

Backpropagation

CS 510

Lecture #12

March 30th, 2020

A Rebirth, one of several



WIKIPEDIA
The Free Encyclopedia

[Main page](#)
[Contents](#)
[Featured content](#)
[Current events](#)
[Random article](#)
[Donate to Wikipedia](#)
[Wikipedia store](#)

Interaction

[Help](#)
[About Wikipedia](#)
[Community portal](#)
[Recent changes](#)
[Contact page](#)

Not logged in [Talk](#) [Contributions](#) [Create account](#) [Log in](#)

Article [Talk](#) [Read](#) [More](#)

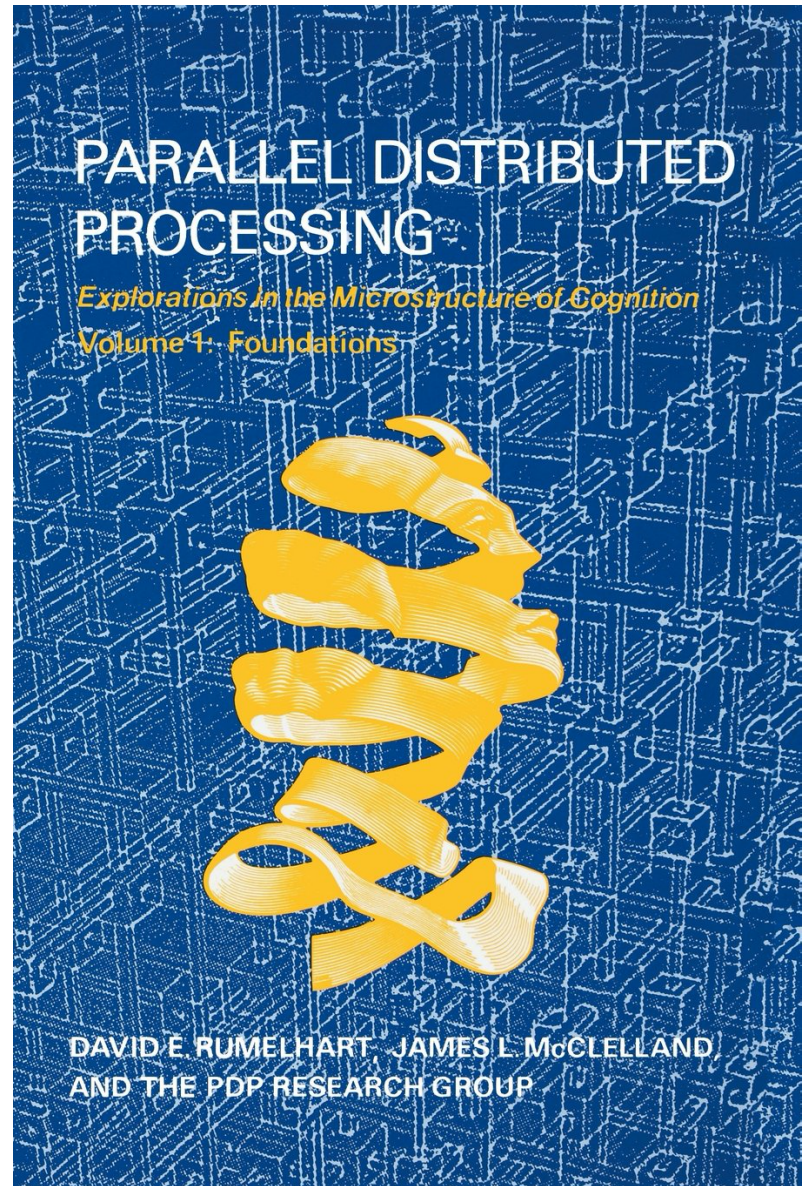
Ronald J. Williams

From Wikipedia, the free encyclopedia

Ronald J. Williams is professor of [computer science](#) at [Northeastern University](#), and one of the pioneers of [neural networks](#). He co-authored a paper on the [backpropagation](#) algorithm which triggered a boom in neural network research.^[1] He also made fundamental contributions to the fields of [recurrent neural networks](#)^{[2][3]} and [reinforcement learning](#).^[4]

References [\[edit \]](#)

- ↑ David E. Rumelhart, Geoffrey E. Hinton und Ronald J. Williams. Learning representations by back-propagating errors., Nature (London) 323, S. 533-536

