

CS-GY 6923: Lecture 9

Neural Nets Introduction, Back propagation

NYU Tandon School of Engineering, Prof. Christopher Musco

Key Concept

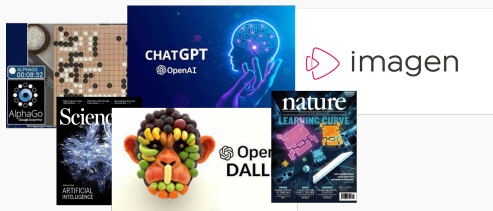
Approach in prior classes:

- 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 hot-topic in machine learning right now. With no sign of slowing down.

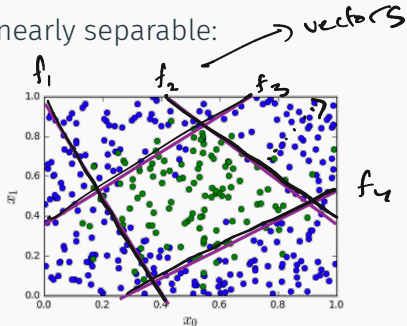
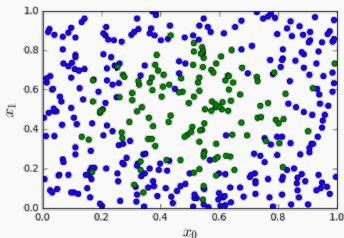


Focus of investment at universities, government research labs, funding agencies, and large tech companies.

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:

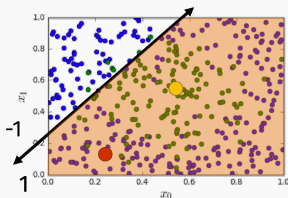
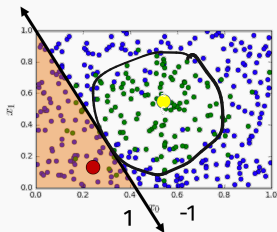


Could use feature transformations or a non-linear kernel.

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

SIMPLE MOTIVATING EXAMPLE

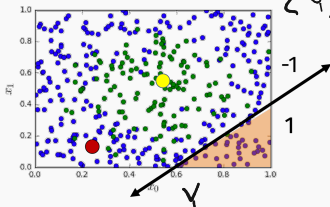
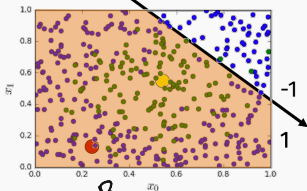
For each linear classifier β , add a new $-1, 1$ feature for every example $\mathbf{x} = [x_0, x_1]$ depending on the sign of $\langle \mathbf{x}, \beta \rangle$.



$$\begin{bmatrix} -1 \\ 1 \\ 1 \\ -1 \end{bmatrix} = \beta_5$$

$$\langle \beta_1, \beta_5 \rangle = 3$$

$$\langle \beta_4, \beta_5 \rangle = 4$$



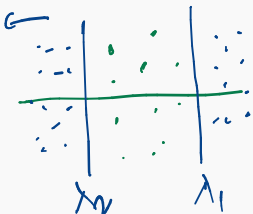
$$\bullet \Rightarrow [1, 1, 1, -1]$$

$$\bullet \Rightarrow [-1, 1, 1, -1]$$

SIMPLE MOTIVATING EXAMPLE

$$\begin{bmatrix} .2, .8, \\ .5, .5 \\ \vdots \\ .5, 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} = \begin{bmatrix} -1, -1, +1, -1 \\ -1, +1, +1, -1 \\ \vdots \\ -1, -1, -1, -1 \end{bmatrix} \begin{matrix} \rightarrow 0_H \\ \\ \\ b_S \end{matrix}$$

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



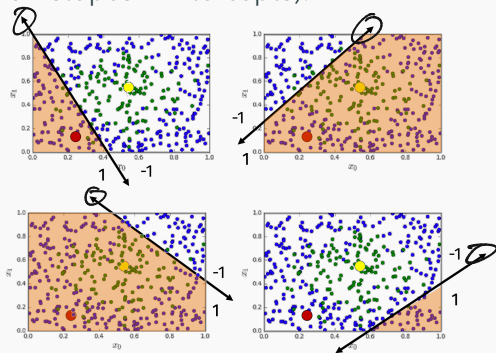
$$\begin{bmatrix} -1, -1, +1, -1 \\ -1, +1, +1, -1 \\ \vdots \\ -1, -1, -1, -1 \end{bmatrix} \xrightarrow{?} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

left of λ_1
 $\rightarrow 1$
 right of λ_2
 $\rightarrow 2$
 otherwise get -1

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

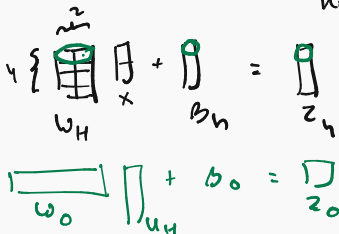
Input: $\underline{x} = x_1, \dots, x_{N_I}$ → number of inputs
Model: $f(x, \Theta)$: → number of hidden features or "neurons"

