

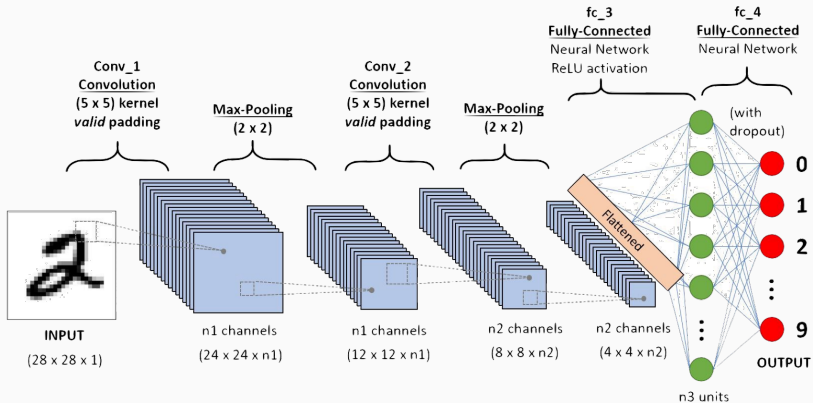
# **CS-GY 6923: Lecture 12**

## **Finishing CNN, Transfer Learning, and Autoencoders**

---

NYU Tandon School of Engineering, Akbar Rafiey

# Overall network architecture



# Pooling and downsampling

Max Pooling

29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2  
pool size

100	184
12	45

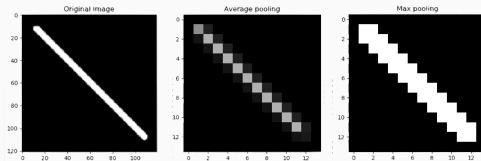
Average Pooling

31	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2  
pool size

36	80
12	15

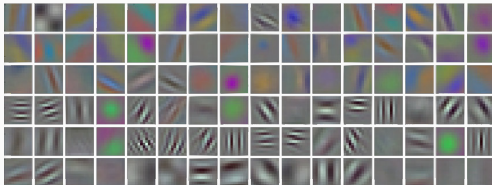
- Reduces number of variables.
- Helps “smooth” result of convolutional filters.
- Improves shift-invariance.



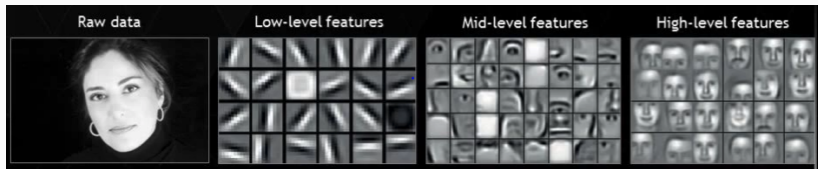
# Understanding layers

What type of convolutional filters do we learn from gradient descent?

**Lots of edge detectors in the first layer!**

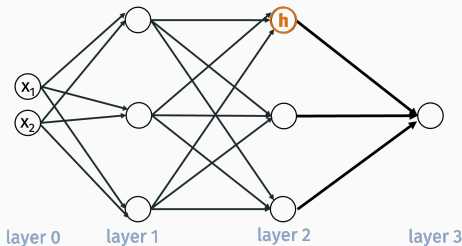


Other layers are harder to understand... but roughly hidden variables later in the network encode for “higher level features”:

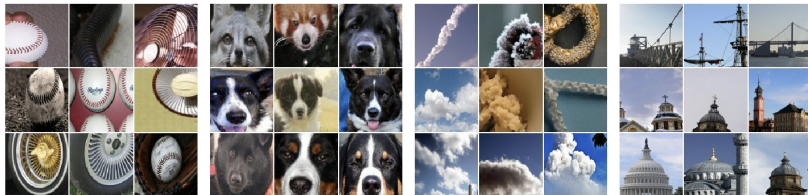


# Understanding layers

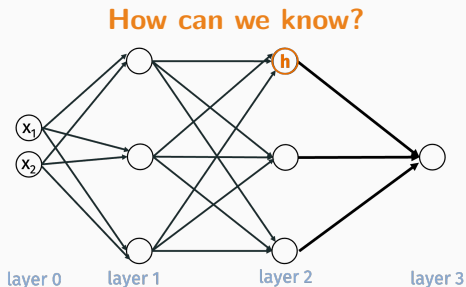
How can we know?



Go through dataset and find the inputs that most “excite” a given neuron  $h$ . I.e. for which  $|h(\mathbf{x})|$  is largest.



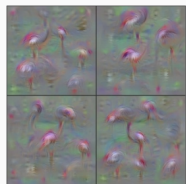
# Understanding layers



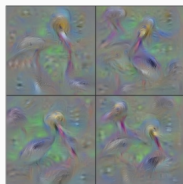
**Alternative approach:** Solve the optimization problem  $\max_{\mathbf{x}} |h(\mathbf{x})|$  e.g. using gradient descent.

# Understanding layers

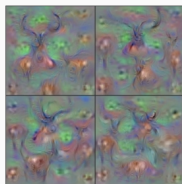
Early work had some interesting results.



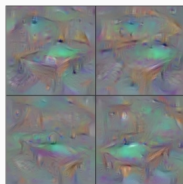
Flamingo



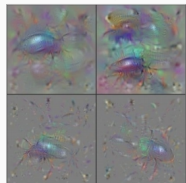
Pelican



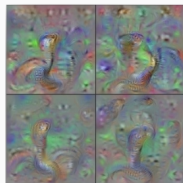
Hartebeest



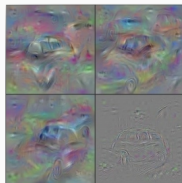
Billiard Table



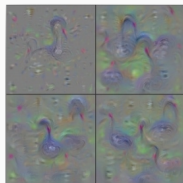
Ground Beetle



Indian Cobra



Station Wagon

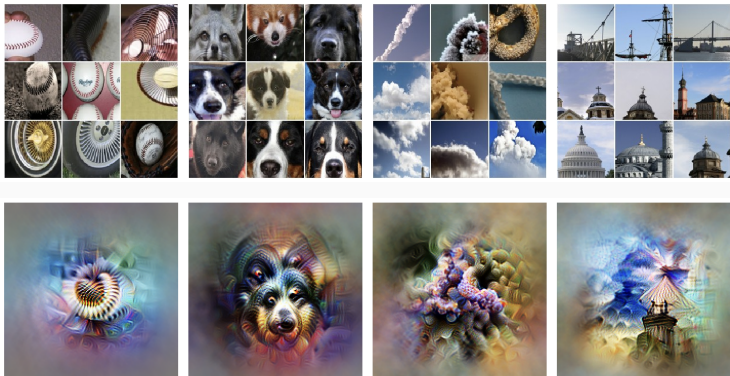


Black Swan

“Understanding Neural Networks Through Deep Visualization”, Yosinski et al.

# Understanding layers

There has been a lot of work on improving these methods by regularization. I.e. solve  $\max_{\mathbf{x}} |h(\mathbf{x})| + g(\mathbf{x})$  where  $g$  constrains  $\mathbf{x}$  to look more like a “natural image”.

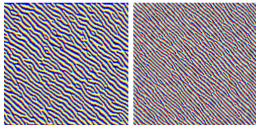
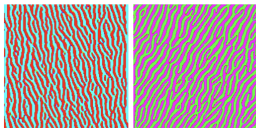
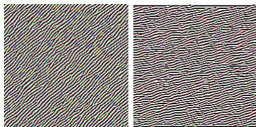


If you are interested in learning more on these techniques, there is a great Distill article at: <https://distill.pub/2017/feature-visualization/>.

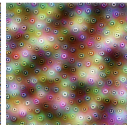
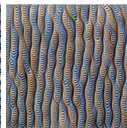
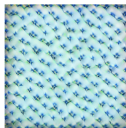
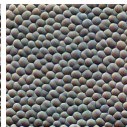
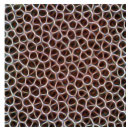


# Understanding layers

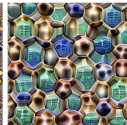
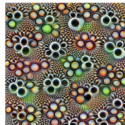
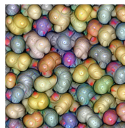
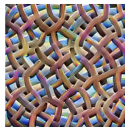
Nodes at different layers have different layers capture increasingly more abstract concepts.



**Edges** (layer conv2d0)



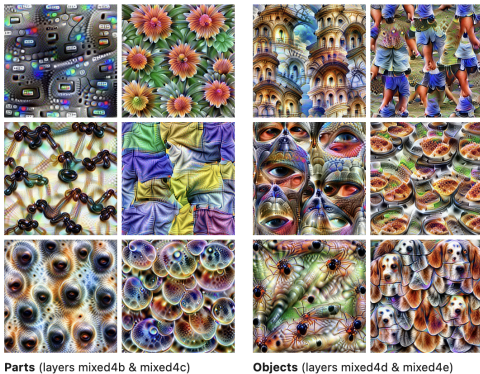
**Textures** (layer mixed3a)



**Patterns** (layer mixed4a)

# Understanding layers

Nodes at different layers have different layers capture increasingly more abstract concepts.

