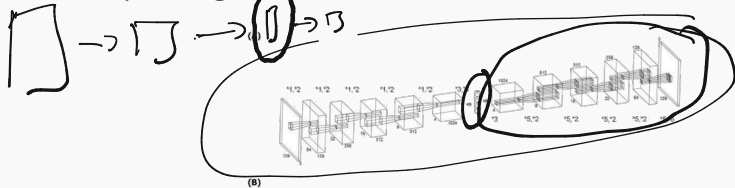
$$\underline{f(x)} = d(e(x))$$



Often symmetric, but does not have to be.

AUTOENCODER RECONSTRUCTION

Example image reconstructions from autoencoder:



<https://www.biorxiv.org/content/10.1101/214247v1.full.pdf>

Input parameters: $d = \underline{49152}$.

Bottleneck "latent" parameters: $k = \underline{1024}$.

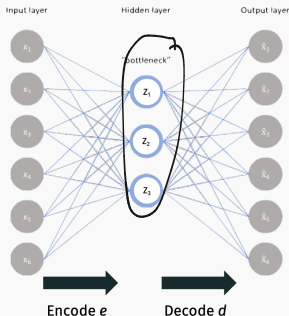
The best autoencoders do not work as well as supervised methods for feature extraction, but they require less labeled data. Recent progress on **self-supervised** learning gets closer to supervised feature learning. Might discuss a bit later.

There are a lot of cool applications of autoencoders beyond feature learning!

- Learned data compression.
- Denoising and in-painting.
- Data/image synthesis.

AUTOENCODERS FOR IMAGE COMPRESSION

Due to their bottleneck design, autoencoders perform **dimensionality reduction** and thus data compression.



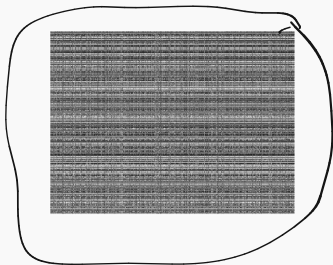
Given input \mathbf{x} , we can completely recover $f(\mathbf{x})$ from $\mathbf{z} = e(\mathbf{x})$. \mathbf{z} typically has many fewer dimensions than \mathbf{x} and for a typical image $f(\mathbf{x})$ will closely approximate \mathbf{x} .

AUTOENCODERS FOR IMAGE COMPRESSION

The best lossy compression algorithms are tailor made for specific types of data:

- JPEG 2000 for images
- MP3 for digital audio.
- MPEG-4 for video.

All of these algorithms take advantage of specific structure in these data sets. E.g. JPEG assumes images are locally “smooth”.



AUTOENCODERS FOR IMAGE COMPRESSION

With enough input data, autoencoders can be trained to find this structure on their own.



“End-to-end optimized image compression”, Ballé, Laparra, Simoncelli

Need to be careful about how you choose loss function, design the network, etc. but can lead to much better image compression than “hand-tuned” algorithms like JPEG.

AUTOENCODERS FOR IMAGE RESTORATION

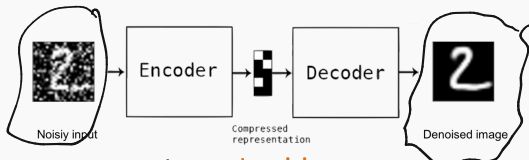


Image denoising



Image inpainting

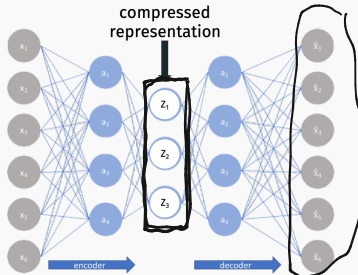
Train autoencoder on uncorrupted images (unsupervised). Pass corrupted image x through autoencoder and return $f(x)$ as repaired result.

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS

$$(256)^{\frac{1}{(128 \times 128 \times 3)}}$$

Why does this work?

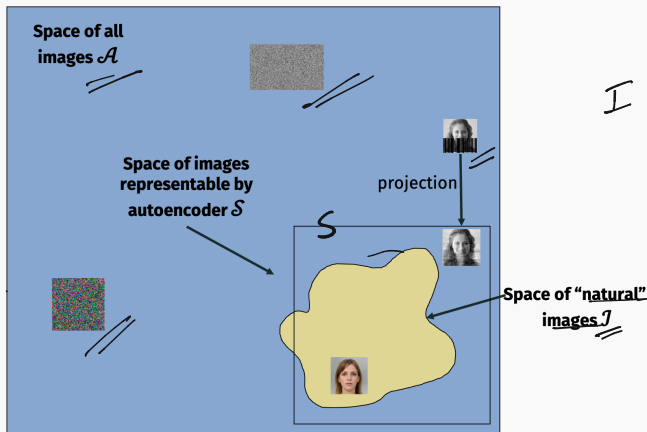
$$(10)^{\frac{1}{13}}$$



Consider 128 \times 128 \times 3 images with pixels values in $0, 1, \dots, 255$.
How many possible images are there?

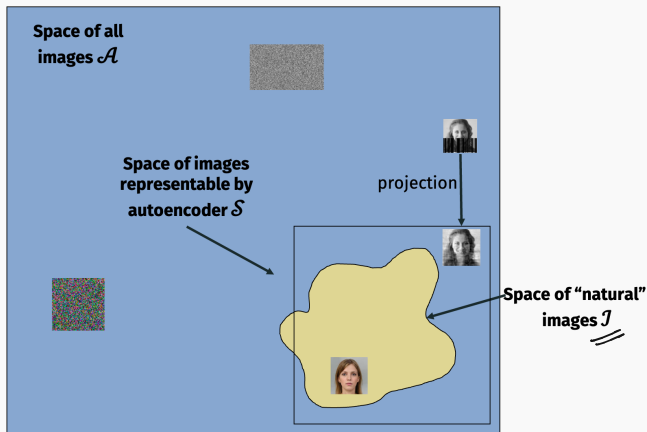
If \mathbf{z} holds k values between in $0, .1, .2, \dots, .1$, how many unique images \mathbf{w} can be output by the autoencoder function f ?

