

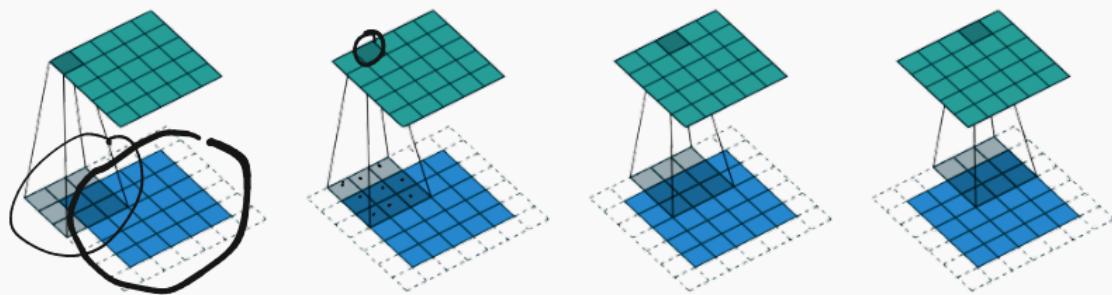
CS-GY 6923: Lecture 12

Finish Convolutional Networks, Adversarial Examples, Autoencoders

NYU Tandon School of Engineering, Prof. Christopher Musco

RECALL FROM LAST LECTURE

Common way of processing images, time series, audio, etc. is via convolution with a filter:



Can perform operations like smoothing, template matching, edge detection, etc.

EDGE DETECTION



I_C

$*$ $\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$



E_C



I_L

$*$ $\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$



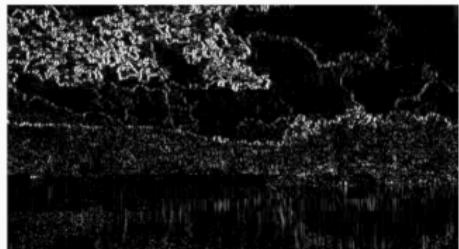
E_L

EDGE DETECTION + PATTERN MATCHING

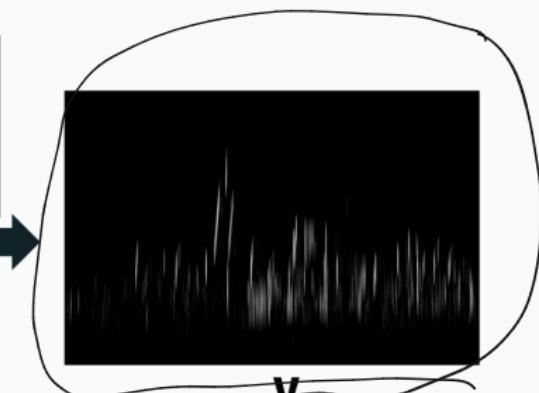
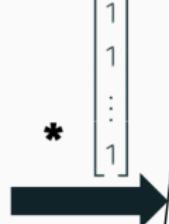
Feed edge detection result into pattern matcher that looks for long vertical lines.



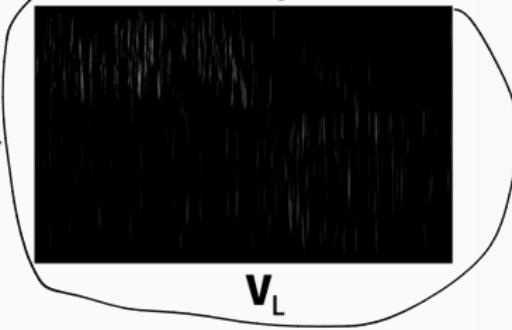
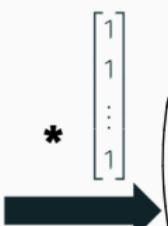
E_C



E_L

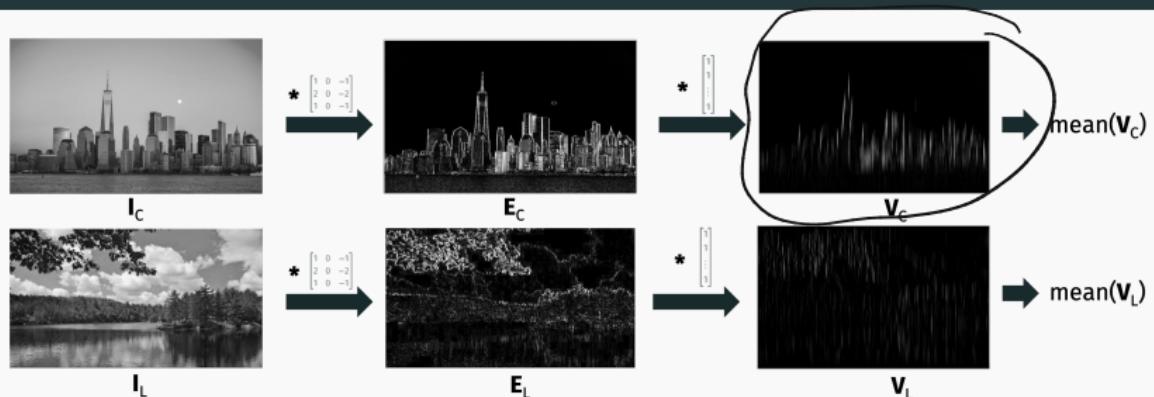


V_C



V_L

HIERARCHICAL CONVOLUTIONAL FEATURES



$$\text{mean}(V_C) = .062 \quad \text{vs.} \quad \text{mean}(V_L) = .054$$

The image with highest average response to (edge detector) + (vertical pattern) is the city scape.

$\text{mean}(V) = \underline{V}^T \underline{\beta}$ where $\underline{\beta} = [1/n, \dots, 1/n]$. So the new features in V could be combined with a simple linear classifier to separate cityscapes from landscapes.

HIERARCHICAL CONVOLUTIONAL FEATURES

Hierarchical combinations of simple convolution filters are
very powerful for understanding images.

Edge detection seems like a critical first step.

Lots of evidence from biology.

VISUAL SYSTEM

Light comes into the eye through the lens and is detected by an array of photosensitive cells in the **retina**.

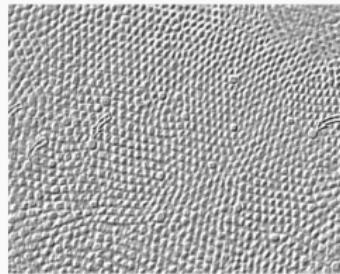
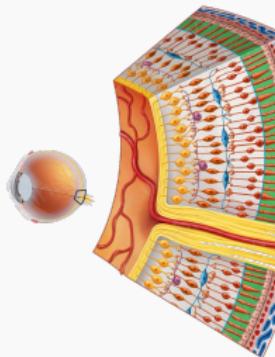
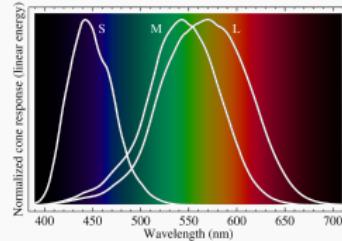


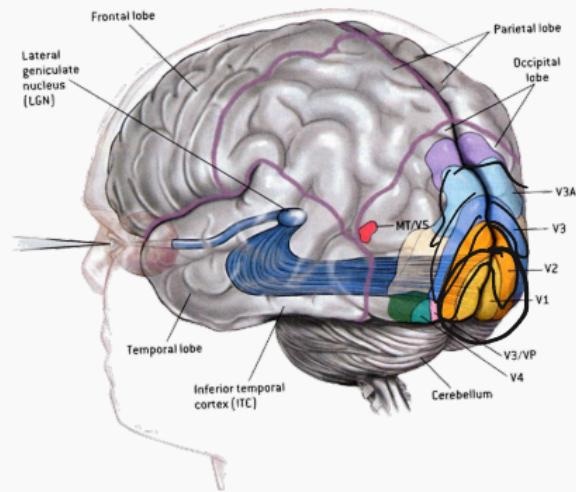
Fig. 13. Tangential section through the human fovea.
Larger cones (arrows) are blue cones. From Ahnelt et al. 1987.

Rod cells are sensitive to all light, larger **cone** cells are sensitive to specific colors. We have three types of cones:



VISUAL SYSTEM

Signal passes from the retina to the primary (V1) visual cortex, which has neurons that connect to higher level parts of the brain.

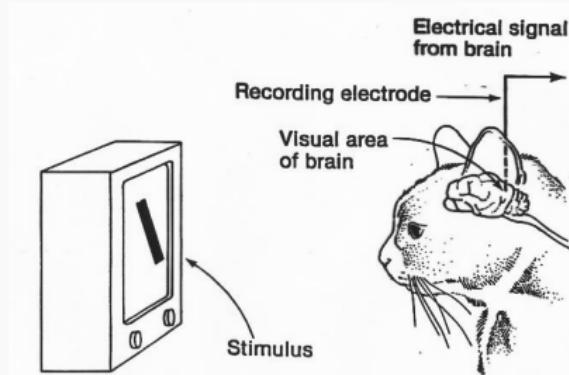


What sort of processing happens in the primary cortex?

Lots of edge detection!

EDGE DETECTORS IN CATS

Huber + Wiesel, 1959: "Receptive fields of single neurones in the cat's striate cortex." Won Nobel prize in 1981.



Different neurons fire when the cat is presented with stimuli at different angles. Cool video at

<https://www.youtube.com/watch?v=0GxVfKJqX5E>.

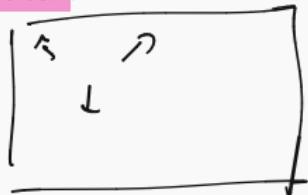
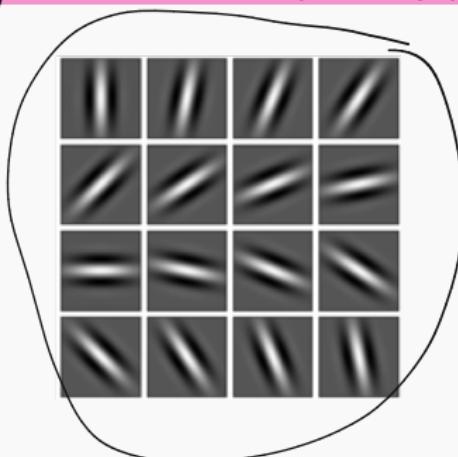
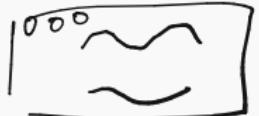
"What the Frog's Eye Tells the Frog's Brain", Lettvin et al. 1959. Found explicit edge detection circuits in a frogs visual cortex.

EXPLICIT FEATURE ENGINEERING

State of the art until 13 years ago:

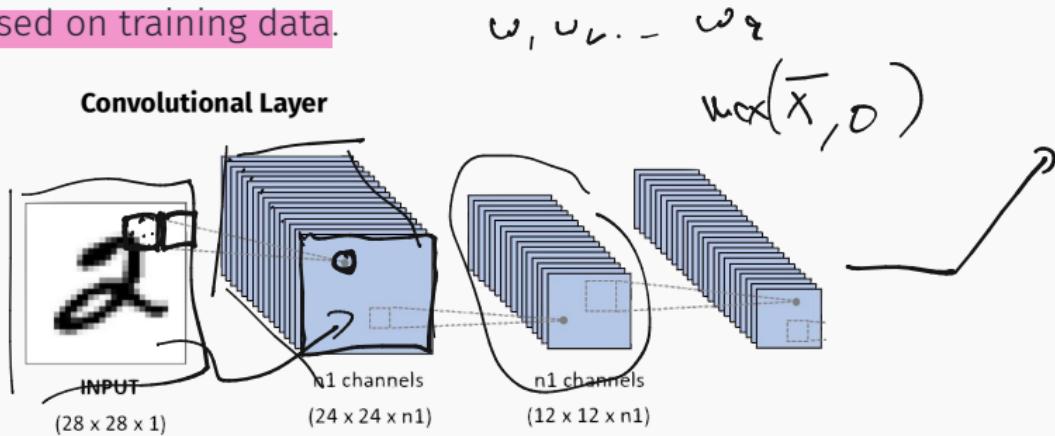
- Convolve image with edge detection filters at many different angles.
- Hand engineer features based on the responses.
- SIFT and HOG features were especially popular.

Histograms of Oriented
Features



CONVOLUTIONAL NEURAL NETWORKS

Neural network approach: Learn the parameters of the convolution filters based on training data.



First convolutional layer involves n convolution filters W_1, \dots, W_q .

