

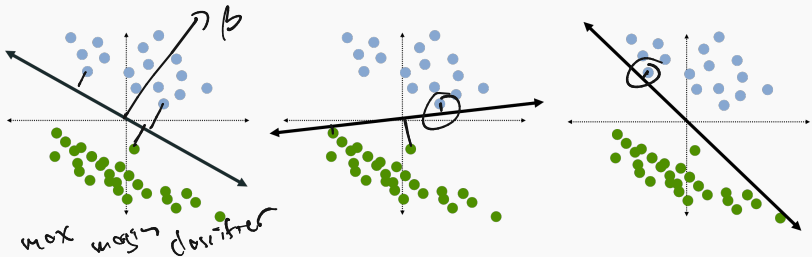
CS-GY 6923: Lecture 10

Finish SVMs, Neural Nets Introduction, Back propagation

NYU Tandon School of Engineering, Prof. Christopher Musco

SUPPORT VECTOR MACHINES

Goal: Find a separating hyperplane for linearly separable classification problem.



Ideally, choose the hyperplane that(maximizes margin.)

OPTIMIZATION FORMULATION

Original problem: $\arg \max_{\beta} \left[\min_{i \in 1, \dots, n} \frac{y_i \cdot \langle x_i, \beta \rangle}{\|\beta\|_2} \right]$.

Equivalent formulation:

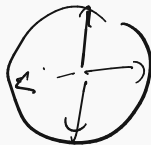
$$\min_{\beta} \|\beta\|_2^2$$

subject to

$$y_i \cdot \langle x_i, \beta \rangle \geq 1 \text{ for all } i.$$

Under this formulation $m = \frac{1}{\|\beta\|_2}$.

- Can be solved using a constrained optimization method.
- Can be combined with any non-linear kernel.
- Classification only requires computing kernel similarity with the support vectors.



$$w = \gamma$$
$$\beta: \|\beta\|_2 = 1$$
$$\hookrightarrow 0.$$

$$y: -\langle x_i, \beta \rangle \geq 1$$

CLASSIFICATION

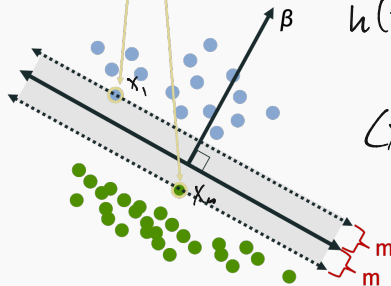
β

$$\mathbb{1}[\langle x_i, \beta \rangle > 0]$$

$$\beta = \sum_{i=1}^n \alpha_i \phi(x_i)$$

$$\hookrightarrow \alpha_1, \dots, \alpha_n$$

support
vectors



$$k(x_i, x_{new})$$

$$k(x_v, x_{new})$$

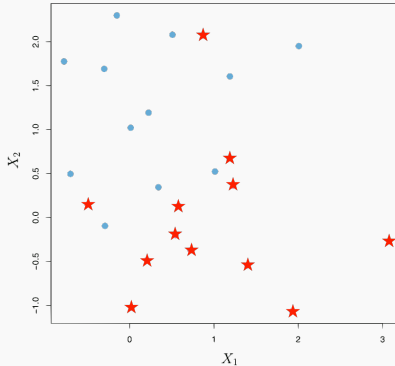
$$\langle x_i, \beta \rangle$$

$$= \sum_{i=1}^n \alpha_i k(x_i, x_{new})$$

- When using a kernel like $k(x_i, x_j) = e^{-\|x_i - x_j\|_2^2}$, classification for a new points x_{new} only requires computing kernel similarity with the support vectors. Logistic regression requires similarity with all training points.

HARD-MARGIN SVM

Hard-margin SVMs have a few other critical issues in practice:

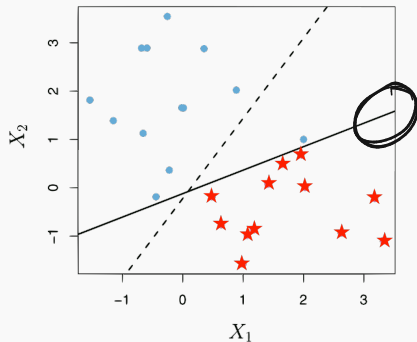
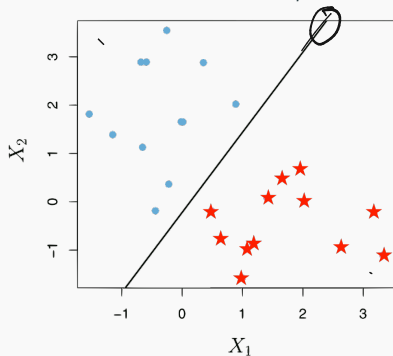


Data might not be linearly separable, in-which case the maximum margin classifier is not even defined.

Less likely to be an issue when using a non-linear kernel. If \mathbf{K} is full rank then perfect separation is always possible. And typically it is, e.g. for an RBF kernel or moderate degree polynomial kernel.

HARD-MARGIN SVM

Another critical issue in practice:



Hard-margin SVM classifiers are not robust.

Solution: Allow the classifier to make some “mistakes”! A mistake can either be a misclassification, or simply a point allowed to be “inside” the margin.

a_1, \dots, a_n

Hard margin objective:

$$\min_{\beta} \|\beta\|_2^2 \quad \text{subject to} \quad y_i \cdot \langle x_i, \beta \rangle \geq 1 \text{ for all } i.$$

Soft margin objective:

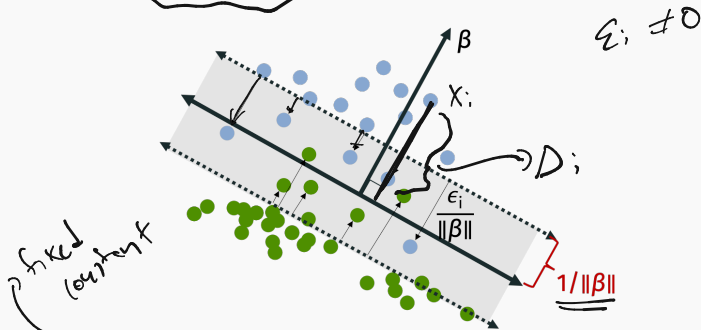
$$\min_{\beta, \epsilon_i} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad y_i \cdot \langle x_i, \beta \rangle \geq 1 - \epsilon_i \text{ for all } i.$$

where $\epsilon_i \geq 0$ is a non-negative “slack variable”.

$\epsilon_i / \|\beta\|_2$ is the magnitude of the “error” (distance past the margin) we allow x_i to travel. Recalling that margin is $1/\|\beta\|_2$, $\epsilon_i \geq 1$ corresponds to a misclassification.

SOFT-MARGIN SVM

Recall that $\Delta_i = \frac{y_i \cdot \langle x_i, \beta \rangle}{\|\beta\|_2}$.

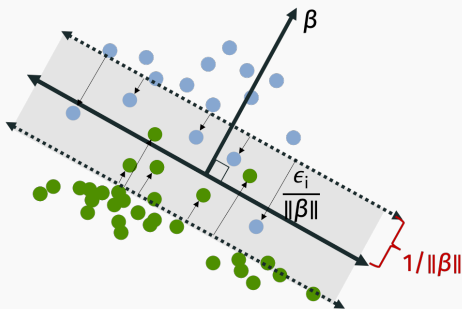


Soft margin objective:

$$\min_{\substack{\beta, \\ \epsilon_1, \dots, \epsilon_n}} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad \left(\frac{y_i \cdot \langle x_i, \beta \rangle}{\|\beta\|_2} \right) \geq \frac{1}{\|\beta\|_2} - \frac{\epsilon_i}{\|\beta\|_2} \text{ for all } i.$$

SOFT-MARGIN SVM

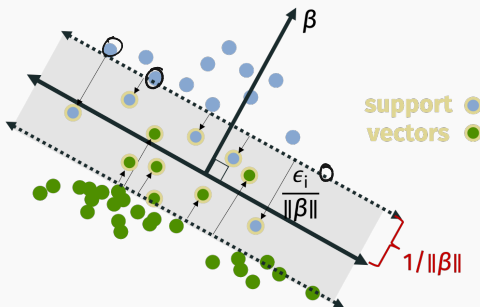
Recall that $\Delta_i = \frac{y_i \cdot \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle}{\|\boldsymbol{\beta}\|_2}$.