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



For a good (accurate, small bottleneck) autoencoder, \mathcal{S} will closely approximate \mathcal{I} . Both will be much smaller than \mathcal{A} .

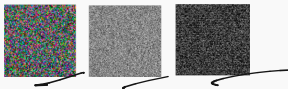
AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



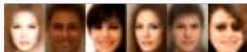
$f(\mathbf{x}) = d(e(\mathbf{x}))$ projects an image \mathbf{x} closer to the space of natural images.

Suppose we want to generate a random natural image. How might we do that?

- **Option 1:** Draw each pixel value in x uniformly at random.
Draws a random image from \mathcal{A} .



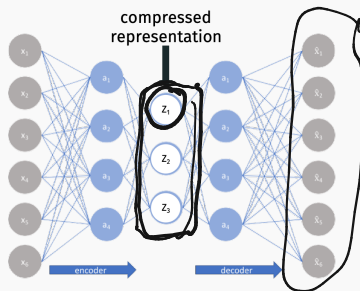
- **Option 2:** Draw x randomly from \mathcal{S} the space of images representable by the autoencoder.



How do we randomly select an image from \mathcal{S} ?

AUTOENCODERS FOR DATA GENERATION

How do we randomly select an image x from S ?



Randomly select code \underline{z} , then set $x = \underline{d(z)}$.¹

¹Lots of details to think about here. In reality, people use “variational autoencoders” (VAEs), which are a natural modification of AEs.

AUTOENCODERS FOR DATA GENERATION



Generative models are a growing area of machine learning, drive by a lot of interesting new ideas. (Generative Adversarial Networks) in particular are now a major competitor with variational autoencoders.

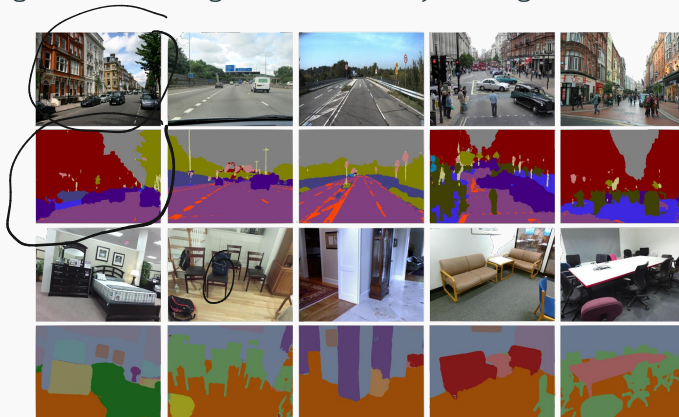
An autoencoder is a model $f: \underline{\mathbb{R}^d} \rightarrow \underline{\mathbb{R}^d}$. In other words, the output is the same dimension as the input:

- Image \rightarrow Image
- Video \rightarrow Video
- Audio clip \rightarrow Audio clip

This structure is also useful for some supervised machine learning problems.

IMAGE SEGMENTATION

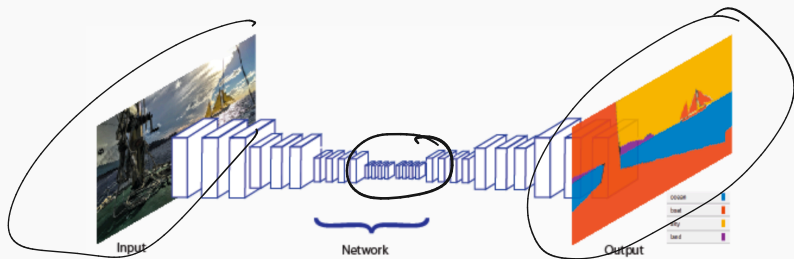
Goal: Learn mask which separates image pixels by what object (foreground or background) that they belong to.



First step in multi-objects classification and scene understanding. Harder than classifying single objects.

END-TO-END IMAGE SEGMENTATION

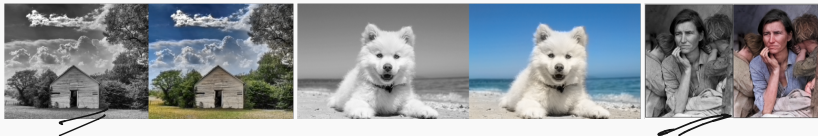
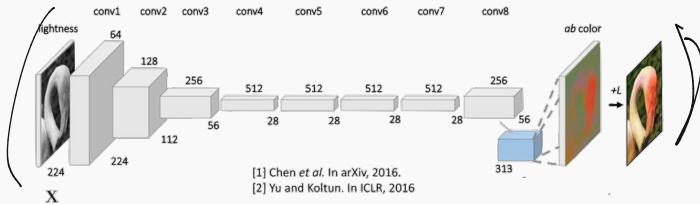
Model: Input is image x , output is image m that has the same size as x , but each pixel value is a label for a segmented region.



Now our training process is actually supervised, but uses the same structure as an autoencoder.

END-TO-END IMAGE COLORIZATION

Model: Input is black and white image x , output is colorized image m .



END-TO-END SUPER RESOLUTION

Model: Input is pixelated or blurred image x , output is full-resolution image m .



PRINCIPAL COMPONENT ANALYSIS

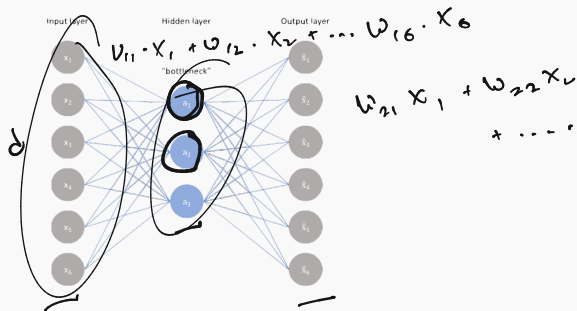
Rest of lecture: Deeper dive into understanding a simple, but powerful autoencoder architecture. Specifically we will view **principal component analysis (PCA)** as a type of autoencoder.