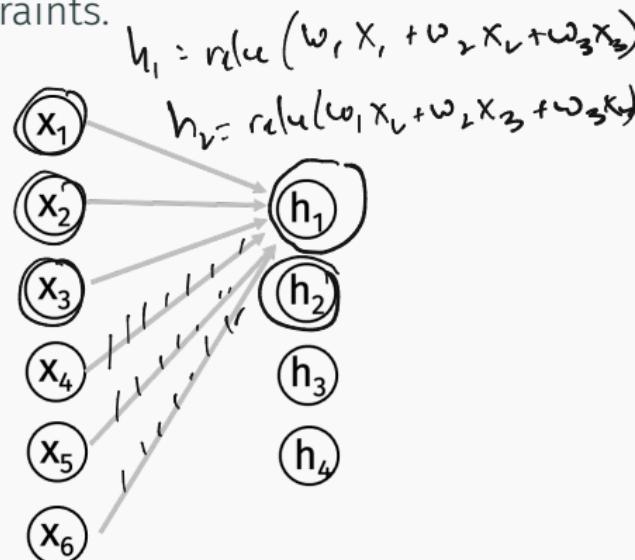
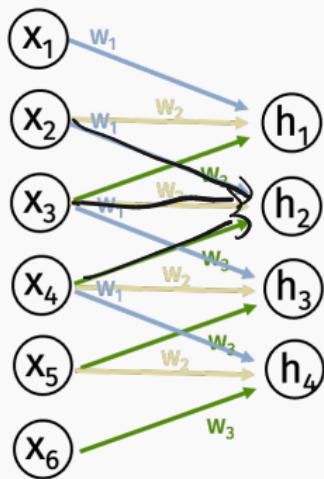
Each is small, e.g. 5×5 . Every entry in W_i is a free parameter: $\sim 25 \cdot q$ parameters to learn.

Produces q matrices of hidden variables: i.e. a tensor with depth q .

Each output in the tensor is processed with a non-linearity. Most commonly a Rectified Linear Unity (ReLU): $x = \max(\bar{x}, 0)$.

WEIGHT SHARING

Convolutional layers can be viewed as fully connected layers with added constraints. Many of the weights are forced to 0 and we have (weight sharing) constraints.



Weight sharing needs to be accounted for when running backprop/gradient descent.

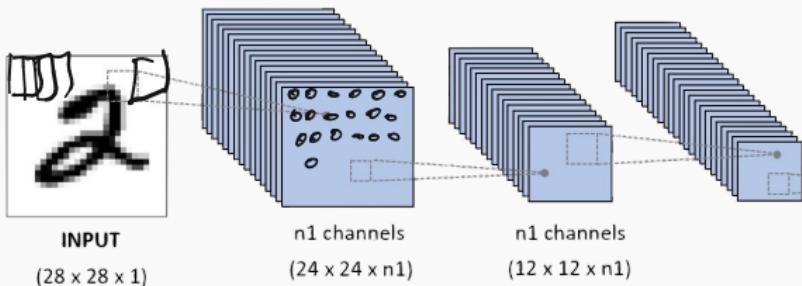
CONVOLUTIONAL NEURAL NETWORKS

For a 28×28 image like MNIST, a fully connected layer that extracts the same features as q , 5×5 filters would require $(28 \cdot 28 \cdot 24 \cdot 24) \cdot q = 451,584 \cdot q$ parameters. Compare to $25q$.

By “baking in” knowledge about what type of features matter, we greatly simplify the network.

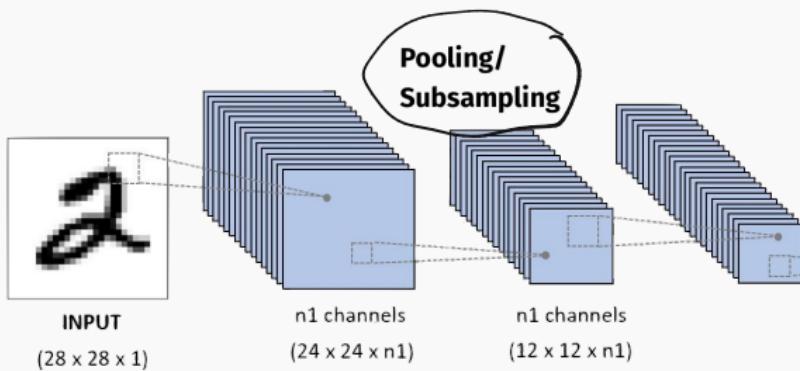
$25q$ 787
 $24 \cdot 27 \cdot 96$ total
hidden neurons

Convolutional Layer



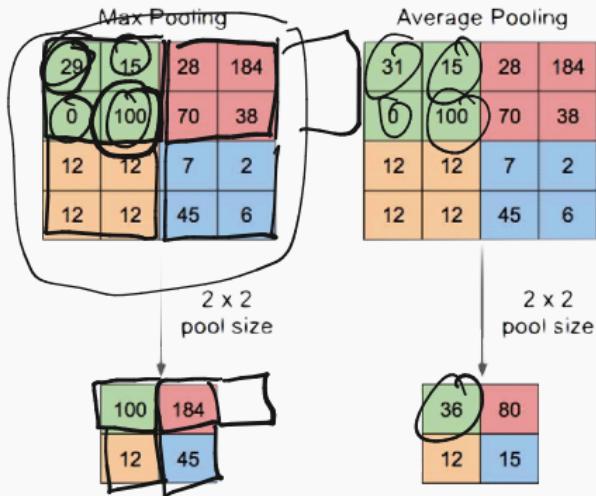
POOLING AND DOWNSAMPLING

Convolution + non-linearity are typically followed by a layer which performs **pooling + down-sampling**.



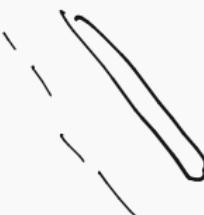
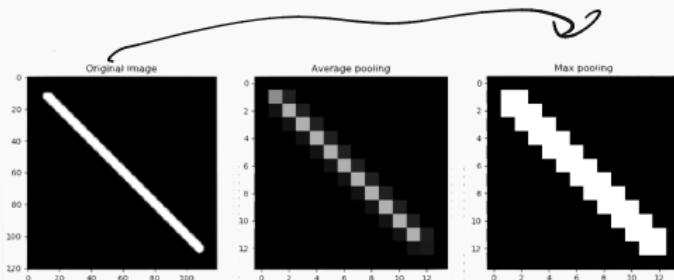
Most common approach is max-pooling.

POOLING AND DOWNSAMPLING

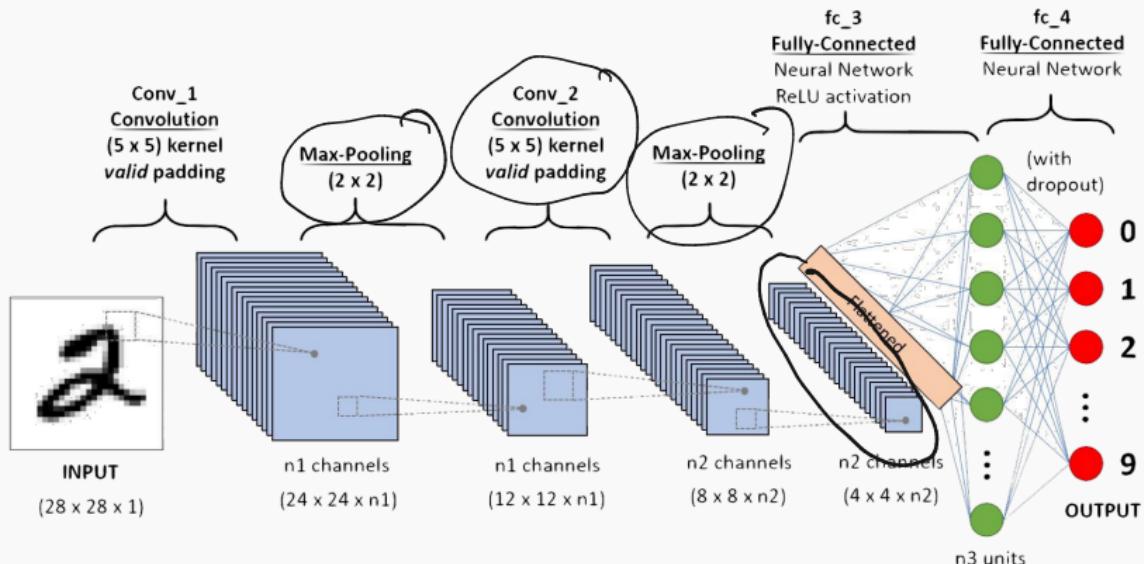


- Reduces number of variables.
- Helps “smooth” result of convolutional filters.
- Improves shift-invariance.

FFF



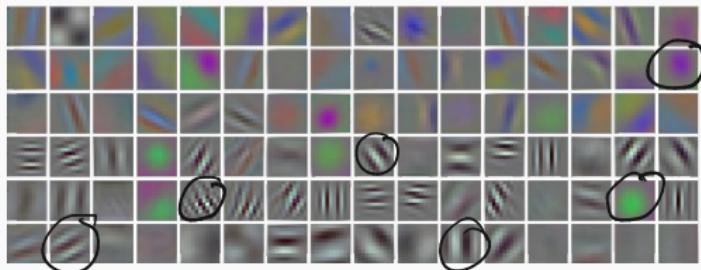
OVERALL NETWORK ARCHITECTURE



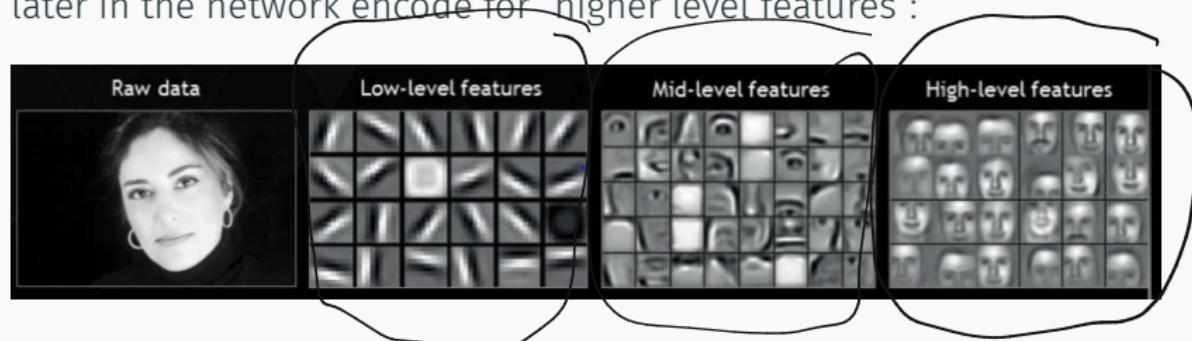
Each layer contains a 3D tensor of variables. Last few layers
are standard fully connected layers.

UNDERSTANDING LAYERS

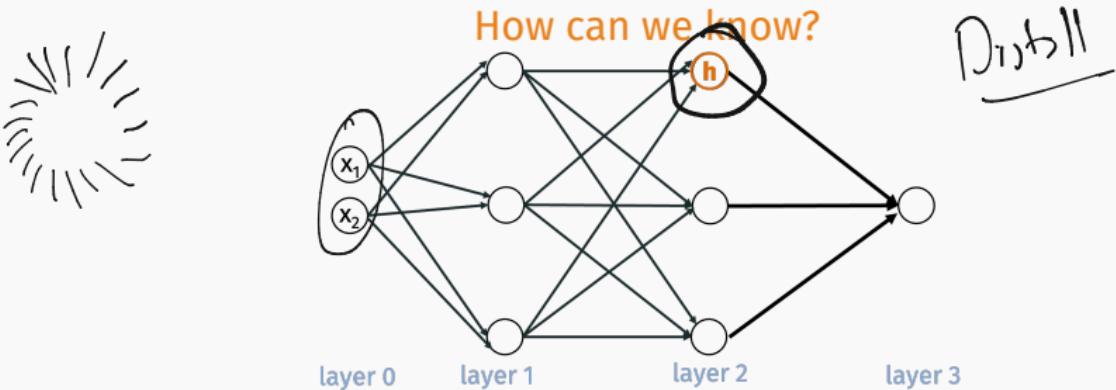
What type of convolutional filters do we learn from gradient descent?
Lots of edge detectors in the first layer!