Soft margin objective:

$$\min_{\boldsymbol{\beta}} \|\boldsymbol{\beta}\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad \frac{y_i \cdot \langle \mathbf{x}_i, \boldsymbol{\beta} \rangle}{\|\boldsymbol{\beta}\|_2} \geq \frac{1}{\|\boldsymbol{\beta}\|_2} - \frac{\epsilon_i}{\|\boldsymbol{\beta}\|_2} \text{ for all } i.$$

SOFT-MARGIN SVM



Any x_i with a non-zero ϵ_i is a support vector. As before, only support vectors are needed for classification in the kernel setting. **Good exercise to prove yourself.**

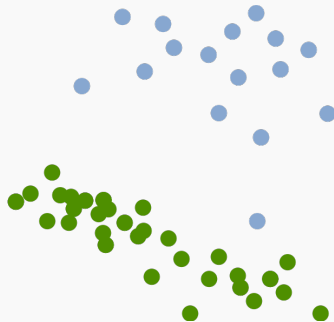
Soft margin objective:

$$\min_{\beta} \underbrace{\|\beta\|_2^2} + C \sum_{i=1}^n \epsilon_i$$

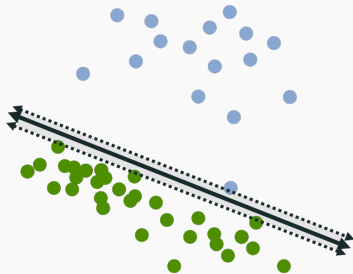
- Large C means penalties are punished more in objective
 \implies smaller margin, less support vectors.
- Small C means penalties are punished less in objective
 \implies larger margin, more support vectors.

When data is linearly separable, as $C \rightarrow \infty$ we will always get a separating hyperplane. A smaller value of C might lead to a more robust solution.

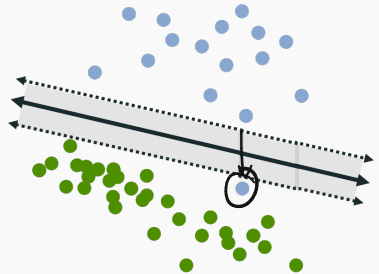
Example dataset:



EFFECT OF C



large C



smaller C

The classifier on the right is intuitively more robust. So for this data, a smaller choice for C might make sense.

Typically the smaller C is, the more support vectors (above image isn't a great example).

COMPARISON TO LOGISTIC REGRESSION

Some basic transformations of the soft-margin objective:

$\epsilon_i = 1$

$$\min_{\beta, \epsilon_1, \dots, \epsilon_n} \|\beta\|_2^2 + C \sum_{i=1}^n \epsilon_i \quad \text{subject to} \quad y_i \cdot \langle x_i, \beta \rangle \geq 1 - \epsilon_i \text{ for all } i.$$

Always have $y_i \cdot \langle x_i, \beta \rangle = 1 - \epsilon_i$ if $\epsilon_i \neq 0$.

$$\min_{\beta} \|\beta\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i \cdot \langle x_i, \beta \rangle).$$

$\epsilon_i \neq 0$

$$\min_{\beta} \frac{1}{C} \|\beta\|_2^2 + \sum_{i=1}^n \max(0, 1 - y_i \cdot \langle x_i, \beta \rangle).$$

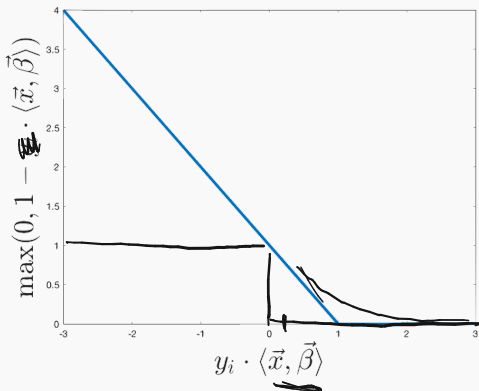
$\lambda = 1/C$

These are all equivalent. $\lambda = 1/C$ is just another scaling parameter. Moved from a constrained problem to a much easier unconstrained optimization problem.

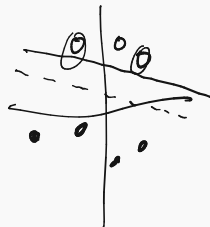
If $y \cdot \langle x_i, \beta \rangle \geq 1$ choose $\epsilon_i = 0$. Otherwise choose $\epsilon_i = 1 - y_i \cdot \langle x_i, \beta \rangle$.

HINGE LOSS

Hinge-loss: $\max(0, 1 - y_i \cdot \langle \vec{x}_i, \vec{\beta} \rangle)$. Recall that $y_i \in \{-1, 1\}$.



$y_i = 1$



Soft-margin SVM:

$$\min_{\vec{\beta}} \left[\sum_{i=1}^n \max(0, 1 - y_i \cdot \langle \vec{x}_i, \vec{\beta} \rangle) + \lambda \|\vec{\beta}\|_2^2 \right]. \quad (1)$$

LOGISTIC LOSS

$$y_i \in \{-1, 1\}$$

Recall the logistic loss for $y_i \in \{0, 1\}$:

$$\begin{aligned} L(\beta) &= - \sum_{i=1}^n y_i \log(h(\langle x_i, \beta \rangle)) + (1 - y_i) \log(1 - h(\langle x_i, \beta \rangle)) \\ &= - \sum_{i=1}^n y_i \log \left(\frac{1}{1 + e^{-\langle x_i, \beta \rangle}} \right) + (1 - y_i) \log \left(\frac{e^{-\langle x_i, \beta \rangle}}{1 + e^{-\langle x_i, \beta \rangle}} \right) \cdot \frac{e^{\langle x_i, \beta \rangle}}{e^{\langle x_i, \beta \rangle}} \\ &= - \sum_{i=1}^n y_i \log \left(\frac{1}{1 + e^{-\langle x_i, \beta \rangle}} \right) + (1 - y_i) \log \left(\frac{1}{1 + e^{\langle x_i, \beta \rangle}} \right) \end{aligned}$$

for $y_i \in \{0, 1\}$

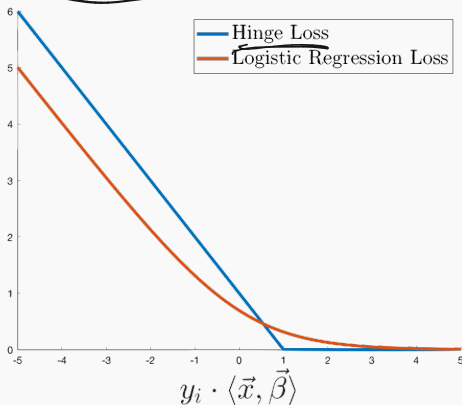
$$= - \sum_{i=1}^n \log \left(\frac{1}{1 + e^{-y_i \langle x_i, \beta \rangle}} \right) \text{ if } y_i \in \{-1, 1\}$$

COMPARISON OF SVM TO LOGISTIC REGRESSION

Compare this to the logistic regression loss reformulated for $y_i \in \{-1, 1\}$:

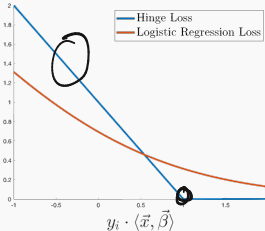
$$\sum_{i=1}^n -\log \left(\frac{1}{1 + e^{-y_i \cdot \langle \vec{x}_i, \vec{\beta} \rangle}} \right)$$

fix x type.



COMPARISON TO LOGISTIC REGRESSION

So, in the end, the function minimized when finding β for the standard **soft-margin SVM** is very similar to the objective function minimized when finding β using **logistic regression** with ℓ_2 regularization.



Both functions can be optimized using first-order methods like gradient descent. This is now a common choice for large problems. **Will explore more on Lab 5.**

NEURAL NETWORKS

Key Concept

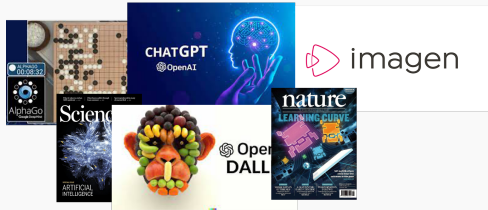
Approach until now:

- Choose good features or a good kernel.
- Use optimization to find best model given those features.

(Neural network approach:

- Learn good features and a good model simultaneously.

The leading method in machine learning right now.

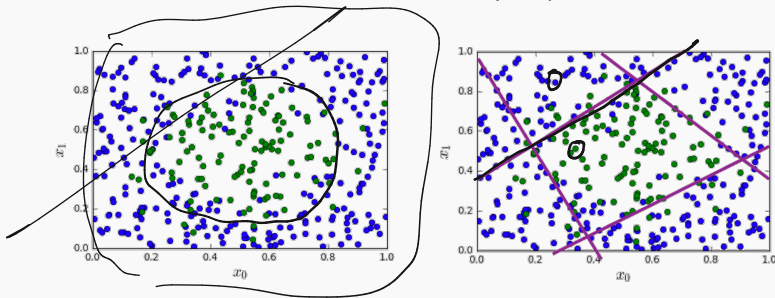


Focus of investment at universities, government research labs, funding agencies, and large tech companies.

Studied since the 1940s/50s. **Why the recent attention?** More on history of neural networks shortly.