PCA is the “linear regression” of unsupervised learning: often the go-to baseline method for feature extraction and dimensionality reduction.

Very important outside machine learning as well.

PRINCIPAL COMPONENT ANALYSIS

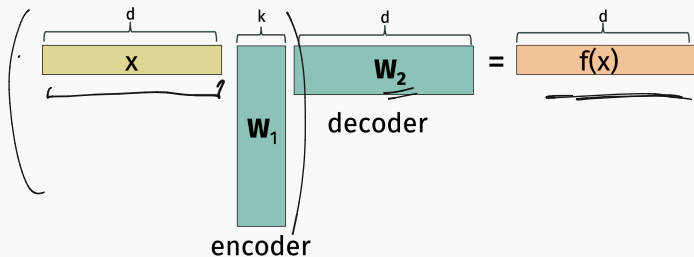
Consider the simplest possible autoencoder:



- One hidden layer. No non-linearity. No biases.
- Latent space of dimension k .
- Weight matrices are $\underline{W}_1 \in \mathbb{R}^{d \times k}$ and $\underline{W}_2 \in \mathbb{R}^{k \times d}$.

PRINCIPAL COMPONENT ANALYSIS

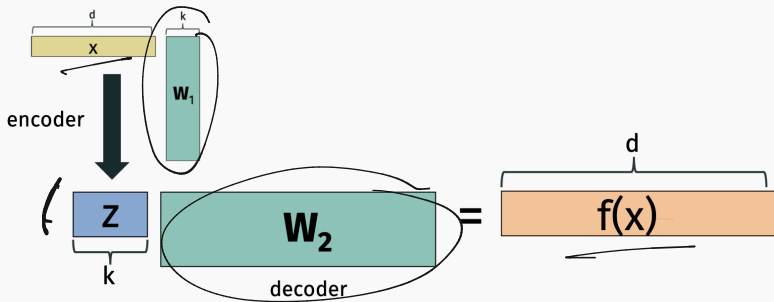
Given input $x \in \mathbb{R}^d$, what is $f(x)$ expressed in linear algebraic terms?



$$(1 \times d) (d \times k) \rightarrow (1 \times k) \quad (1 \times k) (k \times d) \rightarrow (1 \times d)$$

$$f(x) = x^T W_1 W_2$$

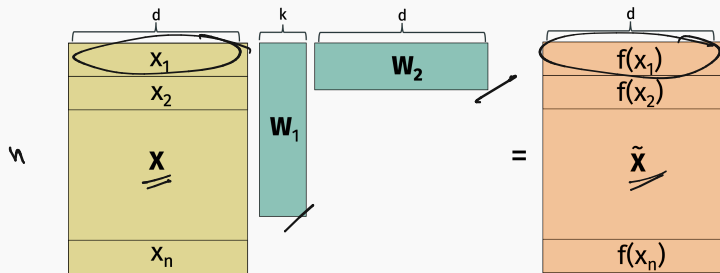
PRINCIPAL COMPONENT ANALYSIS



$$\text{Encoder: } \underline{z} = \underline{x}W_1 \quad \text{Decoder: } \underline{d(z)} = \underline{z}W_2$$

PRINCIPAL COMPONENT ANALYSIS

Given training data set x_1, \dots, x_n , let X denote our data matrix.
Let $\tilde{X} = XW_1W_2$.



$$\tilde{X} \approx X$$

Natural squared autoencoder loss: Minimize $L(\mathbf{X}, \tilde{\mathbf{X}})$ where:

$\mathbf{X} - \tilde{\mathbf{X}}$

$$\begin{aligned} \underline{L(\mathbf{X}, \tilde{\mathbf{X}})} &= \sum_{i=1}^n \|\underline{\mathbf{x}_i} - \underline{f(\mathbf{x}_i)}\|_2^2 \\ &= \sum_{i=1}^n \sum_{j=1}^d (\mathbf{x}_i[j] - \underline{f(\mathbf{x}_i)[j]})^2 \\ &= \underline{\|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2} \end{aligned}$$

$\|\mathbf{M}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^d M_{i,j}^2$

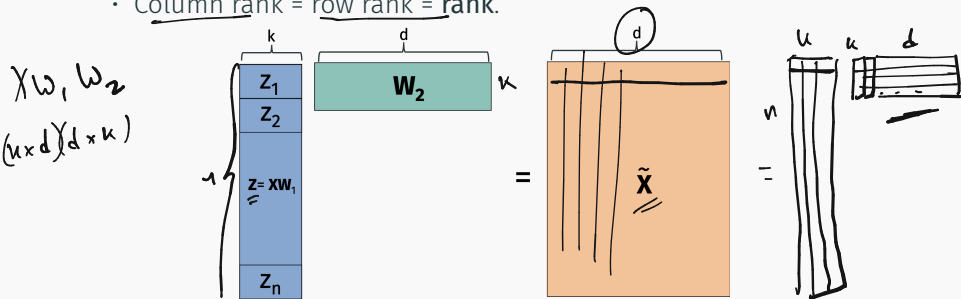
Goal: Find $\underline{\mathbf{W}_1}, \underline{\mathbf{W}_2}$ to minimize the Frobenius norm loss

$\underline{\|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2} = \underline{\|\mathbf{X} - \mathbf{X}\mathbf{W}_1\mathbf{W}_2\|_F^2}$ (sum of squared entries).

LOW-RANK APPROXIMATION

Rank in linear algebra:

- The columns of a matrix with column rank k can all be written as linear combinations of just k columns.
- The rows of a matrix with row rank k can all be written as linear combinations of k rows.
- Column rank = row rank = rank.



\tilde{X} is a **low-rank matrix**. It only has rank k for $k \ll d$.