- $z_H \in \mathbb{R}^{N_H} = \underline{W}_H \underline{x} + \beta_h$

- $\underline{u}_H = \text{sign}(z_H)$

- $z_0 \in \mathbb{R} = \underline{W}_0 \underline{u}_H + \beta_0$

- $\underline{u}_0 = \mathbb{1}[z_0 > \lambda]$

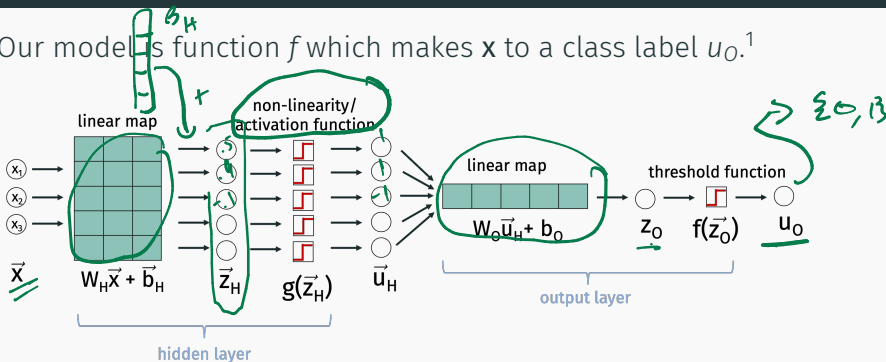


Parameters: $\Theta = [\underline{W}_H \in \mathbb{R}^{N_H \times N_I}, \underline{\beta}_H \in \mathbb{R}^{N_H}, \underline{W}_0 \in \mathbb{R}^{1 \times N_H}, \underline{\beta}_0 \in \mathbb{R}]$.

| W_H, W_0 are weight matrices and β_H, β_0 are bias terms that account for the intercepts of our linear functions.

POSSIBLE MODEL

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



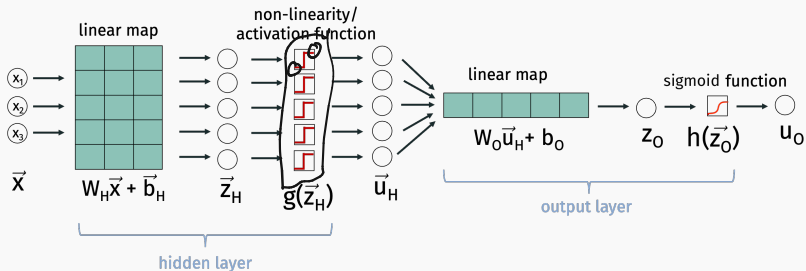
This is called a “multilayer perceptron”: one of the oldest types of neural nets. Dates back to Frank Rosenblatt from 1958

- 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 .



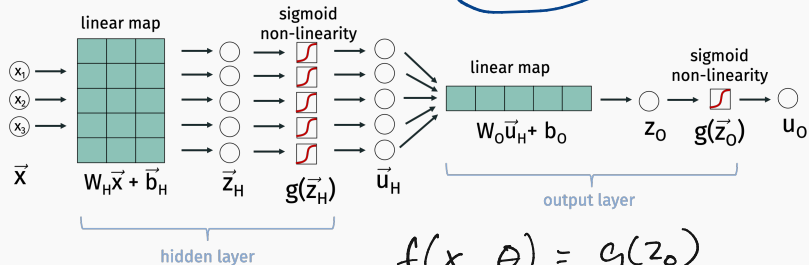
Training the model:

- 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** $g(x) = \frac{1}{1+e^{-x}}$.



- Use cross-entropy loss:

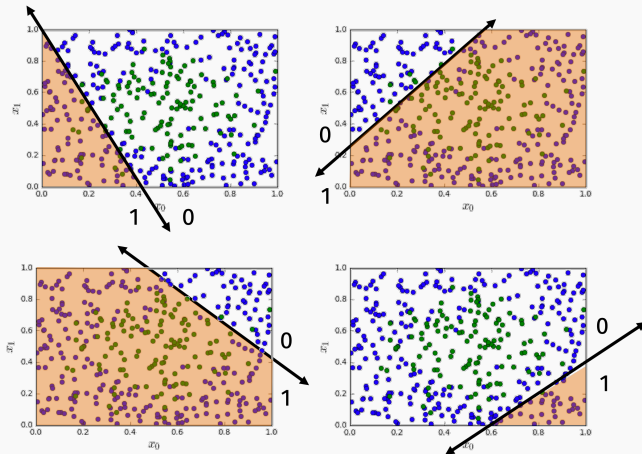
$$L(f(x_i, \theta), y_i) = -y_i \log(f(x_i, \theta)) - (1 - y_i) \log(1 - f(x_i, \theta))$$

- We will discuss soon how to compute gradients.

