

# CPSC 540: Machine Learning

## Neural Networks

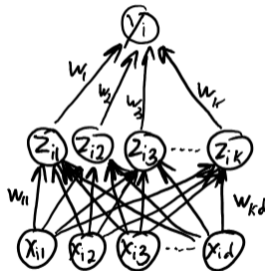
Mark Schmidt

University of British Columbia

Winter 2020

## Feedforward Neural Networks

- In 340 we discussed **feedforward neural networks** for supervised learning.
- With 1 hidden layer the classic model has this structure:



- Motivation:
  - For some problems it's **hard to find good features**.
  - This **learns features  $z$**  that are good for particular supervised learning problem.

## Neural Network Notation

- We'll continue using our supervised learning notation:

$$X = \begin{bmatrix} \text{---} & (x^1)^T & \text{---} \\ \text{---} & (x^2)^T & \text{---} \\ & \vdots & \\ \text{---} & (x^n)^T & \text{---} \end{bmatrix}, \quad y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix},$$

- For the **latent features** and **one hidden layer** we'll use

$$Z = \begin{bmatrix} \text{---} & (z^1)^T & \text{---} \\ \text{---} & (z^2)^T & \text{---} \\ & \vdots & \\ \text{---} & (z^n)^T & \text{---} \end{bmatrix}, \quad v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_k \end{bmatrix}, \quad W = \begin{bmatrix} \text{---} & w_1 & \text{---} \\ \text{---} & w_2 & \text{---} \\ & \vdots & \\ \text{---} & w_k & \text{---} \end{bmatrix},$$

where  $Z$  is  $n$  by  $k$  and  $W$  is  $k$  by  $d$ .

## Introducing Non-Linearity

- We discussed how the “linear-linear” model,

$$z^i = Wx^i, \quad \hat{y}^i = v^T z^i,$$

is **degenerate** since it's still a linear model.

- The classic solution is to introduce a **non-linearity**,

$$z^i = h(Wx^i), \quad \hat{y}^i = v^T z^i,$$

where a common-choice is applying **sigmoid** element-wise,

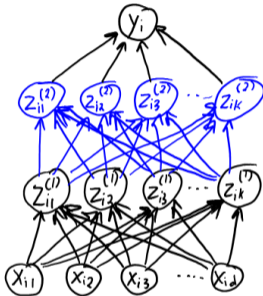
$$z_c^i = \frac{1}{1 + \exp(-w_c^T x^i)},$$

which is said to be the “activation” of neuron  $c$  on example  $i$ .

- A **universal approximator** with  $k$  a function of  $n$  (also true for tanh, ReLU, etc.)

## Deep Neural Networks

- In deep neural networks we add multiple hidden layers,

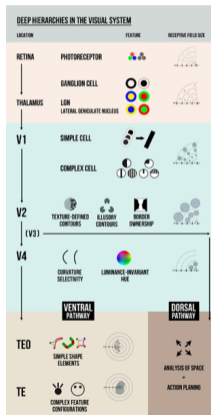


- Mathematically, with 3 hidden layers the classic model uses

$$\hat{y}^i = v^T h(W^3 \underbrace{h(W^2 \underbrace{h(W^1 x^i)}_{z^{i1}})}_{z^{i2}}))_{z^{i3}}).$$

## Biological Motivation

- Deep learning is motivated by theories of deep hierarchies in the brain.



[https://en.wikibooks.org/wiki/Sensory\\_Systems/Visual\\_Signal\\_Processing](https://en.wikibooks.org/wiki/Sensory_Systems/Visual_Signal_Processing)

- But most research is about making models work better, not be more brain-like.

## Deep Neural Network History

- Popularity of deep learning has come in waves over the years.
  - Currently, it is one of the **hottest topics in science**.
- Recent popularity is due to **unprecedented performance** on some difficult tasks:
  - Speech recognition.
  - Computer vision.
  - Machine translation.
- These are mainly due to **big datasets**, **deep models**, and **tons of computation**.
  - Plus tweaks to classic models and focus on structured networks (CNNs, LSTMs).
- For a NY Times article discussing some of the history/successes/issues, see:

<https://mobile.nytimes.com/2016/12/14/magazine/the-great-ai-awakening.html>

## Training Deep Neural Networks

- If we're training a network with 3 hidden layers and squared error, our objective is

$$f(v, W^1, W^2, W^3) = \frac{1}{2} \sum_{i=1}^n \underbrace{(v^T h(W^3 h(W^2 h(W^1 x^i))))}_{\hat{y}^i} - y^i)^2.$$

- Usual training procedure is **stochastic gradient**.
  - **Highly non-convex and notoriously difficult to tune.**
  - But we're discovering sets of **tricks to make things easier** to tune.
- Recent empirical/theoretical work indicates non-convexity may not be an issue:
  - **All local minima may be good** for “large enough” networks.



# Training Deep Neural Networks

- Some common data/optimization tricks we discussed in 340:
  - **Data transformations.**
    - For images, translate/rotate/scale/crop each  $x^i$  to make more data.
  - **Data standardization:** centering and whitening.
  - Adding **bias variables.**
  - **Parameter initialization:** “small but different”, standardizing within layers.
  - **Step-size selection:** “babysitting”, Bottou trick, Adam.
  - **Momentum:** heavy-ball and Nesterov-style modifications.
  - **Batch normalization:** adaptive standardizing within layers.
  - **ReLU:** replacing sigmoid with  $\max\{0, w_c^T x^i\}$ .
    - Avoids gradients extremely-close to zero.

# Training Deep Neural Networks

- Common forms tricks to fight overfitting:
  - Standard **L2-regularization** or **L1-regularization** “weight decay” .
    - Sometimes with different  $\lambda$  for each layer.
    - Recent work shows this **introduces bad local optima**.
  - **Early stopping** of the optimization based on validation accuracy.
  - **Dropout** randomly zeroes  $z$  values to discourage dependence.
  - **Implicit regularization** from using SGD.
  - **Hyper-parameter optimization** to choose various tuning parameters.
    - “Neural architecture search”: recent methods include search over graph structures.
  - **Special architectures** like **convolutional neural networks**:
    - Yields  $W^m$  that are **very sparse** and have many **tied parameters**.