SIMPLE MOTIVATING EXAMPLE

Classification when data is not linearly separable:

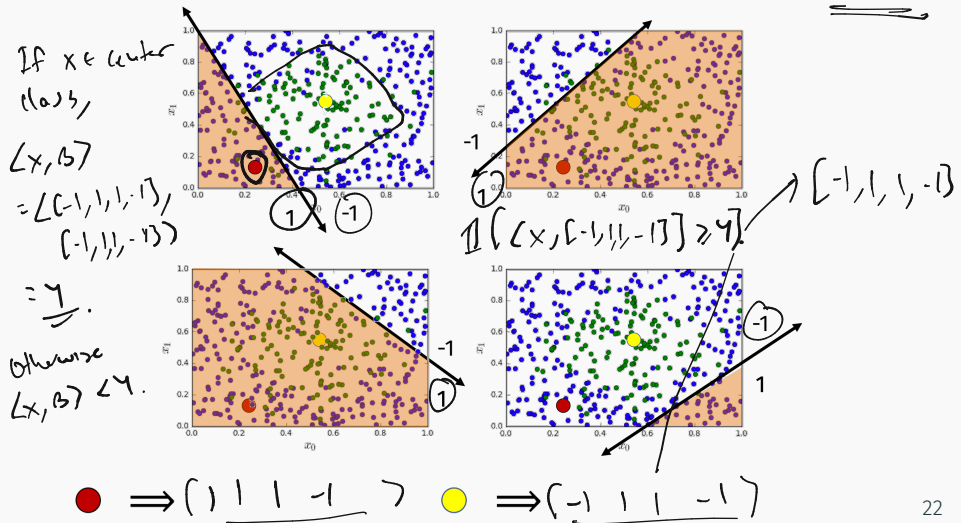


Could use feature transformations or a non-linear kernel.

Alternative approach: Divide the space up into regions using multiple linear classifiers.

SIMPLE MOTIVATING EXAMPLE

For each linear classifier β , add a new $-1, 1$ feature for every example $\mathbf{x} = [x_0, x_1]$ depending on the sign of $\langle \mathbf{x}, \beta \rangle$. $\beta = (-1, 1, 1, -1)$



SIMPLE MOTIVATING EXAMPLE

$$\begin{bmatrix} .2, .8 \\ .5, .5 \\ \vdots \\ .5, 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} = \begin{bmatrix} -1, -1, +1, -1 \\ -1, +1, +1, -1 \\ \vdots \\ -1, -1, -1, -1 \end{bmatrix} \begin{bmatrix} .1 \\ \vdots \\ -1 \end{bmatrix}$$

$\nearrow T$

Question: After data transformation, how should we map each new vector u_i to a class label?

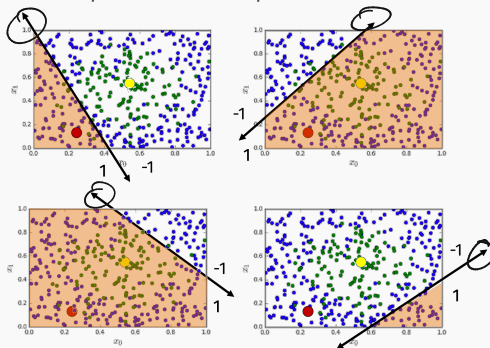
$$1(T \geq .4)$$

$$\begin{bmatrix} -1, -1, +1, -1 \\ -1, +1, +1, -1 \\ \vdots \\ -1, -1, -1, -1 \end{bmatrix} \xrightarrow{?} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

SIMPLE MOTIVATING EXAMPLE

Our machine learning algorithms needs to **learn two things**:

- The original linear functions which divide our data set into regions (their slopes + intercepts).



- Another linear function which maps our new features to an output class probability.

POSSIBLE MODEL

Input: $\underline{x} = x_1, \dots, x_{N_I} \in \mathbb{R}^d$

$$N_H \begin{bmatrix} \circ \\ \circ \\ \circ \\ \circ \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} z_u \end{bmatrix}$$

x_i

Model: $f(\underline{x}, \Theta)$:

$$\circ \underline{z}_H \in \mathbb{R}^{N_H} = \underline{W}_H \underline{x} + \underline{\beta}_H \xrightarrow{\mathbb{R}^{N_H}} \mathbb{1}[(\underline{x}, \underline{w}) > \tau, \lambda]$$

$$\circ \underline{u}_H = \text{sign}(\underline{z}_H)$$

$$\cdot \underline{z}_O \in \mathbb{R} = \underline{W}_O \underline{u}_H + \beta_O$$

$$\cdot \underline{u}_O = \mathbb{1}[\underline{z}_O \text{ ~~is not~~}]$$

$$\begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + \beta_O$$

u_H

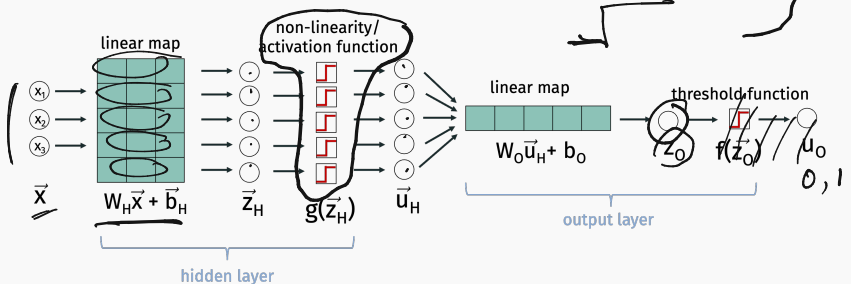
Parameters: $\Theta = [\underbrace{W_H \in \mathbb{R}^{N_H \times N_I}, \beta_H \in \mathbb{R}^{N_H}, W_O \in \mathbb{R}^{1 \times N_H}, \beta_O \in \mathbb{R}}]$.

W_H, W_O are weight matrices and β_H, β_O are bias terms that account for the intercepts of our linear functions.

N_H : number of hidden neurons = 4

POSSIBLE MODEL

Our model is function f which makes x to a class label u_0 .¹



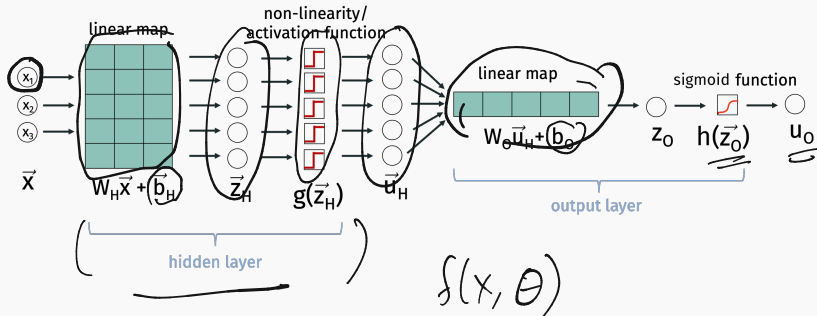
This is called a “multilayer perceptron”: one of the oldest types of neural nets. Dates back to Frank Rosenblatt from 1958

- Number of input variables $N_I = 3$
- Number of hidden variables $N_H = 5$
- Number of output variables $N_O = 1$

¹For regression, would cut off at z_0 to get continuous output.

POSSIBLE MODEL

Our model is function f which maps \mathbf{x} to a class label u_0 .

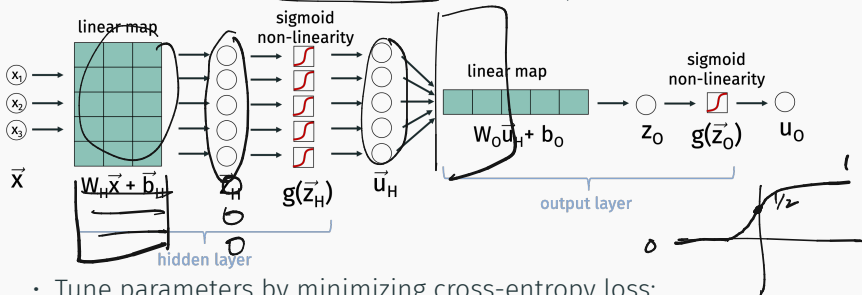


Training the model:

- Choose a loss function $L(f(\mathbf{x}, \Theta), y)$.
- Find optimal parameters: $\Theta^* = \arg \min_{\Theta} \sum_{i=1}^n L(f(\mathbf{x}_i, \Theta), y_i)$
using gradient descent.

FINAL MODEL

A more typical model uses smoother activation functions aka non-linearities, which are more amenable to computing gradients. E.g. we might use the sigmoid function $g(x) = \frac{1}{1+e^{-x}}$.