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

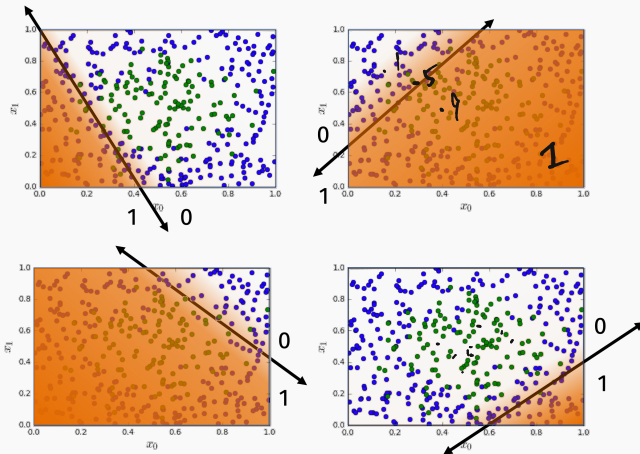
FEATURE EXTRACTION

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. \mathcal{D}

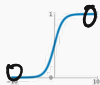


Things we can change in this basic classification network:

- 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

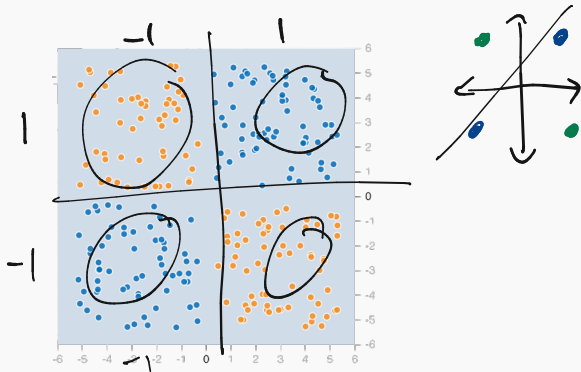
$$\max(0, x)$$



$$\text{ReLU}(x) = \max(0, x)$$

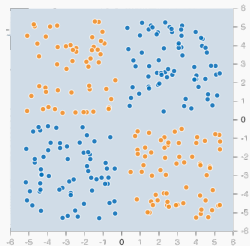
TEST YOUR INTUITION

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

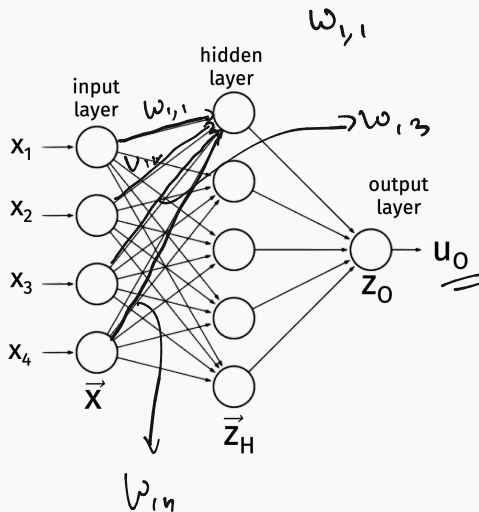


<https://playground.tensorflow.org/>

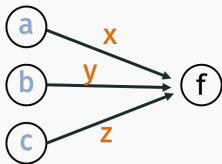
TEST YOUR INTUITION



Another common diagram for a 2-layered network:

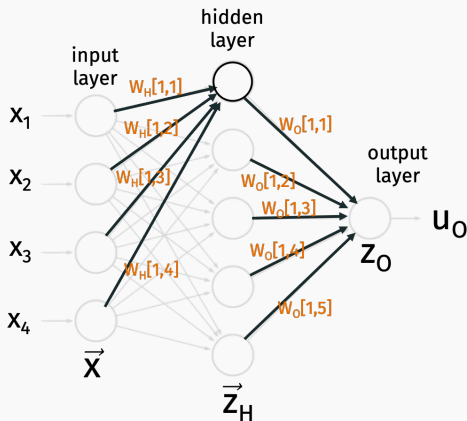


Neural network math:



$$\underline{f} = ax + by + cz$$

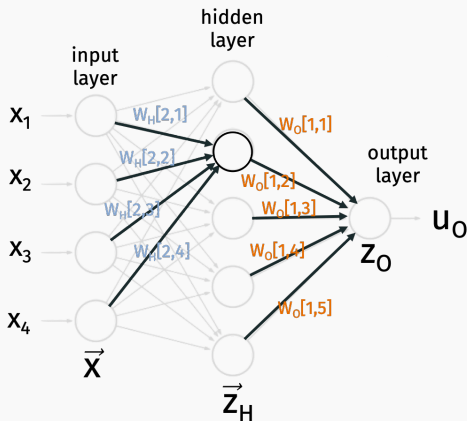
How to interpret:



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.

