

PCA

CMSC 422

SOHEIL FEIZI

sfeizi@cs.umd.edu

Today's topics

- SGD with momentum
- Improved NN architectures
- PCA

Try different architectures and training parameters here:

<http://playground.tensorflow.org>



Epoch
000,000

Learning rate
0.03

Activation
Tanh

Regularization
None

Regularization rate
0

Problem type
Classification

DATA

Which dataset do you want to use?



Ratio of training to test data: 50%



Noise: 0



Batch size: 10



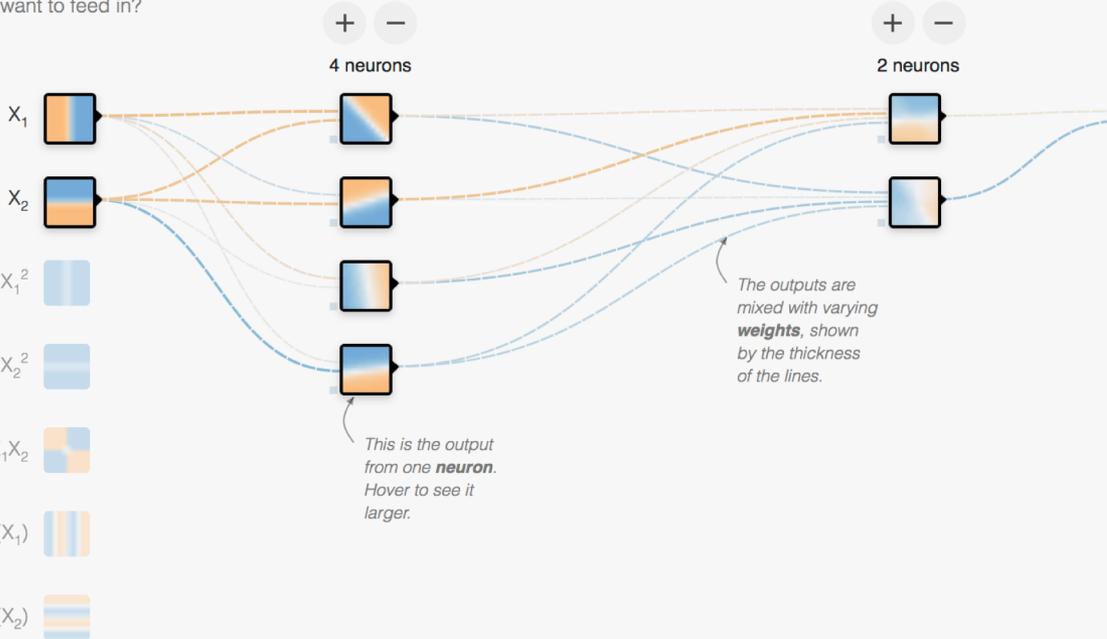
REGENERATE

FEATURES

Which properties do you want to feed in?

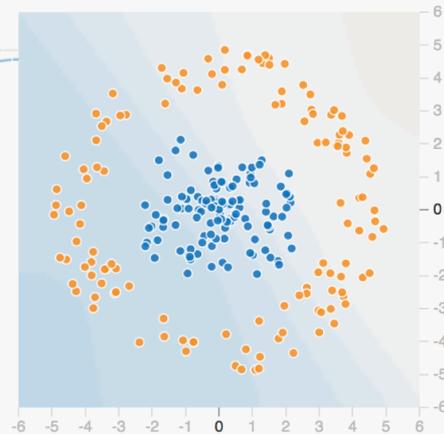
- X_1
- X_2
- X_1^2
- X_2^2
- X_1X_2
- $\sin(X_1)$
- $\sin(X_2)$

+ - 2 HIDDEN LAYERS



OUTPUT

Test loss 0.507
Training loss 0.505



Colors shows data, neuron and weight values.

Show test data Discretize output

Tricky issues with neural network training

- Sensitive to initialization
 - Objective is non-convex, many local optima
 - In practice: start with random values rather than zeros
- Many other hyper-parameters
 - Number of hidden units (and potentially hidden layers)
 - Gradient descent learning rate
 - Stopping criterion

Neural networks vs. linear classifiers

Advantages of Neural Networks:

- More expressive
- Less feature engineering

Challenges using Neural Networks:

- Harder to train
- Harder to interpret

Neural Network Architectures

- We focused on a **multi-layer feedforward** network
- Many other deeper architectures
 - Convolutional networks
 - Recurrent networks (LSTMs)
 - Dense Nets, ResNets, etc

Issues in Deep Neural Networks

- Long training time
 - There are sometimes a lot of training data
 - Many iterations (epochs) are typically required for optimization
 - Computing gradients in each iteration takes too much time

Improving on Gradient Descent: Stochastic Gradient Descent (SGD)

- Update weights for each example

$$L = \frac{1}{2} (y^n - \hat{y}^n)^2 \quad \mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \epsilon \frac{\partial L^n}{\partial \mathbf{w}_i}$$

+ Fast, online

– Sensitive to noise

- Mini-batch SGD: Update weights for a small set of examples

$$L = \frac{1}{2} \sum_{n \in B} (y^n - \hat{y}^n)^2 \quad \mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \epsilon \frac{\partial L^B}{\partial \mathbf{w}_i}$$

+ Fast, online

+ Robust to noise

Improving on Gradient Descent: SGD with Momentum

- Update based on gradients + previous direction

$$v_i(t) = \alpha v_i(t-1) - (1-\alpha) \frac{\partial L}{\partial w_i}(t)$$

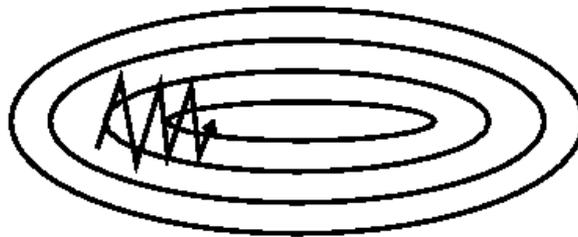
$$\mathbf{w}(t+1) = \mathbf{w}(t) + \epsilon \mathbf{v}(t)$$

+ **Converge faster**

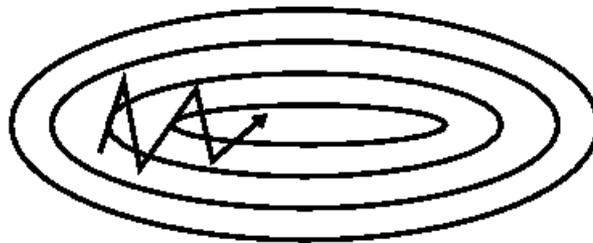
+ **Avoid oscillation**

Improving on Gradient Descent: SGD with Momentum

SGD w/o momentum



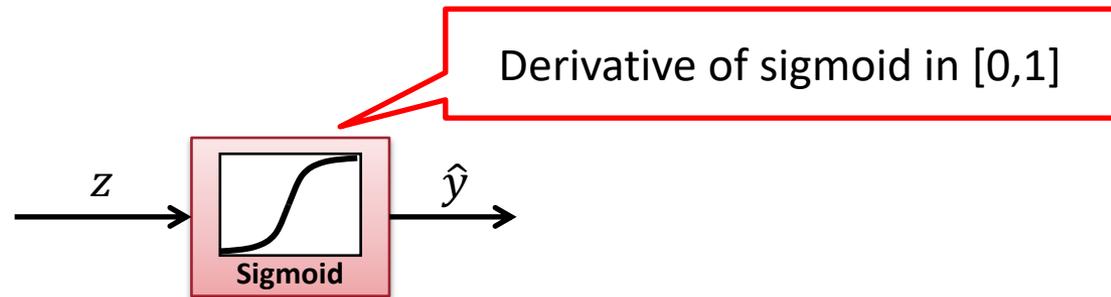
SGD with momentum
helps dampen
oscillations



Vanishing Gradient Problem

In deep networks

- Gradients in the lower layers are typically extremely small
- Optimizing multi-layer neural networks takes huge amount of time

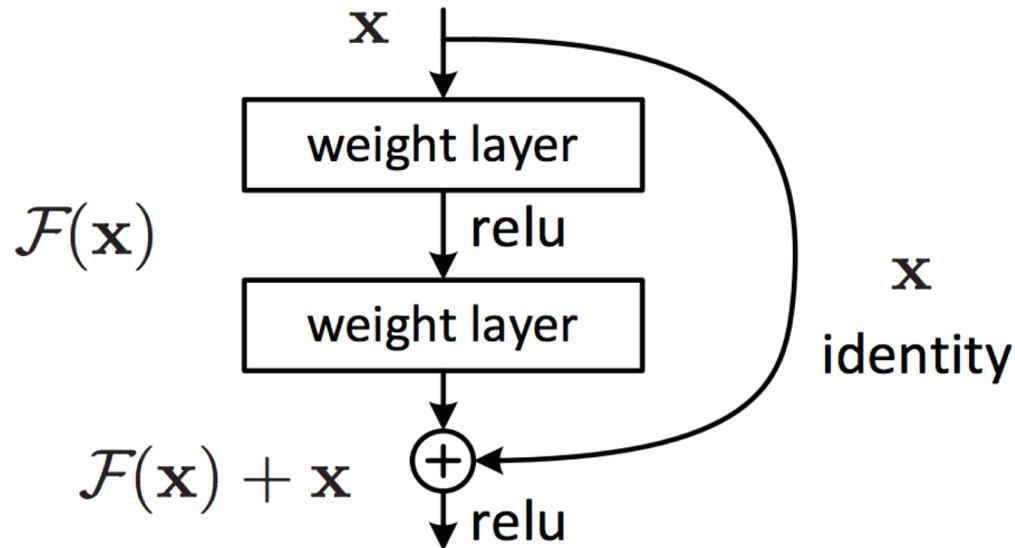


$$\frac{\partial E}{\partial w_{ki}} = \sum_n \frac{\partial z_i^n}{\partial w_{ki}} \frac{d\hat{y}_i^n}{dz_i^n} \frac{\partial E}{\partial \hat{y}_i^n} = \sum_n \frac{\partial z_i^n}{\partial w_{ki}} \frac{d\hat{y}_i^n}{dz_i^n} \sum_j w_{ij} \frac{d\hat{y}_j^n}{dz_j^n} \frac{\partial E}{\partial \hat{y}_j^n}$$

Vanishing Gradient Problem

- Vanishing gradient problem can be mitigated
 - Using custom neural network architectures
 - Using other non-linearities
 - E.g., Rectifier: $f(x) = \max(0, x)$

ResNet



Since last lecture : 20266

[Deep Residual Learning for Image Recognition](https://arxiv.org/abs/1512.03571)

[https://arxiv.org > cs](https://arxiv.org/abs/1512.03571) ▼

by K He - 2015 - [Cited by 19999](#) - [Related articles](#)

Dec 10, 2015 - Abstract: Deeper neural networks are more difficult to train.

We present a residual learning framework to ease the training of networks that are ...

Why Neural Networks?

Perceptron

- Proposed by Frank Rosenblatt in 1957
- Real inputs/outputs, threshold activation function

Revival in the 1980's

Backpropagation discovered in 1970's but popularized in 1986

- David E. Rumelhart, Geoffrey E. Hinton, Ronald J. Williams. "Learning representations by back-propagating errors." In Nature, 1986.

MLP is a universal approximator

- Can approximate any non-linear function in theory, given enough neurons, data
- Kurt Hornik, Maxwell Stinchcombe, Halbert White. "Multilayer feedforward networks are universal approximators." Neural Networks, 1989

Generated lots of excitement and applications

Neural Networks Applied to Vision

LeNet – vision application

- LeCun, Y; Boser, B; Denker, J; Henderson, D; Howard, R; Hubbard, W; Jackel, L, “Backpropagation Applied to Handwritten Zip Code Recognition,” in Neural Computation, 1989
- USPS digit recognition, later check reading

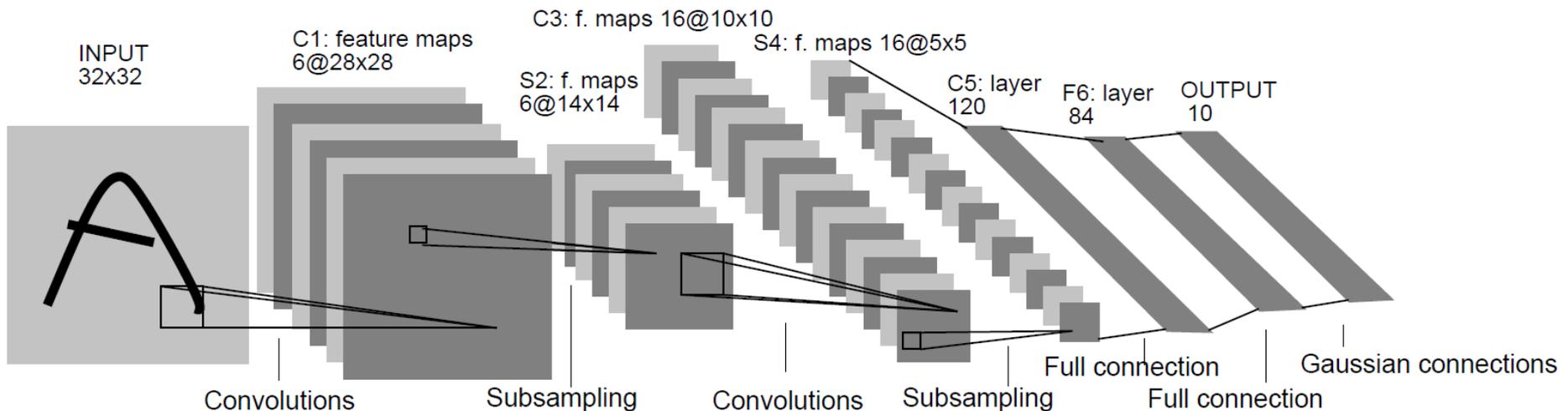


Image credit: LeCun, Y., Bottou, L., Bengio, Y., Haffner, P. “Gradient-based learning applied to document recognition.” Proceedings of the IEEE, 1998.

New “winter” and revival in early 2000’s

New “winter” in the early 2000’s due to

- problems with training NNs
- Support Vector Machines (SVMs), Random Forests (RF) – easy to train, nice theory

Revival again by 2011-2012

- Name change (“neural networks” -> “deep learning”)
- + Algorithmic developments that made training somewhat easier
- + Big data + GPU computing
- = performance gains on many tasks (esp Computer Vision)

Big Data

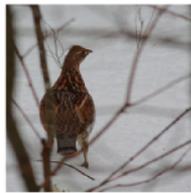
- ImageNet Large Scale Visual Recognition Challenge
 - 1000 categories w/ 1000 images per category
 - 1.2 million training images, 50,000 validation, 150,000 testing



flamingo



cock



ruffed grouse



quail



partridge

...



Egyptian cat



Persian cat



Siamese cat

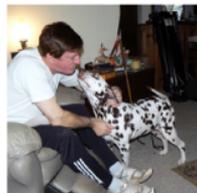


tabby



lynx

...



dalmatian



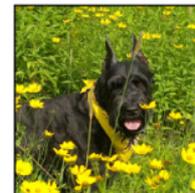
keeshond



miniature schnauzer



standard schnauzer



giant schnauzer

...

AlexNet Architecture

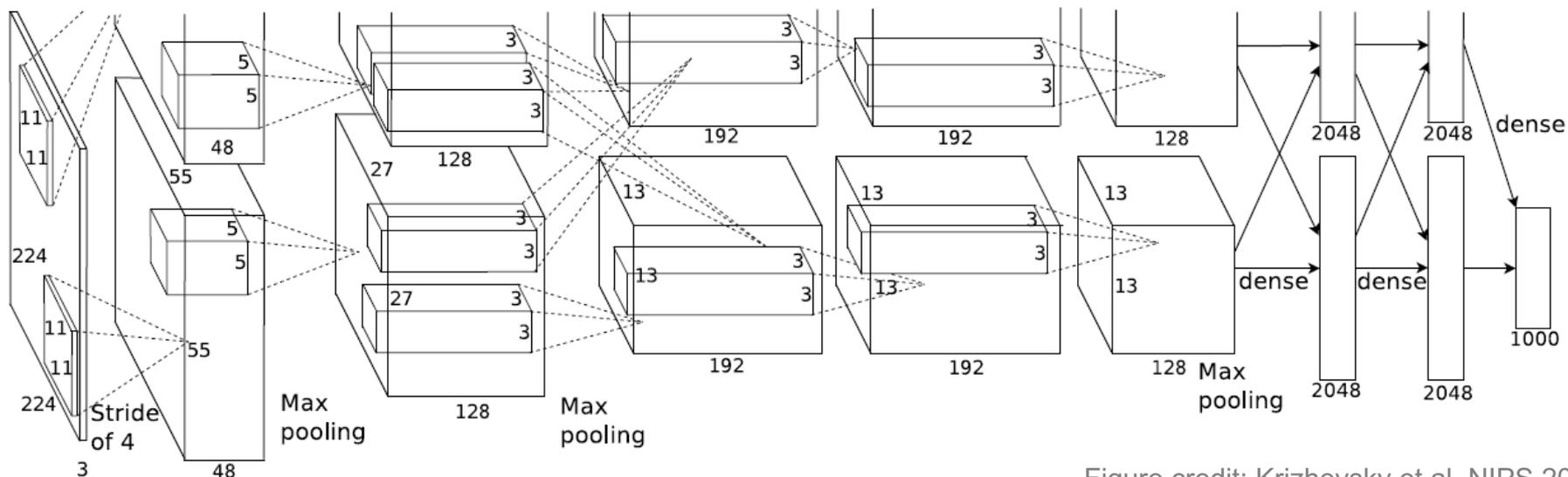


Figure credit: Krizhevsky et al, NIPS 2012.

60 million parameters!

Various tricks

- ReLU nonlinearity
- Overlapping pooling
- Local response normalization
- Dropout – set hidden neuron output to 0 with probability .5
- Data augmentation
- Training on GPUs

GPU Computing

- **Big data** and **big models** require lots of computational power
- GPUs
 - thousands of cores for parallel operations
 - multiple GPUs
 - still took about 5-6 days to train AlexNet on two NVIDIA GTX 580 3GB GPUs (much faster today)

Image Classification Performance

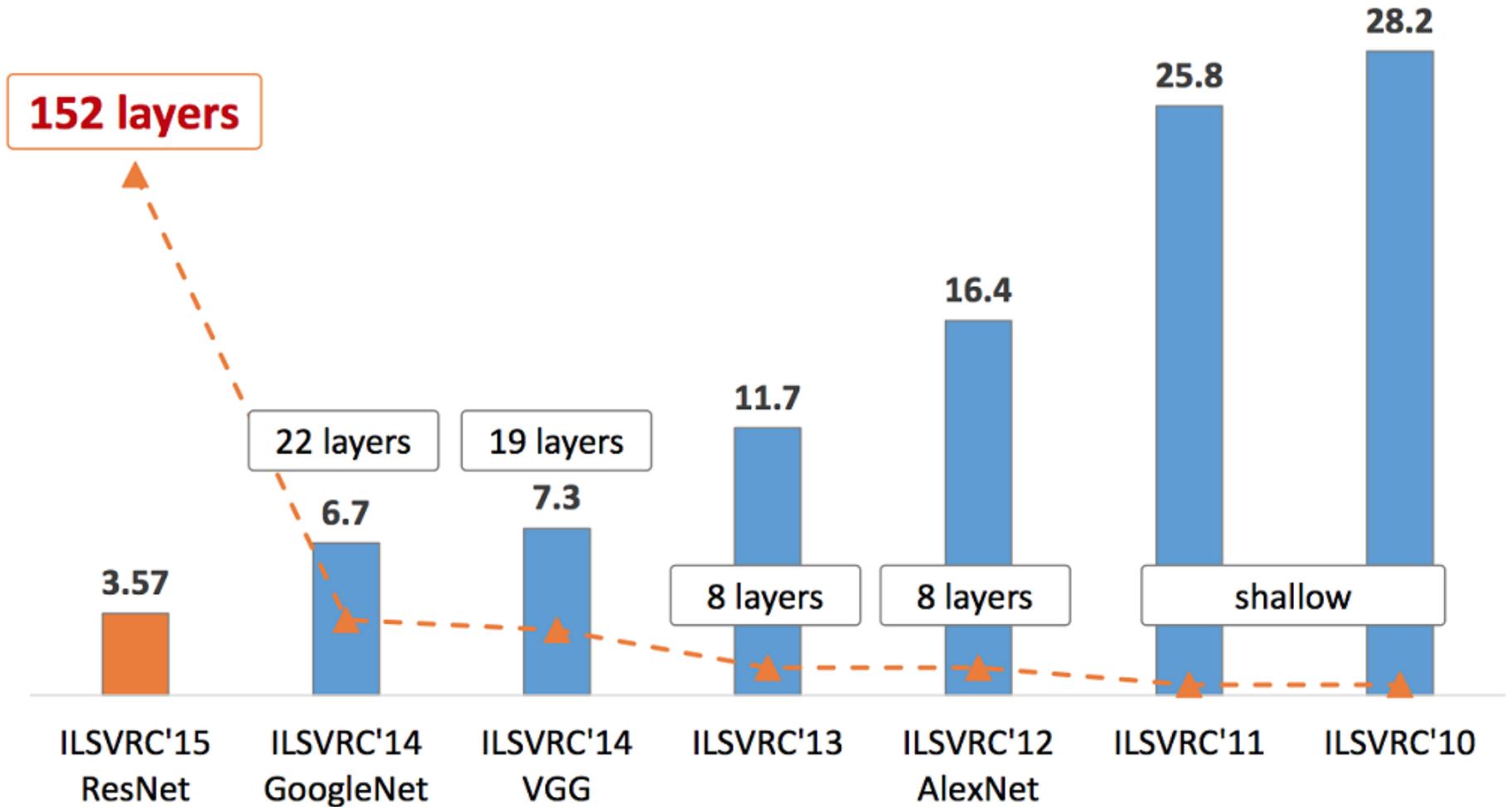
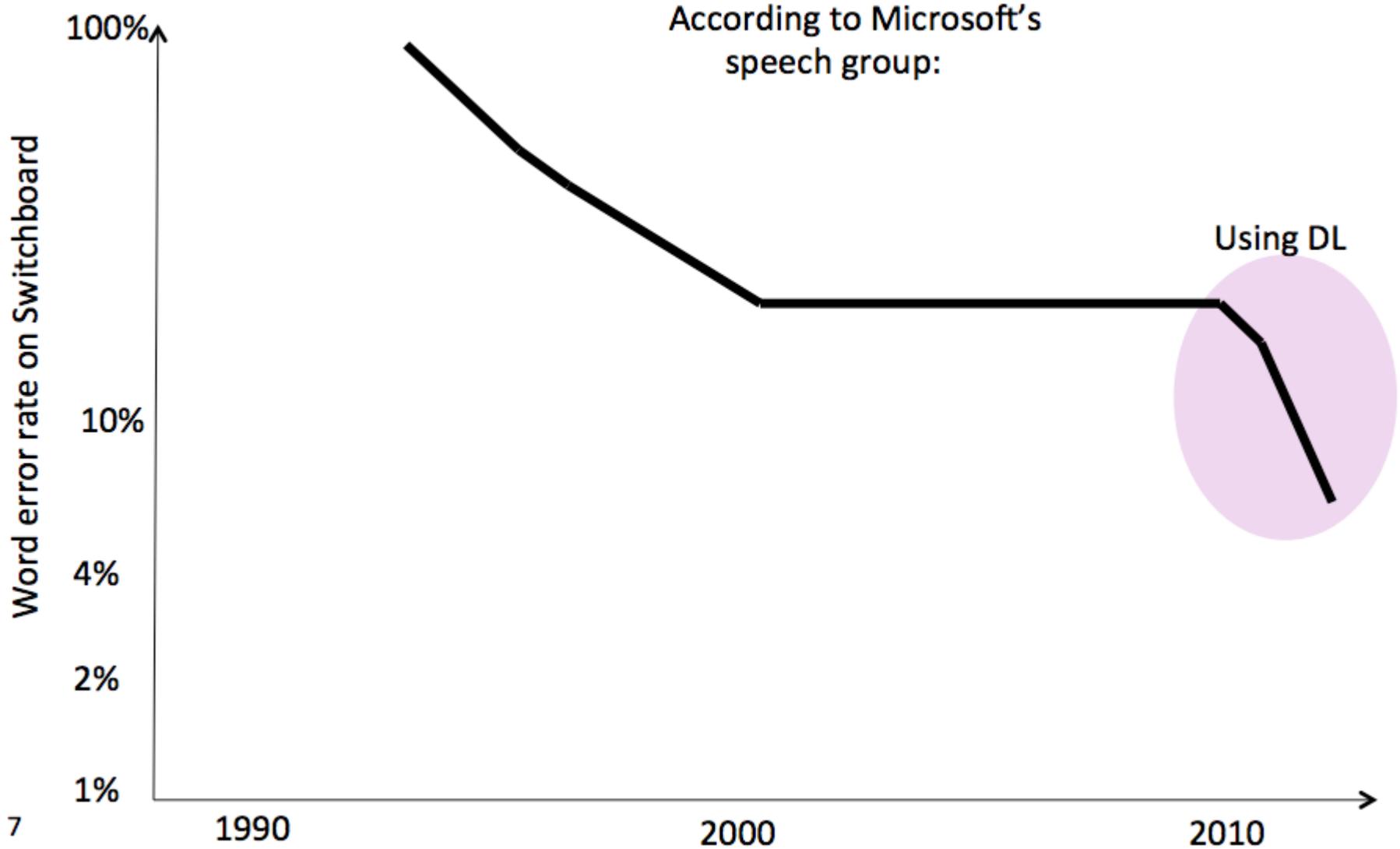


Image Classification Top-5 Errors (%)

Figure from: K. He, X. Zhang, S. Ren, J. Sun. "Deep Residual Learning for Image Recognition". arXiv 2015. (slides)

Speech Recognition



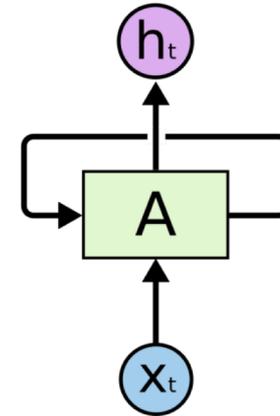
Recurrent Neural Networks for Language Modeling

- Speech recognition is difficult due to ambiguity
 - “how to recognize speech”
 - or “how to wreck a nice beach”?
- Language model gives probability of next word given history
 - $P(\text{“speech”} \mid \text{“how to recognize”})?$

Recurrent Neural Networks

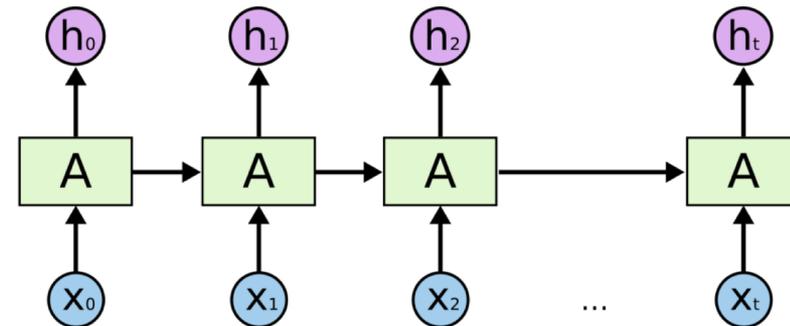
Networks with loops

- The output of a layer is used as input for the same (or lower) layer
- Can model dynamics (e.g. in space or time)



Loops are unrolled

- Now a standard feed-forward network with many layers
- Suffers from vanishing gradient problem
- In theory, can learn long term memory, in practice not (Bengio et al, 1994)



Long Short Term Memory (LSTM)

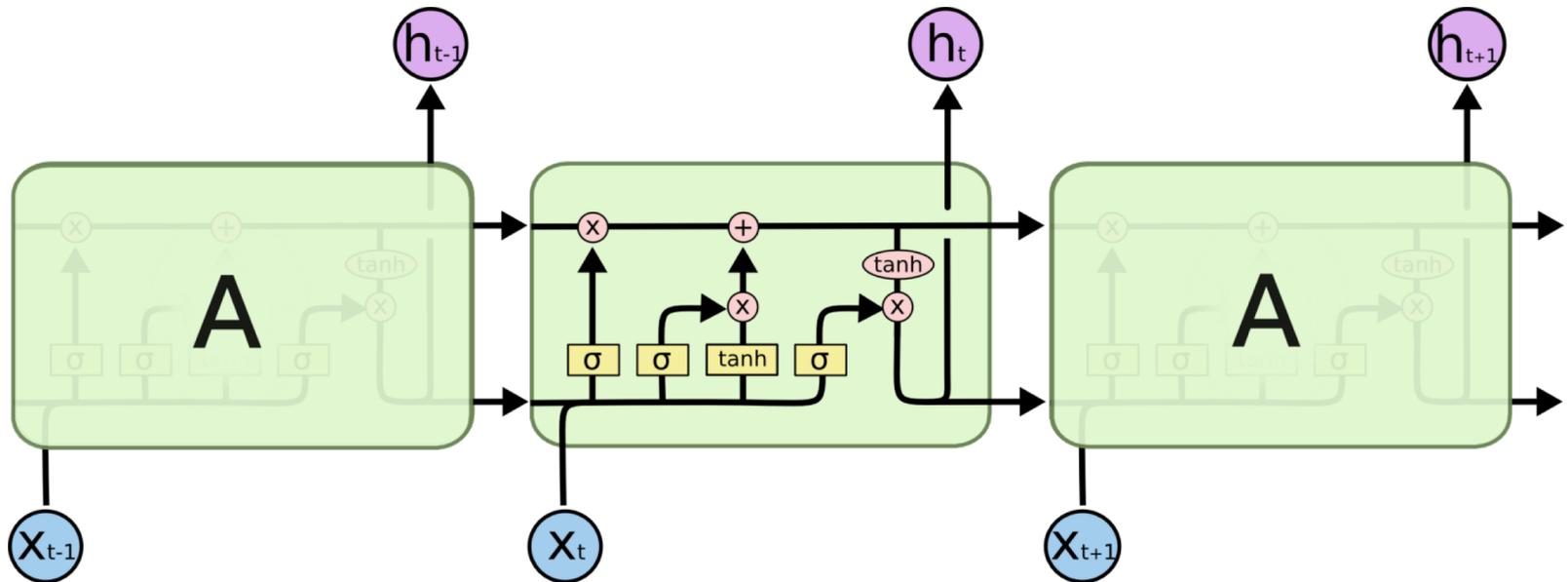
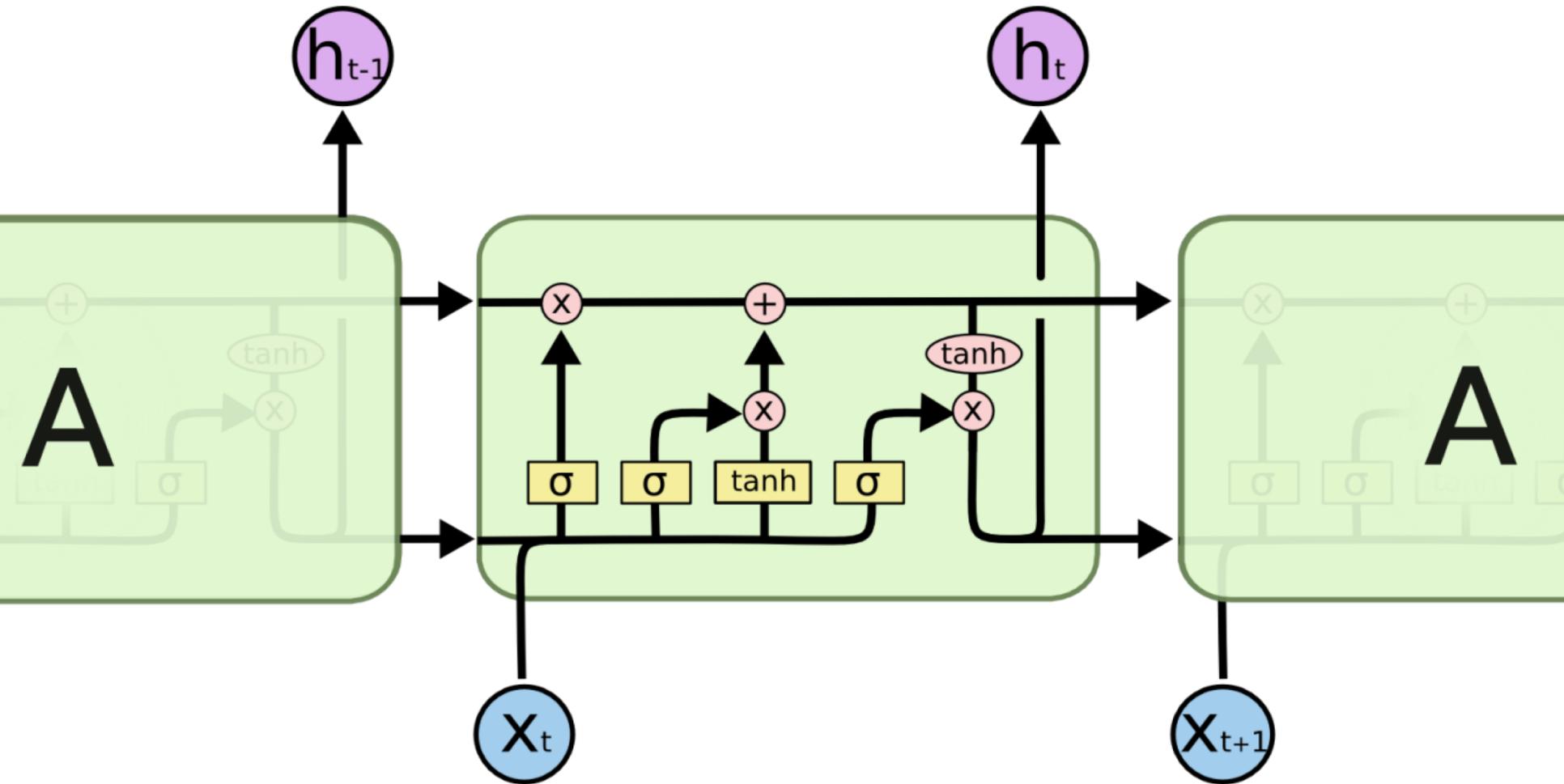


Image credit: Christopher Colah's blog, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

- A type of RNN explicitly designed not to have the vanishing or exploding gradient problem
- Models long-term dependencies
- Memory is propagated and accessed by gates
- Used for speech recognition, language modeling ...

Long Short Term Memory (LSTM)



What you should know about deep neural networks

- Why they are difficult to train
 - Initialization
 - Overfitting
 - Vanishing gradient
 - Require large number of training examples
- What can be done about it
 - Improvements to gradient descent
 - Stochastic gradient descent
 - Momentum
 - Weight decay
 - Alternate non-linearities and new architectures

References (& great tutorials) if you want to explore further:

<http://www.andreykurenkov.com/writing/a-brief-history-of-neural-nets-and-deep-learning-part-1/>

<http://cs231n.github.io/neural-networks-1/>

<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Keeping things in perspective...

In 1958, the New York Times reported the perceptron to be "the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."

Unsupervised Learning

Principal Component Analysis

Unsupervised Learning

- Discovering hidden structure in data
- What algorithms do we know for unsupervised learning?
 - K-Means Clustering
- Today: how can we learn better representations of our data points?

Dimensionality Reduction

- Goal: extract hidden lower-dimensional structure from high dimensional datasets
- Why?
 - To visualize data more easily
 - To remove noise in data
 - To lower resource requirements for storing/processing data
 - To improve classification/clustering

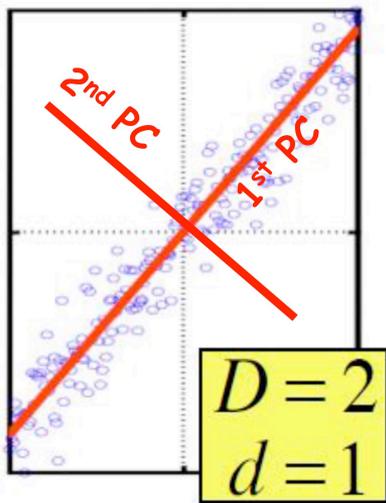
- Linear algebra review:
 - Matrix decomposition with eigenvectors and eigenvalues

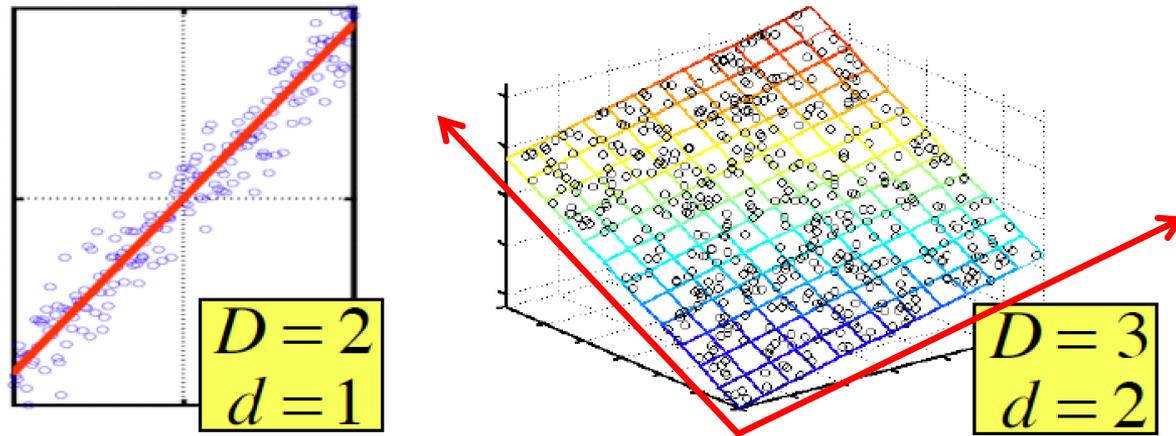
Principal Component Analysis

- Goal: Find a **projection** of the data onto directions that **maximize variance** of the original data set
 - Intuition: those are directions in which most information is encoded
- Definition: **Principal Components** are orthogonal directions that capture most of the variance in the data

PCA: finding principal components

- 1st PC
 - Projection of data points along 1st PC discriminates data most along any one direction
- 2nd PC
 - next orthogonal direction of greatest variability
- And so on...





Examples of data points in D dimensional space that can be effectively represented in a d -dimensional subspace ($d < D$)

PCA: notation

- Data points
 - Represented by matrix X of size $N \times D$
 - Let's assume data is centered
- Principal components are d vectors: v_1, v_2, \dots, v_d
 $v_i \cdot v_j = 0, i \neq j$ and $v_i \cdot v_i = 1$
- The sample variance data projected on vector v is
$$\sum_{i=1}^n (x_i^T v)^2 = (Xv)^T (Xv)$$

PCA formally

- Finding vector that maximizes sample variance of projected data:

$$\operatorname{argmax}_v v^T X^T X v \text{ such that } v^T v = 1$$

- A constrained optimization problem
 - Lagrangian folds constraint into objective:
$$\operatorname{argmax}_v v^T X^T X v - \lambda(v^T v - 1)$$
 - Solutions are vectors v such that $X^T X v = \lambda v$
 - i.e. eigenvectors of $X^T X$ (sample covariance matrix)

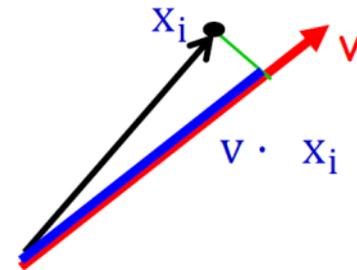
PCA formally

- The eigenvalue λ denotes the amount of variability captured along dimension v
 - Sample variance of projection $v^T X^T X v = \lambda$
- If we rank eigenvalues from large to small
 - The 1st PC is the eigenvector of $X^T X$ associated with largest eigenvalue
 - The 2nd PC is the eigenvector of $X^T X$ associated with 2nd largest eigenvalue
 - ...

Alternative interpretation of PCA

- PCA finds vectors v such that projection on to these vectors minimizes reconstruction error

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - (\mathbf{v}^T \mathbf{x}_i) \mathbf{v}\|^2$$



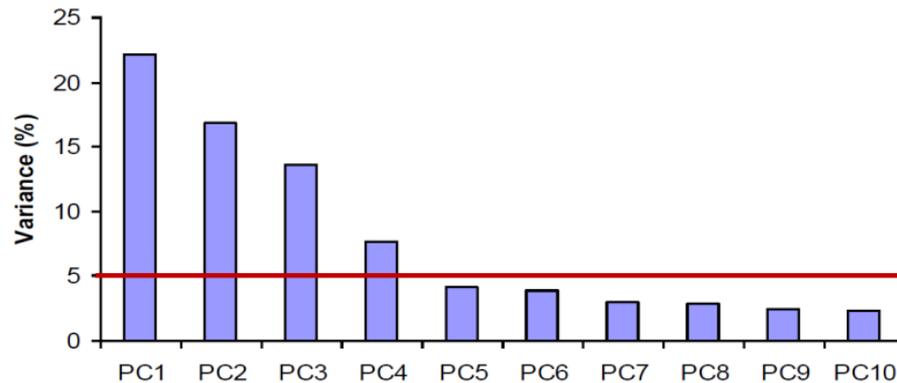
Resulting PCA algorithm

Algorithm 36 PCA(\mathbf{D} , K)

- 1: $\boldsymbol{\mu} \leftarrow \text{MEAN}(\mathbf{X})$ // compute data mean for centering
 - 2: $\mathbf{D} \leftarrow (\mathbf{X} - \boldsymbol{\mu}\mathbf{1}^\top)^\top (\mathbf{X} - \boldsymbol{\mu}\mathbf{1}^\top)$ // compute covariance, $\mathbf{1}$ is a vector of ones
 - 3: $\{\lambda_k, \mathbf{u}_k\} \leftarrow$ top K eigenvalues/eigenvectors of \mathbf{D}
 - 4: **return** $(\mathbf{X} - \boldsymbol{\mu}\mathbf{1}) \mathbf{U}$ // project data using \mathbf{U}
-

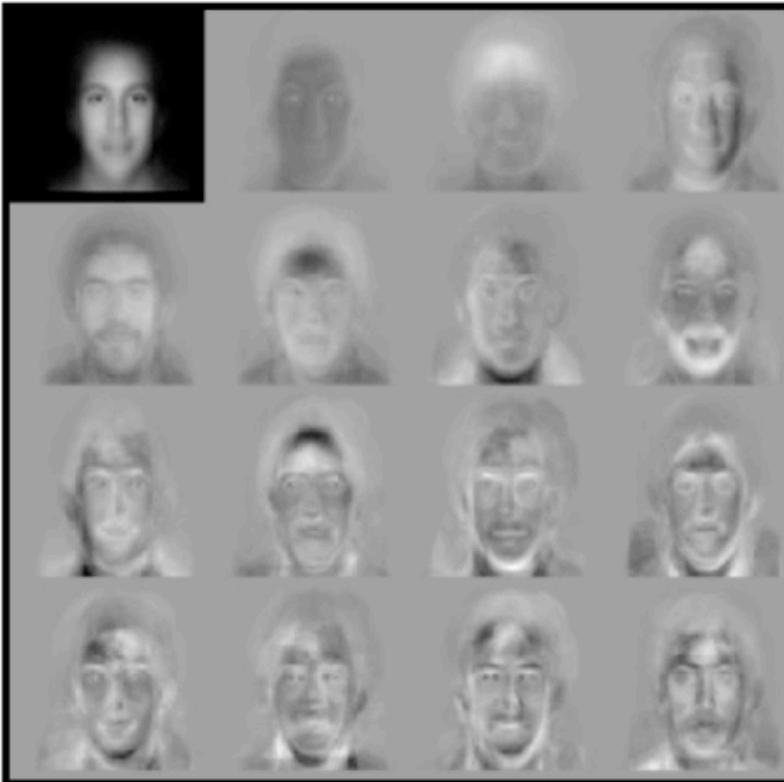
How to choose the hyperparameter K?

- i.e. the number of dimensions



- We can ignore the components of smaller significance

An example: Eigenfaces



Eigenfaces
from 7562
images:

**top left image
is linear
combination
of rest.**

Sirovich & Kirby (1987)
Turk & Pentland (1991)

PCA pros and cons

- Pros
 - Eigenvector method
 - No tuning of the parameters
 - No local optima
- Cons
 - Only based on covariance (2nd order statistics)
 - Limited to linear projections

What you should know

- Principal Components Analysis
 - Goal: Find a **projection** of the data onto directions that **maximize variance** of the original data set
 - PCA **optimization objectives** and resulting **algorithm**
 - Why this is useful!