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



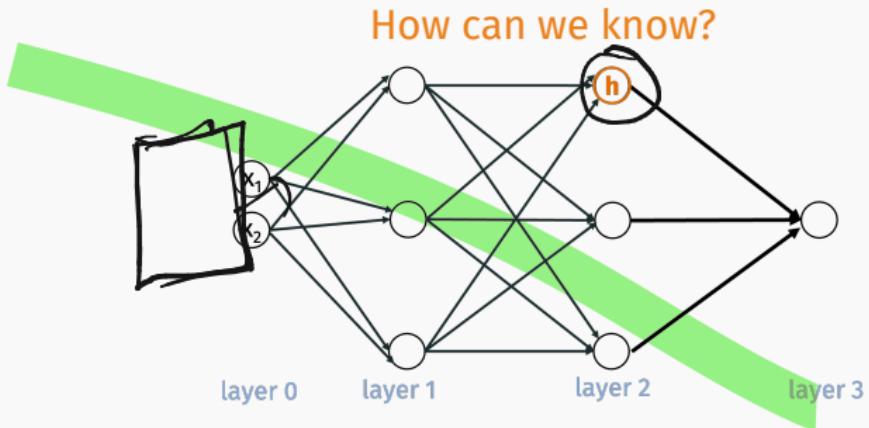
UNDERSTANDING LAYERS



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



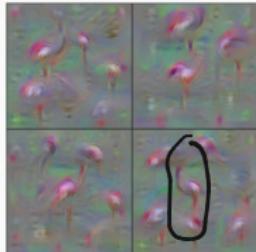
UNDERSTANDING LAYERS



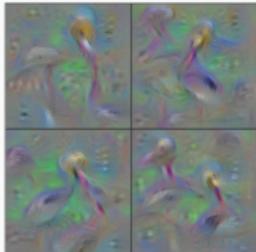
Alternative approach: Solve the optimization problem
 $\max_x |h(x)|$ e.g. using gradient descent.

UNDERSTANDING LAYERS

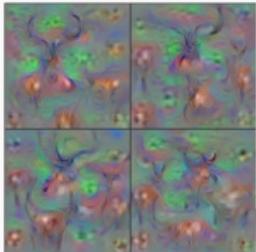
Early work had some interesting results.



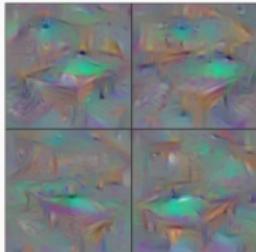
Flamingo



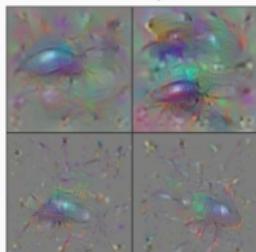
Pelican



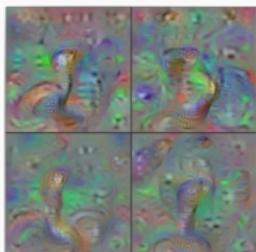
Hartebeest



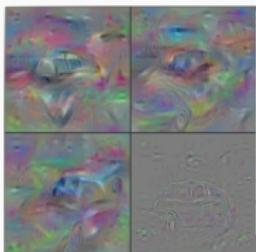
Billiard Table



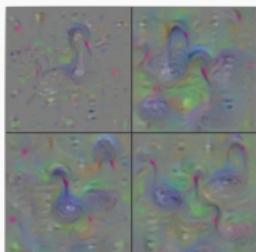
Ground Beetle



Indian Cobra



Station Wagon

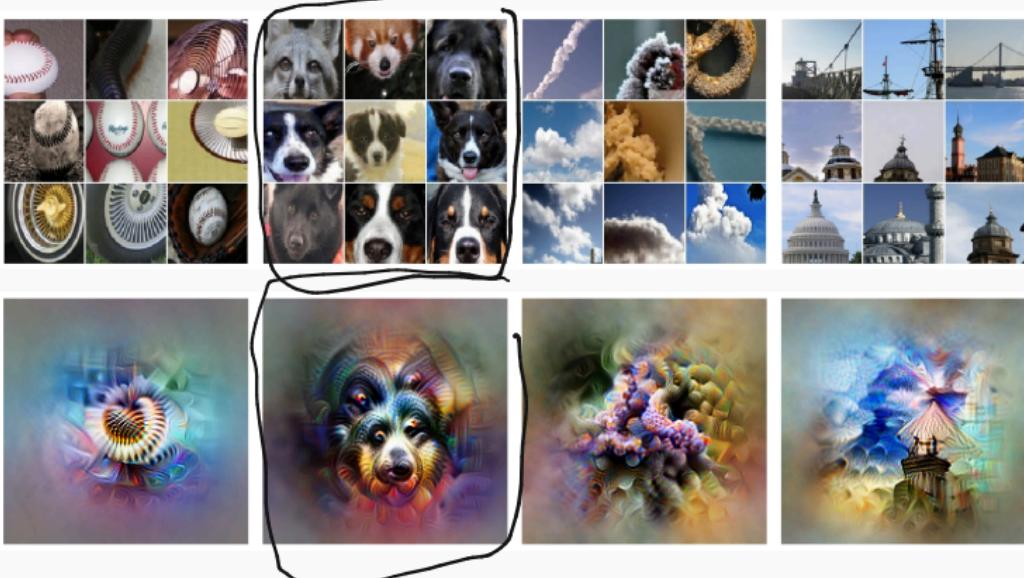


Black Swan

“Understanding Neural Networks Through Deep Visualization”, Yosinski, Clune, Nguyen, Fuchs, Lipson.

UNDERSTANDING LAYERS

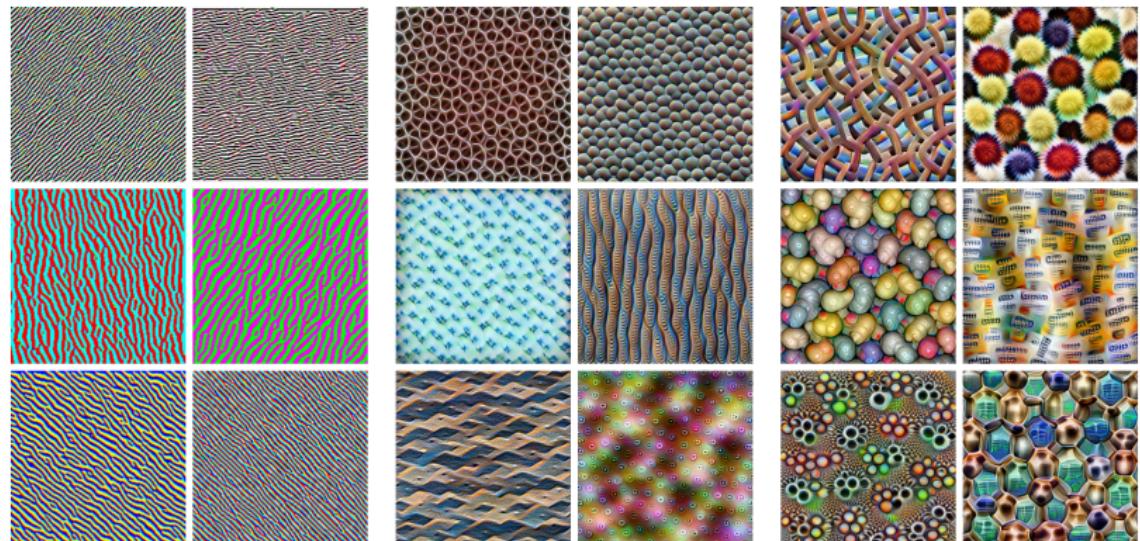
There has been a lot of work on improving these methods by regularization. I.e. solve $\max_x |h(x)| + \underline{g(x)}$ where g constrains x to look more like a “natural image”.



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.



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.

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”.

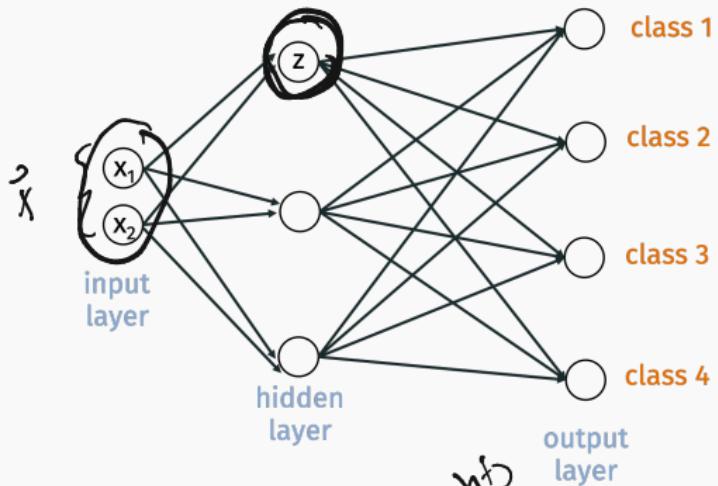
- Batch normalization (accelerate training).
- Dropout (prevent over-fitting)
- Residual connections (accelerate training, allow for more depth – 100s of layers).
- Data augmentation.



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

BATCH NORMALIZATION

Start with any neural network architecture:



For input 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 x . We can write it as $\bar{z}(x)$. Consider the mean and standard deviation of the hidden variable over our entire dataset $x_1 \dots, x_n$:

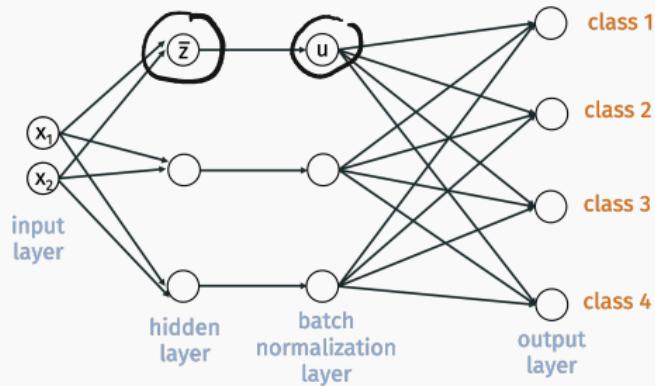
$$\mu = \frac{1}{n} \sum_{j=1}^n \bar{z}(x_j)$$

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

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 γ and c .

BATCH NORMALIZATION

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

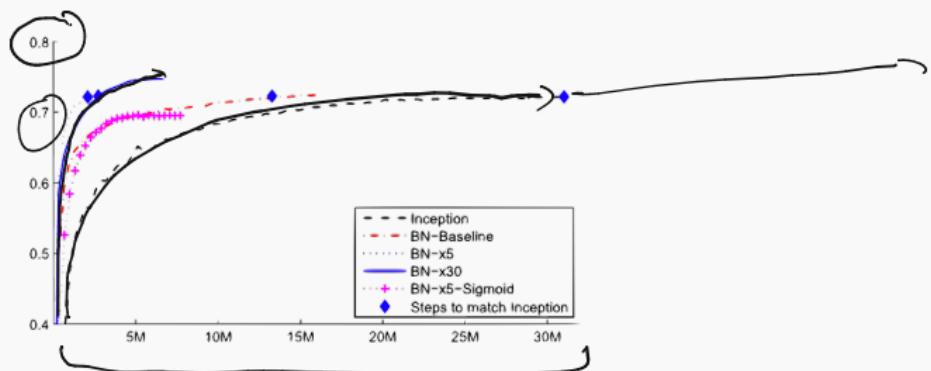
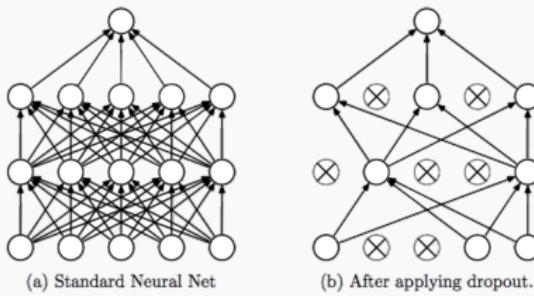


Figure 2: Single crop validation accuracy of Inception and its batch-normalized variants, vs. the number of training steps.

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.

DROPOUT

Proposed in 2012: “Dropout: A Simple Way to Prevent Neural Networks from Overfitting”, Srivastava, Hinton, Krizhevsky, Sutskever, Salakhutdinov:



During training, ignore a random subset of neurons during each gradient step. Select each neuron to be included independently with probability p (typically $p \approx .5$). During testing, no dropout is used.

DROPOUT

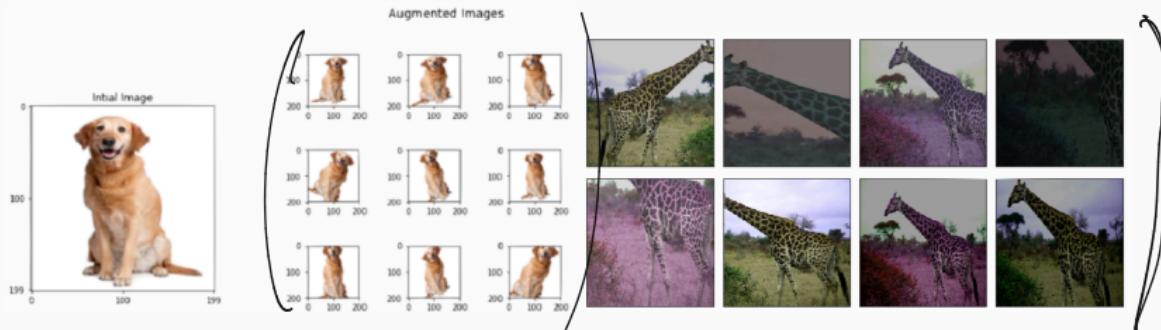
- Only used on fully connected layers.
- Simultaneously performs model regularization (model simplification) and model averaging.
- Has become less important in modern CNNs (convolutional neural nets) as the final fully connected layers become less important. But still a very helpful technique to know about!

For example, will be very helpful in avoiding overfitting in the demo on convolutional nets, since we train a pretty shallow network with the last layer doing a lot of the heavy lifting.

DATA AUGMENTATION

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

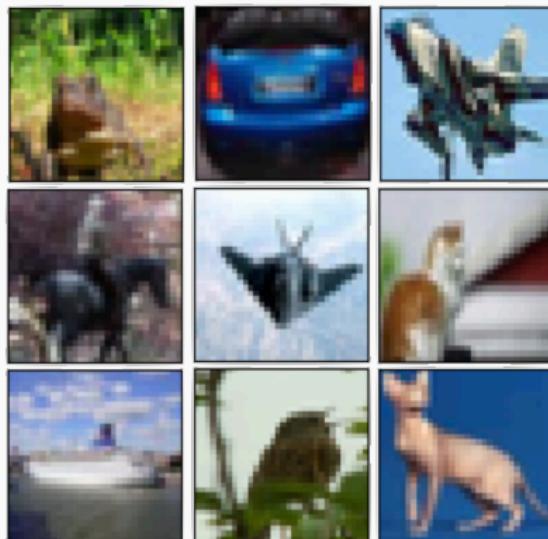
- More training data typically leads to a more accurate model.
- Artificially enlarge training data with simple transformations.



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.

DEEP LEARNING TRICKS

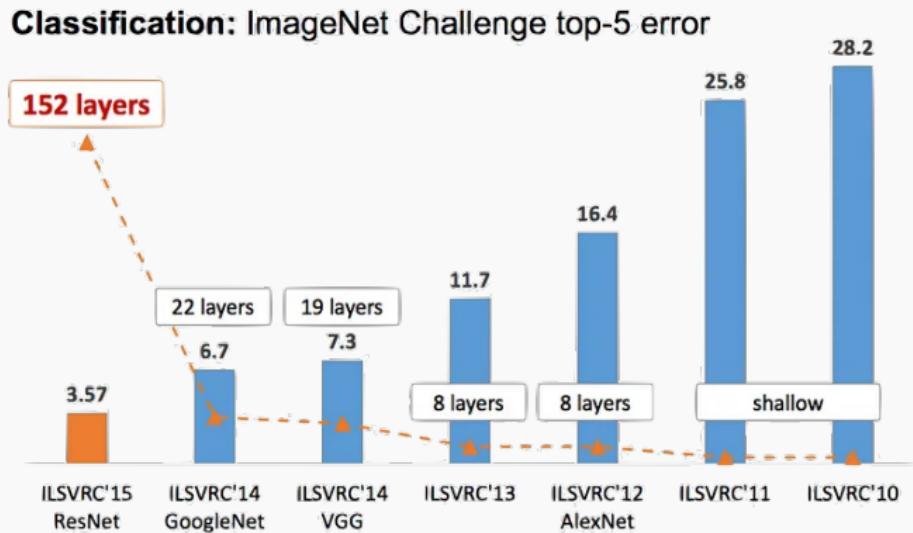
Try these techniques out in `demo_cnn_classifier.ipynb` on CIFAR-10 dataset.



airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck

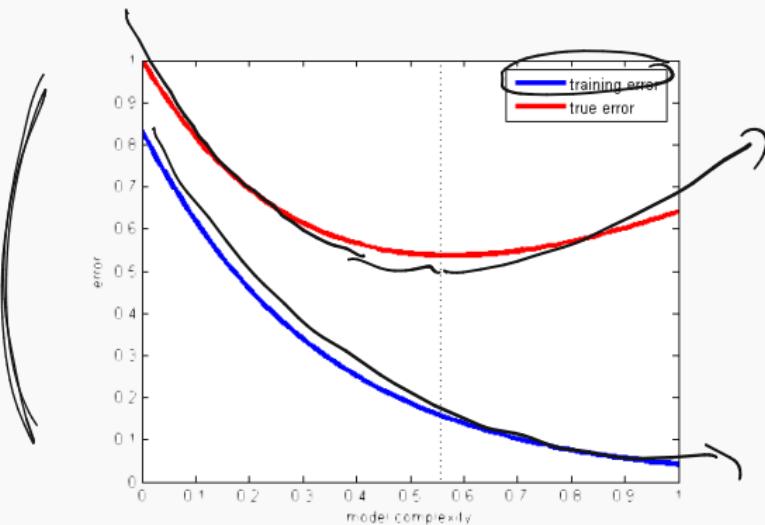
DEEPER AND DEEPER, BIGGER AND BIGGER

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



GENERALIZATION FOR NEURAL NETWORKS

Even with weight sharing, convolution, etc. modern neural networks typically have 100s of millions or billions of parameters. And we often 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 easily overfit to training data. Papers have shown that they can fit large vision data sets with random class labels to 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

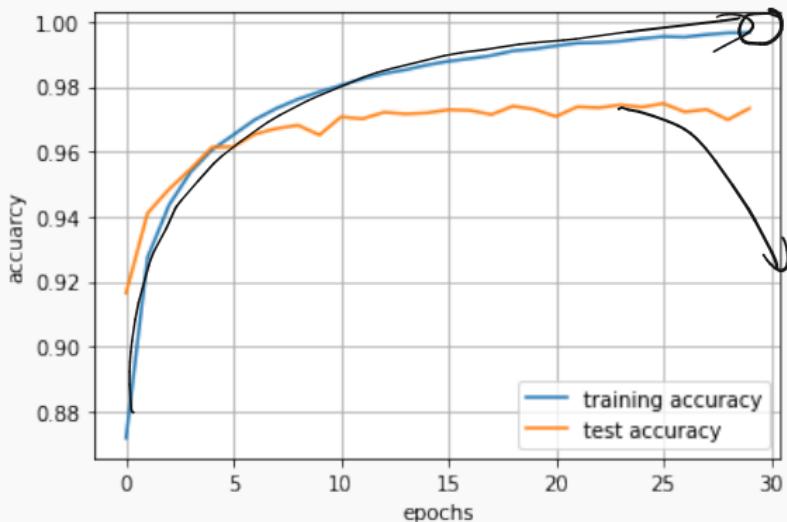
Oriol Vinyals

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. For example, overfitting is common when using fully connected networks.

GENERALIZATION FOR NEURAL NETWORKS

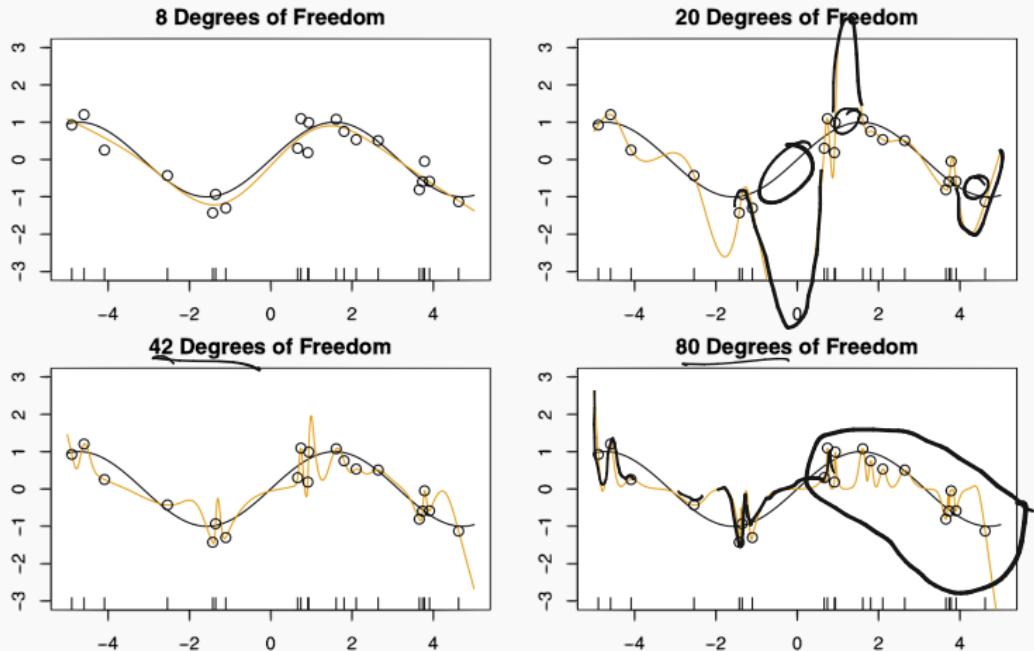
We even see this lack of overfitting for MNIST data. See `keras_demo_mnist.ipynb` that I posted on the website:



Overparameterization seems to be part of the story.

GENERALIZATION FOR NEURAL NETWORKS

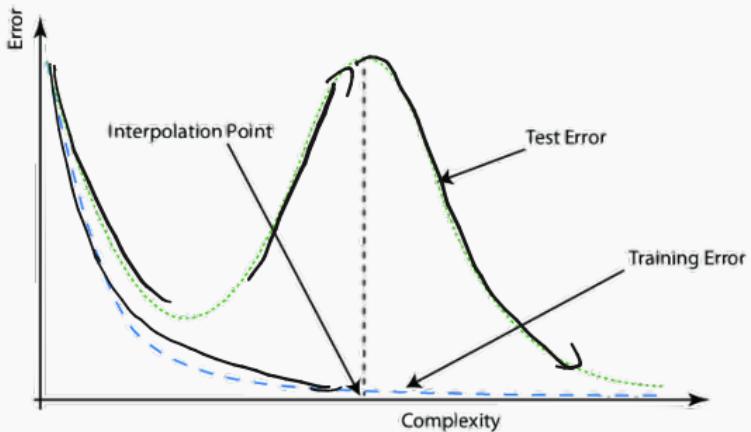
Growing realization is that this phenomena doesn't only apply to neural nets – it can also be true for overparameterized polynomials.



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.

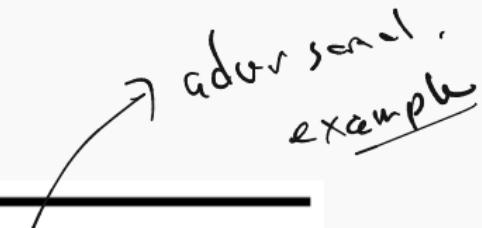


Caveat: We don't usually see this same curve for neural networks, but maybe gives some hint about what is going on.

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

Dumitru Erhan

Google Inc.

Ian Goodfellow

University of Montreal

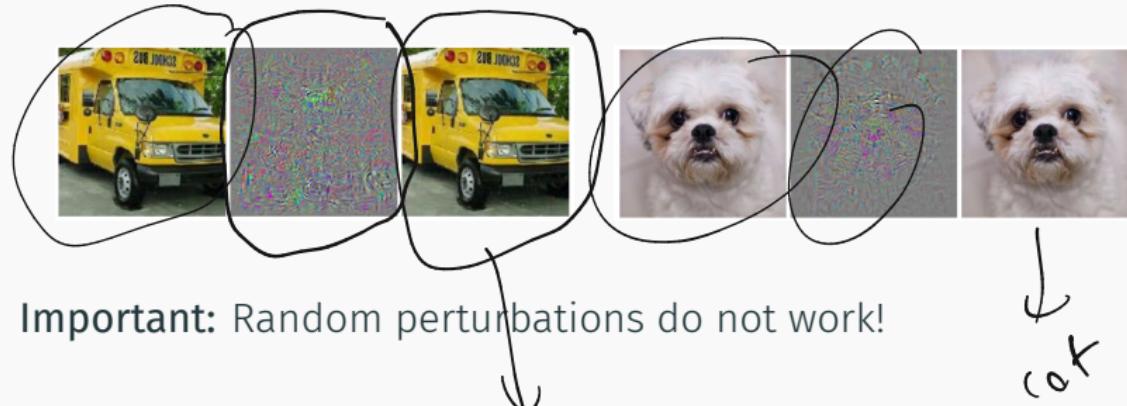
Rob Fergus

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!

horse

cat

ADVERSARIAL EXAMPLES

$$\max_{\delta} L(\theta, x + \delta, y) + \lambda \|\delta\|_2^2$$

How to find “good” perturbations:

Fix model f_θ input x , correct label y . Consider the loss $\ell(\theta, x, y)$.

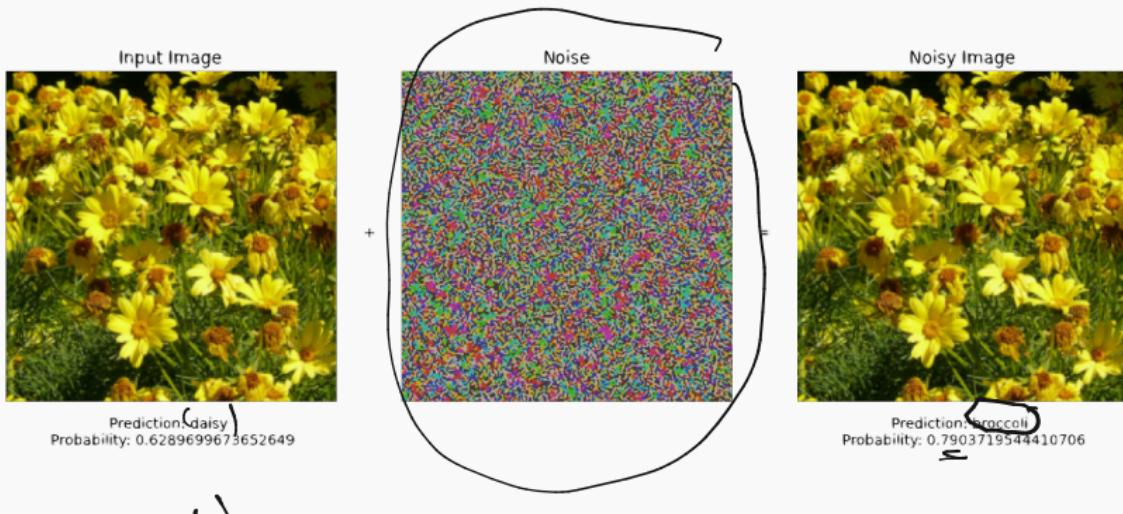
Solve the optimization problem:

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

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

We will post a lab where you can find your own adversarial examples for a model called Resnet18. The entire model + weights are available pretrained through PyTorch, so we do not need to train it ourselves.



Break until 4.05.

TRANSFER LEARNING

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 unlabeled 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,¹ Ameya Joshi,¹) Minsu Cho,¹ Kewal Jani,¹ Shivani Chiranjeevi,² Zi Kang Deng,³ Aditya Balu,² Asheesh K. Singh,² Soumik Sarkar,² Nirav Merchant,³ Arti Singh,² Baskar Ganapathysubramanian,² Chinmay Hegde¹

¹ New York University

² Iowa State University

³ University of Arizona

Aedes Vexans



Creatonotos Gangis



Daphnis Neril



Hypena Deceptalis

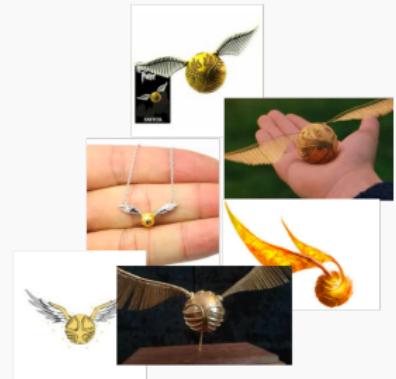


Pyralis Farinalis



ONE-SHOT LEARNING

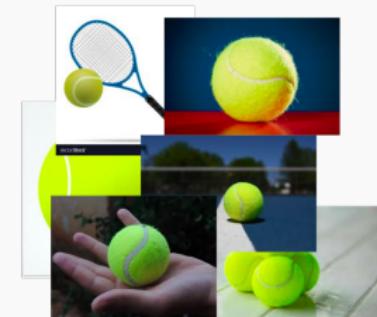
A human could probably achieve near perfect classification accuracy even given access to a single labeled example from each class:



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

TRANSFER LEARNING

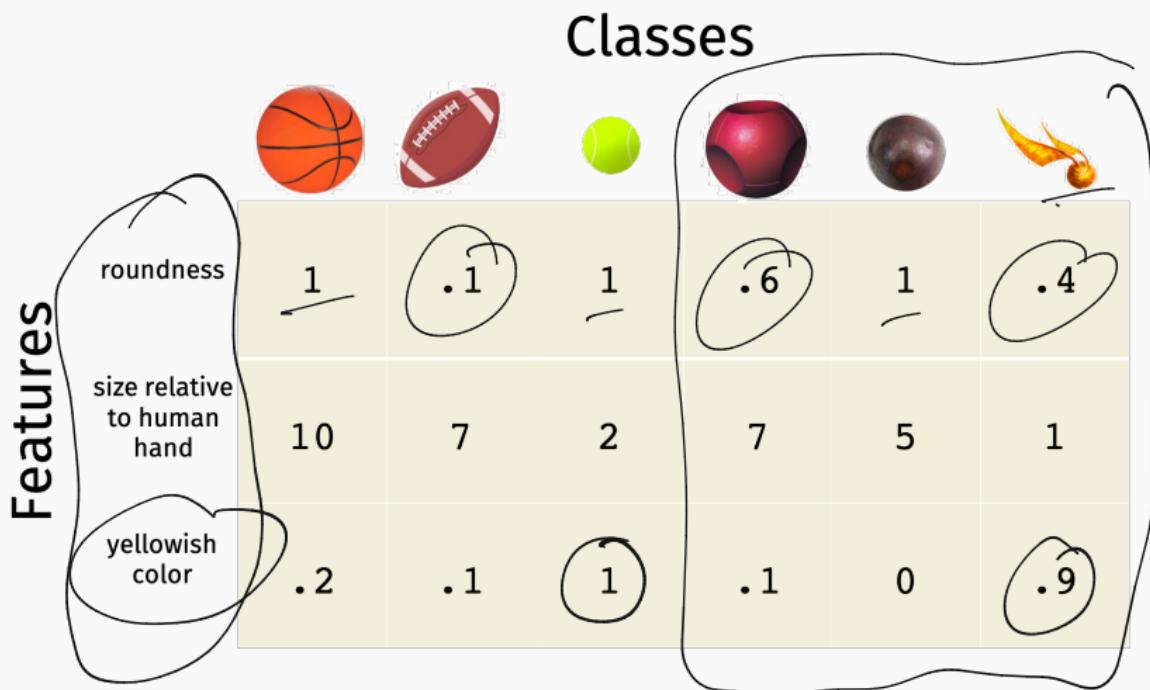
Transfer knowledge from one task we already know how to solve to another.



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.

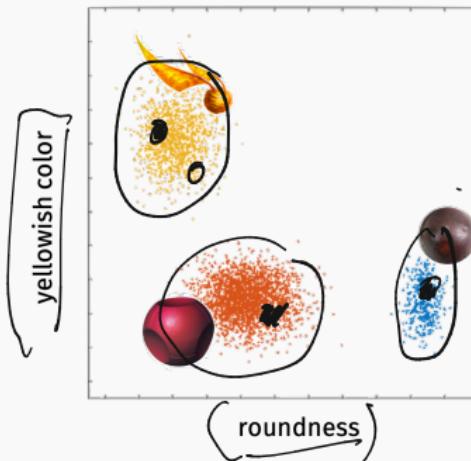
FEATURE LEARNING

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



FEATURE LEARNING

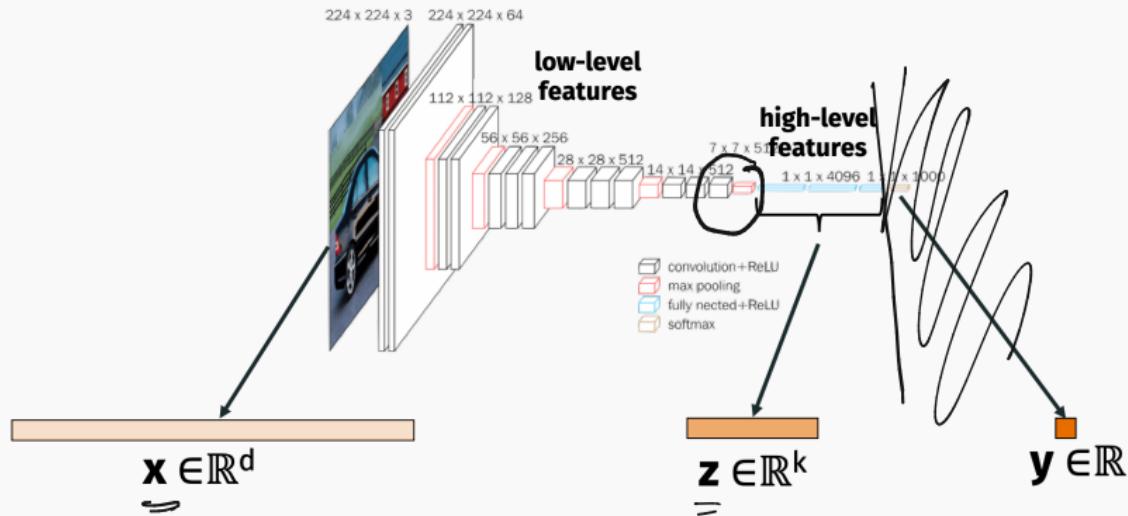
If these features are highly informative (i.e. lead to highly separable data) few training examples are needed to learn.



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 z for any new image 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)

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 10 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.

Next few lectures: How do we learn interesting features
without access to labels?

First of many simple but clever ideas: 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.

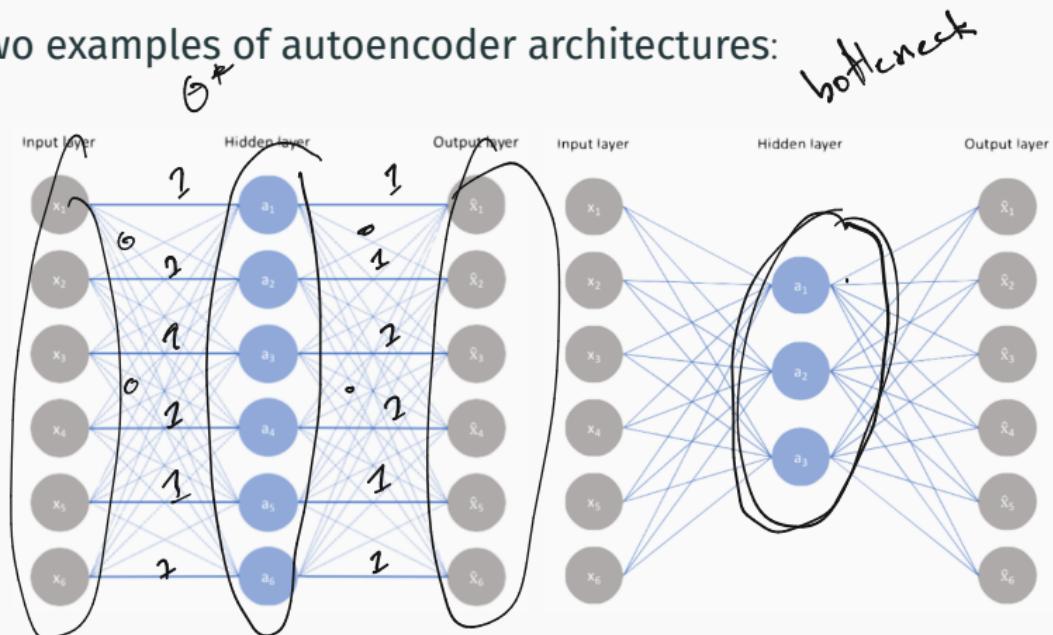
- Let $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be our model.
- Let L_{θ} be a loss function. E.g. squared loss:
$$L_{\theta}(\mathbf{x}) = \|\mathbf{x} - f_{\theta}(\mathbf{x})\|_2^2.$$
- Train model: $\theta^* = \underbrace{\min_{\theta} \sum_{i=1}^n L_{\theta}(\mathbf{x}_i)}_{\text{If } f_{\theta} \text{ is a model that incorporates feature learning, then these features can be used for supervised tasks.}}$

If f_{θ} is a model that incorporates feature learning, then these features can be used for supervised tasks.

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).

AUTOENCODER

Two examples of autoencoder architectures:

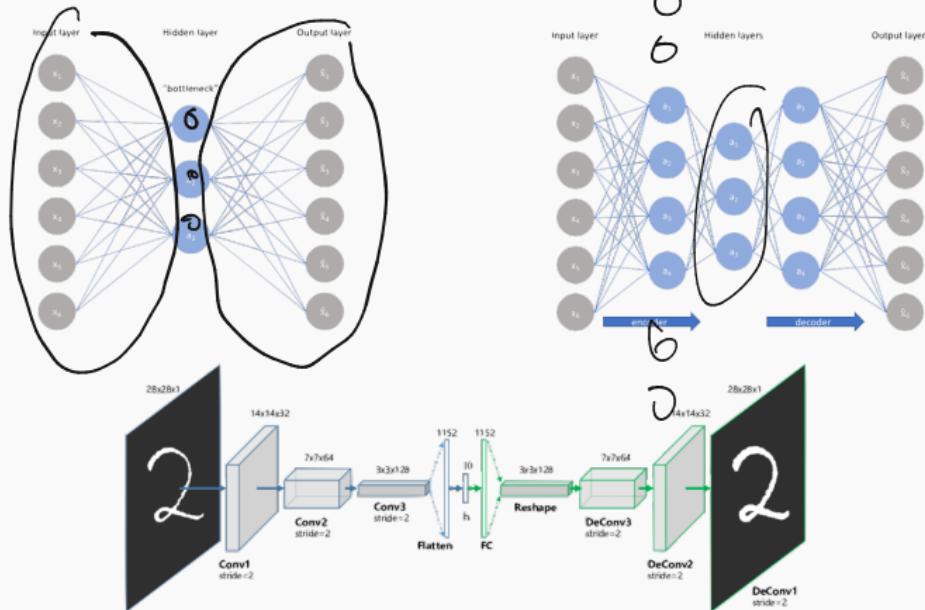


Which would lead to better feature learning?

W^{*})

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.



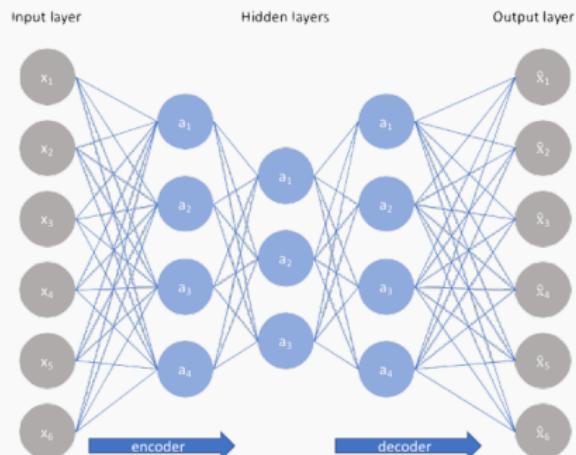
AUTOENCODER

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

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

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

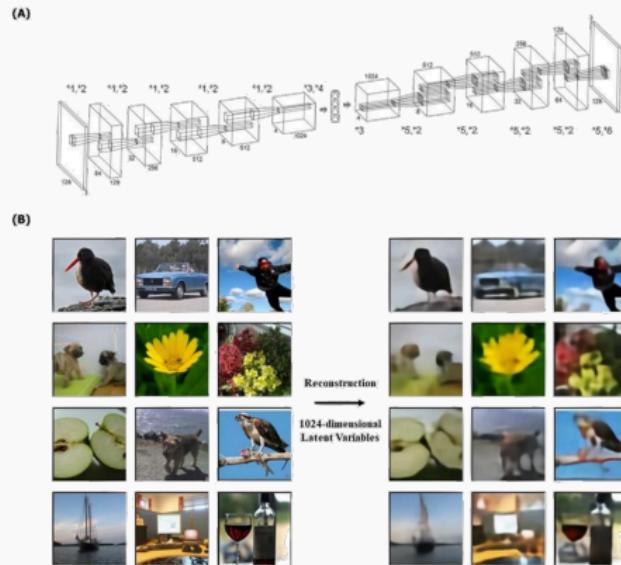
$$f(x) =$$



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: $d = 49152$.

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.¹

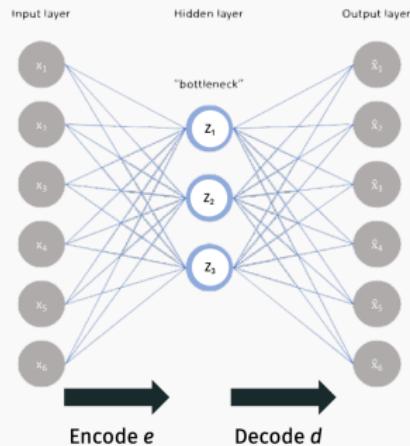
There are a lot of cool applications of autoencoders beyond feature learning!

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

¹Recent 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.



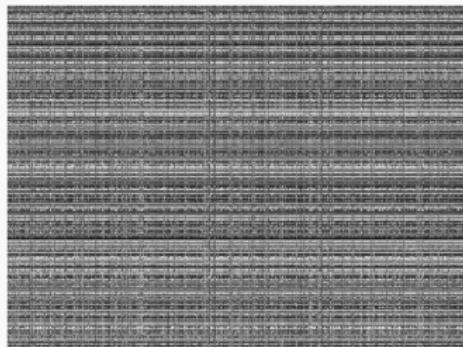
Given input x , we can completely recover $f(x)$ from $z = e(x)$. z typically has many fewer dimensions than x and for a typical image $f(x)$ will closely approximate x .

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

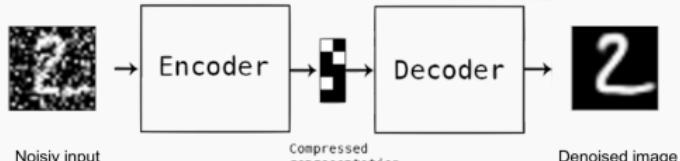


Image **denoising**

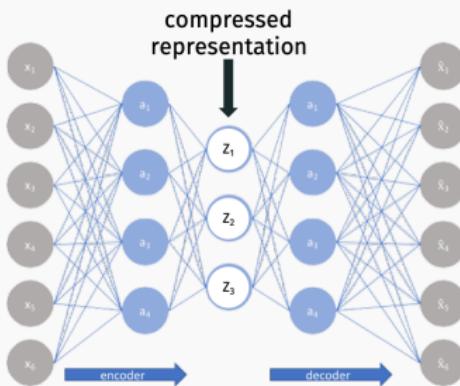


Image **inpainting**

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

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS

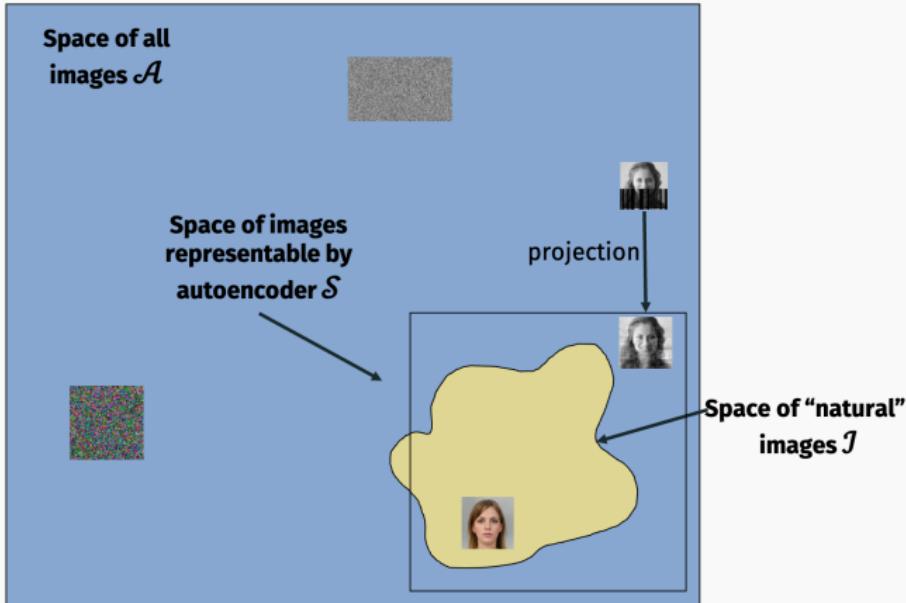
Why does this work?



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

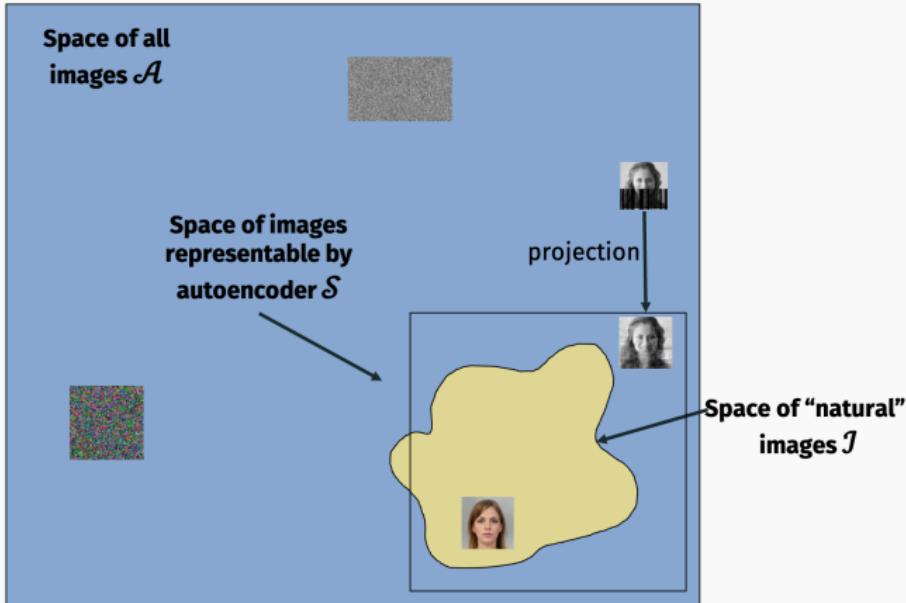
If z holds k values in $0, .1, .2, \dots, 1$, how many unique images w can be output by the autoencoder function f ?

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

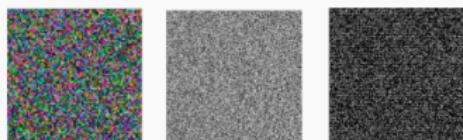


$f(x) = d(e(x))$ projects an image x closer to the space of natural images.

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} .



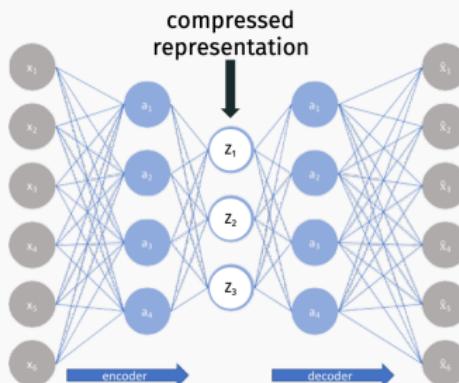
- **Option 2:** Draw \mathbf{x} randomly from \mathcal{S} , the space of images representable by the autoencoder.



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

AUTOENCODERS FOR DATA GENERATION

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

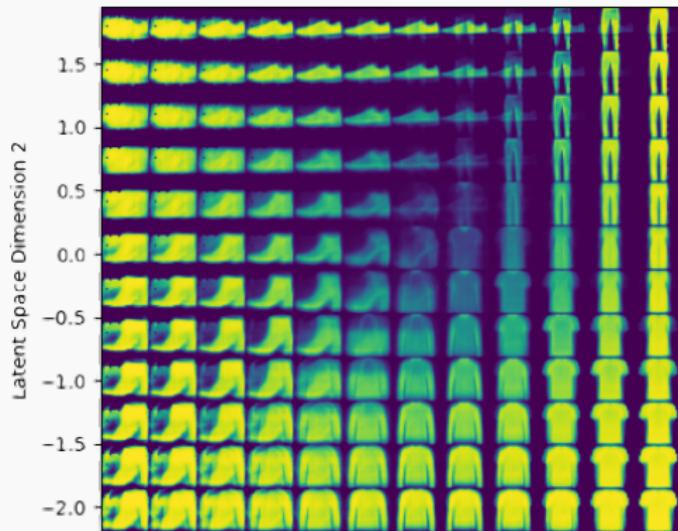
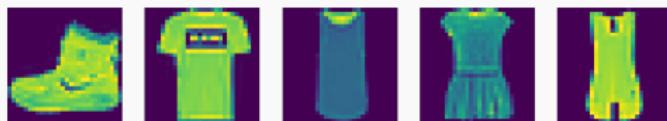


Randomly select code z , then set $x = d(z)$.²

²Lots 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

We will upload a demo on autoencoder based image generation for the "Fashion MNIST" data set:



PRINCIPAL COMPONENT ANALYSIS

PRINCIPAL COMPONENT ANALYSIS

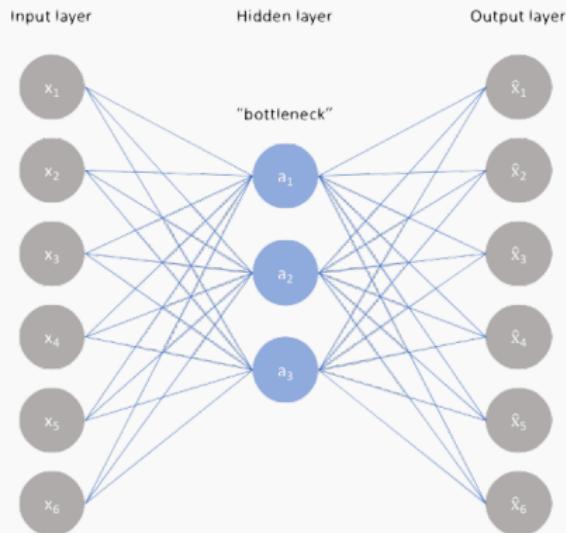
Deeper dive into understanding a simple, but powerful autoencoder architecture. Specifically we will view **principal component analysis (PCA)** as a type of autoencoder.

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

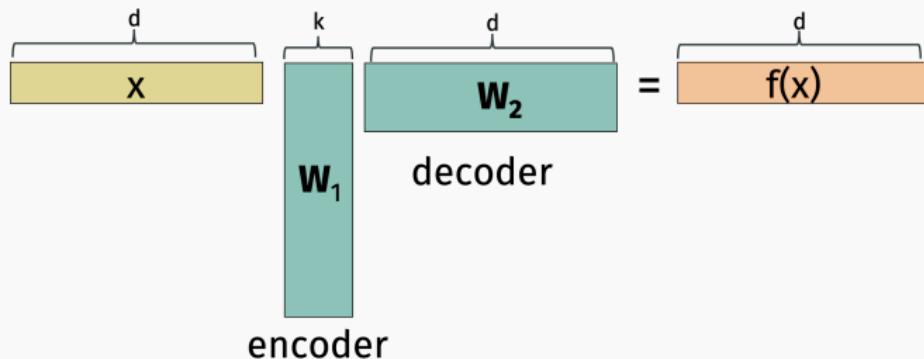
Consider the simplest possible autoencoder:



- One hidden layer. No non-linearity. No biases.
- Latent space of dimension k .
- Weight matrices are $W_1 \in \mathbb{R}^{d \times k}$ and $W_2 \in \mathbb{R}^{k \times d}$.

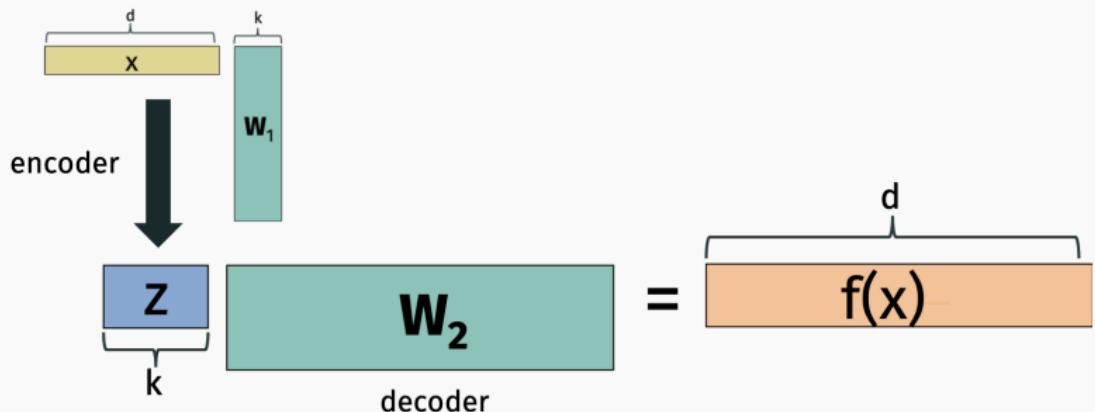
PRINCIPAL COMPONENT ANALYSIS

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



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

PRINCIPAL COMPONENT ANALYSIS

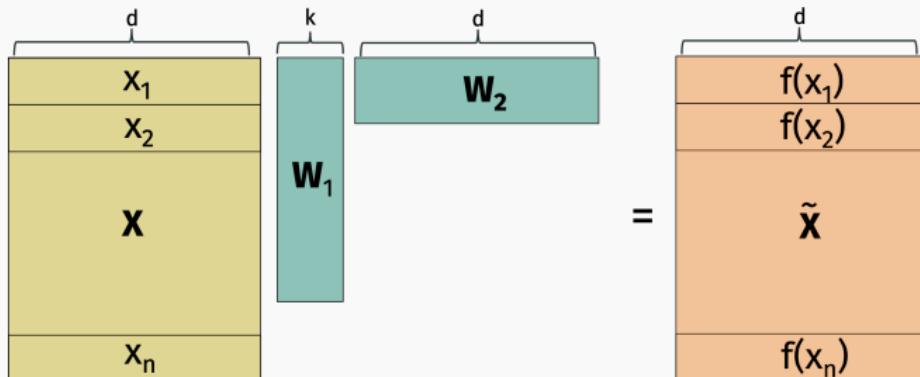


$$\text{Encoder: } e(x) = x^T w_1.$$

$$\text{Decoder: } d(z) = z w_2$$

PRINCIPAL COMPONENT ANALYSIS

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



FROBENIUS NORM

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

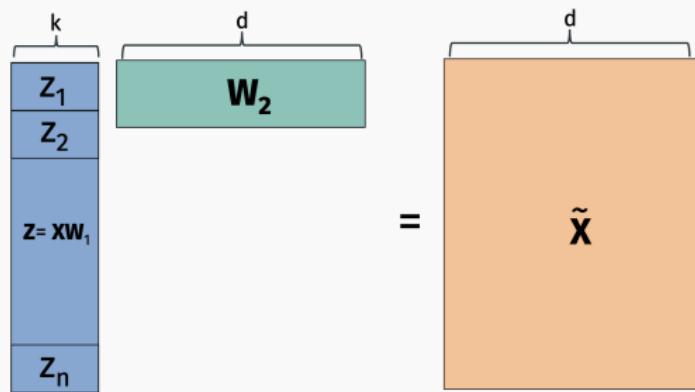
$$\begin{aligned} L(\mathbf{X}, \tilde{\mathbf{X}}) &= \sum_{i=1}^n \|\mathbf{x}_i - f(\mathbf{x}_i)\|_2^2 \\ &= \sum_{i=1}^n \sum_{j=1}^d (\mathbf{x}_i[j] - f(\mathbf{x}_i)[j])^2 \\ &= \|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2 \end{aligned}$$

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**.



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

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:

$$\begin{matrix} d \\ \text{---} \\ n \end{matrix} \quad X = \begin{matrix} \text{left singular vectors} \\ \text{---} \\ U \end{matrix} \quad \begin{matrix} \text{singular values} \\ \text{---} \\ \Sigma \end{matrix} \quad \begin{matrix} \text{right singular vectors} \\ \text{---} \\ V^T \end{matrix}$$

Where $U^T U = I$, $V^T V = I$, 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} .

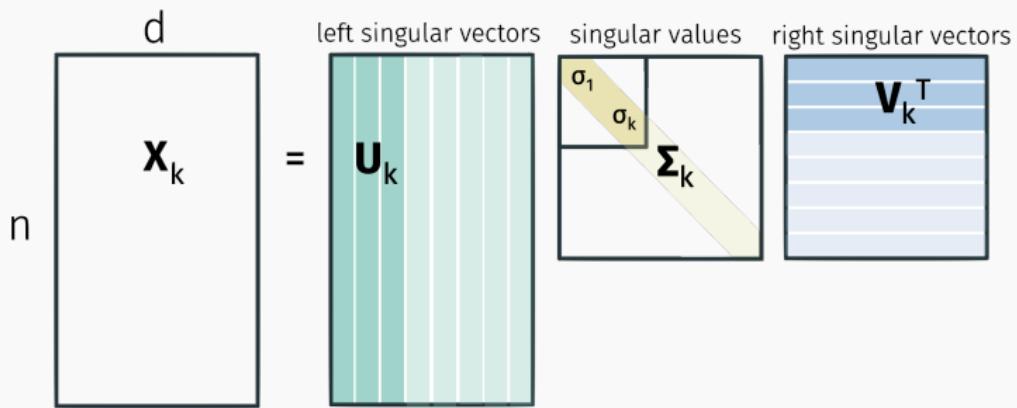
$$\mathbf{U}^\top \mathbf{U} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}$$

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

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

SINGULAR VALUE DECOMPOSITION

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



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

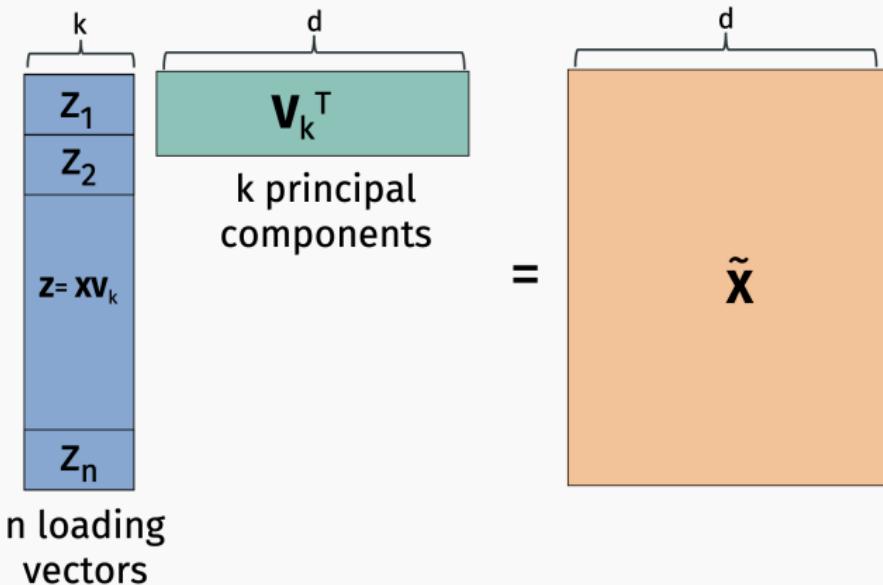
$$X_k = \arg \min_{\tilde{X} \text{ with rank } \leq k} \|X - \tilde{X}\|_F^2.$$

SINGULAR VALUE DECOMPOSITION

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

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

PRINCIPAL COMPONENT ANALYSIS



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.svds(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 $\lambda\mathbf{q} = \mathbf{M}\mathbf{q}$ for any scalar λ .

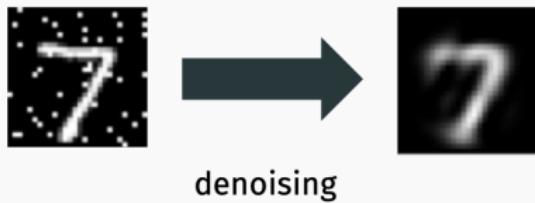
- \mathbf{U} 's columns (the left singular vectors) are the orthonormal eigenvectors of $\mathbf{X}\mathbf{X}^T$.
- \mathbf{V} 's columns (the right singular vectors) are the orthonormal eigenvectors of $\mathbf{X}^T\mathbf{X}$.
- $\sigma_i^2 = \lambda_i(\mathbf{X}\mathbf{X}^T) = \lambda_i(\mathbf{X}^T\mathbf{X})$

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



LOW-RANK APPROXIMATION

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

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

original data

rank 1 approx.

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

rank 2 approx.

0	1	0	0	1	0	0	0	0	0
0	0	0	0	0	1	0	0	0	0
0	0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	0	1
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

rank 3 approx.

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

rank 4 approx.

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

rank 5 approx.

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

rank 6 approx.

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

rank 7 approx.

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

rank 8 approx.

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

rank 9 approx.

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

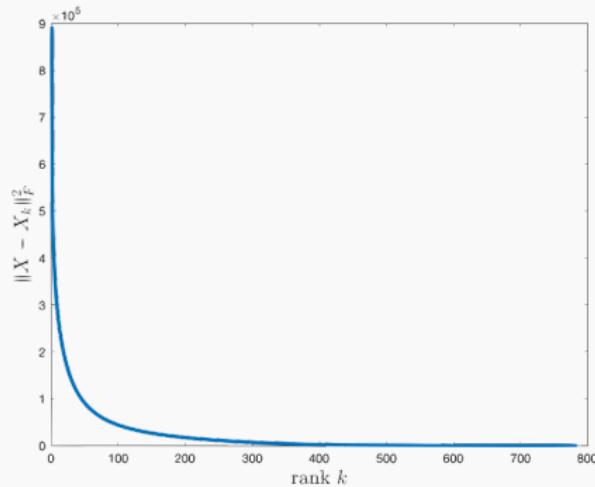
rank 50 approx.

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

LOW RANK APPROXIMATION

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

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



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

$\text{sale price} \approx 1.05 \cdot \text{list price}$.

$\text{property tax} \approx .01 \cdot \text{list price}$.

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.

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

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

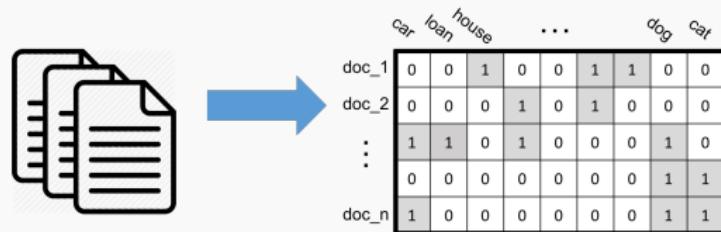
LOW-RANK INTUITION

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:

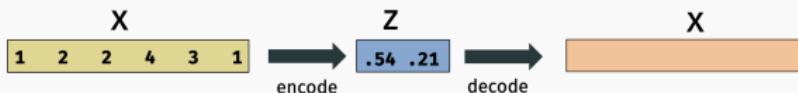


EXAMPLES OF LOW-RANK STRUCTURE

SNPs matrices tend to be very low-rank.

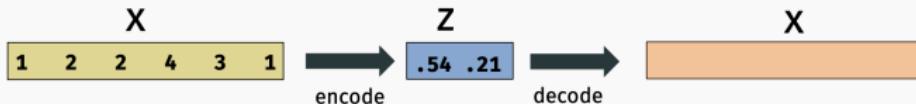
	single nucleotide polymorphisms (SNPs) loci				
	144	312	436	800	943
individual 1	A	T	T	C	G
individual 2	T	G	G	C	C
...					
individual n	C	A	T	A	G

Most of the information in 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.

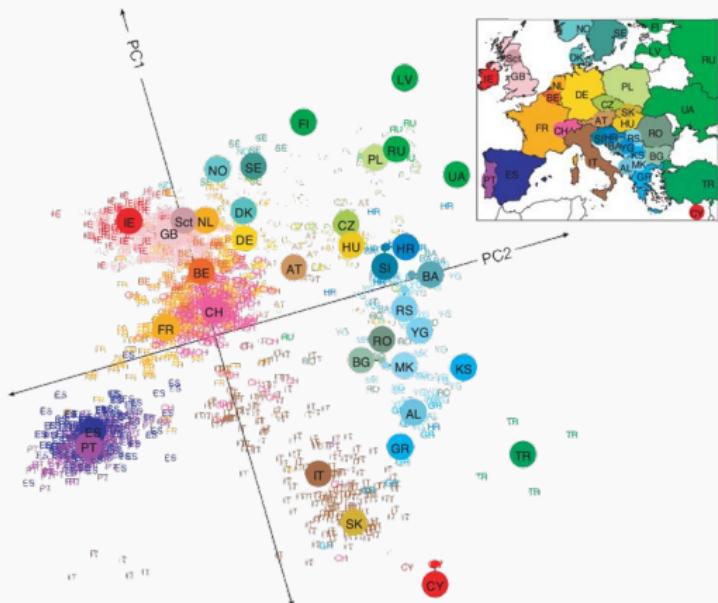
$z[1] \approx$ relative north-south position of birth place

$z[2] \approx$ 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.



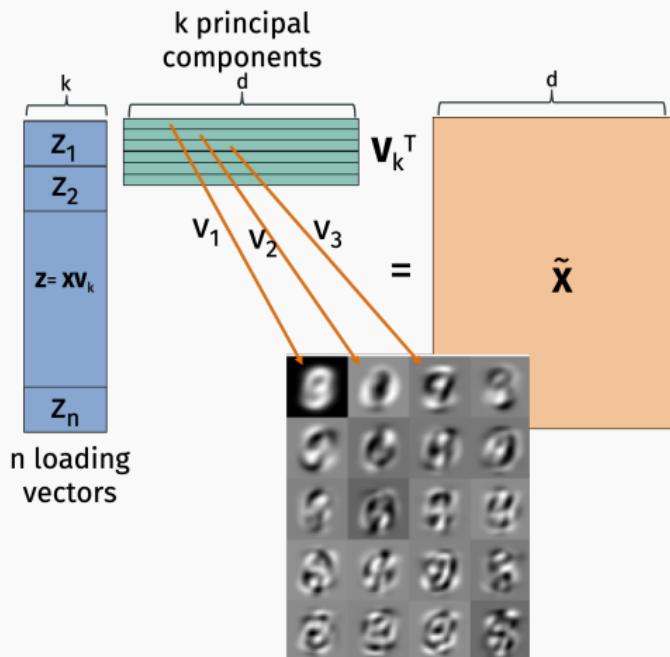
Genetic data can be nicely visualized using PCA! Plot each data example x using two loading variables in z .

PRINCIPAL COMPONENTS

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

PRINCIPAL COMPONENTS

MNIST principal components:

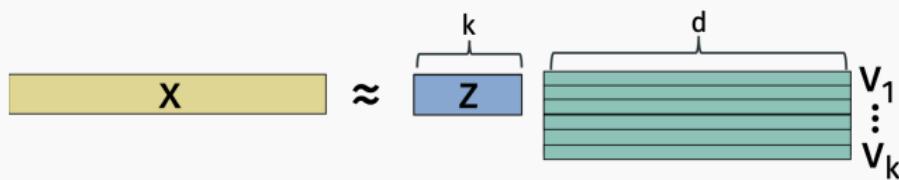


Often principal components are difficult to interpret.

LOADING VECTORS

What do the **loading vectors** looks 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} .

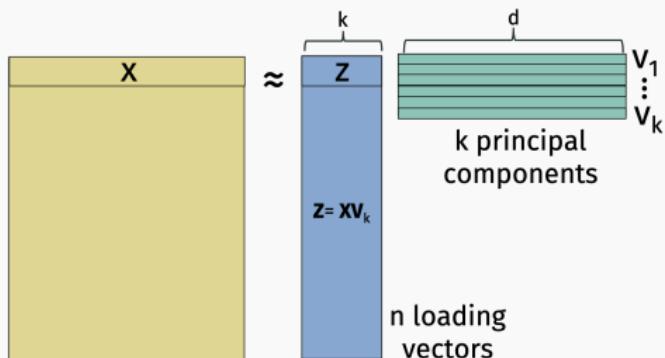


$$\mathbf{x} \approx z_1 \cdot \mathbf{v}_1 + z_2 \cdot \mathbf{v}_2 + z_3 \cdot \mathbf{v}_3 + z_4 \cdot \mathbf{v}_4 + \dots$$

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

LOADING VECTORS: SIMILARITY VIEW

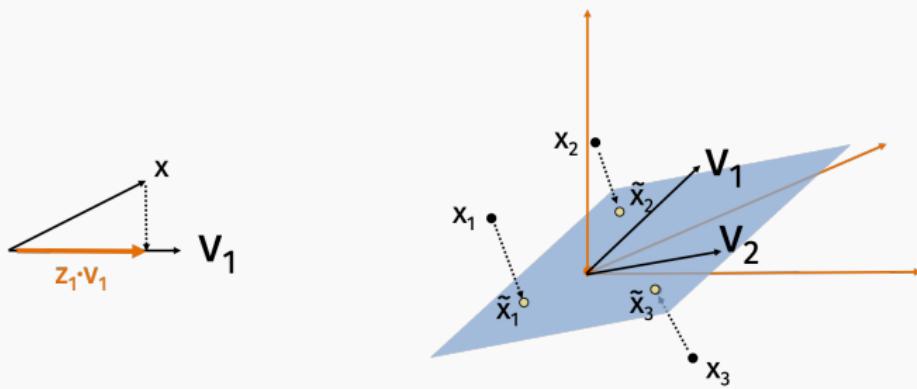
For any x with loading vector z , z_i is the inner product similarity between x and the i^{th} principal component v_i .



$$z_1 = \langle \begin{matrix} 0 \\ x \end{matrix}, \begin{matrix} 0 \\ v_1 \end{matrix} \rangle \quad z_2 = \langle \begin{matrix} 0 \\ x \end{matrix}, \begin{matrix} 0 \\ v_2 \end{matrix} \rangle \quad z_3 = \langle \begin{matrix} 0 \\ x \end{matrix}, \begin{matrix} 0 \\ v_3 \end{matrix} \rangle \dots$$

LOADING VECTORS: PROJECTION VIEW

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

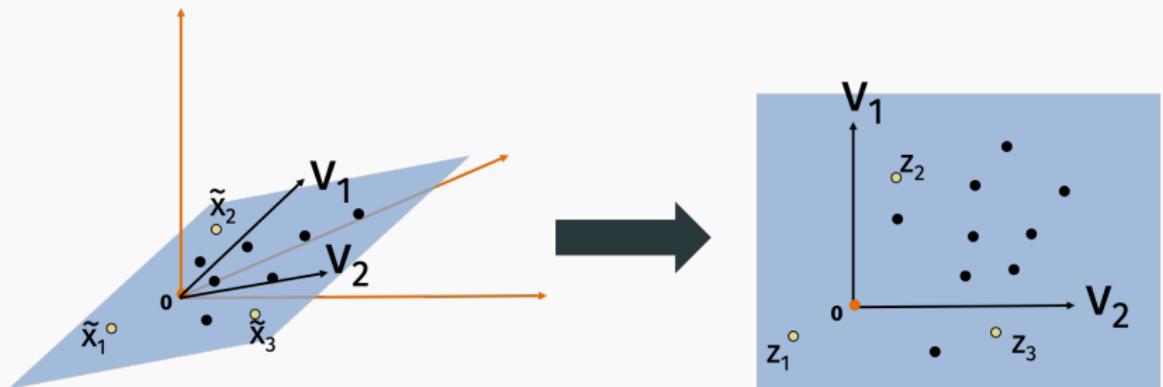


Since $\mathbf{v}_1, \dots, \mathbf{v}_k$ are orthonormal, this operation is a **projection** onto first k principal components.

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 x_i , the loading vector z_i contains the coordinates in the projection space:



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:

$$\|\tilde{\mathbf{x}}_i\|_2^2 = \|\mathbf{z}_i\|_2^2$$

$$\langle \tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j \rangle = \langle \mathbf{z}_i, \mathbf{z}_j \rangle$$

$$\|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j\|_2^2 = \|\mathbf{z}_i - \mathbf{z}_j\|_2^2$$

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

$$\|\mathbf{x}_i\|_2^2 \approx \|\mathbf{z}_i\|_2^2$$

$$\langle \mathbf{x}_i, \mathbf{x}_j \rangle \approx \langle \mathbf{z}_i, \mathbf{z}_j \rangle$$

$$\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \approx \|\mathbf{z}_i - \mathbf{z}_j\|_2^2$$

When we come back from break, will use this to motivate
semantic embeddings.