How to interpret:

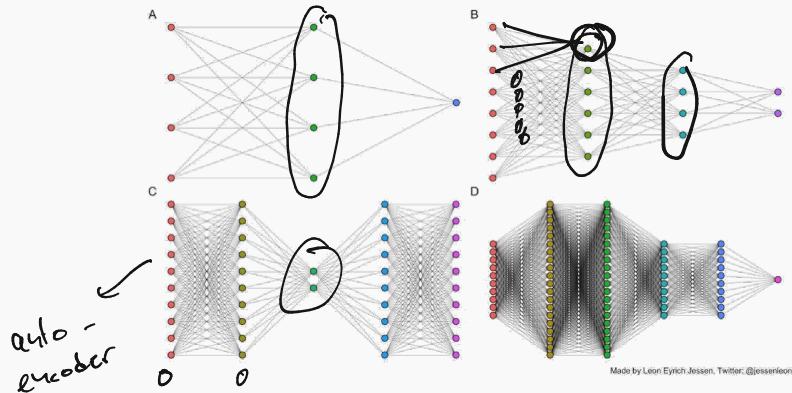


W_H and W_O are our weight matrices from before.

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:



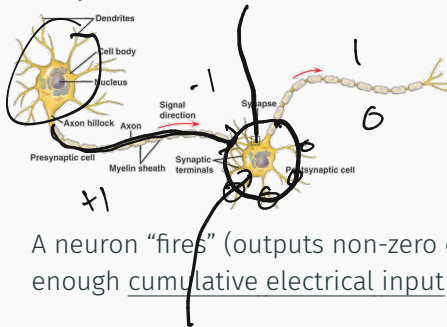
Visualize number of variables, types of connections, number of layers and their relative sizes.

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

SOME HISTORY AND MOTIVATION

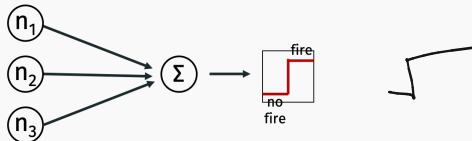
CONNECTION TO BIOLOGY

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 Perceptron is one of the first "artificial" neural networks.
- Continued work throughout the 1960s.

Main issue with neural network methods: They are hard to train. Generally require a lot of computation power. 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:

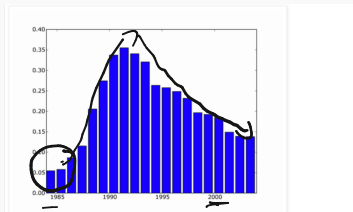
- Work well “out of the box”.
- Relatively easy to understand theoretically.
- Not too computationally expensive for moderately sized datasets.

Fun blog post to check out from 2005:

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

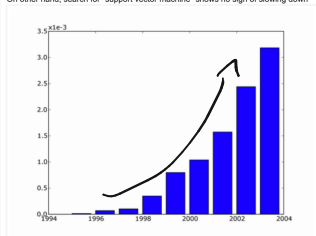
NEURAL NETWORK DECLINE

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



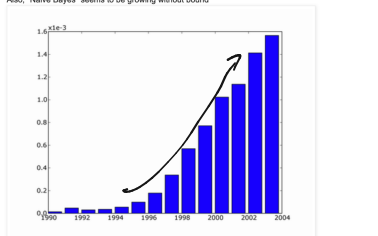
You can see a major upward trend starting around 1985 (that's when Yann LeCun and several others independently rediscovered backpropagation algorithm), peaking in 1992, and going downwards from then.

On other hand, search for "support vector machine" shows no sign of slowing down



(1995 is when Vapnik and Cortez proposed the algorithm)

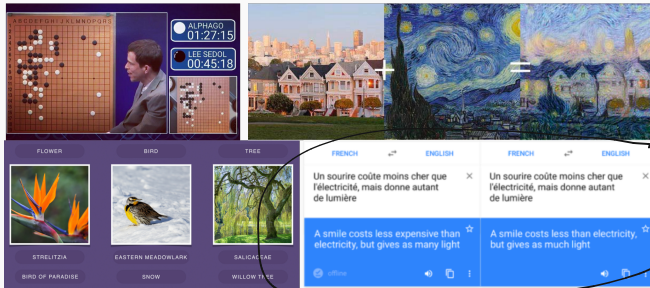
Also, "Naive Bayes" seems to be growing without bound



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...

IMAGENET 14,197,122 images, 21841 synsets indexed
 Explore Download Challenges Publications Updates About
 Not logged in. Login | Signup

ImageNet is an image database organized according to a hierarchical structure (currently only the nouns), in which each node of the hierarchy is depicted by hundreds of images. Currently we have an average of over five hundred images per node. We hope that this database will become a useful resource for researchers, educators, students and all of you who share an interest in pictures.

Click here to learn more about ImageNet, Click here to join the ImageNet mailing list.