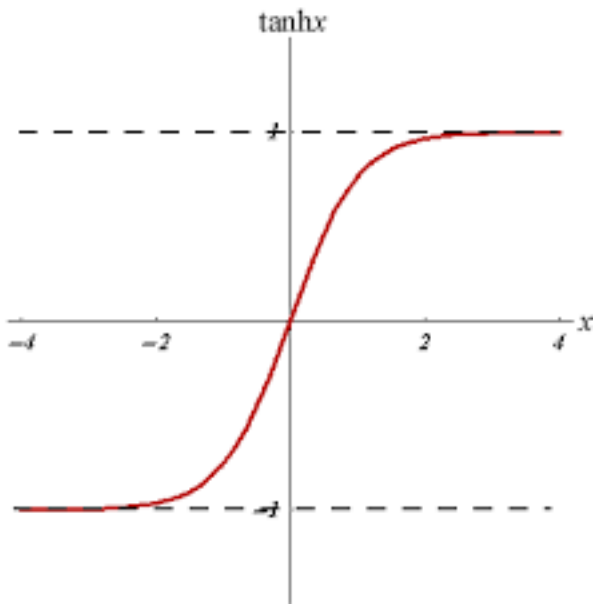


Multi-layer Perceptrons II

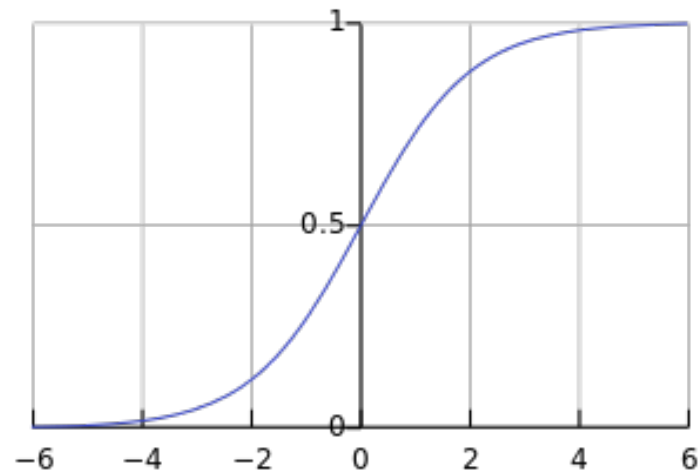
- Classic perceptrons threshold linear functions
 - $f(x) = h(w \cdot x + b)$
 - $h()$ is a threshold-based *activation function*
 - Converts activations into decisions
- But if we want to combine perceptrons?
 - Thresholding individual perceptrons is not useful
 - Replacing $h()$ with identity would allow us to sum linear responses
 - But a sum of linear responses is just another linear response

Sigmoid Activation Functions

- $f(x) = s(w \cdot x + b)$



$$y = \tanh(x)$$



$$y = (1 + e^{-x})^{-1}$$

Activation Function Properties

- Activation functions must
 - Be non-linear
- Activation functions may
 - map an infinite domain to a finite range
 - Like $[-1, 1]$ (for tanh) or $[0, 1]$ (for logistic)
 - Keeps values from growing too large/small
 - Sometimes called “squashing”
 - Have non-zero derivatives everywhere
 - Useful for training

Backpropagation

- Backpropagation is the algorithm that describes how we update weights in a network, given
 - Training samples
 - Training labels
 - A cost function
- Its used for (almost) all networks
- Network nodes may be
 - Non-linear perceptrons (the most common)
 - Convolutional units
 - Pooling units
 - Batch normalization units
 - ...