LOW-RANK APPROXIMATION



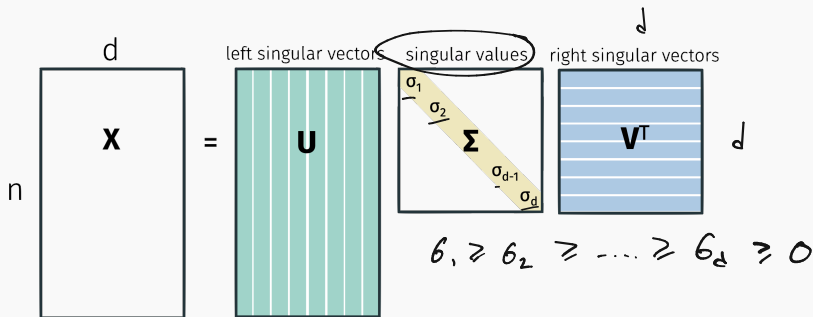
Principal component analysis is the task of finding \mathbf{W}_1 , \mathbf{W}_2 , which amounts to finding a rank k matrix $\tilde{\mathbf{X}}$ which approximates the data matrix \mathbf{X} as closely as possible.

Finding the best \mathbf{W}_1 and \mathbf{W}_2 is a non-convex problem. We could try running an iterative method like gradient descent anyway. But there is also a direct algorithm!

SINGULAR VALUE DECOMPOSITION

Any matrix X can be written:

$$X = U \Sigma V^T$$



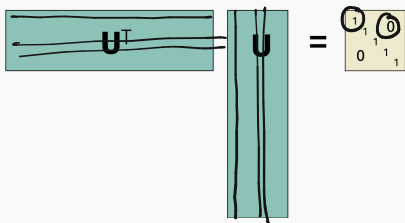
Where $\underline{U^T U = I}$, $\underline{V^T V = I}$, and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$. I.e. U and V are orthogonal matrices.

This is called the **singular value decomposition**.

Can be computed in $O(nd^2)$ time (faster with approximation algos).

ORTHOGONAL MATRICES

Let $\mathbf{u}_1, \dots, \mathbf{u}_n \in \mathbb{R}^n$ denote the columns of \mathbf{U} . I.e. the left singular vectors of \mathbf{X} .



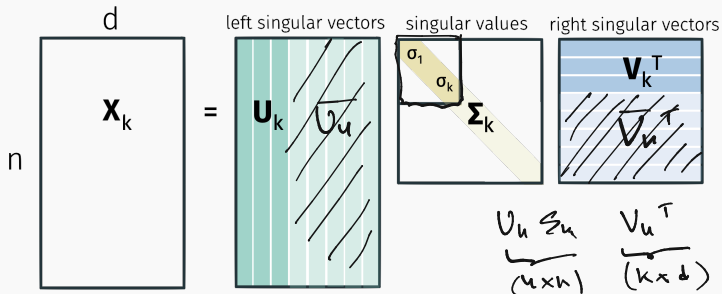
$$\|u_i\|_2^2 = 1$$

$$u_i^T u_j = 0$$

$$v_i^T v_i = \|v_i\|_2^2 = 1$$

SINGULAR VALUE DECOMPOSITION

Can read off optimal low-rank approximations from the SVD:



Eckart-Young-Mirsky Theorem: For any $k \leq d$, $(X_k = U_k \Sigma_k V_k^T)$ is the optimal k rank approximation to X :

$$\underline{X_k} = \arg \min_{\tilde{X} \text{ with rank } \leq k} \underline{\|X - \tilde{X}\|_F^2}$$

SINGULAR VALUE DECOMPOSITION

Claim: $X_k = U_k \Sigma_k V_k^T = X V_k V_k^T$

$$X V_k = \begin{bmatrix} U_k & \text{---} \\ \hline \text{---} & \text{---} \end{bmatrix} \begin{bmatrix} \Sigma_k & \text{---} \\ \hline \text{---} & \text{---} \end{bmatrix} \begin{bmatrix} V_k^T & \text{---} \\ \hline \text{---} & \text{---} \end{bmatrix} \begin{bmatrix} V_k & \text{---} \\ \hline \text{---} & \text{---} \end{bmatrix}$$

$$V^T V = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 0 & \text{---} \\ & & & & \ddots & \\ & & & & & 1 & \\ & & & & & & 0 & \text{---} \\ & & & & & & & \ddots & \end{bmatrix}$$

$$X V_k = U_k \Sigma_k$$

$$X U_k V_k^T = U_k \Sigma_k V_k^T$$

$$U \leq V^T V u$$

$$= \begin{bmatrix} \sigma_1 & \dots & \sigma_u \\ \hline 0 \end{bmatrix}$$

$$U_k \Sigma_k$$

$$\begin{bmatrix} 1 & & \\ \hline & \ddots & \\ & & 1 \\ \hline & & & 0 \end{bmatrix} = V^T V_k$$

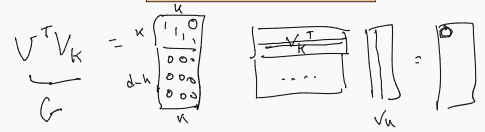
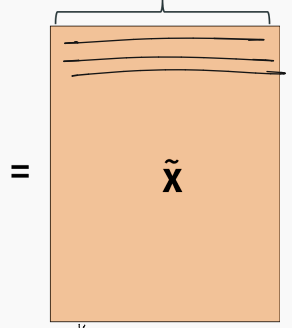
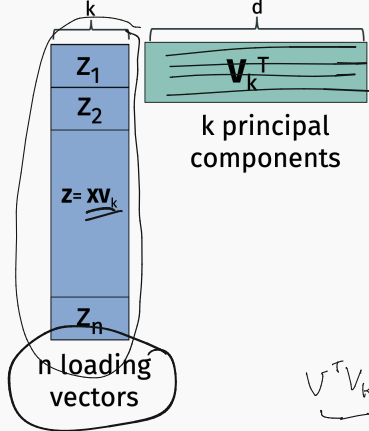
$X W_1 W_2$ is the best rank- k approximation to X .

So for a model with k hidden variables, we obtain an optimal autoencoder by setting $\underline{W}_1 = \underline{V}_k$, $\underline{W}_2 = \underline{V}_k^T$. $f(x) = x V_k V_k^T$.