ILSVRC 2017
 ILSVRC 2016
 ILSVRC 2015
 ILSVRC 2014
 ILSVRC 2013
 ILSVRC 2012
 ILSVRC 2011
 ILSVRC 2010

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	0.28191	2.1144	using sift and lbp feature with two non-linear coding representations and stochastic SVM, optimized for top-5 hit rate

2010 Results

Team name	Filename	Error (5 guesses)	Description
SuperVision	test-preds-141-146.2009-131-137-145-146.2011-145f.	0.15315	Using extra training data from ImageNet Fall 2011 release
SuperVision	test-preds-131-137-145-135-145f.txt	0.16422	Using only supplied training data
ISI	pred_FVs_wLACs_weighted.txt	0.26172	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

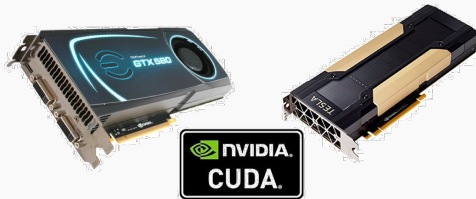
We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. 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).

12:45

Two main algorithmic tools for training neural network models:

1. Stochastic gradient descent.)
2. Backpropagation.)

TRAINING NEURAL NETWORKS

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

$$g_l(W_l(\dots W_3 \cdot g_2(W_2 \cdot g_1(W_1 x + \beta_1) + \beta_2) + \beta_3 \dots) + \beta_l).$$

W_i and β_i are the weight matrix and bias vector for layer i and g_i is the non-linearity (e.g. sigmoid). $\theta = [W_0, \beta_0, \dots, W_l, \beta_l]$ is a vector of all entries in these matrices.

Goal: Given training data $(x_1, y_1), \dots, (x_n, y_n)$ minimize the loss

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

Handwritten annotations: "true label" with an arrow pointing to y_i ; "prediction" with an arrow pointing to $f(\theta, x_i)$.

Example: We might use the binary cross-entropy loss for binary classification:


$$L(y_i, f(\theta, x_i)) = -y_i \log(f(\theta, x_i)) - (1 - y_i) \log(1 - f(\theta, x_i))$$

Handwritten annotations: "outputs probabilities" with an arrow pointing to $f(\theta, x_i)$.

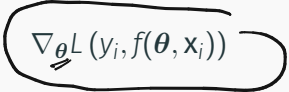
GRADIENT OF THE LOSS

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 $\mathbf{W}_0, \beta_0, \dots, \mathbf{W}_\ell, \beta_\ell$.

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

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


So we can focus on computing:

$$\nabla_{\boldsymbol{\theta}} L(y_i, f(\boldsymbol{\theta}, \mathbf{x}_i))$$


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

CHAIN RULE REVIEW

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

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

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

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

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:

$$\frac{df}{dx} = \frac{df}{dy} \frac{dy}{dx} = \frac{df}{dx}$$

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

where $c = y(x+t) - y(x)$. 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 $y(x), z(x), w(x)$ be functions of x and let $f(y, z, w)$ be a function of y, z, w .

$$\frac{df}{dx} = \frac{df}{dy} \cdot \frac{dy}{dx} + \frac{df}{dz} \cdot \frac{dz}{dx} + \frac{df}{dw} \cdot \frac{dw}{dx}$$

$$\rightarrow x^3 \cdot x^2 = x^5$$

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

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

$$= \underset{\downarrow}{z} \cdot 3x^2 + y \cdot 2x$$

$$= x^2 \cdot 3x^2 + x^3 \cdot 2x = 5x^4$$

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

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

GRADIENT OF THE LOSS

Applying chain rule each partial derivative of the loss:

$$\nabla_{\theta} L(y, f(\theta, x)) = \frac{\partial L}{\partial f(\theta, x)} \nabla_{\theta} f(\theta, x)$$

Handwritten annotations:
An arrow points from $\in \{0, 1\}$ to y .
A bracket under y is labeled $\in \mathbb{R}$.
The terms $\frac{\partial L}{\partial f(\theta, x)}$ and $\nabla_{\theta} f(\theta, x)$ are circled in the original image.

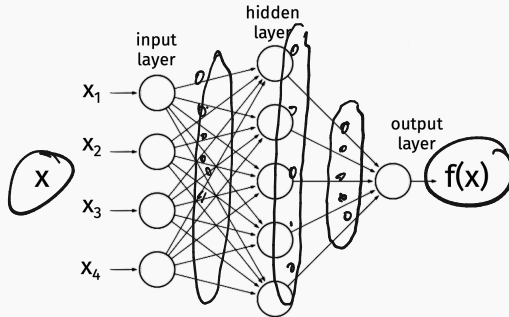
Binary cross-entropy example:

$$L(y, f(\theta, x)) = -y \log(f(\theta, x)) - (1 - y) \log(1 - f(\theta, x))$$

$$\frac{\partial L}{\partial f(\theta, x)} = -y \cdot \frac{1}{f(\theta, x)} - (1 - y) \frac{1}{1 - f(\theta, x)} \cdot (-1)$$