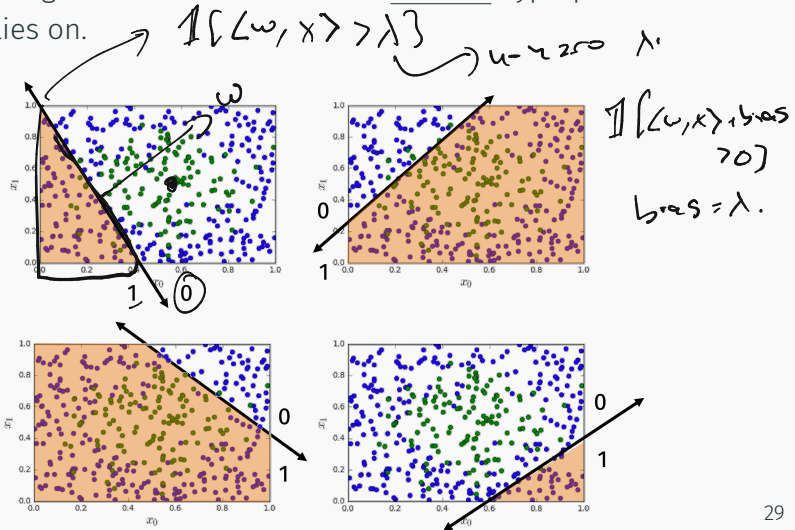
- Tune parameters by minimizing cross-entropy loss:

$$\left(\sum_{i=1}^n L(f(\mathbf{x}_i, \Theta), y_i) = \sum_{i=1}^n -y_i \log(\underline{f(\mathbf{x}_i, \Theta)}) - (1 - y_i) \log(1 - f(\mathbf{x}_i, \Theta)) \right)$$

- We will discuss soon how to compute gradients.

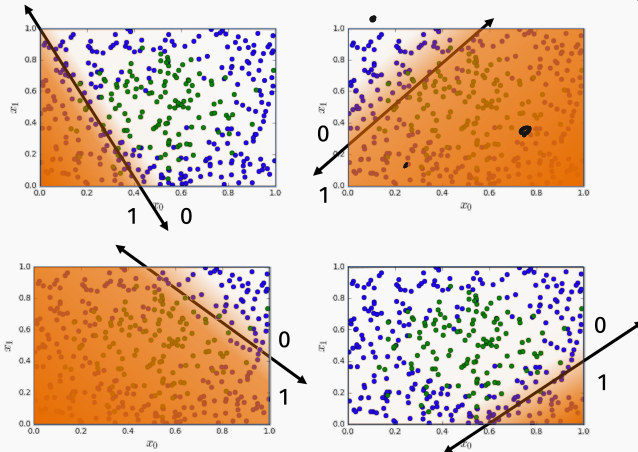
FEATURE EXTRACTION

Features learned using step-function activation are binary, depending on which side of a set of learned hyperplanes each point lies on.



FEATURE EXTRACTION

Features learned using sigmoid activation are real valued in $[0, 1]$. Mimic binary features.

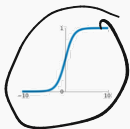


Things we can change in this basic classification network:

- More or less hidden variables.
- We could add more layers.
- Different non-linearity/activation function.)
- Different loss function.

Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

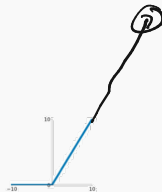


tanh

$$\tanh(x)$$



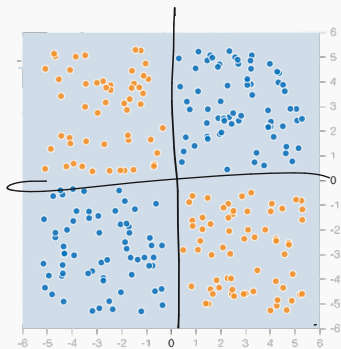
ReLU
 $\max(0, x)$



$$\text{relu}(100) = 100$$
$$\text{relu}(-5) = 0$$

TEST YOUR INTUITION

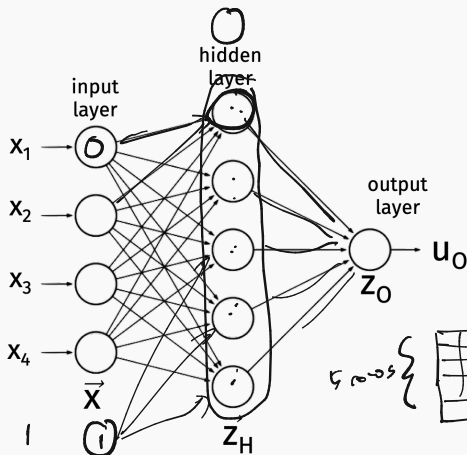
How many hidden variables (e.g. splitting hyperplanes) would be needed to classify this dataset correctly?



<https://playground.tensorflow.org/>)

NOTATION

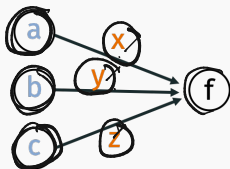
Another common diagram for a 2-layered network:



$x \in \mathbb{R}^5$
 Break to
 3:45

$4 \times 5 = 20 \rightarrow 20 \text{ weights}$

Neural network math:

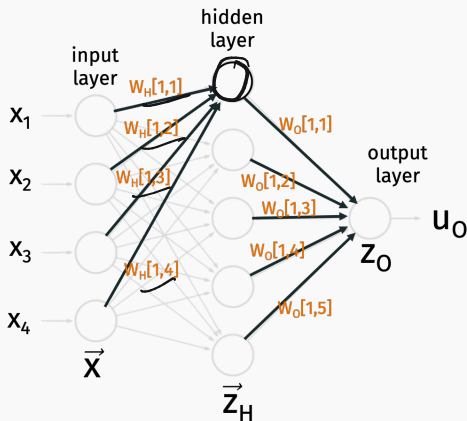


~~$$f = ax + by + cz$$~~

$$f = S(\underline{a}x + \underline{b}y + \underline{c}z + \underline{bias})$$

↓
activation

How to interpret:

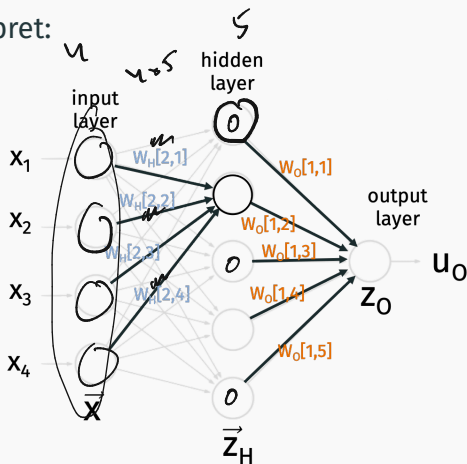


W_H and W_O are our weight matrices from before.

Note: This diagram does not explicitly show the bias terms or the non-linear activation functions.

NOTATION

How to interpret:

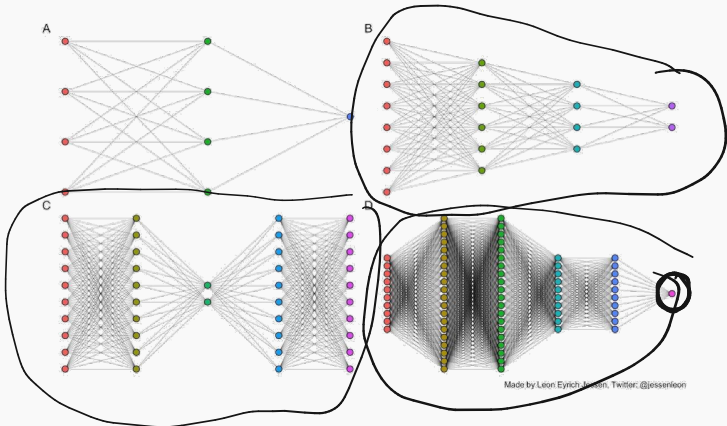


W_H and W_O are our weight matrices from before.

Note: This diagram depicts a network with “fully-connected” layers. Every variable in layer i is connected to every variable in layer $i + 1$.

ARCHITECTURE VISUALIZATION

Effective way of visualize “architecture” of a neural network:

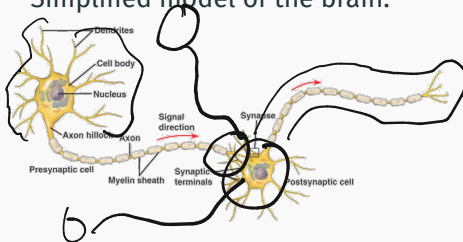


Visualize number of variables, types of connections, number of layers and their relative sizes.

These are all **feedforward** neural networks. No backwards (**recurrent**) connections.

SOME HISTORY AND MOTIVATION

Simplified model of the brain:

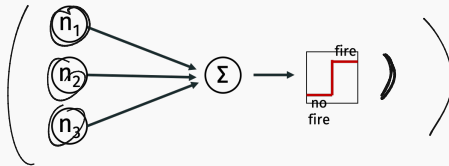


Dendrites: Input electrical current from other neurons.

Axon: Output electrical current to other neurons.

Synapse: Where these two connect.

A neuron “fires” (outputs non-zero electric charge) if it receives enough cumulative electrical input from all neurons connected to it.