PRINCIPAL COMPONENT ANALYSIS

$$U_k \Sigma_k = \underbrace{X}_{k} U_k = U \Sigma V^T V_k$$

$$U \Sigma \underbrace{G}_{d} = U_k \Sigma_k$$



Usually x 's columns (features) are mean centered and normalized to variance 1 before computing principal components.

Computing the SVD.

- Full SVD:

`U, S, V = scipy.linalg.svd(X).`

Runs in $O(nd^2)$ time.

- Just the top k components:

`U, S, V = scipy.sparse.linalg.svd(X, k).`

Runs in roughly $O(ndk)$ time.

CONNECTION TO EIGENDECOMPOSITION

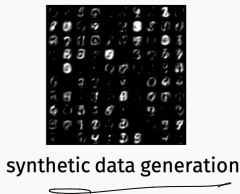
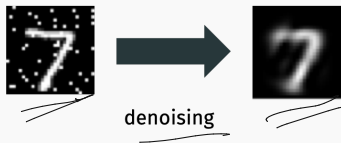
Recall that for a matrix $\mathbf{M} \in \mathbb{R}^{p \times p}$, \mathbf{q} is an eigenvector of \mathbf{M} if $\lambda \mathbf{q} = \mathbf{Mq}$ for any scalar λ .

- \mathbf{U} 's columns (the left singular vectors) are the orthonormal eigenvectors of \mathbf{XX}^T . $n \times n$ $n \times n$
- \mathbf{V} 's columns (the right singular vectors) are the orthonormal eigenvectors of $\mathbf{X}^T\mathbf{X}$. $d \times d$ $n \times n$
- $\sigma_i^2 = \lambda_i(\mathbf{XX}^T) = \lambda_i(\mathbf{X}^T\mathbf{X})$

Exercise: Verify this directly. This means you can use any (symmetric) eigensolver for computing the SVD.

Like any autoencoder, PCA can be used for:

- Feature extraction
- Denoising and rectification
- Data generation
- Compression
- Visualization



LOW-RANK APPROXIMATION

The larger we set k , the better approximation we get.

$$d = 28 \times 28 \\ = 784$$

7	2	1	0	4	1	9	9	0	9
0	6	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

original data

$$u = 10$$



rank 1 approx.

8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8	8	8

rank 2 approx.

0	0	7	0	1	0	8	0	0	9
0	0	9	0	1	0	8	0	0	9
7	0	0	0	1	0	8	0	0	9
7	1	9	0	7	0	9	1	0	1
7	1	7	3	3	3	9	9	9	9
8	8	8	8	0	0	1	8	1	8
8	8	8	8	7	0	8	3	0	8
8	0	0	7	1	7	3	9	1	9
8	0	0	7	1	7	3	9	1	9
8	8	7	7	8	1	6	0	1	8

rank 3 approx.

9	3	7	0	7	1	9	9	0	9
0	0	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 4 approx.

9	3	7	0	7	1	9	9	0	9
0	0	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 5 approx.

9	3	7	0	7	1	9	9	0	9
0	0	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 6 approx.

7	2	1	0	4	1	9	9	0	9
0	6	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 7 approx.

7	2	1	0	4	1	9	9	0	9
0	6	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 8 approx.

7	2	1	0	4	1	9	9	0	9
0	6	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 9 approx.

7	2	1	0	4	1	9	9	0	9
0	6	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

rank 50 approx.

7	2	1	0	4	1	9	9	0	9
0	6	9	0	1	5	9	7	8	4
7	6	6	5	4	0	7	4	0	1
3	1	3	4	7	2	7	1	2	1
1	7	4	2	3	5	1	2	4	4
6	3	5	5	6	0	4	1	9	5
7	8	5	3	7	4	6	4	3	0
7	0	2	7	1	7	3	2	9	7
1	6	2	7	8	4	7	3	6	1
3	6	2	3	1	4	1	7	6	9

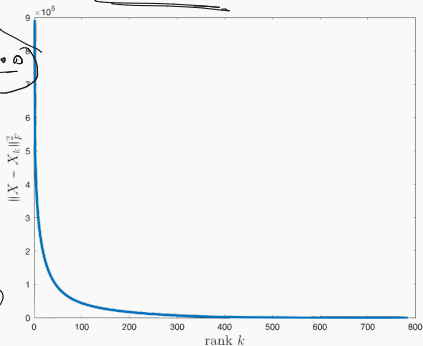
LOW RANK APPROXIMATION

Error vs. k is dictated by X 's singular values. The singular values are often called the **spectrum** of X .

$\sigma_1, \dots, \sigma_d$

$$\|X - X_k\|_F^2 = \sum_{i=k+1}^d \sigma_i^2$$

$$\|X - X U_k V_k^T\|_F^2$$



LOW RANK INTUITION

sale price \approx list price

Which of these data sets would you expect to have a good low-rank approximation? Why?

\bullet 0.05 \bullet list price

1. House data:

~~1~~ ~~3~~ 2

$\mathbf{x} = [\# \text{ bedrooms}, \# \text{ bathrooms}, \text{list price}, \text{sale price}, \text{property tax}]$

2. Student data:

$\mathbf{x} = [\text{gender}, \text{year}, \text{age}, \text{GPA}, \text{engineering major}]$

list price \approx $c + a(\# \text{ bedrooms}) + b(\# \text{ bathrooms})$

COLUMN REDUNDANCY

Colinearity of data features leads to an approximately low-rank data matrix.

	bedrooms	bathrooms	sq.ft.	floors	list price	sale price
home 1	2	2	1800	2	200,000	195,000
home 2	4	2.5	2700	1	300,000	310,000
.
.
.
home n	5	3.5	3600	3	450,000	450,000

sale price $\approx 1.05 \cdot$ list price.

property tax $\approx .01 \cdot$ list price.