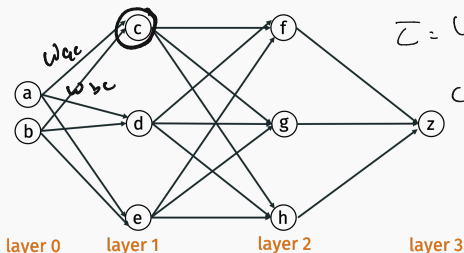
GRADIENT OF THE LOSS

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



Back-propagation is an efficient way to compute $\nabla_{\theta} f(\theta, \mathbf{x})$. 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



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

↓
non-linearity

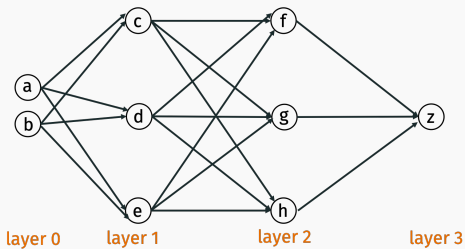
Notation for few slides:

- $\underline{a}, \underline{b}, \dots, \underline{z}$ 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.
- $\underline{\beta}_j$ 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

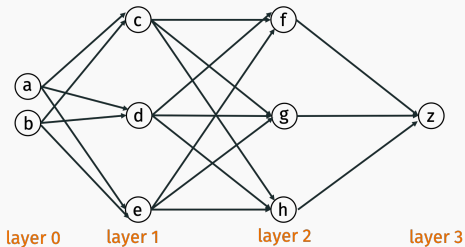
Goal: Compute the gradient $\nabla f(\boldsymbol{\theta}, \mathbf{x})$, which contains the partial derivatives with respect to every parameter:

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



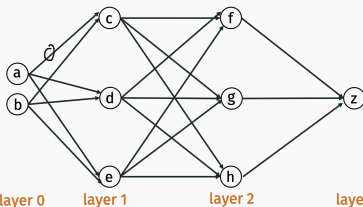
- Using current parameters, compute the output z by moving from left to right.
- Store all intermediate results:

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

BACKPROP EXAMPLE

$O(P)$

$O(P)$



Step 1: Forward pass.

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

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

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

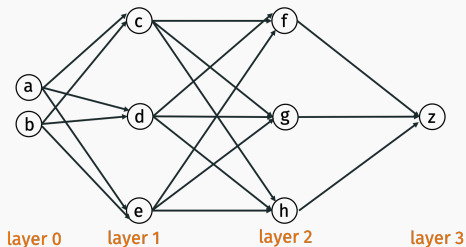
$$\bar{f} = W_{c,f} \cdot c + W_{d,f} \cdot d + W_{e,f} \cdot e + \beta_f \qquad f = s(\bar{f})$$

\vdots

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

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

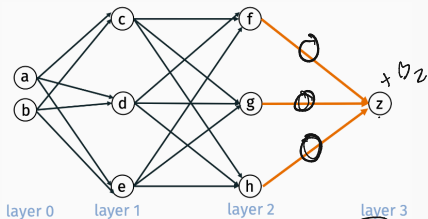
Step 2: Backward pass.



- 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.



$$\bar{z} = w_{f,z} \cdot f + \dots + \underline{b_z}$$

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

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

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

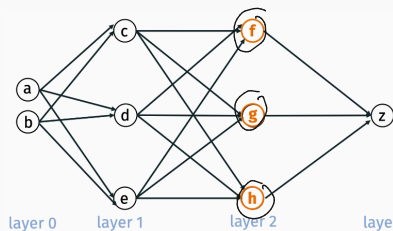
$$\frac{\partial z}{\partial w_{f,z}} = \frac{\partial \bar{z}}{\partial w_{f,z}} \cdot \frac{\partial z}{\partial \bar{z}} = f \cdot s'(\bar{z})$$

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

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

BACKPROP EXAMPLE

Step 2: Backward pass.

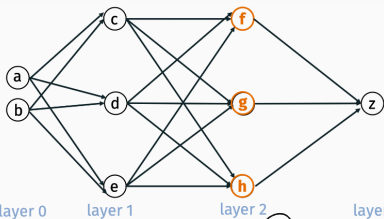


$$\frac{\partial z}{\partial f} = \frac{\partial \bar{z}}{\partial f} \cdot \frac{\partial z}{\partial \bar{z}} = W_{f,z} \cdot s'(\bar{z})$$
$$\frac{\partial z}{\partial g} = \frac{\partial \bar{z}}{\partial g} \cdot \frac{\partial z}{\partial \bar{z}} = W_{g,z} \cdot s'(\bar{z})$$
$$\frac{\partial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

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

BACKPROP EXAMPLE

Step 2: Backward pass.



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

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

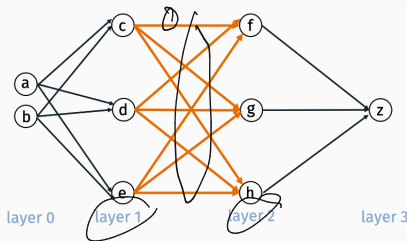
$$\begin{aligned} \frac{\partial z}{\partial \bar{f}} &= \frac{\partial z}{\partial f} \cdot \frac{\partial f}{\partial \bar{f}} = \frac{\partial z}{\partial f} \cdot s'(\bar{f}) \\ \frac{\partial z}{\partial \bar{g}} &= \frac{\partial z}{\partial g} \cdot \frac{\partial g}{\partial \bar{g}} = \frac{\partial z}{\partial g} \cdot s'(\bar{g}) \\ \frac{\partial z}{\partial \bar{h}} &= \frac{\partial z}{\partial h} \cdot \frac{\partial h}{\partial \bar{h}} = \frac{\partial z}{\partial h} \cdot s'(\bar{h}) \end{aligned}$$

And for “pre-nonlinearity” nodes.

BACKPROP EXAMPLE

Step 2: Backward pass. Next layer.

$$\bar{f}, \bar{g}, \bar{h}$$



$$\bar{f} = \underline{w_{cf}} \cdot c + \dots$$

$$f = \dots + b_f$$

$$\frac{\partial z}{\partial \beta_f} = \frac{\partial z}{\partial \bar{f}} \cdot \frac{\partial \bar{f}}{\partial \beta_f} = \left(\frac{\partial z}{\partial \bar{f}} \right) \cdot 1 =$$

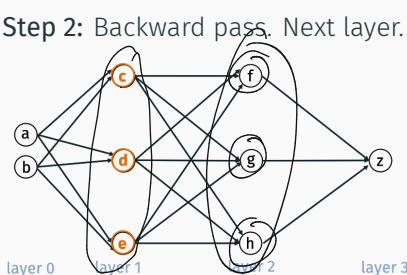
$$\frac{\partial z}{\partial W_{c,f}} = \left(\frac{\partial z}{\partial \bar{f}} \right) \left(\frac{\partial \bar{f}}{\partial W_{c,f}} \right) = \frac{\partial z}{\partial \bar{f}} \cdot c$$

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

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

BACKPROP EXAMPLE

Step 2: Backward pass. Next layer. $z(f, g, h)$



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

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

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

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

δ_2

δ_1

Linear algebraic view.

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

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

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

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

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

Linear algebraic view.

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

$$\underline{\delta}_3 = [1] \quad \underline{\delta}_2 = \begin{bmatrix} \partial z/\partial f \\ \partial z/\partial g \\ \partial z/\partial h \end{bmatrix} \quad \delta_1 = \begin{bmatrix} \partial z/\partial c \\ \partial z/\partial d \\ \partial z/\partial e \end{bmatrix}$$

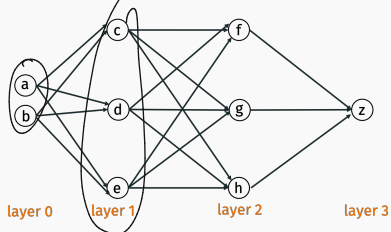
Let $\bar{\delta}_j$ be a vector containing $\partial z/\partial \bar{j}$ for all nodes j in layer i .

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

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

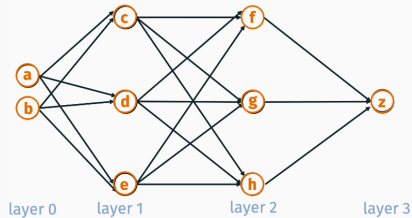
BACKPROP LINEAR ALGEBRA

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



$$W_2 = \begin{bmatrix} W_{f,z} & W_{g,z} & W_{h,z} \end{bmatrix} \quad W_1 = \begin{bmatrix} W_{c,f} & W_{d,f} & W_{e,f} \\ W_{c,g} & W_{d,g} & W_{e,g} \\ W_{c,h} & W_{d,h} & W_{e,h} \end{bmatrix} \quad W_0 = \begin{bmatrix} W_{a,c} & W_{b,c} \\ W_{a,d} & W_{b,d} \\ W_{a,e} & W_{b,e} \end{bmatrix}$$

BACKPROP LINEAR ALGEBRA

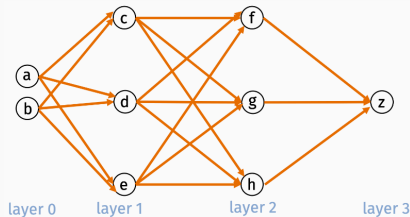


Claim 1: Node derivative computation is matrix multiplication.

$$\underline{\delta}_i = \underline{W}_i^T \underline{\delta}_{i+1}$$

BACKPROP LINEAR ALGEBRA

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



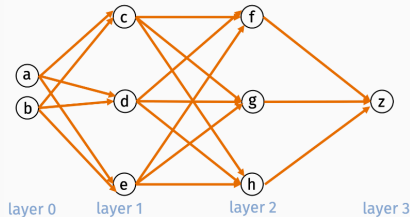
$$\Delta_2 = \begin{bmatrix} \partial z / \partial W_{f,z} & \partial z / \partial W_{g,z} & \partial z / \partial W_{h,z} \end{bmatrix}$$

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

$$\Delta_0 = \dots$$



BACKPROP LINEAR ALGEBRA



Claim 2: Weight derivative computation is an outer-product.

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

$$O(P)$$

Takeaways:

- Backpropagation can be used to compute derivatives for all weights and biases for any feedforward neural network.
- Final computation boils down to linear algebra operations (matrix multiplication and vector operations) which can be performed quickly on a GPU.

Backpropagation allows us to compute $\nabla L(y_i, f(\boldsymbol{\theta}, \mathbf{x}_i))$ for a single training example (\mathbf{x}_i, y_i) . Computing entire gradient requires computing:

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

Computing the entire sum would be very expensive.

$O((\text{time for backprop}) \cdot n)$ time.

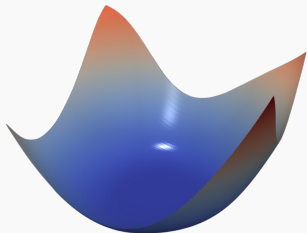
SGD iteration:

- Initialize θ_0 (typically randomly).
- For $t = 1, \dots, T$:
 - Choose j uniformly at random.
 - Compute stochastic gradient $\mathbf{g} = \nabla L_j(\theta_t)$.
 - For neural networks this is done using backprop with training example (\mathbf{x}_j, y_j) .
 - Update $\theta_{t+1} = \theta_t - \eta \mathbf{g}$

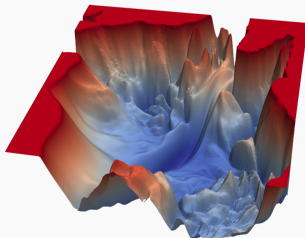
Move in direction of steepest descent in expectation.

CONVERGENCE

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



convex loss



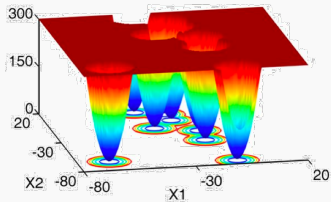
cross-entropy loss for
neural net

Neural networks very much do not...

CONVERGENCE

But SGD still performs remarkably well in practice. Understanding this phenomenon is a major open research question in machine learning and optimization.

- Initialization seems important (random uniform vs. random Gaussian vs. Xavier initialization vs. He initialization vs. etc.)
- SGD finds “good” local minima?



We already discussed a few practical modifications of SGD:

- Using “mini-batch” gradients. $\sum_{i=1}^B \nabla L_{j_i}(\boldsymbol{\theta})$.
- Shuffling then cycling through training data instead of picking a training data point at random each time.

Practical Modification: Per-parameter adaptive learning rate.

Let $\mathbf{g} = \begin{bmatrix} g_1 \\ \vdots \\ g_p \end{bmatrix}$ be a stochastic or batch stochastic gradient. Our typical parameter update looks like:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathbf{g}.$$

We've already seen a simple method for adaptively choosing the learning rate/step size η . Worked well for convex functions.

Practical Modification: Per-parameter adaptive learning rate.

In practice, neural networks can often be optimized much faster by using “adaptive gradient methods” like Adagrad, Adadelta, RMSProp, and ADAM. These methods make updates of the form:

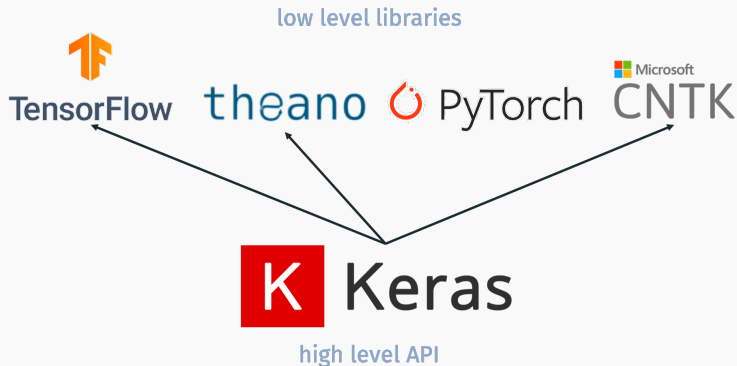
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \begin{bmatrix} \eta_1 \cdot g_1 \\ \vdots \\ \eta_p \cdot g_p \end{bmatrix}$$

So we have a separate learning rate for each entry in the gradient (e.g. parameter in the model). And each η_1, \dots, η_p is chosen adaptively.

Two demos on neural networks:

- `keras_demo_synthetic.ipynb`
- `keras_demo_mnist.ipynb`

Please spend some time working through these!



Low-level libraries have built in optimizers (SGD and improvements) and can automatically perform backpropagation for arbitrary network structures. Also optimize code for any available GPUs.

Keras has high level functions for defining and training a neural network architecture.

Define model:

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

Compile model:

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

Train model:

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

CONVOLUTIONAL NEURAL NETWORKS (CNNs)