Output charge can be positive or negative (excitatory vs. inhibitory)

Inspired early work on neural networks:

- 1940s Donald Hebb proposed a Hebbian learning rule for how brains neurons change over time to allow learning.
- 1950s Frank Rosenblatt's single-layer Perceptron is one of the first attempts to create an "artificial" neural networks.
- Continued work throughout the 1960s.

Main issue with neural network methods: They are hard to train. Gradient descent converges very slowly. Also pretty finicky: user needs to be careful with initialization, regularization, etc. when training. We have gotten a lot better at resolving these issues though!

EARLY NEURAL NETWORK EXPLOSION

Around 1985 several groups (re)-discovered the backpropagation algorithm which allows for efficient training of neural nets via (stochastic) gradient descent. Along with increased computational power this led to a resurgence of interest in neural network models.

Backpropagation Applied to Handwritten Zip Code Recognition

Y. LeCun

B. Boser

J. S. Denker

D. Henderson

R. E. Howard

W. Hubbard

L. D. Jackel

AT&T Bell Laboratories Holmdel, NJ 07733 USA

The ability of learning networks to generalize can be greatly enhanced by providing constraints from the task domain. This paper demonstrates how such constraints can be integrated into a backpropagation network through the architecture of the network. This approach has been successfully applied to the recognition of handwritten zip code digits provided by the U.S. Postal Service. A single network learns the entire recognition operation, going from the normalized image of the character to the final classification.

Very good performance on problems like digit recognition.

From (1990s - 2010,) kernel methods, SVMs, and probabilistic methods began to dominate the literature in machine learning:

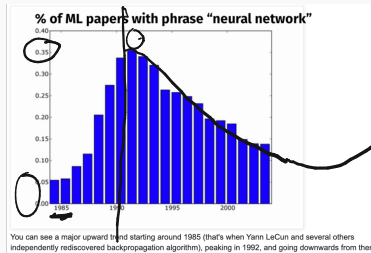
- Work well “out of the box”.
- Relatively easy to understand theoretically.
- Not too computationally expensive for moderately sized datasets.

(Fun blog post to check out from 2005:)

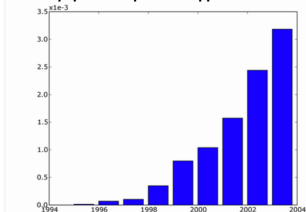
<http://yaroslavvb.blogspot.com/2005/12/trends-in-machine-learning-according.html>

NEURAL NETWORK DECLINE

Finding trends in machine learning by search papers in Google Scholar that match a certain keyword:

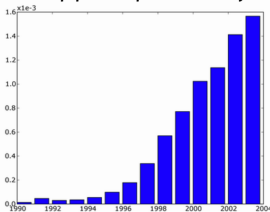


% of ML papers with phrase "support vector machine"



(1995 is when Vapnik and Cortez proposed the algorithm)

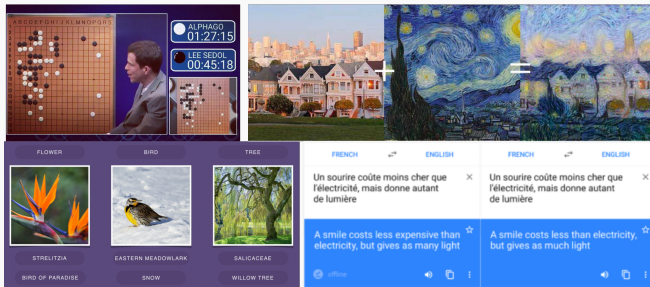
% of ML papers with phrase "naive bayes"



(If I were to trust this, I would say that Naive Bayes research the hottest machine learning area right now)

MODERN NEURAL NETWORK RESURGENCE

In recent years this trend completely turned around:



State-of-the-art results in game playing, image recognition, content generation, natural language processing, machine translation, many other areas.

2019 TURING AWARD WINNERS

“For conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing.”



Yann LeCun




Geoff Hinton



Yoshua Bengio

What were these breakthroughs? What made training large neural networks computationally feasible?

All changed with the introduction of AlexNet and the 2012 ImageNet Challenge...



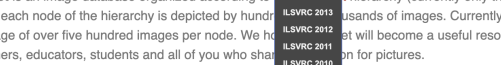
[Explore](#)
[Download](#)
[Challenges](#)
[Publications](#)
[Updates](#)
[About](#)

14,197,122 images, 21841 synsets indexed

Not logged in. [Login](#) | [Signup](#)

ImageNet is an image database organized according to a hierarchy in which each node of the hierarchy is depicted by hundreds of images. We have an average of over five hundred images per node. We have images from researchers, educators, students and all of you who share your images.

Click [here](#) to learn more about ImageNet, Click [here](#) to join the ImageNet mailing list.



What do these images have in common? *Find out!*

Very general image classification task.

All changed with AlexNet and the 2012 ImageNet Challenge...

team name	team members	filename	flat cost	hie cost	description
NEC-UIUC	NEC: Yuanqing Lin, Fengjun Lv, Shenghuo Zhu, Ming Yang, Timothee Cour, Kai Yu UIUC: LiangLiang Cao, Zhen Li, Min-Hsuan Tsai, Xi Zhou, Thomas Huang Rutgers: Tong Zhang	flat_opt.txt	<u>0.28191</u>	2.1144	using <u>sift</u> and <u>lbp</u> feature with two non-linear coding representations and stochastic <u>SVM</u> optimized for top-5 hit rate

2010 Results

Team name	Filename	Error (5 guesses)	Description
<u>SuperVision</u>	test-preds-141-146.2009-131- 137-145-146.2011-145f.	<u>0.15315</u>	Using extra training data from ImageNet Fall 2011 release
<u>SuperVision</u>	test-preds-131-137-145-135- 145f.txt	0.16422	Using only supplied training data
<u>ISI</u>	pred_FVs_wLACs_weighted.txt	<u>0.26172</u>	Weighted sum of scores from each classifier with SIFT+FV, LBP+FV, GIST+FV, and CSIFT+FV, respectively.

2012 Results

ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky
University of Toronto
kriz@cs.utoronto.ca

Ilya Sutskever
University of Toronto
ilya@cs.utoronto.ca

Geoffrey E. Hinton
University of Toronto
hinton@cs.utoronto.ca

Abstract

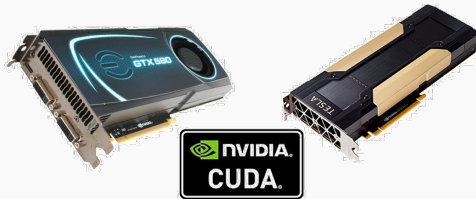
We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

Why 2012?

- Clever ideas in changing neural network architecture and training. E.g. ReLU non-linearities, dropout regularization, batch normalization, data augmentation.
- Wide-spread access to GPU computing power.)

Hardware innovation: Widely available, inexpensive GPUs allowing for cheap, highly parallel linear algebra operations.

- 2007: Nvidia released CUDA platform, which allows GPUs to be easily programmed for general purposed computation.



AlexNet architecture used 60 million parameters. Could not have been trained using CPUs alone (except maybe on a government super computer).

Two main algorithmic tools for training neural network models:

1. Stochastic gradient descent.

2. Backpropagation. → algorithm for computing stochastic gradients.

TRAINING NEURAL NETWORKS

Let $f(\underline{\theta}, \underline{x})$ be our neural network. A typical ℓ -layer feed forward model has the form:

$$g_\ell(\underline{W}_\ell(\dots \underline{W}_3 \cdot \underline{g}_2(\underline{W}_2(\underline{g}_1(\underline{W}_1 \underline{x} + \underline{\beta}_1) + \underline{\beta}_2) + \underline{\beta}_3 \dots) + \underline{\beta}_\ell).$$

\underline{W}_i and $\underline{\beta}_i$ are the weight matrix and bias vector for layer i and g_i is the non-linearity (e.g. sigmoid). $\underline{\theta} = [\underline{W}_0, \underline{\beta}_0, \dots, \underline{W}_\ell, \underline{\beta}_\ell]$ is a vector of all entries in these matrices.

Goal: Given training data $(\underline{x}_1, y_1), \dots, (\underline{x}_n, y_n)$ minimize the loss

$$\mathcal{L}(\underline{\theta}) = \sum_{i=1}^n L(y_i, \underline{f}(\underline{\theta}, \underline{x}_i)),$$

where L is, e.g., binary cross-entropy (logistic) loss:

$$L(y_i, \underline{f}(\underline{\theta}, \underline{x}_i)) = -y_i \log(f(\underline{\theta}, \underline{x}_i)) - (1 - y_i) \log(1 - f(\underline{\theta}, \underline{x}_i)).$$

GRADIENT OF THE LOSS

Approach: minimize the loss by using gradient descent. Which requires us to compute the gradient of the loss function, $\nabla \mathcal{L}$. Note that this gradient has an entry for every value in $\mathbf{W}_0, \boldsymbol{\beta}_0, \dots, \mathbf{W}_\ell, \boldsymbol{\beta}_\ell$.