COLUMN REDUNDANCY

Sometimes these relationships are simple, other times more complex. But as long as there exists linear relationships between features, we will have a lower rank matrix.

$$\underline{\text{yard size}} \approx \underline{\text{lot size}} - \cancel{\text{1}} \cdot \underline{\text{square footage}}.$$

$$\underline{\text{cumulative GPA}} \approx \frac{1}{4} \cdot \text{year 1 GPA} + \frac{1}{4} \cdot \text{year 2 GPA} \\ + \frac{1}{4} \cdot \text{year 3 GPA} + \frac{1}{4} \cdot \text{year 4 GPA}.$$

LOW-RANK INTUITION

Two other examples of data with good low-rank approximations:

1. Genetic data:

single nucleotide polymorphisms (SNPs) loci

	144	312	436	800	943
individual 1	A	T	T	C	G
individual 2	T	G	G	C	C
...					
individual n	C	A	T	A	G

2. “Term-document” matrix with bag-of-words data:

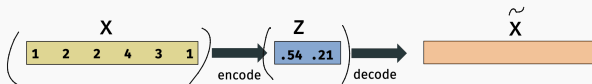
	car	loan	house	...	dog	cat			
doc_1	0	0	1	0	0	1	1	0	0
doc_2	0	0	0	1	0	1	0	0	0
...	1	1	0	1	0	0	0	1	0
...	0	0	0	0	0	0	0	1	1
doc_n	1	0	0	0	0	0	0	1	1

EXAMPLES OF LOW-RANK STRUCTURE

SNPs matrices tend to be very low-rank.

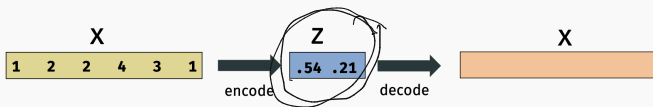
	single nucleotide polymorphisms (SNPs) loci				
	144	312	436	800	943
individual 1	A	T	T	C	G
individual 2	T	G	G	C	C
...					
individual n	C	A	T	A	G

Most of the information in x is explained by just a few **latent variable**.



EXAMPLES OF LOW-RANK STRUCTURE

“Genes Mirror Geography Within Europe” – Nature, 2008.



In data collected from European populations, latent variables capture information about geography.

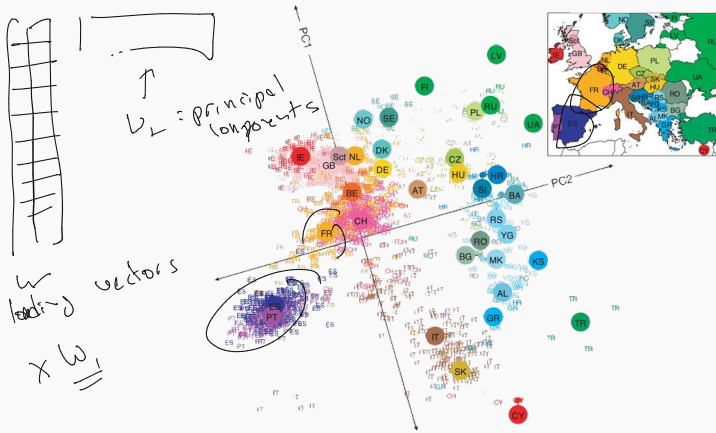
$z[1] \approx$ relative north-south position of birth place

$z[2] \approx$ relative east-west position of birth place

Individuals born in similar places tend to have similar genes.

PCA FOR DATA VISUALIZATION

“Genes Mirror Geography Within Europe” – Nature, 2008.



Genetic data can be nicely visualized using PCA! Plot each data example x using two loading variables in z .

SIMILARITY PRESERVATION

$$\|\tilde{x}_i - \tilde{x}_j\|_2^2 = \|\tilde{x}_i\|_2^2 - 2\langle \tilde{x}_i, \tilde{x}_j \rangle + \|\tilde{x}_j\|_2^2 = \|z_i\|_2^2 - 2\langle z_i, z_j \rangle + \|z_j\|_2^2$$

Important note for data visualization and more: Latent feature vectors preserve similarity and distance information in the original data.

$$\|x_i\|_2^2 \approx \|z_i\|_2^2 \quad \|x_i - x_j\|_2^2 \approx \|z_i - z_j\|_2^2$$

Let $x_1, \dots, x_n \in \mathbb{R}^d$ be our original data vectors, $z_1, \dots, z_n \in \mathbb{R}^k$ be our loading vectors (encoding), and $\tilde{x}_1, \dots, \tilde{x}_n \in \mathbb{R}^d$ be our low-rank approximated data.

$$\|a\|_2^k = a^T a$$

Recall that $\tilde{x}_i = V_k z_i$. Because V_k is orthogonal, we have:

$$X \approx X V_u V_u^T$$

$$f(x_i) = \tilde{x}_i$$

$$(z_i^T V_k^T)^T = (V_k z_i)$$

$$\begin{aligned} \|\tilde{x}_i\|_2^2 &= \|z_i\|_2^2 \\ \langle \tilde{x}_i, \tilde{x}_j \rangle &= \langle z_i, z_j \rangle \\ \|\tilde{x}_i - \tilde{x}_j\|_2^2 &= \|z_i - z_j\|_2^2 \\ \tilde{x}_i^T \tilde{x}_j &= z_i^T V_u^T V_u z_j \\ &= z_i^T z_j \end{aligned}$$

$$\begin{aligned} z_i &= V_k^T x_i \\ \|\tilde{x}_i\|_2^2 &= \|V_k z_i\|_2^2 \\ &= (z_i^T V_k^T) V_k z_i \\ &= z_i^T \underbrace{V_k^T V_k}_{I} z_i \\ &= z_i^T z_i = \|z_i\|_2^2 \end{aligned}$$

SIMILARITY PRESERVATION

$$\|x_i - x_j\|_2^2 \approx \|\tilde{x}_i - \tilde{x}_j\|_2^2$$

Conclusion: If our data had a good low rank approximation, we expect that:

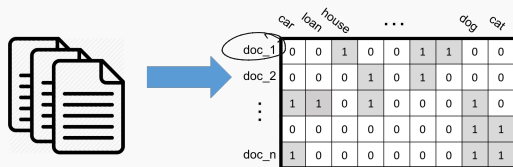
$$\|x_i\|_2^2 \approx \|z_i\|_2^2$$

$$\langle x_i, x_j \rangle \approx \langle z_i, z_j \rangle$$

$$\|x_i - x_j\|_2^2 \approx \|z_i - z_j\|_2^2$$

TERM DOCUMENT MATRIX

Word-document matrices tend to be low rank.