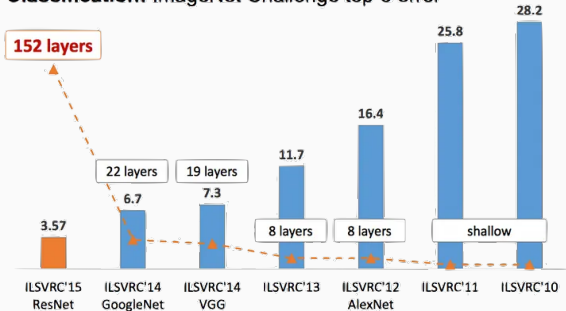


**General observation:** Depth more important than width. Alexnet 2012 had 8 layers, modern convolutional nets can have 100s.

# Deeper and deeper, bigger and bigger

After AlexNet (8 layers, 60 million parameters) achieved state of the art performance on ImageNet, progress proceeded rapidly:

**Classification: ImageNet Challenge top-5 error**



## Tricks of the trade

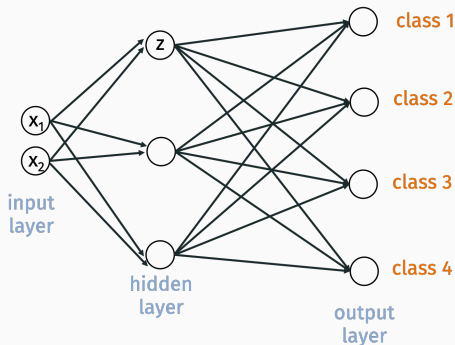
Beyond techniques discussed for general neural nets (back-prop, batch gradient descent, adaptive learning rates) training deep networks requires a lot of “tricks”.

- Batch normalization (accelerate training).
- Dropout (prevent over-fitting)
- Residual connections (accelerate training, allow for more depth – 100s of layers).
- Data augmentation.

**And deep networks require lots of training data and lots of time.**

# Batch normalization

Start with any neural network architecture:



For input  $\mathbf{x}$ ,

$$\bar{z} = \mathbf{w}^T \mathbf{x} + b$$

$$z = s(\bar{z})$$

where  $\mathbf{w}$ ,  $b$ , and  $s$  are weights, bias, and non-linearity.

## Batch normalization

$\bar{z}$  is a function of the input  $\mathbf{x}$ . We can write it as  $\bar{z}(\mathbf{x})$ . Consider the mean and standard deviation of the hidden variable over our entire dataset  $\mathbf{x}_1 \dots, \mathbf{x}_n$ :

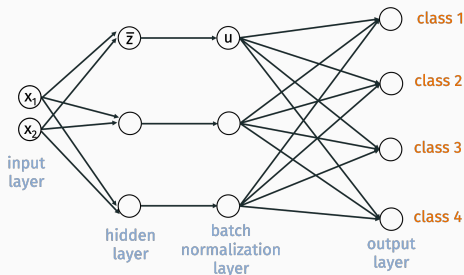
$$\mu = \frac{1}{n} \sum_{j=1}^n \bar{z}(\mathbf{x}_j)$$

$$\sigma^2 = \frac{1}{n} \sum_{j=1}^n (\bar{z}(\mathbf{x}_j) - \mu)^2$$

Just as normalization (mean centering, scaling to unit variance) is sometimes used for input features, batch-norm applies normalization to learned features.

# Batch normalization

Can add a batch normalization layer after any layer:



$$\bar{u} = \frac{\bar{z} - \mu}{\sigma}$$

$$u = s(\bar{u}).$$

Has the effect of mean-centering/normalizing  $\bar{z}$ . Typically we actually allow  $u = s(\gamma \cdot \bar{u} + c)$  for learned parameters  $\gamma$  and  $c$ .

# Batch normalization

Proposed in 2015: “Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift”, Ioffe, Szegedy.

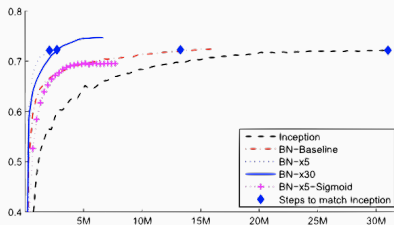


Figure 2: *Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.*

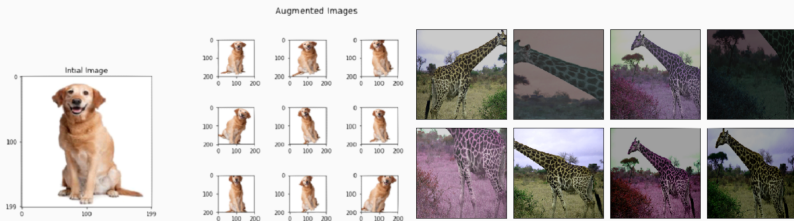
Doesn't change the expressive power of the network, but allows for significant convergence acceleration. It is not yet well understood why batch normalization speeds up training.



# Data augmentation

Great general tool to know about. **Main idea:**

- More training data typically leads to a more accurate model.
- Artificially enlarge training data with simple transformations.



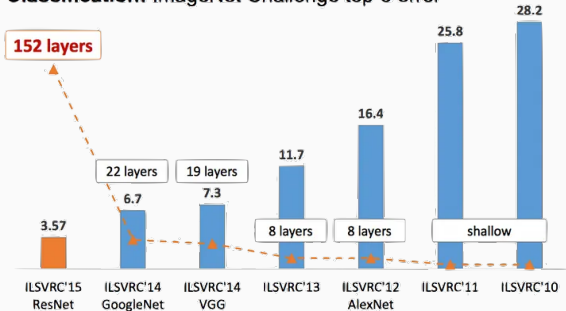
Take training images and randomly shift, flip, rotate, skew, darken, lighten, shift colors, etc. to create new training images. **Final classifier will be more robust to these transformations.**

Need to take a full course on neural networks/deep learning to learn more! State-of-the-art techniques are constantly evolving.

# Deeper and deeper, bigger and bigger

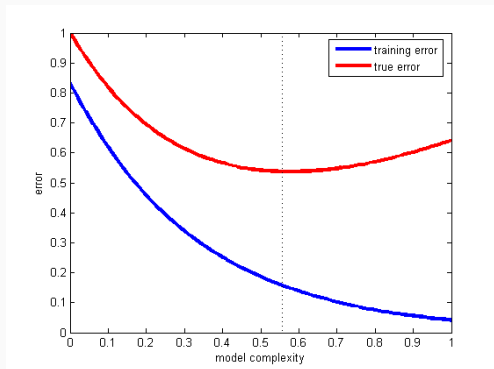
After AlexNet (8 layers, 60 million parameters) achieved state of the art performance on ImageNet, progress proceeded rapidly:

**Classification: ImageNet Challenge top-5 error**



# Generalization for neural networks

Even with weight sharing, convolution, etc. modern neural networks typically have 100s of millions of parameters. And we don't train them with regularization. Intuitively we might expect them to overfit to training data.



# Generalization for neural networks

In fact, we now know that modern neural nets can easily overfit to training data. This work showed that we can fit large vision data sets with random class labels to essentially perfect accuracy.

## UNDERSTANDING DEEP LEARNING REQUIRES RE-THINKING GENERALIZATION

**Chiyuan Zhang\***  
Massachusetts Institute of Technology  
chiyuan@mit.edu

**Samy Bengio**  
Google Brain  
bengio@google.com

**Moritz Hardt**  
Google Brain  
mrtz@google.com

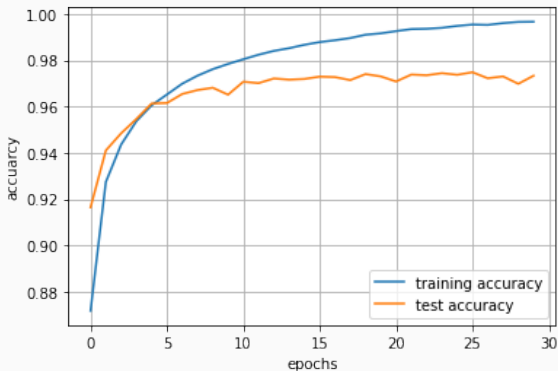
**Benjamin Recht†**  
University of California, Berkeley  
brecht@berkeley.edu

**Oriol Vinyals**  
Google DeepMind  
vinyals@google.com

But we don't always see a large gap between training and test error. **Don't take this to mean overfitting isn't a problem when using neural nets!** It's just not always a problem.

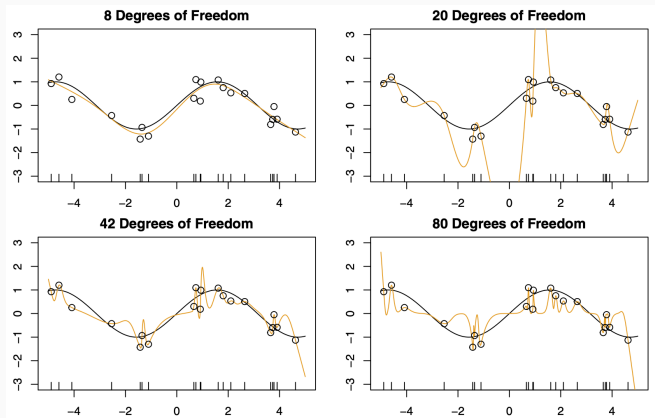
# Generalization for neural networks

We even see this lack of overfitting for MNIST data. Check `keras_demo_mnist.ipynb` I posted on the website:



# Generalization for neural networks

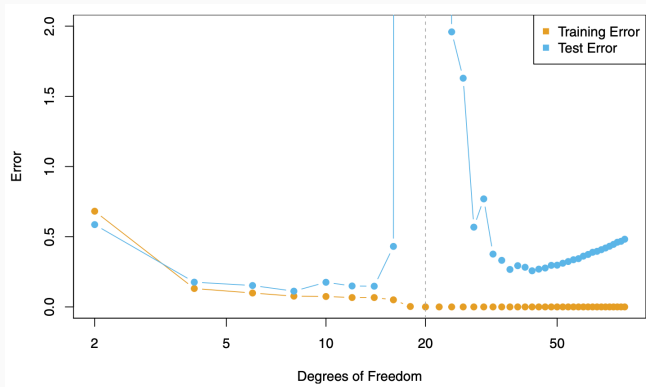
One growing realization is that this phenomena doesn't only apply to neural networks – it can also be true for fitting highly-overparameterized polynomials.



The choice of training algo (e.g. gradient descent) seems important.

# Double descent

We sometimes see a “double descent curve” for these models. Test error is worst for “just barely” overparameterized models, but gets better with lots of overparameterization.



We don't usually see this same curve for neural networks.



# Overfitting in neural nets

**Take away:** Modern neural network overfit, but still seem fairly robust. Perform well on any new test data we throw at them.

Or do they?

---

## Intriguing properties of neural networks

---

**Christian Szegedy**  
Google Inc.

**Wojciech Zaremba**  
New York University

**Ilya Sutskever**  
Google Inc.

**Joan Bruna**  
New York University

**Dimitru Erhan**  
Google Inc.

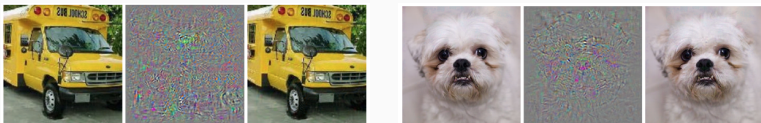
**Ian Goodfellow**  
University of Montreal

**Rob Fergus**  
New York University  
Facebook Inc.

# Adversarial examples

# Adversarial examples

**Main discovery:** It is possible to find imperceptibly small perturbations of input images that will fool deep neural networks. This seems to be a universal phenomenon.



**Important:** Random perturbations do not work!

## Adversarial examples

How to find “good” perturbations:

Fix model  $f_\theta$ , input  $\mathbf{x}$ , correct label  $y$ . Consider the loss  $\ell(\theta, \mathbf{x}, y)$ .

Solve the optimization problem:

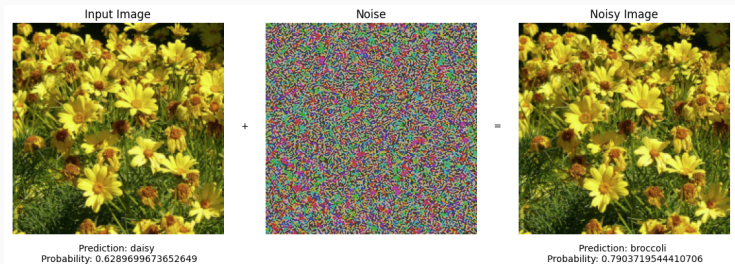
$$\max_{\delta, \|\delta\| \leq \epsilon} \ell(\theta, \mathbf{x} + \delta, y)$$

Can be solved using gradient descent! We just need to compute the derivative of the loss with respect to the image pixels.

Backprop can do this easily.

# Adversarial examples

Teal put together a really cool lab where you can find your own adversarial examples for a model called Resnet18. The entire model + weights are available through PyTorch, so we do not need to train it ourselves (i.e. this is a pre-trained model).



# Transfer Learning and Autoencoders

# Transfer learning

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

But they require lots and lots of data. ImageNet has 14 million labeled images. Mostly of everyday objects.

# One-shot learning

What if you want to apply deep convolutional networks to a problem where you do not have a lot of **labeled data** in the first place?



quaffle



bludger



snitch

**Example:** Classify images of different Quidditch balls.



**Real example:** Classify images of insects for use in agricultural applications in new localities.

## Zero-Shot Insect Detection via Weak Language Supervision

**Benjamin Feuer,<sup>1</sup> Ameya Joshi,<sup>1</sup> Minsu Cho,<sup>1</sup> Kewal Jani,<sup>1</sup> Shivani Chiranjeevi,<sup>2</sup> Zi Kang Deng,<sup>3</sup>  
Aditya Balu,<sup>2</sup> Asheesh K. Singh,<sup>2</sup> Soumik Sarkar,<sup>2</sup> Nirav Merchant,<sup>3</sup> Arti Singh,<sup>2</sup>  
Baskar Ganapathysubramanian,<sup>2</sup> Chinmay Hegde<sup>1</sup>**

<sup>1</sup> New York University

<sup>2</sup> Iowa State University

<sup>3</sup> University of Arizona

Aedes Vexans



Cretonotos Gangis



Daphnis Neril



Hypena Deceptalis

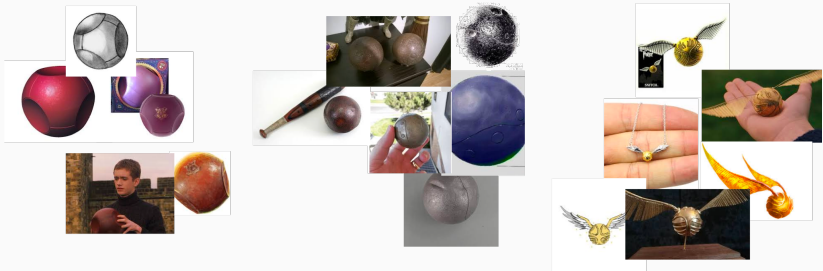


Pyralis Farinalis



# One-shot learning

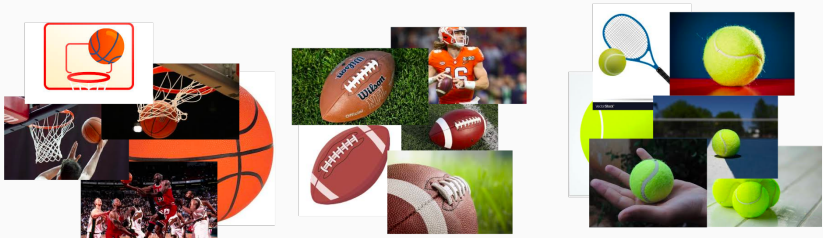
A human could probably achieve near perfect classification accuracy even given access to a **single labeled example** from each class:



**Major question in ML:** How? Can we design ML algorithms which can do the same?

# Transfer learning

Transfer knowledge from one task we already know how to solve to another.



For example, we have learned from past experience that balls used in sports have consistent shapes, colors, and sizes. These features can be used to distinguish balls of different type.

# Feature learning

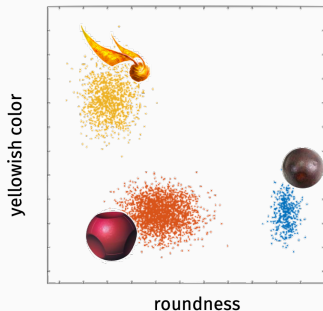
Examples of possible high-level features a human would learn:

**Classes**

							
Features	roundness	1	.1	1	.6	1	.4
	size relative to human hand	10	7	2	7	5	1
	yellowish color	.2	.1	1	.1	0	.9

## Feature learning

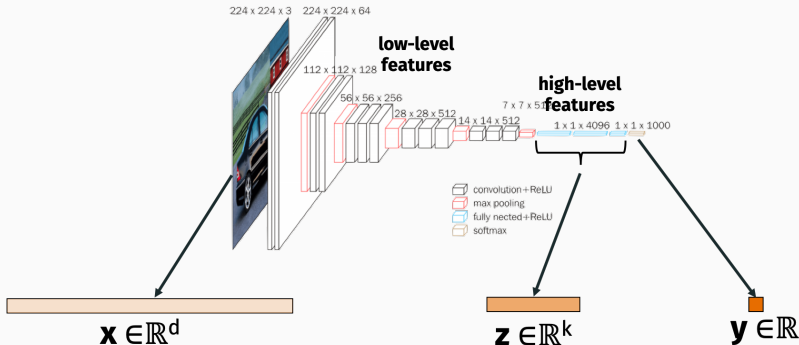
If these features are highly informative (i.e. lead to highly separable data) few training examples are needed to learn.



Might suffice to classify ball using nearest training example in feature space, even if just a handful of training examples.

# Transfer learning

**Empirical observation:** Features learned when training models like deep neural nets seem to capture exactly these sorts of high-level properties.



Even if we can't put into words what each feature in  $z$  means...

# Transfer learning

This is now a common technique in computer vision:

1. Download network trained on 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.).

This approach has even been used on the quidditch problem:

[github.com/thatbrguy/Object-Detection-Quidditch](https://github.com/thatbrguy/Object-Detection-Quidditch)

# Unsupervised feature 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**.



## Supervised vs. unsupervised learning

- **Supervised learning:** All input data examples come with targets/labels. What machines have been really good at for the past 8 years.
- **Unsupervised learning:** No input data examples come with targets/labels. Interesting problems to solve include clustering, anomaly detection, semantic embedding, etc.
- **Semi-supervised learning:** Some (typically very few) input data examples come with targets/labels. What human babies are really good at, and we have recently made machines a lot better at.

# Autoencoder

**Simple but clever idea:** If we have inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  but few or no targets  $y_1, \dots, y_n$ , just make the inputs the targets.

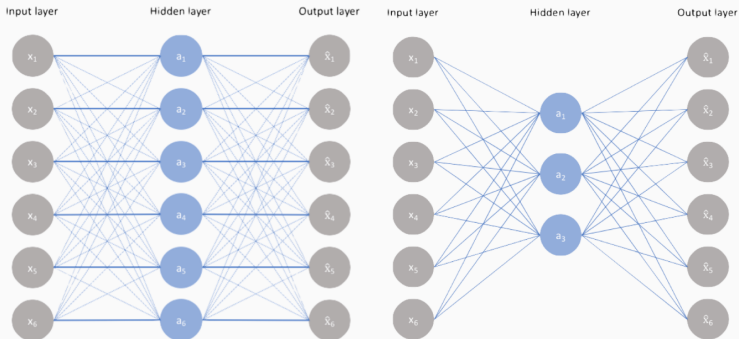
- Let  $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be our model.
- Let  $L_\theta$  be a loss function. E.g. squared loss:  
$$L_\theta(\mathbf{x}) = \|\mathbf{x} - f_\theta(\mathbf{x})\|_2^2.$$
- Train model:  $\theta^* = \min_\theta \sum_{i=1}^n L_\theta(\mathbf{x}_i)$ .

If  $f_\theta$  is a model that incorporates feature learning, then these features can be used for supervised tasks.

$f_\theta$  is called an **autoencoder**. It maps input space to input space (e.g. images to images, french to french, PDE solutions to PDE solutions).

# Autoencoder

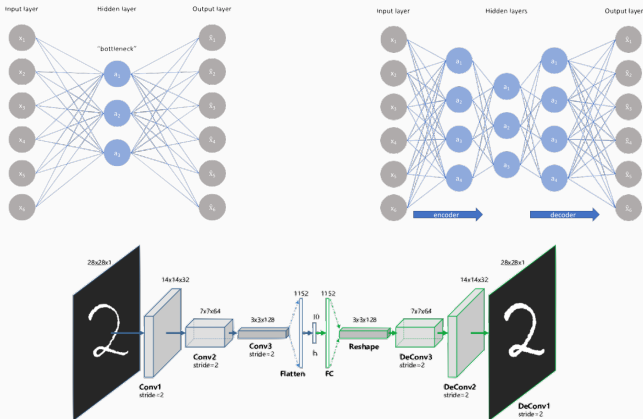
Two examples of autoencoder architectures:



Which would lead to better feature learning?

# Autoencoder

**Important property of autoencoders:** no matter the architecture, 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.



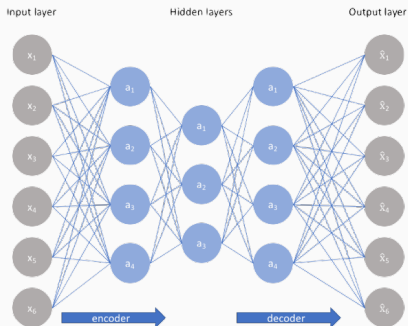
# Autoencoder

Separately name the mapping from input to bottleneck and from bottleneck to output.

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

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

$$f(\mathbf{x}) =$$

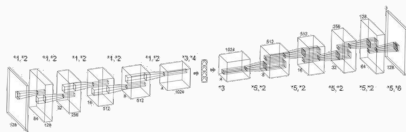


Often symmetric, but does not have to be.

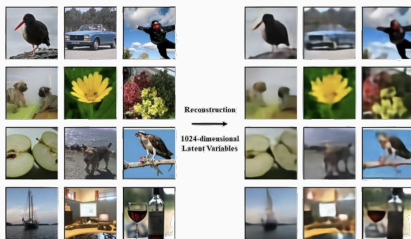
# Autoencoder reconstruction

## Example image reconstructions from autoencoder:

(A)



(B)



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

Input parameters:  $d = 49152$ .

Bottleneck “latent” parameters:  $k = 1024$ .

# Autoencoders for feature extraction

The best autoencoders do not work as well as supervised methods for feature extraction, but they require no labeled data.<sup>1</sup>

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

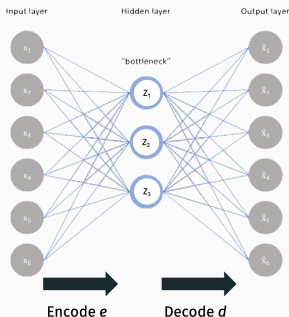
- Learned data compression.
- Denoising and in-painting.
- Data/image synthesis.

---

<sup>1</sup>Recent progress on **self-supervised** learning achieves the best of both worlds – state-of-the-art feature learning with no labeled data.

# Autoencoders for data 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  $\mathbf{x}$ .

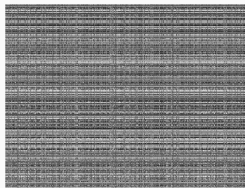


# Autoencoders for image compression

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.



**Proposed method**, 5908 bytes (0.167 bit/px), PSNR: luma 23.38 dB/chroma 31.86 dB, MS-SSIM: 0.9219



**JPEG 2000**, 5908 bytes (0.167 bit/px), PSNR: luma 23.24 dB/chroma 31.04 dB, MS-SSIM: 0.8803



**Proposed method**, 6021 bytes (0.170 bit/px), PSNR: 24.12 dB, MS-SSIM: 0.9292

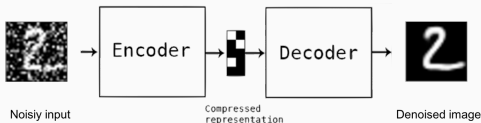


**JPEG 2000**, 6037 bytes (0.171 bit/px), PSNR: 23.47 dB, MS-SSIM: 0.9036

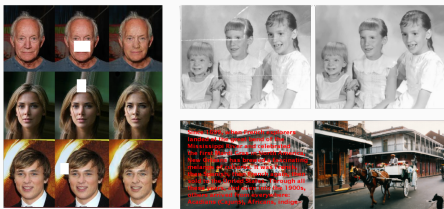
“End-to-end optimized image compression”, Ballé, Laparra, Simoncelli

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 correction



## Image denoising

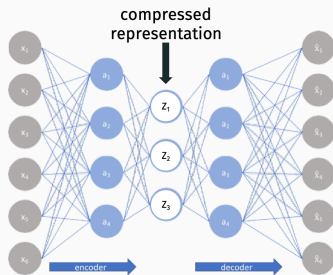


## Image inpainting

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

# Autoencoders learn compressed representations

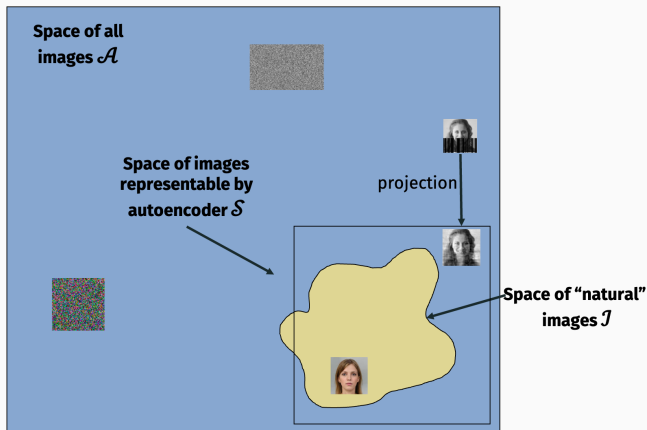
## Why does this work?



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

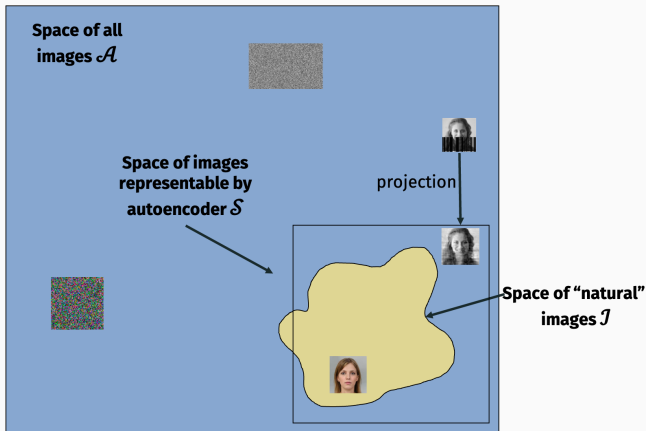
If  $\mathbf{z}$  holds  $k$  values in  $0, .1, .2, \dots, 1$ , how many unique images  $\mathbf{w}$  can be output by the autoencoder function  $f$ ?

# 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

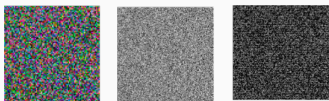


$f(\mathbf{x}) = d(e(\mathbf{x}))$  projects an image  $\mathbf{x}$  closer to the space of natural images.

# Autoencoders for data generation

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

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



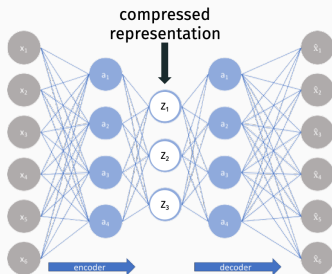
- **Option 2:** Draw  $\mathbf{x}$  randomly from  $\mathcal{S}$ , the space of images representable by the autoencoder.



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

# Autoencoders for data generation

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



Randomly select code  $\mathbf{z}$ , then set  $\mathbf{x} = d(\mathbf{z})$ .<sup>2</sup>

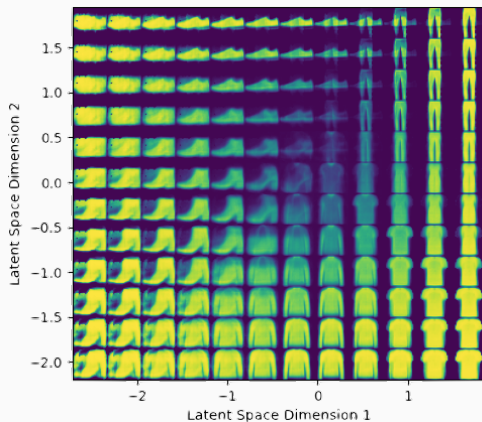
---

<sup>2</sup>Lots of details to think about here. In reality, people use “variational autoencoders” (VAEs), which are a natural modification of AEs.



# Autoencoders for data generation demo

Demo for the "Fashion MNIST" data set:



# Principal Component Analysis (PCA)

# Principal Component Analysis (PCA)

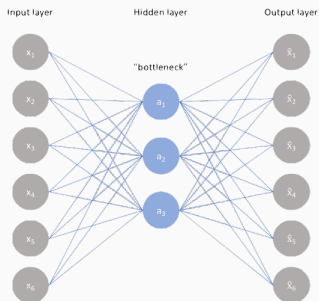
**Rest of lecture:** Deeper dive into understanding a simple, but powerful autoencoder architecture. Specifically we will view **Principal Component Analysis (PCA)** as a type of autoencoder.

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 (PCA)

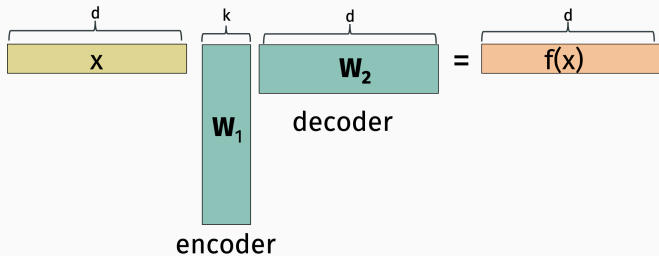
Consider the simplest possible autoencoder:



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

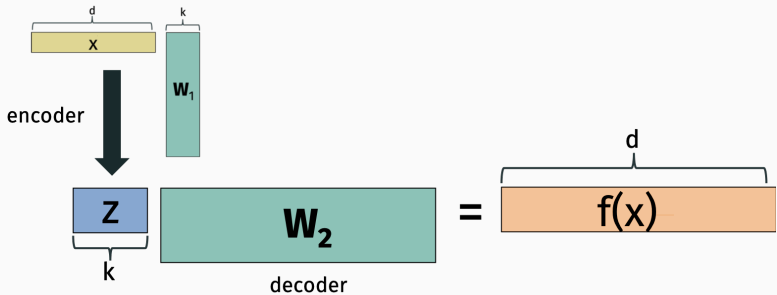
# Principal Component Analysis (PCA)

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



$$f(\mathbf{x})^T = \mathbf{x}^T \mathbf{W}_1 \mathbf{W}_2$$

# Principal Component Analysis (PCA)

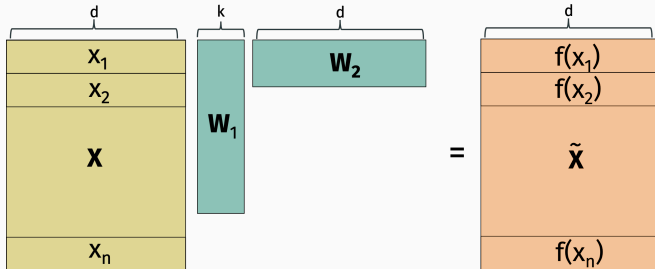


**Encoder:**  $e(x) = x^T W_1$ .

**Decoder:**  $d(z) = zW_2$

# Principal Component Analysis (PCA)

Given training data set  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , let  $\mathbf{X}$  denote our data matrix.  
Let  $\tilde{\mathbf{X}} = \mathbf{X}\mathbf{W}_1\mathbf{W}_2$ .



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

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

**Goal:** Find  $\mathbf{W}_1, \mathbf{W}_2$  to minimize the Frobenius norm loss  $\|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2 = \|\mathbf{X} - \mathbf{X}\mathbf{W}_1\mathbf{W}_2\|_F^2$  (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**.

The diagram shows the equation  $Z = X\tilde{X}$ . On the left, a blue vertical rectangle represents matrix  $Z$ , with a bracket above it labeled  $k$ . It is divided into four sections:  $Z_1$ ,  $Z_2$ ,  $Z = XW_1$ , and  $Z_n$ . To its right is a green horizontal rectangle representing matrix  $W_2$ , with a bracket above it labeled  $d$ . An equals sign is placed between the two sides. On the right side of the equals sign is a large orange vertical rectangle representing matrix  $\tilde{X}$ , with a bracket above it labeled  $d$ .

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

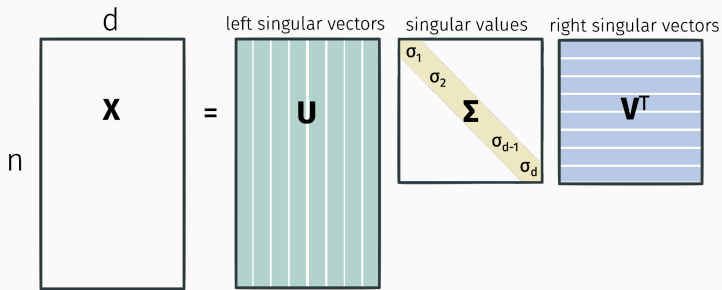
## Low-rank approximation

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  $\mathbf{X}$  can be written:



Where  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$ ,  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ , and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$ . I.e.  $\mathbf{U}$  and  $\mathbf{V}$  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  $\mathbf{X}$ .

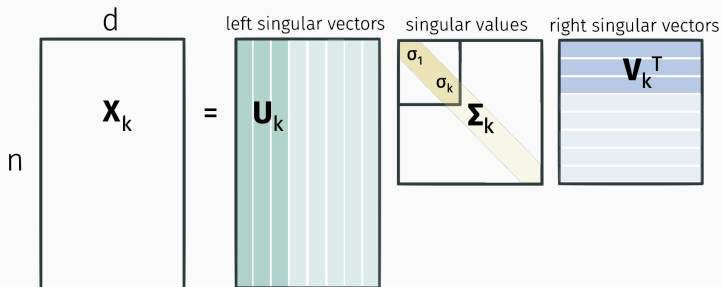
$$\mathbf{U}^T \quad \mathbf{U} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 0 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\|\mathbf{u}_i\|_2^2 =$$

$$\mathbf{u}_i^T \mathbf{u}_j =$$

# Singular value decomposition

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



**Eckart–Young–Mirsky Theorem:** For any  $k \leq d$ ,  
 $\mathbf{X}_k = \mathbf{U}_k \Sigma_k \mathbf{V}_k^T$  is the optimal  $k$  rank approximation to  $\mathbf{X}$ :

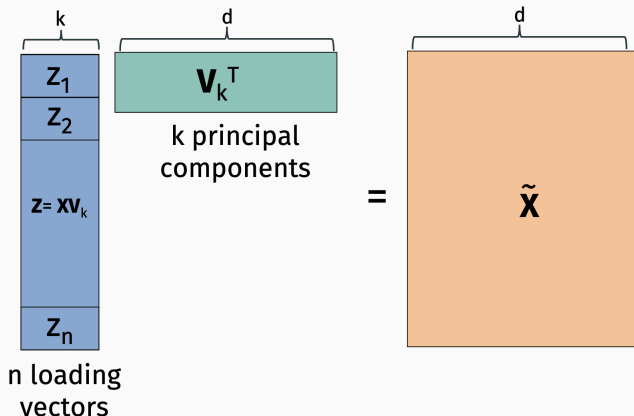
$$\mathbf{X}_k = \underset{\tilde{\mathbf{X}} \text{ with rank } \leq k}{\operatorname{arg\,min}} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2.$$

## Singular value decomposition

**Claim:**  $\mathbf{X}_k = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T = \mathbf{X} \mathbf{V}_k \mathbf{V}_k^T$ .

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

# Principal Component Analysis (PCA)



Usually  $\mathbf{x}$ '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.svds(X, k).
```

Runs in roughly  $O(ndk)$  time.



## Connection to eigen-decomposition

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

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

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

# PCA applications

Like any autoencoder, PCA can be used for:

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



denoising



synthetic data generation

# Low-rank approximation

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

```
72104194959
0690159784
7665407401
3134727121
1742351244
6355604195
7853746430
7009173297
1627847361
3683141769
```

original data

rank 1 approx.

```
8888888888888
8888888888888
8888888888888
8888888888888
8888888888888
8888888888888
8888888888888
8888888888888
8888888888888
8888888888888
```

rank 2 approx.

```
888100799809
00890189784
7665407401
779097097097
7773307880
8880000181
88833710880
007178279
8887781607
8817781709
```

rank 3 approx.

```
92709719909
0090189784
7665407401
31340727131
1792351244
895604195
9097950430
9007173297
9027847361
3683141769
```

rank 4 approx.

```
92709719909
0090189784
7665407401
31340727131
1792351244
895604195
9097950430
9007173297
9027847361
3683141769
```

rank 5 approx.

```
72104194959
0090189784
7665407401
31340727131
1792351244
895604195
9097950430
9007173297
9027847361
3683141769
```

```
72104194959
0090189784
7665407401
31340727131
1792351244
6355604195
9097950430
7009173297
1627847361
3683141769
```

rank 6 approx.

```
72104194959
0090189784
7665407401
31340727131
1792351244
6355604195
9097950430
7009173297
1627847361
3683141769
```

rank 7 approx.

```
72104194959
0090189784
7665407401
31340727131
1792351244
6355604195
9097950430
7009173297
1627847361
3683141769
```

rank 8 approx.

```
72104194959
0090189784
7665407401
31340727131
1792351244
6355604195
9097950430
7009173297
1627847361
3683141769
```

rank 9 approx.

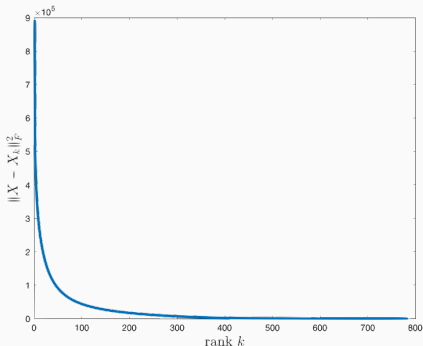
```
72104194959
0090189784
7665407401
31340727131
1792351244
6355604195
9097950430
7009173297
1627847361
3683141769
```

rank 50 approx.

# Low rank approximation

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

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



## 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.

$$\text{yard size} \approx \text{lot size} - \frac{1}{2} \cdot \text{square footage.}$$

$$\begin{aligned} \text{cumulative GPA} &\approx \frac{1}{4} \cdot \text{year 1 GPA} + \frac{1}{4} \cdot \text{year 2 GPA} \\ &\quad + \frac{1}{4} \cdot \text{year 3 GPA} + \frac{1}{4} \cdot \text{year 4 GPA.} \end{aligned}$$

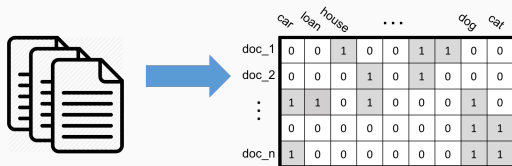
# Low-rank intuition

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:

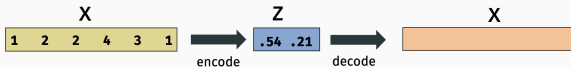


## Examples of low-rank structure

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

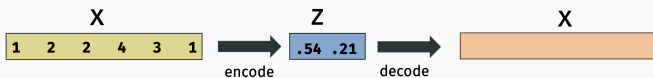
Most of the information in  $\mathbf{x}$  is explained by just a few **latent variable**.





## Examples of low-rank structure

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



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

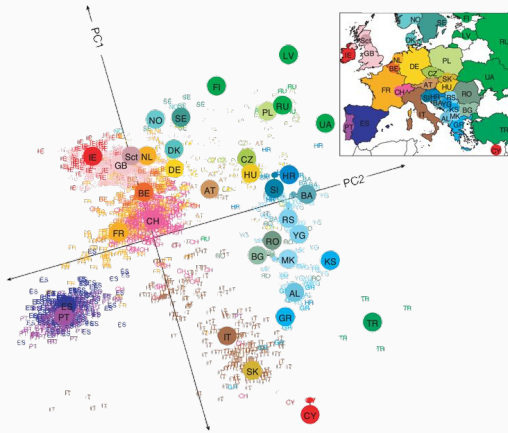
$z[1] \approx$  relative north-south position of birth place

$z[2] \approx$  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.

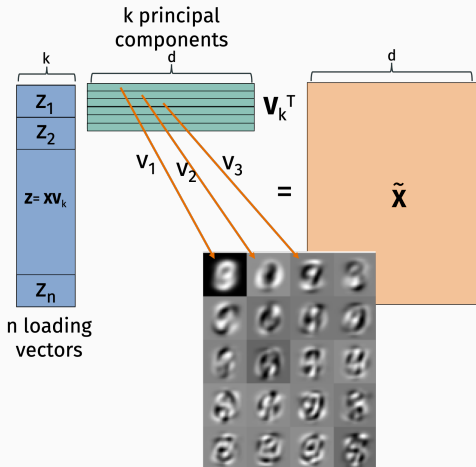


Genetic data can be nicely visualized using PCA! Plot each data example  $x$  using two loading variables in  $z$ .

For more complex data, what do principal components and loading vectors look like?

# Principal components

MNIST principal components:

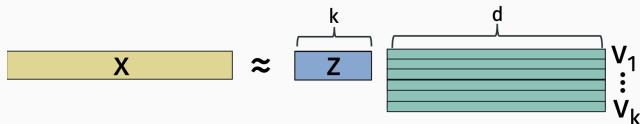


Often principal components are difficult to interpret.

# Loading vectors

What do the **loading vectors** look like?

The loading vector  $\mathbf{z}$  for an example  $\mathbf{x}$  contains coefficients which recombine the top  $k$  principal components  $\mathbf{v}_1, \dots, \mathbf{v}_k$  to approximately reconstruct  $\mathbf{x}$ .

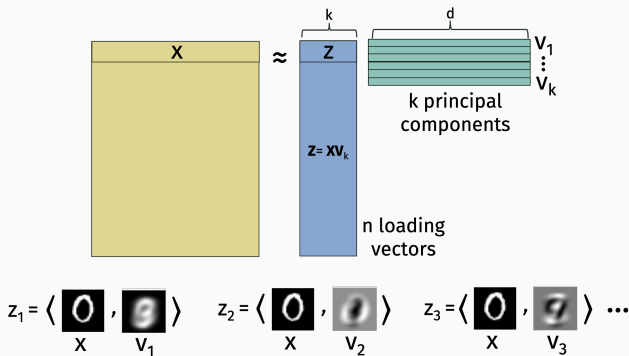


An equation showing the reconstruction of a handwritten digit '0'. On the left is a black square containing a white '0', labeled  $\mathbf{x}$  below it. To its right is an approximation symbol  $\approx$ . This is followed by a sum of terms:  $z_1 \cdot \mathbf{v}_1 + z_2 \cdot \mathbf{v}_2 + z_3 \cdot \mathbf{v}_3 + z_4 \cdot \mathbf{v}_4 + \dots$ . Each term consists of a coefficient  $z_i$  multiplied by a grayscale image of a digit '0' labeled  $\mathbf{v}_i$  below it.

Provide a short “finger print” for any image  $\mathbf{x}$  which can be used to reconstruct that image.

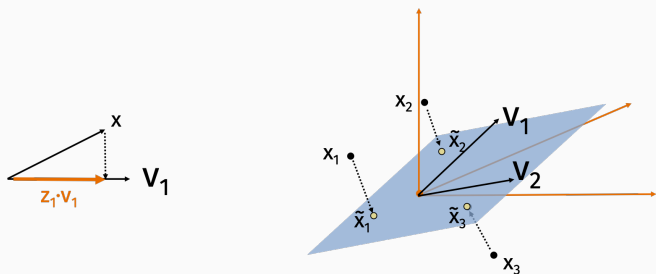
## Loading vectors: similarity view

For any  $\mathbf{x}$  with loading vector  $\mathbf{z}$ ,  $z_i$  is the inner product similarity between  $\mathbf{x}$  and the  $i^{\text{th}}$  principal component  $\mathbf{v}_i$ .



## Loading vectors: projection view

So we approximate  $\mathbf{x} \approx \tilde{\mathbf{x}} = \langle \mathbf{x}, \mathbf{v}_1 \rangle \cdot \mathbf{v}_1 + \dots + \langle \mathbf{x}, \mathbf{v}_k \rangle \cdot \mathbf{v}_k$ .

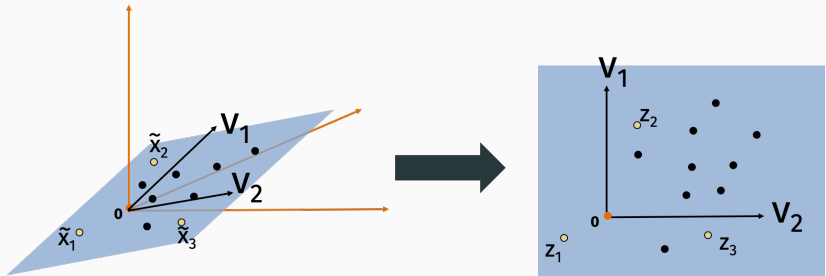


Since  $\mathbf{v}_1, \dots, \mathbf{v}_k$  are orthonormal, this operation is a **projection** onto first  $k$  principal components.

I.e. we are projecting  $\mathbf{x}$  onto the  $k$ -dimensional subspace spanned by  $\mathbf{v}_1, \dots, \mathbf{v}_k$ .

## Loading vectors: projection view

For an example  $\tilde{x}_i$ , the loading vector  $z_i$  contains the coordinates in the projection space:





## Similarity preservation

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

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

We have:

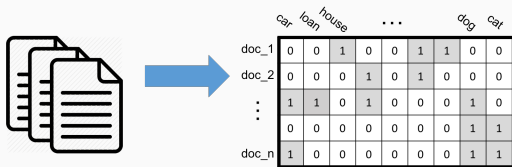
$$\begin{aligned}\|\tilde{\mathbf{x}}_i\|_2^2 &= \|\mathbf{z}_i\|_2^2 \\ \langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle &= \langle \mathbf{z}_i, \mathbf{z}_j \rangle \\ \|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2^2 &= \|\mathbf{z}_i - \mathbf{z}_j\|_2^2\end{aligned}$$

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

$$\begin{aligned}\|\mathbf{x}_i\|_2^2 &\approx \|\mathbf{z}_i\|_2^2 \\ \langle \mathbf{x}_i, \mathbf{x}_j \rangle &\approx \langle \mathbf{z}_i, \mathbf{z}_j \rangle \\ \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 &\approx \|\mathbf{z}_i - \mathbf{z}_j\|_2^2\end{aligned}$$

# Term document matrix

Word-document matrices tend to be low rank.

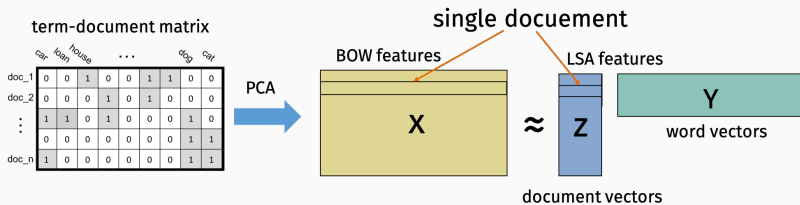


Documents tend to fall into a relatively small number of different categories, which use similar sets of words:

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

# 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 **z** corresponds to a latent “category” or “topic”. Corresponding row in **Y** corresponds to the “frequency” with which different words appear in documents on that topic.

## Latent semantic analysis

Similar documents have similar LSA document vectors. I.e.  $\langle \mathbf{z}_i, \mathbf{z}_j \rangle$  is large.

- $\mathbf{z}_i$  provides a more compact “finger print” for documents than the long bag-of-words vectors. Useful for e.g search engines.
- Comparing document vectors is often more effective than comparing raw BOW features. Two documents can have  $\langle \mathbf{z}_i, \mathbf{z}_j \rangle$  large even if they have no overlap in words. E.g. because both share a lot of words with words with another document  $k$ , or with a bunch of other documents.

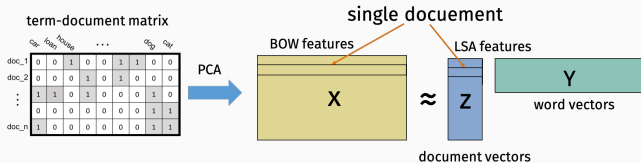
# Eigenfaces

Same fingerprinting idea was also important in early facial recognition systems based on “eigenfaces”:

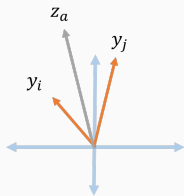


Each image above is one of the principal components of a dataset containing images of faces.

# Word embeddings



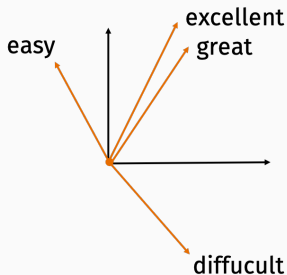
- $\langle \mathbf{y}_i, \mathbf{z}_a \rangle \approx 1$  when  $doc_a$  contains  $word_i$ .
- If  $word_i$  and  $word_j$  both appear in  $doc_a$ , then  $\langle \mathbf{y}_i, \mathbf{z}_a \rangle \approx \langle \mathbf{y}_j, \mathbf{z}_a \rangle \approx 1$ , so we expect  $\langle \mathbf{y}_i, \mathbf{y}_j \rangle$  to be large.



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.



Extremely useful “side-effect” of LSA.

Capture e.g. the fact that “great” and “excellent” are near synonyms. Or that “difficult” and “easy” are antonyms.



## Word embeddings: motivating problem

**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.*

[ , , , , , , , ]

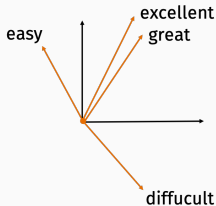
**Review 3:** *Not entirely sure this was worth \$20. Mom couldn't figure out how to use it and it's fairly difficult to turn for someone with arthritis.*

[ , , , , , , , ]

# 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).

# Using word embeddings

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

Very small and handy for traveling or camping. **remove "stop words"** → [ small, handy, traveling, camping ]

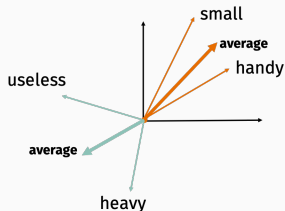
[ small, handy, traveling, camping ] **word embedding** →   
 $y_1 y_2 \dots y_q$

  
 $y_1 y_2 \dots y_q$  **???** →   
feature vector

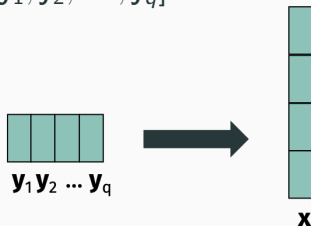
# Using word embeddings

A few simple options:

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

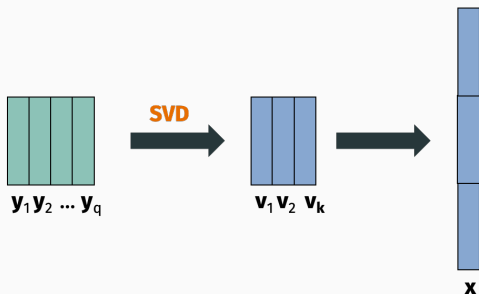


Feature vector  $\mathbf{x} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_q]$ .



## Using word embeddings

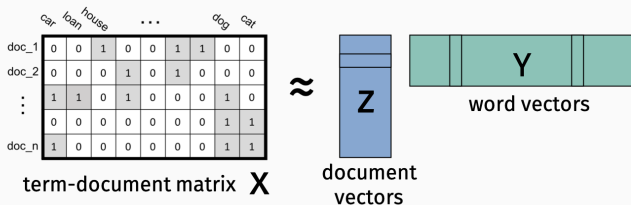
To avoid issues with inconsistent sentence length, word ordering, etc., can concatenate a fixed number of top principal components of the matrix of word vectors:



There are much more complicated approaches that account for word position in a sentence. Lots of pretrained libraries available (e.g. Facebook's InferSent).

# Word embeddings

Another view on word embeddings from LSA:

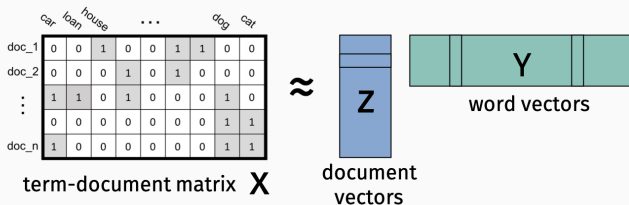


We chose  $\mathbf{Z}$  to equal  $\mathbf{X}\mathbf{V}_k = \mathbf{U}_k\mathbf{\Sigma}_k$  and  $\mathbf{Y} = \mathbf{V}_k^T$ .

Could have just as easily set  $\mathbf{Z} = \mathbf{U}_k$  and  $\mathbf{Y} = \mathbf{\Sigma}_k\mathbf{V}_k^T$ , so  $\mathbf{Z}$  has orthonormal columns.

# Word embeddings

Another view on word embeddings from LSA:

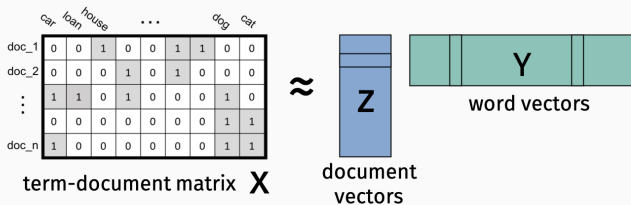


- $X \approx ZY$
- $X^T X \approx Y^T Z^T Z Y = Y^T Y$
- So for  $word_i$  and  $word_j$ ,  $\langle \mathbf{y}_i, \mathbf{y}_j \rangle \approx [X^T X]_{i,j}$ .

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



# Word embeddings



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

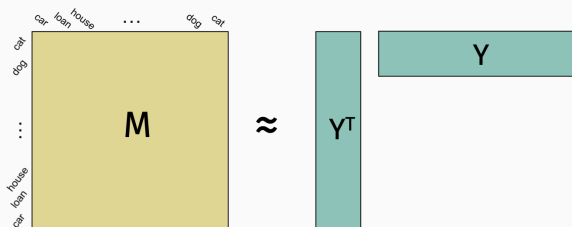
# Word embeddings

$\langle \mathbf{y}_i, \mathbf{y}_j \rangle$  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 similarity matrix  $\mathbf{M} \in \mathbb{R}^{n \times n}$  with  $\mathbf{M}_{i,j} = k(word_i, word_j)$ .
3. Find low rank approximation  $\mathbf{M} \approx \mathbf{Y}^T \mathbf{Y}$  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 3, 4, or 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.

# Modern word embeddings

Computing word similarities for “window size” 4:

The girl walks to her **dog to the park.**  
It can take a long time to park your car in NYC.  
**The dog park is** always crowded on Saturdays.

The girl walks to her dog to the park.  
It can take a long time to park your car in NYC.  
The dog **park is always crowded** on Saturdays.

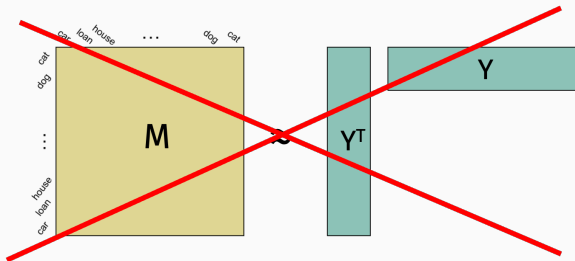
The girl walks to **her dog to the park.**  
It can take a long time to park your car in NYC.  
**The dog park is** always crowded on Saturdays.

	dog	park	crowded	the
dog	0	2	0	3
park	2	0	1	2
crowded	0	1	0	0
the	3	2	0	0

**Current state of the art models:** GloVe, word2vec.

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

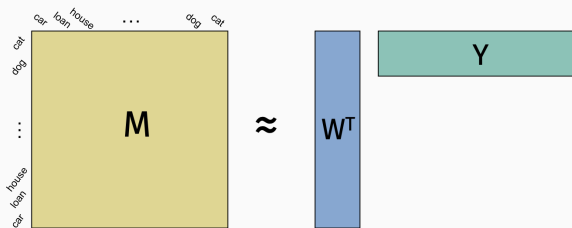
## Caveat about factorization



SVD will not return a symmetric factorization in general. In fact, if  $M$  is not positive semidefinite<sup>3</sup> then the optimal low-rank approximation does not have this form.

<sup>3</sup>I.e.,  $k(\text{word}_i, \text{word}_j)$  is not a positive semidefinite kernel.

## Caveat about factorization



- For each word  $i$  we get a left and right embedding vector  $\mathbf{w}_i$  and  $\mathbf{y}_i$ . It's reasonable to just use one or the other.
- If  $\langle \mathbf{y}_i, \mathbf{y}_j \rangle$  is large and positive, we expect that  $\mathbf{y}_i$  and  $\mathbf{y}_j$  have similar similarity scores with other words, so they typically are still related words.
- Another option is to use as your features for a word the concatenation  $[\mathbf{w}_i, \mathbf{y}_i]$

## Easiest way to use word embeddings

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

- **Original gloVe website:**

`https://nlp.stanford.edu/projects/glove/`

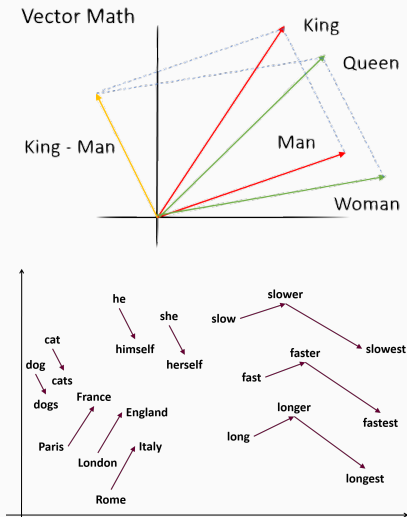
- **Compilation of many sources:**

`https://github.com/3Top/word2vec-api`



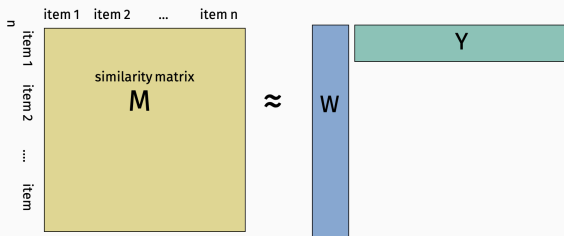
# Word embeddings math

Lots of cool demos online for what can be done with these embeddings. E.g. “vector math” to solve analogies.



## Semantic embeddings

The same approach used for word embeddings can be used to obtain meaningful numerical features for any other data where there is a natural notion of similarity.

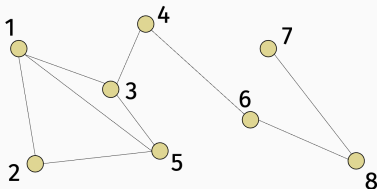


For example, the items could be nodes in a social network graph. Maybe we want to predict an individual's age, level of interest in a particular topic, political leaning, etc.

# Node embeddings



Generate random walks (e.g. “sentences” of nodes) and measure similarity by node co-occurrence frequency.



1, 3, 4, 4, 5, 2, 1, 2, 5  
6, 8, 6, 4, 3, 1, 5, 3, 4  
7, 8, 6, 8, 7, 8, 6, 8, 6  
⋮  
4, 6, 8, 6, 4, 3, 1, 2, 5

## Node embeddings

Again typically normalized and apply a non-linearity (e.g. log) as in word embeddings.

1, 3, 4, 4, 5, 2, 1, 2, 5  
6, 8, 6, 4, 3, 1, 5, 3, 4  
7, 8, 6, 8, 7, 8, 6, 8, 6  
⋮  
4, 6, 8, 6, 4, 3, 1, 2, 5

	node 1	node 2	...	node 8
node 1	0	2		1
node 2	2	0		0
⋮				
node 8	1	0		0

Popular implementations: DeepWalk, Node2Vec. Again initially derived as simple neural network models, but are equivalent to matrix-factorization (Qiu et al. 2018).