As usual, our loss function has finite sum structure, so:

$$\nabla \mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^n \nabla L(\underline{y_i}, \underline{f(\boldsymbol{\theta}, \underline{x_i})})$$

So we can focus on computing:

$$\left(\nabla L(\underline{y_i}, \underline{f(\boldsymbol{\theta}, \underline{x_i})}) \right)$$

for a single training example (\mathbf{x}_i, y_i) .

CHAIN RULE REVIEW

$$f(g(x))$$

For a scalar function $f(x)$, we write the derivative with respect to x as:

$$\underline{f'(x)} = \frac{df}{dx} = \lim_{\underline{t \rightarrow 0}} \frac{f(x+t) - f(x)}{t}$$

For a multivariate function $f(\underline{x}, \underline{y}, \underline{z})$ we write the partial derivative with respect to x as:

$$\left(\frac{df}{dx} = \lim_{t \rightarrow 0} \frac{f(x+t, y, z) - f(x, y, z)}{t} \right)$$

$$\frac{df}{dx} \quad \frac{df}{dy} \quad \frac{df}{dz}$$

CHAIN RULE REVIEW

Let $y(x)$ be a function of x and let $f(y)$ be a function of y . The chain rule says that:

$$f(y(x))$$

$$\left(\frac{df}{dx} = \frac{df}{dy} \cdot \frac{dy}{dx} \right)$$

$$\begin{aligned} \left(\frac{df}{dx} \right) &= \lim_{t \rightarrow 0} \frac{f(y(x+t)) - f(y(x))}{t} \\ &= \lim_{t \rightarrow 0} \frac{f(y(x+t)) - f(y(x))}{y(x+t) - y(x)} \cdot \underbrace{\left(\frac{y(x+t) - y(x)}{t} \right)}_{dy/dx} \\ &= \lim_{t \rightarrow 0} \frac{f(y(x) + c) - f(y(x))}{c} \cdot \underbrace{\frac{y(x+t) - y(x)}{t}}_{dy/dx} \end{aligned}$$

where $c = \underline{y(x+t)} - \underline{y(x)}$.

As long as $\lim_{t \rightarrow 0} y(x+t) - y(x) = 0$ then the first term equals

$\frac{df}{dy}$. The second term equals $\frac{dy}{dx}$.

MULTIVARIABLE CHAIN RULE

Let $\underline{y(x)}$, $\underline{z(x)}$, $\underline{w(x)}$ be functions of x and let $\underline{f(y, z, w)}$ be a function of y, z, w .

$$\left(\frac{df}{dx} \right) = \left(\frac{df}{dy} \right) \left(\frac{dy}{dx} \right) + \frac{df}{dz} \cdot \frac{dz}{dx} + \frac{df}{dw} \cdot \frac{dw}{dx} \dots$$

Example: Let $\underline{y(x)} = \underline{x^3}$ and $\underline{z(x)} = \underline{x^2}$. Let $f(y, z) = \underline{y} \cdot \underline{z}$. Then:

$$\frac{df}{dx} = \left(\frac{df}{dy} \cdot \frac{dy}{dx} \right) + \left(\frac{df}{dz} \cdot \frac{dz}{dx} \right) \quad \frac{df}{dx}$$

$$= z \cdot 3x^2 + y \cdot 2x$$

$$= x^2 \cdot 3x^2 + x^3 \cdot 2x$$

$$= 5x^4$$

$$f(y, z) = y \cdot z$$

$$f(x) = x^3 \cdot x^2 = \underline{\underline{x^5}}$$

$$\frac{df}{dx} = 5x^4$$

Applying chain rule each partial derivative of the loss:

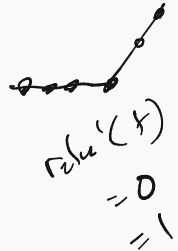
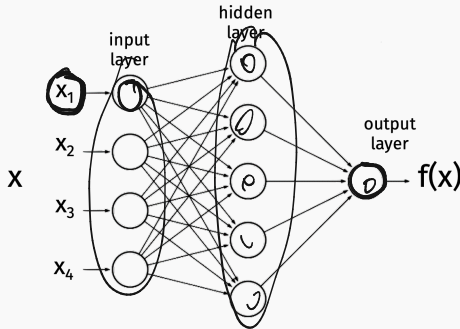
$$\nabla L(y, f(\boldsymbol{\theta}, \mathbf{x})) = \frac{\partial L}{\partial f(\boldsymbol{\theta}, \mathbf{x})} \cdot \nabla f(\boldsymbol{\theta}, \mathbf{x})$$

Binary cross-entropy example:

$$L(y, f(\boldsymbol{\theta}, \mathbf{x})) = -y \log(f(\boldsymbol{\theta}, \mathbf{x})) - (1 - y) \log(1 - f(\boldsymbol{\theta}, \mathbf{x}))$$

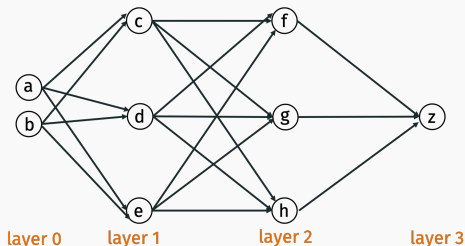
GRADIENT OF THE LOSS

We have reduced our goal to computing $\nabla f(\theta, x)$, where the gradient is with respect to the parameters θ .



Backpropagation is an efficient way to compute $\nabla f(\theta, x)$. It derives its name because we compute gradient from back to front: starting with the parameters closest to the output of the neural net.

BACKPROP EXAMPLE



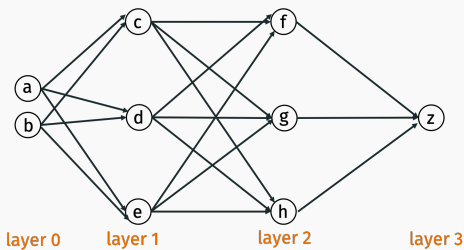
Notation for few slides:

- a, b, \dots, z are the node names, and denote values at the nodes after applying non-linearity.
- $\bar{a}, \bar{b}, \dots, \bar{z}$ denote values before applying non-linearity.
- $W_{i,j}$ is the weight of edge from node i to node j .
- $s(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is the non-linear activation function.
- β_j is the bias for node j .

Example: $h = s(\bar{h}) = s(c \cdot W_{c,h} + d \cdot W_{d,h} + e \cdot W_{e,h} + \beta_h)$

BACKPROP EXAMPLE

For any node j , let \bar{j} denote the value obtained before applying the non-linearity g .



So if $h = s(c \cdot W_{c,h} + d \cdot W_{d,h} + e \cdot W_{e,h} + \beta_h)$ then we use \bar{h} to denote:

$$\bar{h} = c \cdot W_{c,h} + d \cdot W_{d,h} + e \cdot W_{e,h} + \beta_h$$

BACKPROP EXAMPLE

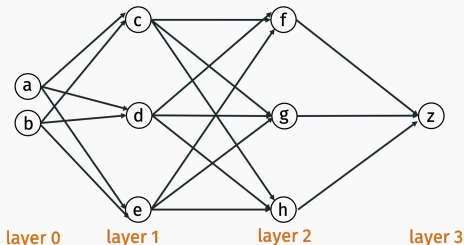
Goal: Compute the gradient $\nabla f(\boldsymbol{\theta}, \mathbf{x})$, which contains the partial derivatives with respect to every parameter:

- $\partial z / \partial \beta_z$
- $\partial z / \partial W_{f,z}, \partial z / \partial W_{g,z}, \partial z / \partial W_{h,z}$
- $\partial z / \partial \beta_f, \partial z / \partial \beta_g, \partial z / \partial \beta_h$
- $\partial z / \partial W_{c,f}, \partial z / \partial W_{c,g}, \partial z / \partial W_{c,h}$
- $\partial z / \partial W_{d,f}, \partial z / \partial W_{d,g}, \partial z / \partial W_{d,h}$
- \vdots
- $\partial z / \partial W_{a,c}, \partial z / \partial W_{a,d}, \partial z / \partial W_{a,e}$

Two steps: Forward pass to compute function value.
Backwards pass to compute gradients.

BACKPROP EXAMPLE

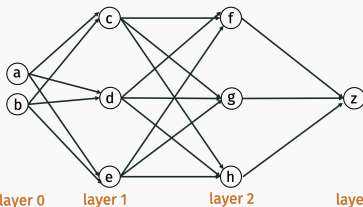
Step 1: Forward pass.



- Using current parameters, compute the output z by moving from left to right.
- Store all intermediate results:

$$\bar{c}, \bar{d}, \bar{e}, c, d, e, \bar{f}, \bar{g}, \bar{h}, f, g, h, \bar{z}, z.$$

BACKPROP EXAMPLE



Step 1: Forward pass.

$$\bar{c} = W_{a,c} \cdot a + W_{b,c} \cdot b + \beta_c \quad c = s(\bar{c})$$

$$\bar{d} = W_{a,d} \cdot a + W_{b,d} \cdot b + \beta_d \quad d = s(\bar{d})$$

$$\bar{e} = W_{a,e} \cdot a + W_{b,e} \cdot b + \beta_e \quad e = s(\bar{e})$$

$$\bar{f} = W_{c,f} \cdot c + W_{d,f} \cdot d + W_{e,f} \cdot e + \beta_f \quad f = s(\bar{f})$$

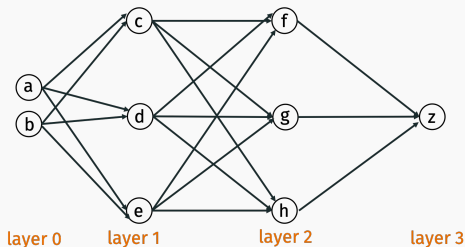
\vdots

$$\bar{z} = W_{f,z} \cdot f + W_{g,z} \cdot g + W_{h,z} \cdot h + \beta_z \quad z = s(\bar{z})$$

Question: What is runtime in terms of # of parameters P ?

BACKPROP EXAMPLE

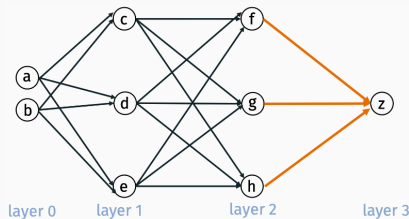
Step 2: Backward pass.



- Using **current parameters** and **computed node values**, compute the partial derivatives of all parameters by moving from right to left.

BACKPROP EXAMPLE

Step 2: Backward pass. Deepest layer.



$$\frac{\partial z}{\partial \beta_z} = \frac{\partial \bar{z}}{\partial \beta_z} \cdot \frac{\partial z}{\partial \bar{z}} = 1 \cdot s'(\bar{z})$$

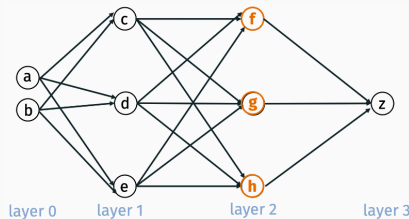
$$\frac{\partial z}{\partial W_{f,z}} = \frac{\partial \bar{z}}{\partial W_{f,z}} \cdot \frac{\partial z}{\partial \bar{z}} = f \cdot s'(\bar{z})$$

$$\frac{\partial z}{\partial W_{g,z}} = \frac{\partial \bar{z}}{\partial W_{g,z}} \cdot \frac{\partial z}{\partial \bar{z}} = g \cdot s'(\bar{z})$$

$$\frac{\partial z}{\partial W_{h,z}} = \frac{\partial \bar{z}}{\partial W_{h,z}} \cdot \frac{\partial z}{\partial \bar{z}} = h \cdot s'(\bar{z})$$

BACKPROP EXAMPLE

Step 2: Backward pass.



$$\frac{\partial z}{\partial f} = \frac{\partial \bar{z}}{\partial f} \cdot \frac{\partial z}{\partial \bar{z}} = W_{f,z} \cdot s'(\bar{z})$$

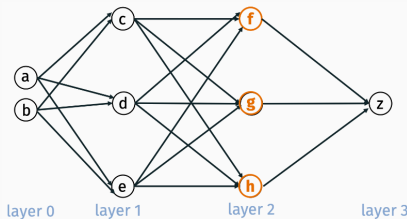
$$\frac{\partial z}{\partial g} = \frac{\partial \bar{z}}{\partial g} \cdot \frac{\partial z}{\partial \bar{z}} = W_{g,z} \cdot s'(\bar{z})$$

$$\frac{\partial z}{\partial h} = \frac{\partial \bar{z}}{\partial h} \cdot \frac{\partial z}{\partial \bar{z}} = W_{h,z} \cdot s'(\bar{z})$$

Compute partial derivatives with respect to nodes, even though these are not used in the gradient.

BACKPROP EXAMPLE

Step 2: Backward pass.



$$\frac{\partial z}{\partial \bar{f}} = \frac{\partial z}{\partial f} \cdot \frac{\partial f}{\partial \bar{f}} = \frac{\partial z}{\partial f} \cdot s'(\bar{f})$$

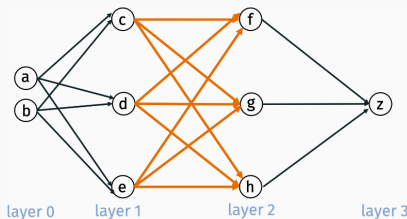
$$\frac{\partial z}{\partial \bar{g}} = \frac{\partial z}{\partial g} \cdot \frac{\partial g}{\partial \bar{g}} = \frac{\partial z}{\partial g} \cdot s'(\bar{g})$$

$$\frac{\partial z}{\partial \bar{h}} = \frac{\partial z}{\partial h} \cdot \frac{\partial h}{\partial \bar{h}} = \frac{\partial z}{\partial h} \cdot s'(\bar{h})$$

And for “pre-nonlinearity” nodes.

BACKPROP EXAMPLE

Step 2: Backward pass. Next layer.



$$\frac{\partial z}{\partial \beta_f} = \frac{\partial z}{\partial \bar{f}} \cdot \frac{\partial \bar{f}}{\partial \beta_f} = \frac{\partial z}{\partial \bar{f}} \cdot 1$$

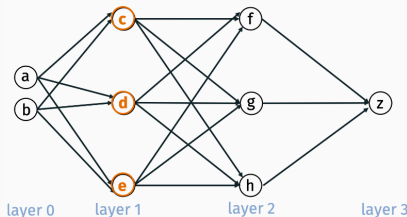
$$\frac{\partial z}{\partial W_{c,f}} = \frac{\partial z}{\partial \bar{f}} \cdot \frac{\partial \bar{f}}{\partial W_{c,f}} = \frac{\partial z}{\partial \bar{f}} \cdot c$$

$$\frac{\partial z}{\partial W_{d,f}} = \frac{\partial z}{\partial \bar{f}} \cdot \frac{\partial \bar{f}}{\partial W_{d,f}} = \frac{\partial z}{\partial \bar{f}} \cdot d$$

$$\frac{\partial z}{\partial W_{e,f}} = \frac{\partial z}{\partial \bar{f}} \cdot \frac{\partial \bar{f}}{\partial W_{e,f}} = \frac{\partial z}{\partial \bar{f}} \cdot e$$

BACKPROP EXAMPLE

Step 2: Backward pass. Next layer. Use multivariate chain rule.



$$\frac{\partial z}{\partial c} = \frac{\partial z}{\partial \bar{f}} \cdot \frac{\partial \bar{f}}{\partial c} + \frac{\partial z}{\partial \bar{g}} \cdot \frac{\partial \bar{g}}{\partial c} + \frac{\partial z}{\partial \bar{h}} \cdot \frac{\partial \bar{h}}{\partial c}$$

$$= \frac{\partial z}{\partial \bar{f}} \cdot W_{c,f} + \frac{\partial z}{\partial \bar{g}} \cdot W_{c,g} + \frac{\partial z}{\partial \bar{h}} \cdot W_{c,h}$$

$$\frac{\partial z}{\partial d} = \frac{\partial z}{\partial \bar{f}} \cdot W_{d,f} + \frac{\partial z}{\partial \bar{g}} \cdot W_{d,g} + \frac{\partial z}{\partial \bar{h}} \cdot W_{d,h}$$

$$\frac{\partial z}{\partial e} = \frac{\partial z}{\partial \bar{f}} \cdot W_{e,f} + \frac{\partial z}{\partial \bar{g}} \cdot W_{e,g} + \frac{\partial z}{\partial \bar{h}} \cdot W_{e,h}$$

Linear algebraic view.

Let \mathbf{v}_i be a vector containing the value of all nodes j in layer i .

$$\mathbf{v}_3 = \begin{bmatrix} z \end{bmatrix} \qquad \mathbf{v}_2 = \begin{bmatrix} f \\ g \\ h \end{bmatrix} \qquad \mathbf{v}_1 = \begin{bmatrix} c \\ d \\ e \end{bmatrix}$$

Let $\bar{\mathbf{v}}_i$ be a vector containing \bar{j} for all nodes j in layer i .

$$\bar{\mathbf{v}}_3 = \begin{bmatrix} \bar{z} \end{bmatrix} \qquad \bar{\mathbf{v}}_2 = \begin{bmatrix} \bar{f} \\ \bar{g} \\ \bar{h} \end{bmatrix} \qquad \bar{\mathbf{v}}_1 = \begin{bmatrix} \bar{c} \\ \bar{d} \\ \bar{e} \end{bmatrix}$$

Note: $\mathbf{v}_i = s(\bar{\mathbf{v}}_i)$, where s is applied entrywise.

Linear algebraic view.

Let δ_i be a vector containing $\partial z / \partial j$ for all nodes j in layer i .

$$\delta_3 = \begin{bmatrix} 1 \end{bmatrix} \quad \delta_2 = \begin{bmatrix} \partial z / \partial f \\ \partial z / \partial g \\ \partial z / \partial h \end{bmatrix} \quad \delta_1 = \begin{bmatrix} \partial z / \partial c \\ \partial z / \partial d \\ \partial z / \partial e \end{bmatrix}$$

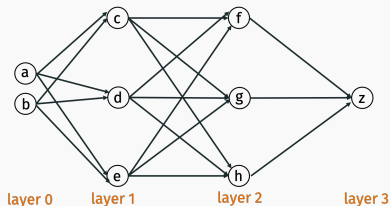
Let $\bar{\delta}_i$ be a vector containing $\partial z / \partial \bar{j}$ for all nodes j in layer i .

$$\bar{\delta}_3 = \begin{bmatrix} \partial z / \partial \bar{z} \end{bmatrix} \quad \bar{\delta}_2 = \begin{bmatrix} \partial z / \partial \bar{f} \\ \partial z / \partial \bar{g} \\ \partial z / \partial \bar{h} \end{bmatrix} \quad \bar{\delta}_1 = \begin{bmatrix} \partial z / \partial \bar{c} \\ \partial z / \partial \bar{d} \\ \partial z / \partial \bar{e} \end{bmatrix}$$

Note: $\bar{\delta}_i = s'(\bar{\mathbf{v}}_i) \times \delta_i$ where \times denotes entrywise multiplication.

BACKPROP LINEAR ALGEBRA

Let W_i be a matrix containing all the weights for edges between layer i and layer $i + 1$.

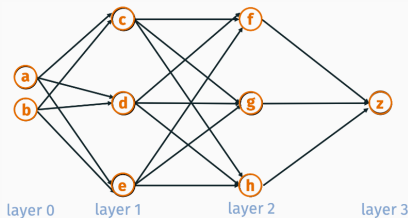


$$W_0 = \begin{bmatrix} W_{a,c} & W_{b,c} \\ W_{a,d} & W_{b,d} \\ W_{a,e} & W_{b,e} \end{bmatrix}$$

$$W_1 = \begin{bmatrix} W_{c,f} & W_{d,f} & W_{e,f} \\ W_{c,g} & W_{d,g} & W_{e,g} \\ W_{c,h} & W_{d,h} & W_{e,h} \end{bmatrix}$$

$$W_2 = \begin{bmatrix} W_{f,z} & W_{g,z} & W_{h,z} \end{bmatrix}$$

BACKPROP LINEAR ALGEBRA



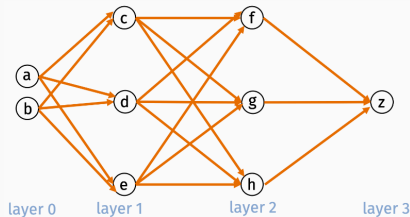
Claim 1: Node derivative computation is matrix multiplication.

$$\delta_i = W_i^T \bar{\delta}_{i+1}$$

What is the computational complexity if $W_i \in \mathbb{R}^{k \times m}$?

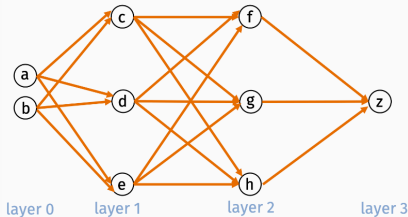
BACKPROP LINEAR ALGEBRA

Let Δ_i be a matrix contain the derivatives for all weights for edges between layer i and layer $i + 1$.



$$\Delta_2 = \begin{bmatrix} \partial z / \partial W_{f,z} & \partial z / \partial W_{g,z} & \partial z / \partial W_{h,z} \end{bmatrix}$$
$$\Delta_1 = \begin{bmatrix} \partial z / \partial W_{c,f} & \partial z / \partial W_{d,f} & \partial z / \partial W_{e,f} \\ \partial z / \partial W_{c,g} & \partial z / \partial W_{d,g} & \partial z / \partial W_{e,g} \\ \partial z / \partial W_{c,h} & \partial z / \partial W_{d,h} & \partial z / \partial W_{e,h} \end{bmatrix}$$
$$\Delta_0 = \dots$$

BACKPROP LINEAR ALGEBRA



Claim 2: Weight derivative computation is an outer-product between the $(i + 1)^{\text{st}}$ derivative vector and the i^{th} value vector.

$$\Delta_i = \mathbf{v}_i \delta_{i+1}^T.$$

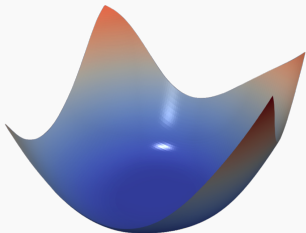
What is the computational complexity of computing the derivatives for a single weight matrix $\mathbf{W}_i \in \mathbb{R}^{k \times m}$?

Takeaways:

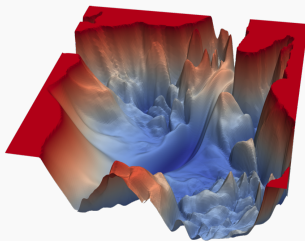
- Backpropagation can be used to compute derivatives for all weights and biases for any feedforward neural network.
- Total computation cost is linear in the number of parameters of the network to compute $f(\boldsymbol{\theta}, \mathbf{x})$ and thus $\nabla L(y, f(\boldsymbol{\theta}, \mathbf{x}))$ for a single training example \mathbf{x}, \mathbf{y} .
- SGD can be run in $O(P)$ time per iteration for a network with P parameters.
- Final computation boils down to linear algebra operations (matrix multiplication and vector operations) which can be performed quickly on a GPU.

CONVERGENCE

Least squares regression, logistic regression, SVMs, even all of these with kernels lead to convex losses.



convex loss



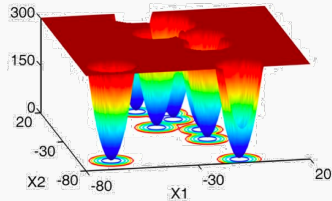
cross-entropy loss for
neural net

Neural networks very much do not...

CONVERGENCE

But SGD still performs remarkably well in practice. Understanding this phenomenon is still an open research question in machine learning and optimization. Current hypotheses include:

- Initialization seems important (random uniform vs. random Gaussian vs. Xavier initialization vs. He initialization vs. etc.)
- Randomization helps in escaping local minima.
- Many local minima are global minima?
- SGD finds “good” local minima?



Issue: Backpropagation + SGD is fast, but tedious to implement.

Typical to use automatic differentiation, which can compute the gradient of pretty much any function you can code up.

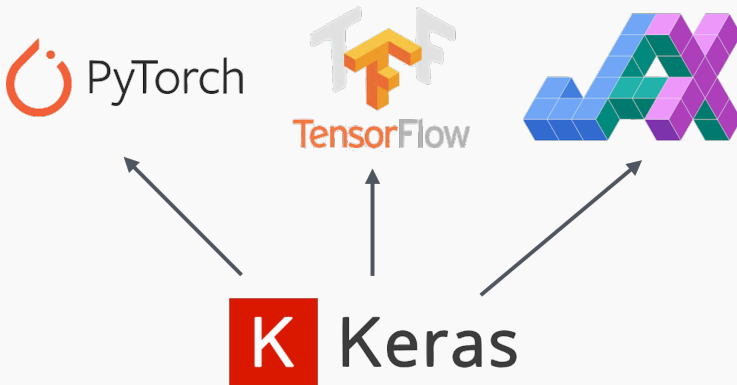
```
def loss(W, b):  
    preds = predict(W, b, inputs)  
    label_probs = preds * targets + (1 - preds)  
    return -np.sum(jnp.log(label_probs))  
  
from jax import grad  
W_grad, b_grad = grad(loss, (0, 1))(W, b)  
print('W_grad', W_grad)  
print('b_grad', b_grad)
```

May mature low-level libraries that handle neural network representation, autodiff, have built in optimizers (SGD, ADAM, etc.), etc.



LIBRARIES

Higher-level libraries like Keras make it even easy to work with this software. Tools for easily defining and building neural networks with specific structure, tracking training, etc.



Define:

```
model = Sequential()  
model.add(Dense(units=nh, input_shape=(nin,), activation='sigmoid', name='hidden'))  
model.add(Dense(units=nout, activation='softmax', name='output'))
```

Compile:

```
opt = optimizers.Adam(lr=0.001) |  
model.compile(optimizer=opt,  
              loss='sparse_categorical_crossentropy',  
              metrics=['accuracy'])
```

Train:

```
hist = model.fit(Xtr, ytr, epochs=30, batch_size=100, validation_data=(Xts,yts))
```

We will release two demos on working with Keras:
keras_demo_synthetic.ipynb and
keras_demo_mnist.ipynb