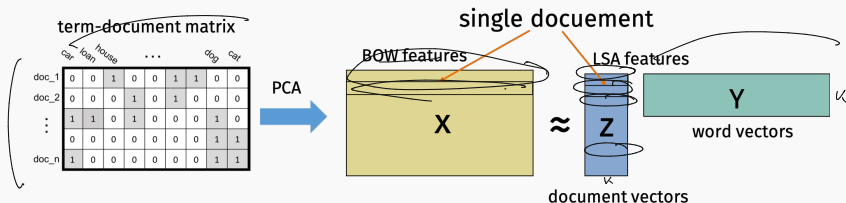
Documents tend to fall into a relatively small number of different categories, which use similar sets of words:

- Financial news: *markets, analysts, dow, rates, stocks*
- US Politics: *president, senate, pass, slams, twitter, media*
- StackOverflow posts: *python, help, convert, javascript*
- Etc.

Intuition that this data should have co-linear **rows**.

LATENT SEMANTIC ANALYSIS

Latent semantic analysis = PCA applied to a word-document matrix (usually from a large corpus). One of the most fundamental techniques in **natural language processing (NLP)**.



Each column of \mathbf{z} corresponds to a latent “category” or “topic”.

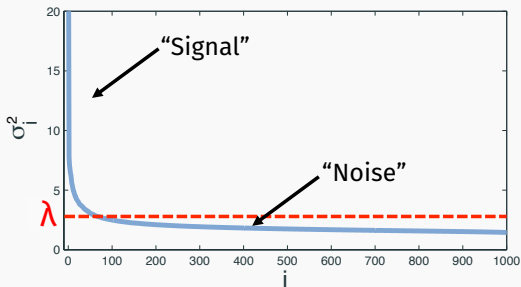
Similar documents have similar LSA document vectors. I.e.

$\langle \mathbf{z}_i, \mathbf{z}_j \rangle$ is large. Provide a more compact “finger print” for documents than the long bag-of-words vectors. Useful for e.g search engines.

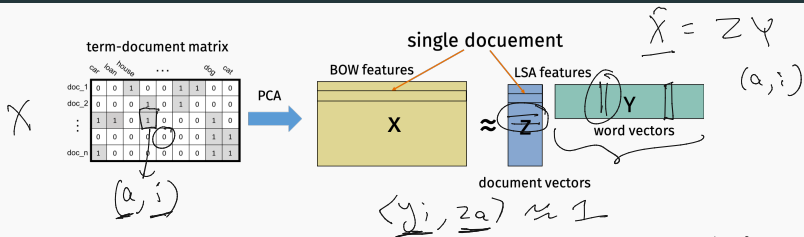
LATENT SEMANTIC ANALYSIS

LSA vectors often provide a more meaningful similarity metric than bag-of-words vectors. Capture high-level categorical information and eliminate document specific quirks.

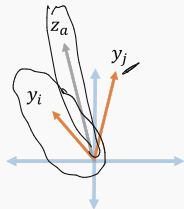
Spectrum of data matrix X .



WORD EMBEDDINGS



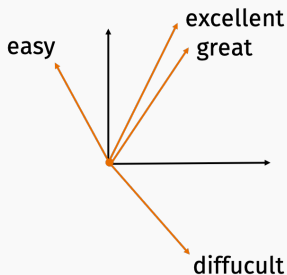
- $\langle y_i, z_a \rangle \approx 1$ when doc_i contains $word_a$. $\langle y_i, z_a \rangle \approx 0$
- If doc_i and doc_j both contain $word_a$, $\langle y_i, z_a \rangle \approx \langle y_j, z_a \rangle \approx 1$.



If two words appear in the same document their, word vectors tend to point more in the same direction.

SEMANTIC EMBEDDINGS

Result: Map words to numerical vectors in a semantically meaningful way. Similar words map to similar vectors. Dissimilar words to dissimilar vectors.



Word embeddings

Extremely useful "side-effect" of LSA.

Review 1: *Very small and handy for traveling or camping. Excellent quality, operation, and appearance.*

Review 2: *So far this thing is great. Well designed, compact, and easy to use. I'll never use another can opener.*

Review 3: *Not entirely sure this was worth \$20. Mom couldn't figure out how to use it and it's fairly difficult to turn for someone with arthritis.*

Goal is to classify reviews as “positive” or “negative”.

BAG-OF-WORDS FEATURES

Vocabulary: Small, handy, excellent, great, quality, compact, easy, difficult.

Review 1: *Very small and handy for traveling or camping. Excellent quality, operation, and appearance.*

[, , , , , , ,]

Review 2: *So far this thing is great. Well designed, compact, and easy to use. I'll never use another can opener.*

[, , , , , , ,]

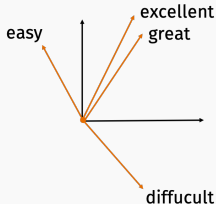
Review 3: *Not entirely sure this was worth \$20. Mom couldn't figure out how to use it and it's fairly difficult to turn for someone with arthritis.*

[, , , , , , ,]

SEMANTIC EMBEDDINGS

Bag-of-words approach typically only works for large data sets.

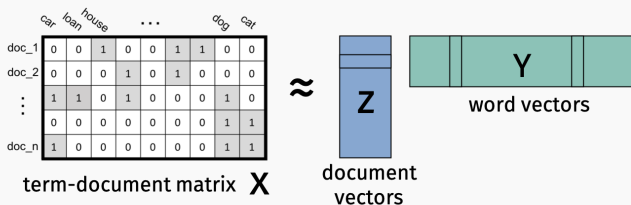
The features do not capture the fact that “great” and “excellent” are near synonyms. Or that “difficult” and “easy” are antonyms.