Goals For Today

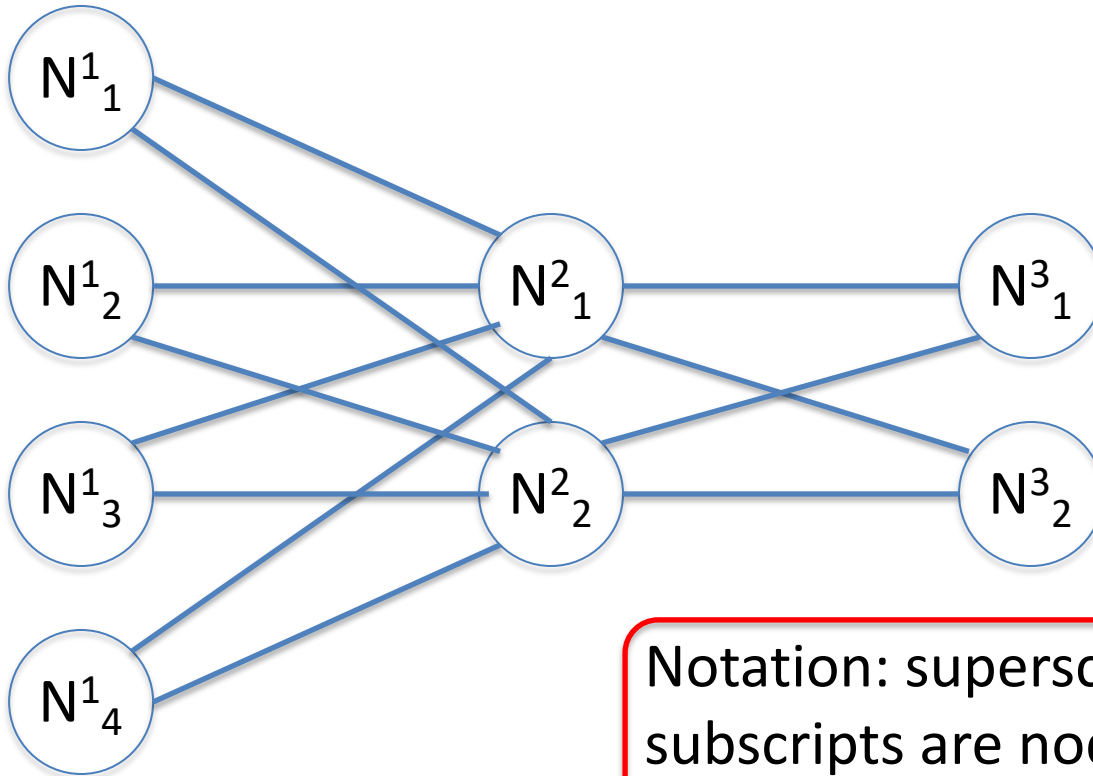
- Walk you through the math of backpropagation
 - Complicated, but just calculus
 - Almost universal : modifiable for different node types (see previous slide)
- Today's derivation assumes multi-layer perceptrons
 - $z(x) = wx + b$
 - $a(x) = h(z(x)) = h(wx + b)$
- Remember the chain rule from calculus:
 - $f(x) = g(h(x)) \rightarrow f'(x) = g'(h(x))h'(x)$

Simple Neural Network

Layer 1

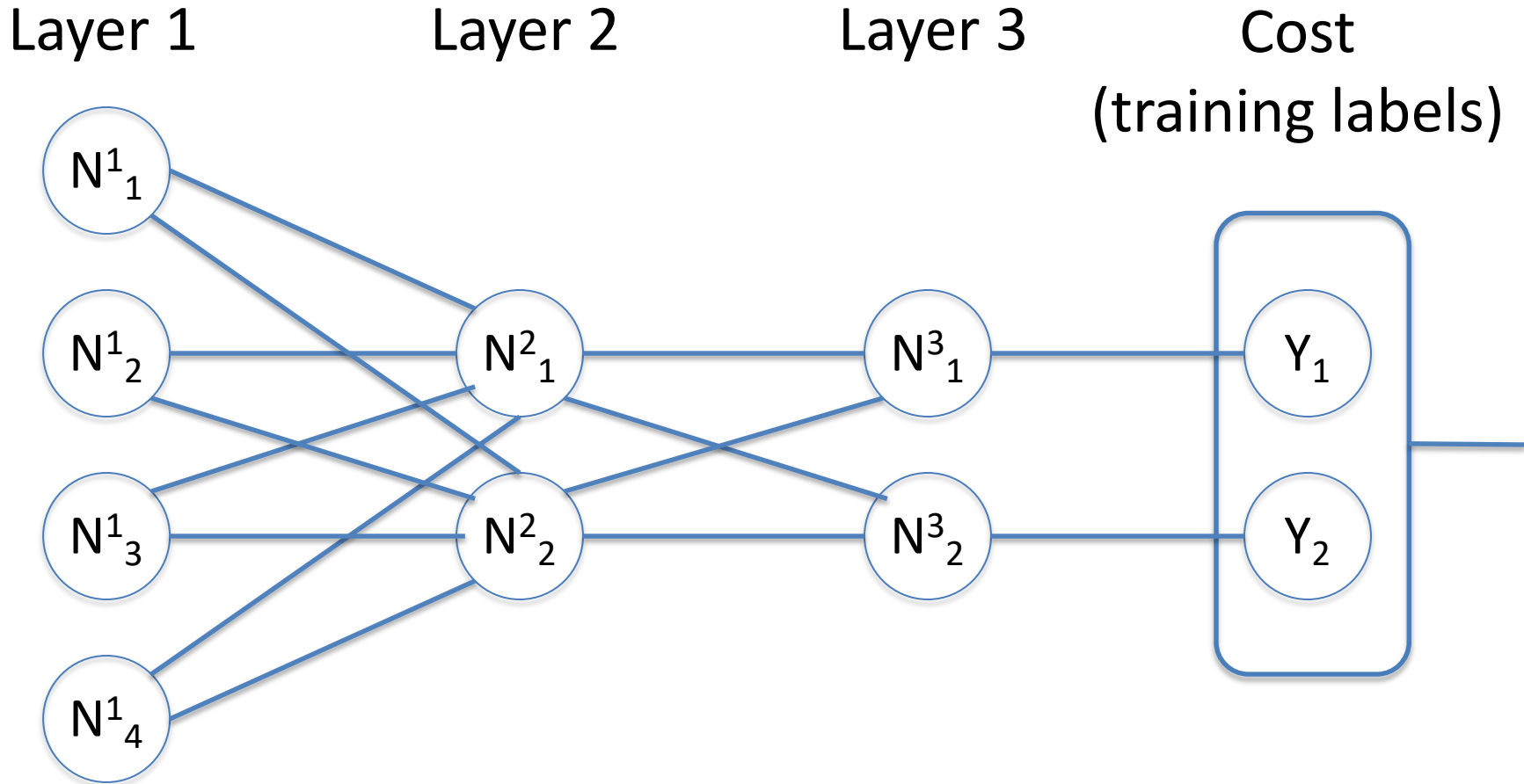
Layer 2

Layer 3



Notation: superscripts are layers,
subscripts are node numbers

Setup for Training: Cost



Cost Functions

- Cost functions measure the gap between the network output and the ideal output

- Two necessary properties

Allows us to optimize per sample

1. An average over samples: $C = \frac{1}{n} \sum_x C_x$

2. Function of output activations: $C = C(a_l)$

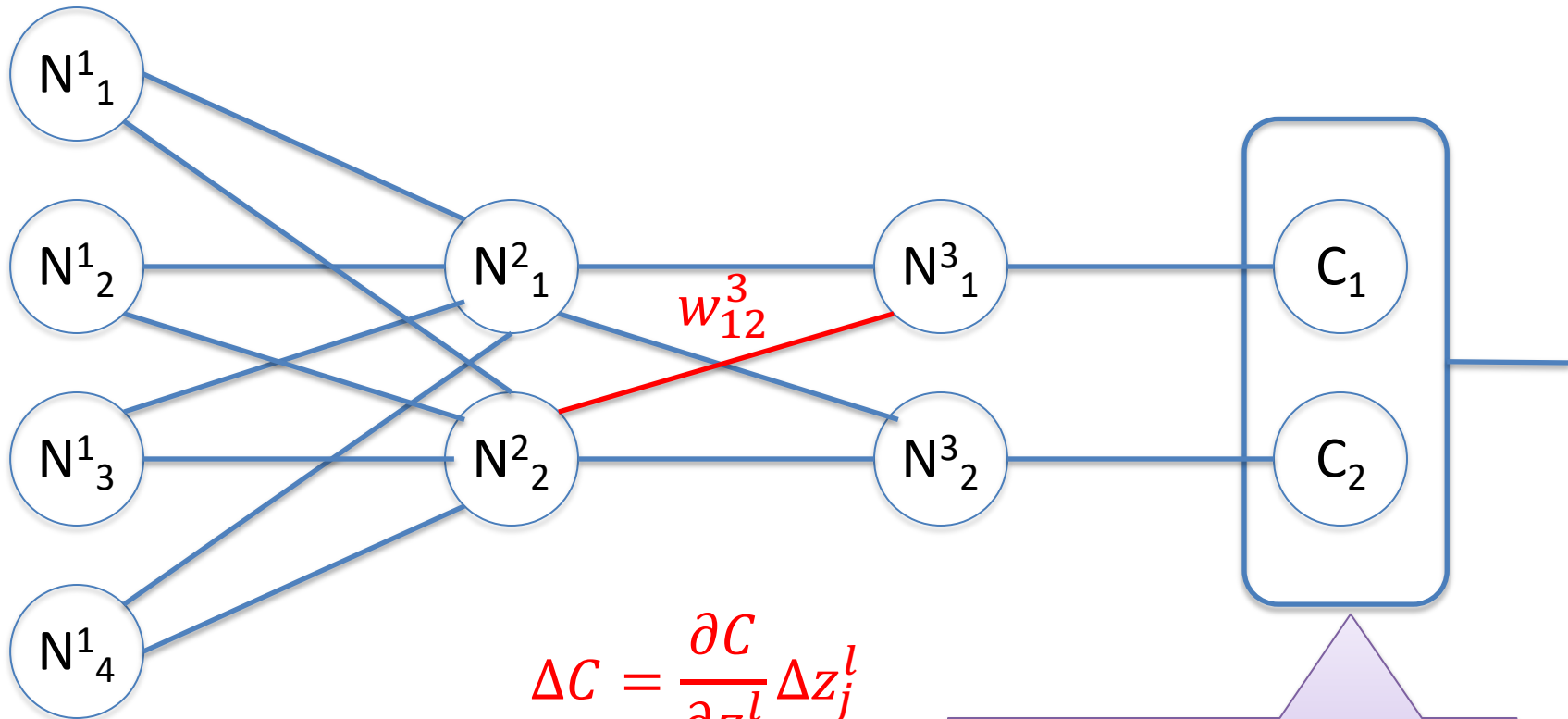
- Example: mean squared error

$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$$

Allows us to initialize the partial derivative computations

δS : local derivatives as error measures

Layer 1 Layer 2 Layer 3 Cost



$$\Delta C = \frac{\partial C}{\partial z_j^l} \Delta z_j^l$$

Subtle change from previous diagram, now showing cost C not training label.

Partial derivatives as error measures

- Imagine you want to change the output z_j^l , by Δz_j^l
- Then $\Delta C = \frac{\partial C}{\partial z_j^l} \Delta z_j^l$
- If $\left| \frac{\partial C}{\partial z_j^l} \right|$ is large, then C becomes smaller by giving Δz_j^l the opposite sign
- But if $\left| \frac{\partial C}{\partial z_j^l} \right|$ is near zero, then Δz_j^l doesn't matter.
 - $\frac{\partial C}{\partial z_j^l}$ is already optimal!
 - $\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}$

Recap - where are we?

- We can optimize on a per-sample basis
 - Because the cost function is an average
- Minimizing the δ s optimizes the net
 - The δ s depend on the data samples
- But how do we minimize the δ s?

- We will assume that nodes have non-linear functions, so $a_j^l = h(z_j^l)$

Output Layer

- $\delta_j^L = \frac{\partial C}{\partial a_j^L} h'(z_j^L)$ by the chain rule
- $\frac{\partial C}{\partial a_j^L}$ is the partial derivative of C with respect to the activation of output unit j
 - If C is LMS (slide #6)
 - $\frac{\partial C}{\partial a_j^L} = a_j^L(x) - y(x)$
 - The difference between the output & desired output

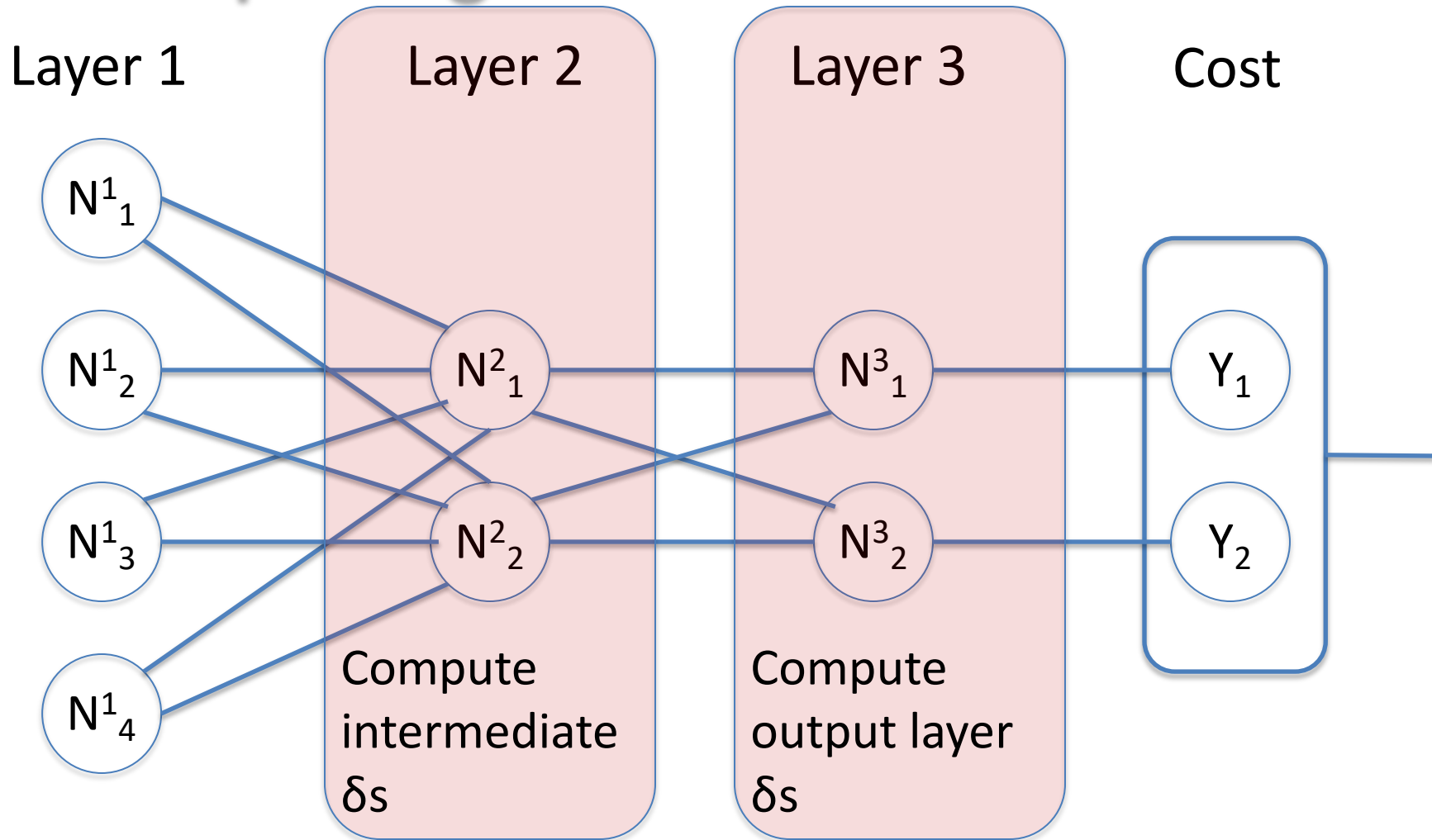
Output Layer (cont.)

- $h'(z_j^L)$ is the derivative of the non-linear transfer function at z_j^L
- If $h(x) = \tanh(x)$, $\sigma'(x) = 1 - \tanh^2(x)$
- If $h(x) = (1 + e^{-x})^{-1}$,
$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$
- $\delta_j^L = \left(a_j^L(x) - y(x) \right) \left(1 - \tanh^2 \left(z_j^L(x) \right) \right)$ or
- $\delta_j^L = \left(a_j^L(x) - y(x) \right) \left(a_j^L(x) \left(1 - a_j^L(x) \right) \right)$

δ^L given δ^{L+1}

- $\delta_j^l = \sigma'(z_j^l) \sum_k w_{kj}^{l+1} \delta_k^{l+1}$
- σ' is computed as on previous slide
- The RHS is just the sum of the impacts
- This is where *backpropagation* comes from
 - Calculate δ s for output layer
 - Then recursively compute δ s for previous layers

Computing δ s...



So...

- Given an input x and output y :
 - We can compute δ_j^l for every node j at every level l
 - Minimizing the δ s will optimize the network
 - Relative to this sample
 - So we need to adjust the weights w_i and b to reduce the δ s
 - But just a little for each input/output pair
 - So we can optimize across all samples

Adjusting b

- Remember that $\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}$ (slide #8)
- And that $z_j^l = w_j^l x + b$
- So $\frac{\partial C}{\partial b_j^l} = \delta_j^l$
- So $b_j^l \leftarrow (1 - \alpha)b_j^l - \alpha\delta_j^l$
 - Where α is a learning rate
 - Regulates how much you react to each sample

Adjusting w's

- $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$
- So $w_{jk}^l \leftarrow (1 - \alpha)w_{jk}^l - \alpha a_k^{l-1} \delta_j^l$
 - Where α is the same learning rate as before
 - We are collectively minimizing the deltas by heading downhill in the $k+1$ dimensional space defined by w & b

Backpropagation (redux)

- Backpropagation updates weights in a network, given
 - Training samples
 - Training labels
 - A cost function
- Network nodes may be
 - Non-linear perceptrons (the most common)
 - Convolutional units
 - Pooling units
 - Batch normalization units
 - ...

Step Back! Other Resources

- Modulo some notation ambiguity the previous formula-based presentation is fine, but for some of us unsatisfying
- It is best to approach the task of understanding backpropagation simultaneously from three angles.
 1. Mathematical formulas (just finished)
 2. Develop an internal visualization
 3. Running code

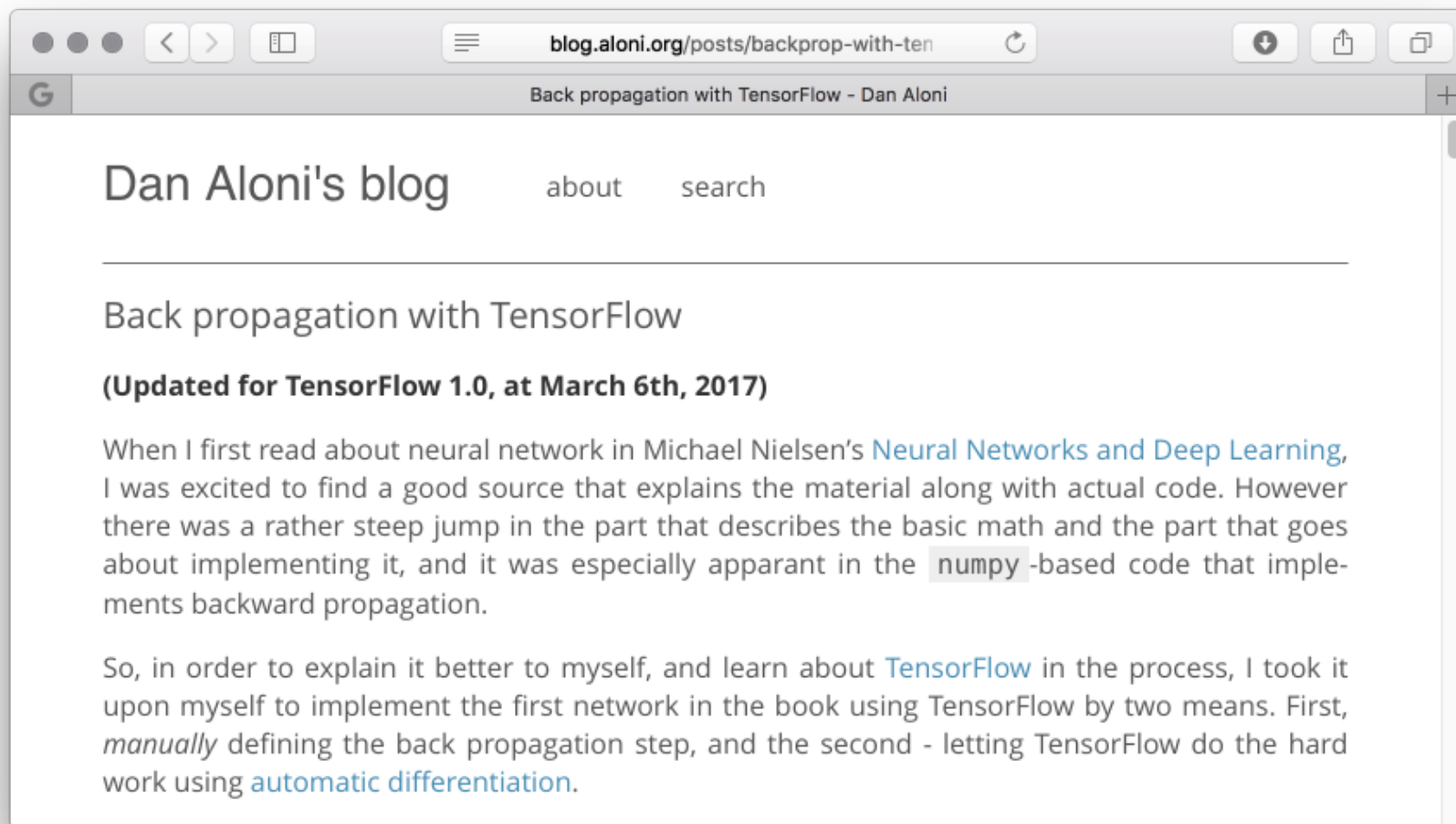
May I Recommend

The screenshot shows a YouTube video player interface. The browser address bar displays `www.youtube.com/watch?v=llg3gGewQ5U`. The YouTube logo and search bar are visible at the top. The video content features a diagram of a neural network on a black background. On the left, a white digit '3' is enclosed in a blue square. To its right, a vertical bracket labeled '784' indicates the input layer, which consists of 784 small white circles. These are connected to three hidden layers of circles. The first hidden layer has 4 white circles, the second has 4 grey circles, and the third has 4 white circles. The output layer on the right has 10 circles labeled 0 through 9, with circle 3 being white and the others grey. A small circular logo is in the bottom right corner of the video frame.

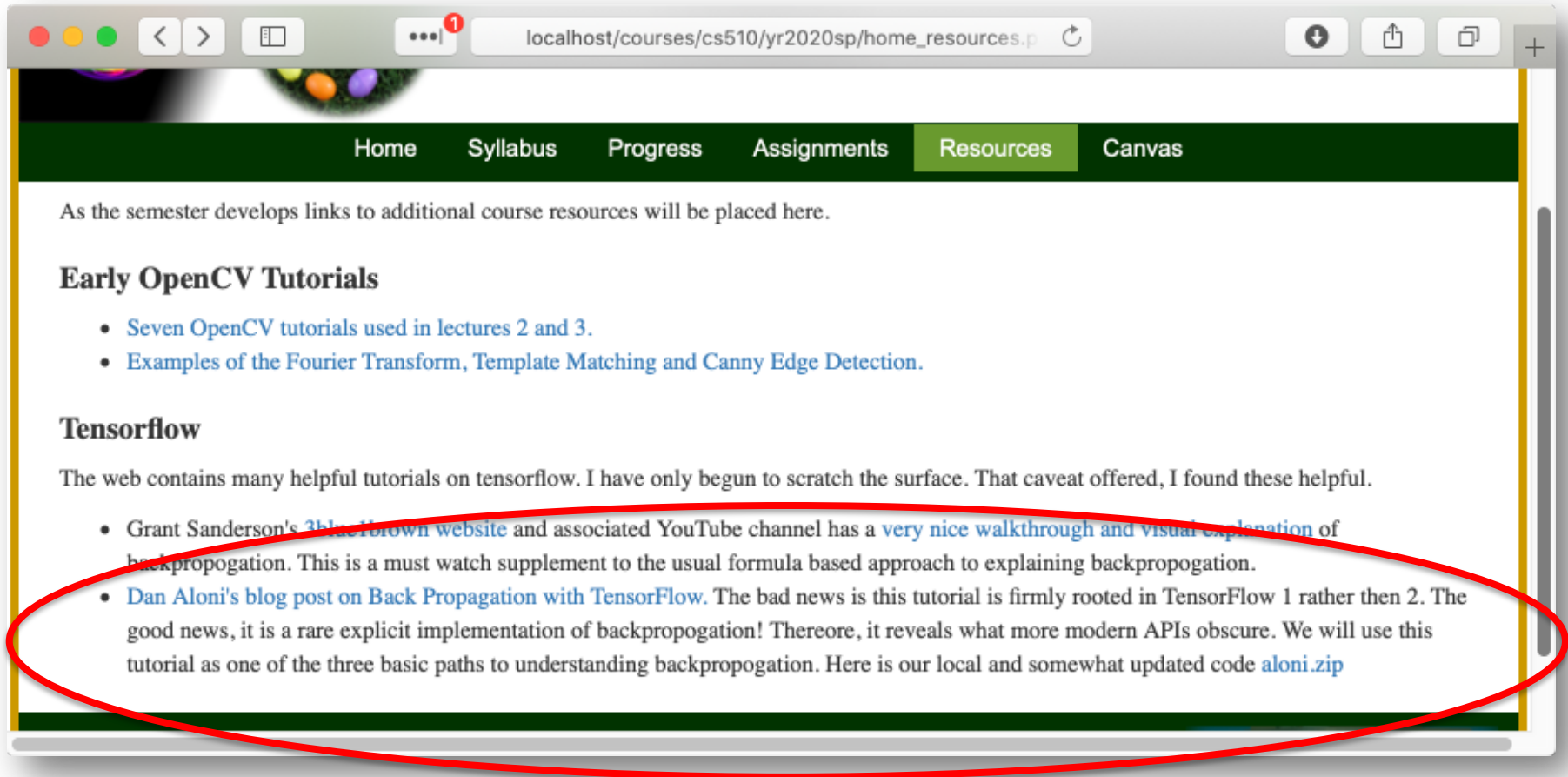
3BLUE1BROWN SERIES S3 • E3
What is backpropagation really doing? | Deep learning, chapter 3
1,698,016 views • Nov 3, 2017

👍 30K 💬 264 ➦ SHARE ⌵ SAVE ⋮

And Also



Please Run and Play With ...



localhost/courses/cs510/yr2020sp/home_resources.p

Home Syllabus Progress Assignments Resources Canvas

As the semester develops links to additional course resources will be placed here.

Early OpenCV Tutorials

- [Seven OpenCV tutorials used in lectures 2 and 3.](#)
- [Examples of the Fourier Transform, Template Matching and Canny Edge Detection.](#)

Tensorflow

The web contains many helpful tutorials on tensorflow. I have only begun to scratch the surface. That caveat offered, I found these helpful.

- [Grant Sanderson's 3blue1brown website](#) and associated YouTube channel has a [very nice walkthrough and visual explanation of backpropagation](#). This is a must watch supplement to the usual formula based approach to explaining backpropagation.
- [Dan Aloni's blog post on Back Propagation with TensorFlow](#). The bad news is this tutorial is firmly rooted in TensorFlow 1 rather than 2. The good news, it is a rare explicit implementation of backpropagation! Therefore, it reveals what more modern APIs obscure. We will use this tutorial as one of the three basic paths to understanding backpropagation. Here is our local and somewhat updated code [aloni.zip](#)

As we will discuss in lecture today, I expect everyone to setup a TF 1.14 environment and play with this code.