This can be addressed by first mapping words to semantically meaningful vectors. That mapping can be trained using a much large corpus of text than the data set you are working with (e.g. Wikipedia, Twitter, news data sets).

WORD EMBEDDINGS

Another view on word embeddings from LSA:

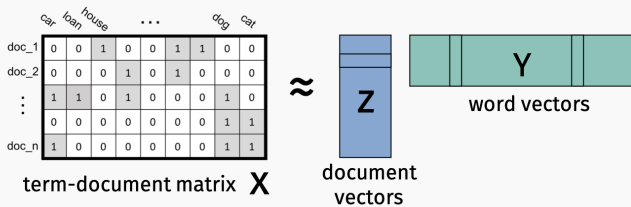


Choose \mathbf{Z} to have orthogonal columns. E.g. $\mathbf{Z} = \mathbf{U}_k$ and $\mathbf{Y} = \mathbf{\Sigma}_k \mathbf{V}_k^T$.

- $\mathbf{X} \approx \mathbf{Z}\mathbf{Y}$
- $\mathbf{X}^T\mathbf{X} \approx \mathbf{Y}^T\mathbf{Z}^T\mathbf{Z}\mathbf{Y} = \mathbf{Y}^T\mathbf{Y}$
- So for $word_i$ and $word_j$, $\langle \mathbf{y}_i, \mathbf{y}_j \rangle \approx [\mathbf{X}^T\mathbf{X}]_{i,j}$.

What does the i, j entry of $\mathbf{X}^T\mathbf{X}$ represent?

WORD EMBEDDINGS



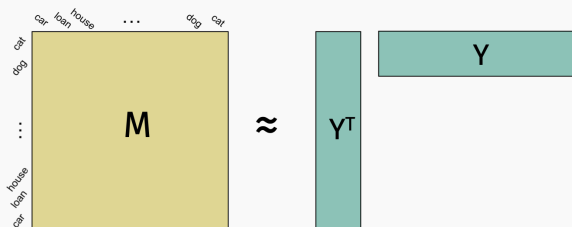
What does the i, j entry of $X^T X$ represent?

$\langle \mathbf{y}_i, \mathbf{y}_j \rangle$ is larger if $word_i$ and $word_j$ appear in more documents together (high value in **word-word co-occurrence matrix**, $\mathbf{X}^T\mathbf{X}$).
Similarity of word embeddings mirrors similarity of word context.

General word embedding recipe:

1. Choose similarity metric $k(word_i, word_j)$ which can be computed for any pair of words.
2. Construct symmetric similarity matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ with $\mathbf{M}_{i,j} = k(word_i, word_j)$.
3. Find symmetric low rank factorization $\mathbf{M} \approx \mathbf{Y}^T\mathbf{Y}$ where $\mathbf{Y} \in \mathbb{R}^{k \times n}$.
4. Columns of \mathbf{Y} are word embedding vectors.

WORD EMBEDDINGS



How do current state-of-the-art methods differ from LSA?

- Similarity based on co-occurrence in smaller chunks of words. E.g. in sentences or in any consecutive sequences of 10 words.
- Usually transformed in non-linear way. E.g.
 $k(\text{word}_i, \text{word}_j) = \frac{p(i,j)}{p(i)p(j)}$ where $p(i,j)$ is the frequency both i, j appeared together, and $p(i), p(j)$ is the frequency either one appeared.

Current state of the art models: GloVe, word2vec.

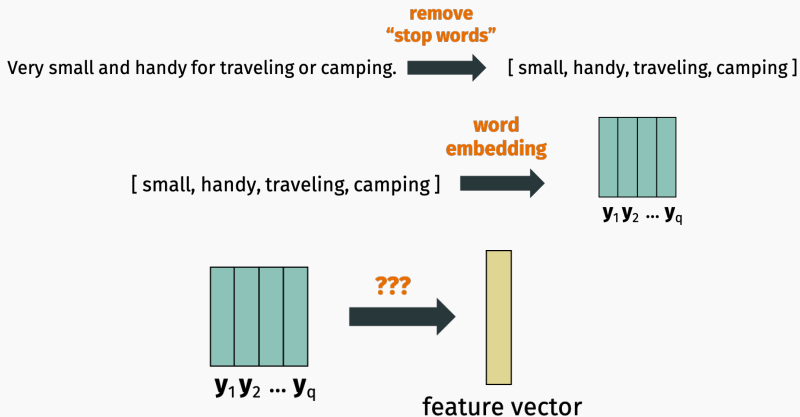
- Based on same principal as LSA.
- **word2vec** was originally presented as a shallow neural network model, but it can be viewed as a matrix factorization method (Levy, Goldberg 2014).
- For **word2vec**, similarity metric is the “point-wise mutual information”: $\log \frac{p(i,j)}{p(i)p(j)}$.

If you want to use word embeddings for your project, the easiest approach is to download pre-trained word vectors:

- Original gloVe website:
`https://nlp.stanford.edu/projects/glove/`
- Compilation of many sources:
`https://github.com/3Top/word2vec-api`

USING WORD EMBEDDINGS

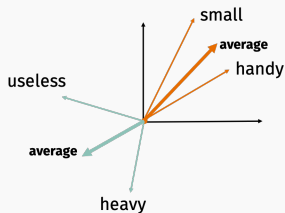
How to go from word embeddings to features for a whole sentence or chunk of text?



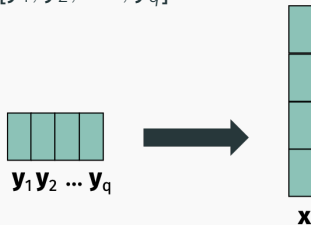
USING WORD EMBEDDINGS

A few simple options:

Feature vector $\mathbf{x} = \frac{1}{q} \sum_{i=1}^q \mathbf{y}_i$.



Feature vector $\mathbf{x} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_q]$.



Better option than concatenation: To avoid issues with inconsistent sentence length, word ordering, etc., can concatenate a fixed number of top principal components of the matrix of word vectors:

