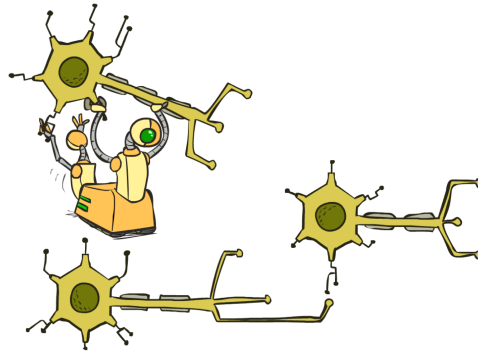


CS 188: Artificial Intelligence

Optimization and Neural Nets

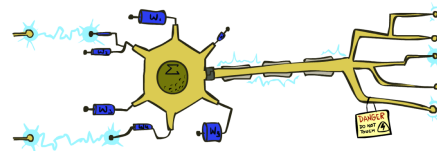


Instructors: Pieter Abbeel and Dan Klein --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

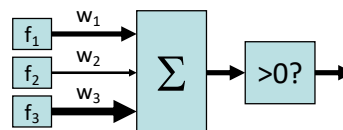
Reminder: Linear Classifiers

- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- If the activation is:
 - Positive, output +1
 - Negative, output -1

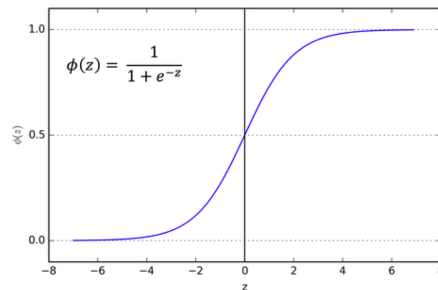


How to get probabilistic decisions?

- Activation: $z = w \cdot f(x)$
- If $z = w \cdot f(x)$ very positive \rightarrow want probability going to 1
- If $z = w \cdot f(x)$ very negative \rightarrow want probability going to 0

- Sigmoid function

$$\phi(z) = \frac{1}{1 + e^{-z}}$$



Best w?

- Maximum likelihood estimation:

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

with: $P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$

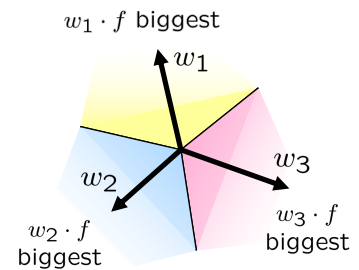
$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

= **Logistic Regression**

Multiclass Logistic Regression

- Multi-class linear classification

- A weight vector for each class: w_y
- Score (activation) of a class y : $w_y \cdot f(x)$
- Prediction w/highest score wins: $y = \arg \max_y w_y \cdot f(x)$



- How to make the scores into probabilities?

$$\underbrace{z_1, z_2, z_3}_{\text{original activations}} \rightarrow \underbrace{\frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}}_{\text{softmax activations}}$$

Best w ?

- Maximum likelihood estimation:

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

with:
$$P(y^{(i)} | x^{(i)}; w) = \frac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_y e^{w_y \cdot f(x^{(i)})}}$$

= Multi-Class Logistic Regression

This Lecture

- Optimization

- i.e., how do we solve:

$$\max_w \ell(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

Hill Climbing

- Recall from CSPs lecture: simple, general idea

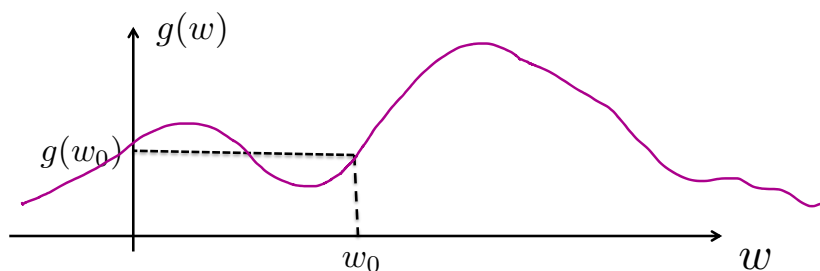
- Start wherever
 - Repeat: move to the best neighboring state
 - If no neighbors better than current, quit



- What's particularly tricky when hill-climbing for multiclass logistic regression?

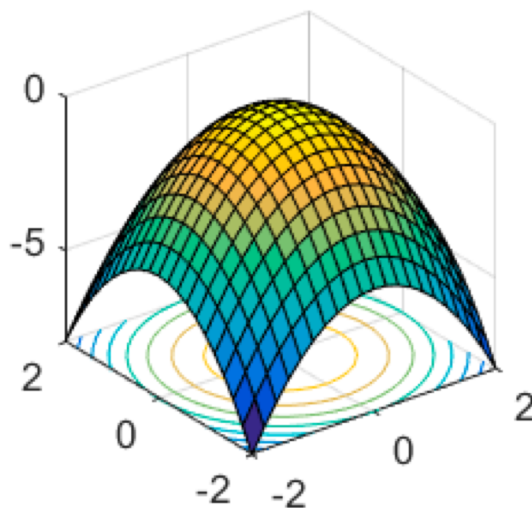
- Optimization over a continuous space
 - Infinitely many neighbors!
 - How to do this efficiently?

1-D Optimization



- Could evaluate $g(w_0 + h)$ and $g(w_0 - h)$
 - Then step in best direction
- Or, evaluate derivative:
$$\frac{\partial g(w_0)}{\partial w} = \lim_{h \rightarrow 0} \frac{g(w_0 + h) - g(w_0 - h)}{2h}$$
 - Tells which direction to step into

2-D Optimization



Gradient Ascent

- Perform update in uphill direction for each coordinate
- The steeper the slope (i.e. the higher the derivative) the bigger the step for that coordinate
- E.g., consider: $g(w_1, w_2)$

- Updates:

$$w_1 \leftarrow w_1 + \alpha * \frac{\partial g}{\partial w_1}(w_1, w_2)$$

$$w_2 \leftarrow w_2 + \alpha * \frac{\partial g}{\partial w_2}(w_1, w_2)$$

- Updates in vector notation:

$$w \leftarrow w + \alpha * \nabla_w g(w)$$

$$\text{with: } \nabla_w g(w) = \begin{bmatrix} \frac{\partial g}{\partial w_1}(w) \\ \frac{\partial g}{\partial w_2}(w) \end{bmatrix} = \text{gradient}$$

Gradient Ascent

- Idea:
 - Start somewhere
 - Repeat: Take a step in the gradient direction

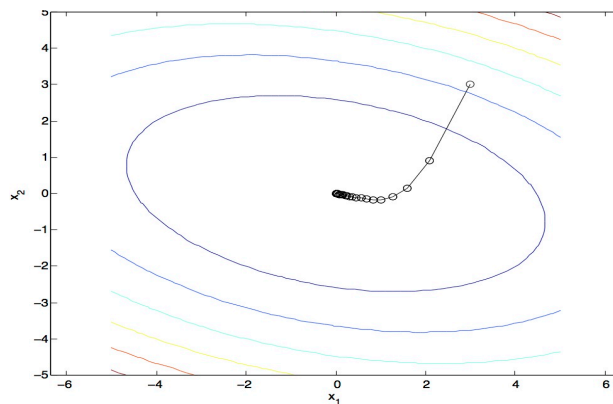


Figure source: Mathworks

What is the Steepest Direction?

$$\max_{\Delta: \Delta_1^2 + \Delta_2^2 \leq \varepsilon} g(w + \Delta)$$



▪ First-Order Taylor Expansion: $g(w + \Delta) \approx g(w) + \frac{\partial g}{\partial w_1} \Delta_1 + \frac{\partial g}{\partial w_2} \Delta_2$

▪ Steepest Descent Direction: $\max_{\Delta: \Delta_1^2 + \Delta_2^2 \leq \varepsilon} g(w) + \frac{\partial g}{\partial w_1} \Delta_1 + \frac{\partial g}{\partial w_2} \Delta_2$

▪ Recall: $\max_{\Delta: \|\Delta\| \leq \varepsilon} \Delta^\top a \rightarrow \Delta = \varepsilon \frac{a}{\|a\|}$

▪ Hence, solution: $\Delta = \varepsilon \frac{\nabla g}{\|\nabla g\|}$ **Gradient direction = steepest direction!**

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \end{bmatrix}$$

Gradient in n dimensions

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \dots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

Optimization Procedure: Gradient Ascent

- `init` w
- `for` $\text{iter} = 1, 2, \dots$
 $w \leftarrow w + \alpha * \nabla g(w)$

- α : learning rate --- tweaking parameter that needs to be chosen carefully
- How? Try multiple choices
 - Crude rule of thumb: update changes w about 0.1 – 1 %

Batch Gradient Ascent on the Log Likelihood Objective

$$\max_w ll(w) = \max_w \underbrace{\sum_i \log P(y^{(i)} | x^{(i)}; w)}_{g(w)}$$

- `init` w
- `for` $\text{iter} = 1, 2, \dots$
 $w \leftarrow w + \alpha * \sum_i \nabla \log P(y^{(i)} | x^{(i)}; w)$

Stochastic Gradient Ascent on the Log Likelihood Objective

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

Observation: once gradient on one training example has been computed, might as well incorporate before computing next one

```
▪ init  $w$ 
▪ for iter = 1, 2, ...
  ▪ pick random  $j$ 
     $w \leftarrow w + \alpha * \nabla \log P(y^{(j)}|x^{(j)}; w)$ 
```

Mini-Batch Gradient Ascent on the Log Likelihood Objective

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

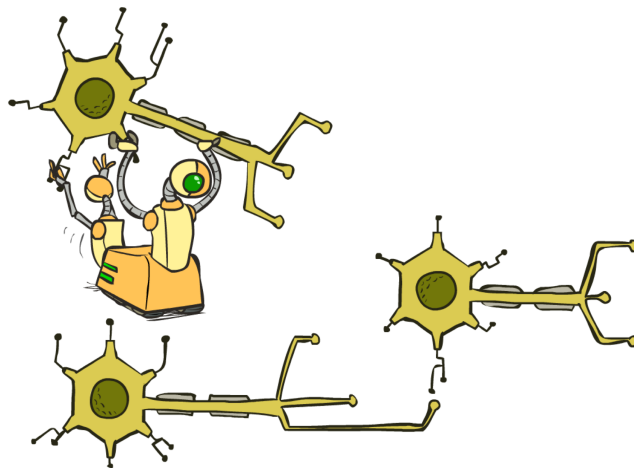
Observation: gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

```
▪ init  $w$ 
▪ for iter = 1, 2, ...
  ▪ pick random subset of training examples  $J$ 
     $w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)}|x^{(j)}; w)$ 
```

How about computing all the derivatives?

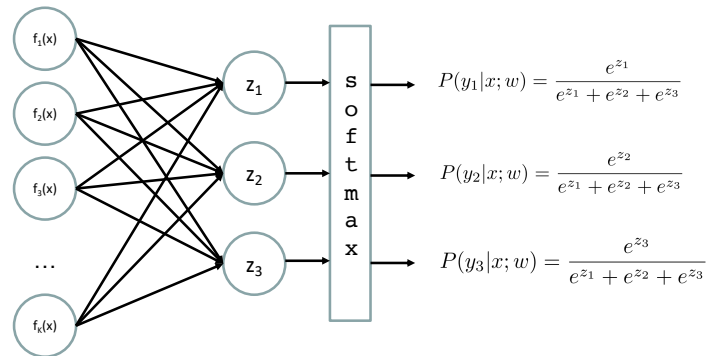
- We'll talk about that once we covered neural networks, which are a generalization of logistic regression

Neural Networks

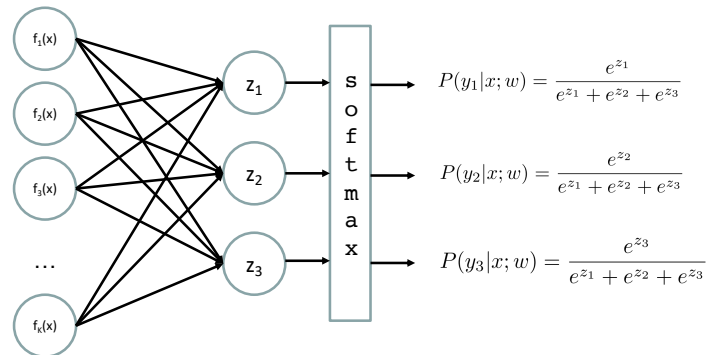


Multi-class Logistic Regression

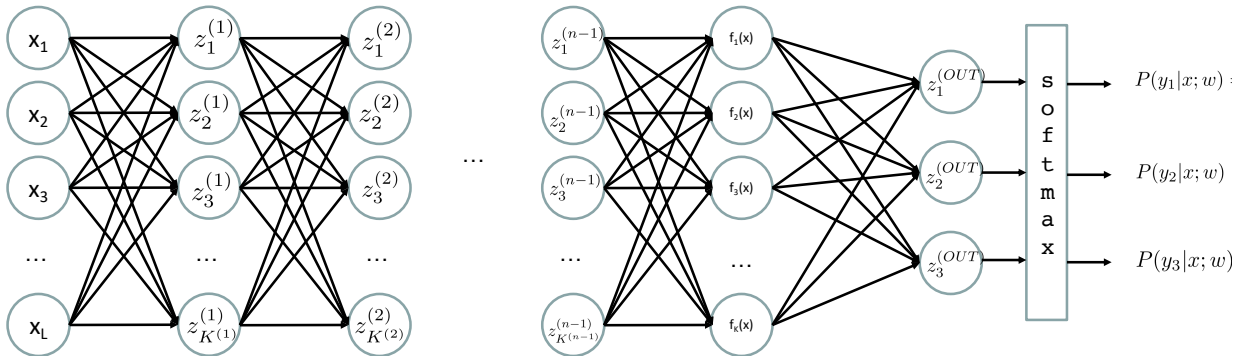
- = special case of neural network



Deep Neural Network = Also learn the features!



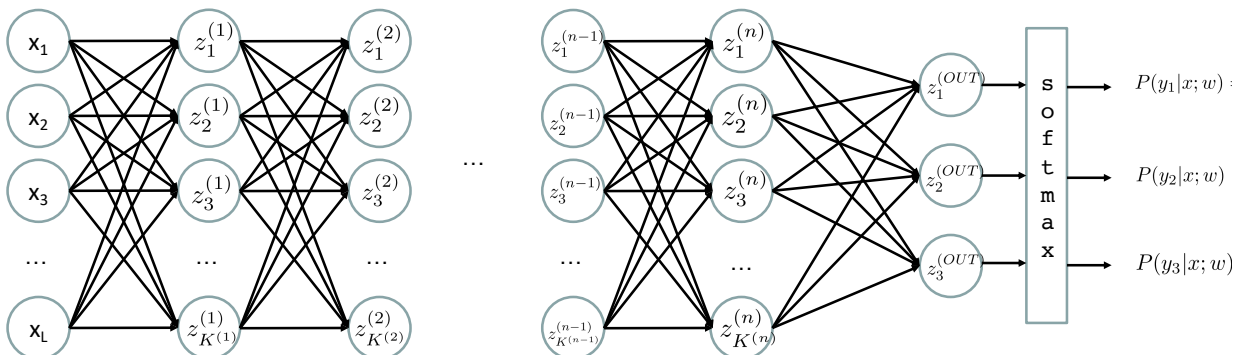
Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

g = nonlinear activation function

Deep Neural Network = Also learn the features!

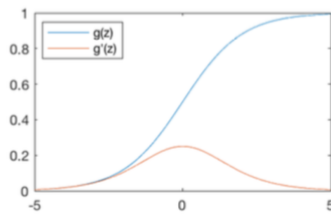


$$z_i^{(k)} = g\left(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)}\right)$$

g = nonlinear activation function

Common Activation Functions

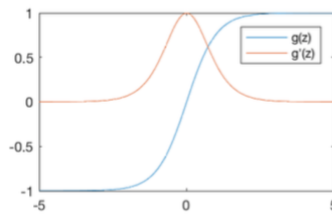
Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

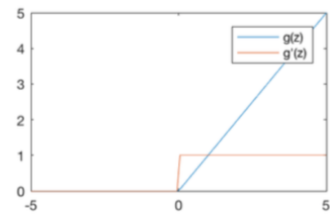
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

[source: MIT 6.S191 introtodeeplearning.com]

Deep Neural Network: Also Learn the Features!

- Training the deep neural network is just like logistic regression:

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

just w tends to be a much, much larger vector ☺

→ just run gradient ascent

+ stop when log likelihood of hold-out data starts to decrease

Neural Networks Properties

- Theorem (Universal Function Approximators). A two-layer neural network with a sufficient number of neurons can approximate any continuous function to any desired accuracy.
- Practical considerations
 - Can be seen as learning the features
 - Large number of neurons
 - Danger for overfitting
 - (hence early stopping!)

Universal Function Approximation Theorem*

Hornik theorem 1: Whenever the activation function is *bounded and nonconstant*, then, for any finite measure μ , standard multilayer feedforward networks can approximate any function in $L^p(\mu)$ (the space of all functions on R^k such that $\int_{R^k} |f(x)|^p d\mu(x) < \infty$) arbitrarily well, provided that sufficiently many hidden units are available.

Hornik theorem 2: Whenever the activation function is *continuous, bounded and non-constant*, then, for arbitrary compact subsets $X \subseteq R^k$, standard multilayer feedforward networks can approximate any continuous function on X arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

- In words: Given any continuous function $f(x)$, if a 2-layer neural network has enough hidden units, then there is a choice of weights that allow it to closely approximate $f(x)$.

Universal Function Approximation Theorem*

<p>Math. Control Signals Systems (1989) 2: 303-314</p> <p>Mathematics of Control, Signals, and Systems © 1989 Springer-Verlag New York Inc.</p> <p>Approximation by Superpositions of a Sigmoidal Function*</p> <p>G. Cybenko[†]</p> <p>Abstract. In this paper we demonstrate that finite linear combinations of continuous functions of a real variable with support in the unit ball, and only such combinations are required to approximate functions. Our main result is an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary discrete signals can be arbitrarily well approximated by continuous feedforward neural networks with only a single hidden layer and any continuous sigmoidal activation function. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.</p> <p>Key words. Neural networks, Approximation, Completeness.</p> <p>1. Introduction</p> <p>A number of diverse application areas are concerned with the representation of general functions of an n-dimensional real variable, $x \in \mathbb{R}^n$, by finite linear combinations of the form</p> $\sum_{j=1}^N a_j \sigma(x_j) + b_j, \quad (1)$ <p>where $x_j \in \mathbb{R}^n$ and $a_j, b_j \in \mathbb{R}$ are fixed. (x^T is the transpose of x so that $x^T x$ is the inner product of x and x.) Here the univariate function σ depends heavily on the context of the application. Our major concern is with so-called sigmoidal σ's:</p> $\sigma(t) = \begin{cases} 1 & \text{as } t \rightarrow +\infty, \\ 0 & \text{as } t \rightarrow -\infty. \end{cases}$ <p>Such functions arise naturally in neural network theory as the activation function of a neural node (or unit as is becoming the preferred term) [L1], [RHIM]. The main result of this paper is a demonstration of the fact that sums of the form (1) are dense in the space of continuous functions on the unit cube if σ is any continuous sigmoidal</p> <p>* Date received: October 21, 1988. Date revised: February 17, 1989. This research was supported in part by NSF Grant DCR-8819103, ONR Contract N00016-86-G-0202 and DOE Grant DE-FG02-88FE22862.</p> <p>[†] Center for Supercomputing Research and Development and Department of Electrical and Computer Engineering, University of Illinois, Urbana, Illinois 61801, U.S.A.</p> <p>303</p>	<p>Neural Networks, Vol. 2, pp. 251-257, 1989. Printed in the U.S.A. All rights reserved.</p> <p>Copyright © 1989 Pergamon Press, Inc.</p> <p>ORIGINAL CONTRIBUTION</p> <p>Approximation Capabilities of Multilayer Feedforward Networks</p> <p>KURT HORNIK</p> <p>Technische Universität Wien, Vienna, Austria</p> <p>(Received 16 January 1989; revised and accepted 16 October 1989)</p> <p>Abstract. We show that multilayer feedforward networks with as few as a single hidden layer and arbitrary bounded and monotonic activation functions are universal approximators with respect to L_2 performance criteria, for arbitrary finite input-output measures provided only that sufficiently many hidden units are available. If the activation function is continuous, bounded and monotonic, then continuous mappings can be learned uniformly over compact input sets. We also give very general conditions ensuring that networks with sufficiently smooth activation functions are capable of arbitrarily accurate approximation to functions and to derivatives.</p> <p>Keywords: Multilayer feedforward networks, Activation functions, Universal approximation capabilities, Input-output measures, L_2 approximation, Uniform approximation, Sobolev spaces, Smooth approximation.</p> <p>1. INTRODUCTION</p> <p>The approximation capabilities of neural network architectures have recently been investigated by many authors, including Carrell and Dickinson (1989), Cybenko (1989), Endershaw (1989), Galland and White (1989), Hornik (1991), Hornik, Stinchcombe, and White (1989, 1990), Lee and Mendel (1990), Lapedis and Fister (1988), Stinchcombe and White (1989, 1990). (This list is by no means complete.)</p> <p>If we think of the network architecture as a rule for computing values at l output units given values at n input units, hence implementing a value of mapping from \mathbb{R}^n to \mathbb{R}^l, we can ask how well arbitrary mappings from \mathbb{R}^n to \mathbb{R}^l can be approximated by the network, in particular, if so many hidden units as required for internal representation and computation may be employed.</p> <p>How to measure the accuracy of approximation depends on how we measure closeness between functions, which in turn varies significantly with the specific problem to be dealt with. In many applications, it is necessary to have the network perform reasonably well on all input samples taken from some compact input set X in \mathbb{R}^n. In this case, closeness is measured by the uniform distance between functions on X, that is,</p> $\ f, g\ _\infty = \sup_{x \in X} f(x) - g(x) .$ <p>In other applications, we think of the inputs as random variables and are interested in the average performance where the average is taken with respect to the input-output measure μ, where $\mu(\mathbb{R}^n) = 1$. In this case, closeness is measured by the L_2 distance</p> $\ f, g\ _2 = \left[\int_X f(x) - g(x) ^2 d\mu(x) \right]^{1/2}.$ <p>$1 \leq p < \infty$ is the most popular choice being $p = 2$, corresponding to mean square error.</p> <p>Of course, there are many more ways of measuring closeness of functions. In particular, in many applications, it is also necessary that the derivatives of the approximating function implemented by the network closely resemble those of the function to be approximated, up to some order. This issue was first taken up by Hornik et al. (1989), who discuss the sources of need of smooth functional approximation in more detail. Typical examples arise in robotics (learning of smooth movements) and signal processing (analysis of chaotic time series). For a recent application in problems of nonparametric inference in statistics and econometrics, see Galland and White (1989).</p> <p>All papers establishing certain approximation ca-</p> <p>251</p>	<p>MULTILAYER FEEDFORWARD NETWORKS WITH NON-POLYNOMIAL ACTIVATION FUNCTIONS CAN APPROXIMATE ANY FUNCTION</p> <p>by</p> <p>Moshe Leshno Faculty of Management Tel Aviv University Tel Aviv, Israel 69978</p> <p>and</p> <p>Shimon Schocken Leonard N. Stern School of Business New York University New York, NY 10003</p> <p>September 1991</p> <p>Center for Research on Information Systems Information Systems Department Leonard N. Stern School of Business New York University Working Paper Series STERN IS-91-26</p> <p>Appeared previously as Working Paper No. 21/91 at The Israel Institute Of Business Research</p>
---	---	--

Cybenko (1989) "Approximations by superpositions of sigmoidal functions"
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"
Leshno and Schocken (1991) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

Fun Neural Net Demo Site

- Demo-site:
 - <http://playground.tensorflow.org/>

How about computing all the derivatives?

Derivatives tables:

$\frac{d}{dx}(a) = 0$	$\frac{d}{dx}[\ln u] = \frac{d}{dx}[\log_e u] = \frac{1}{u} \frac{du}{dx}$
$\frac{d}{dx}(x) = 1$	$\frac{d}{dx}[\log_a u] = \log_a e \frac{1}{u} \frac{du}{dx}$
$\frac{d}{dx}(au) = a \frac{du}{dx}$	$\frac{d}{dx}e^u = e^u \frac{du}{dx}$
$\frac{d}{dx}(u + v - w) = \frac{du}{dx} + \frac{dv}{dx} - \frac{dw}{dx}$	$\frac{d}{dx}a^u = a^u \ln a \frac{du}{dx}$
$\frac{d}{dx}(uv) = u \frac{dv}{dx} + v \frac{du}{dx}$	$\frac{d}{dx}(u^v) = v u^{v-1} \frac{du}{dx} + \ln u \cdot u^v \frac{dv}{dx}$
$\frac{d}{dx}\left(\frac{u}{v}\right) = \frac{1}{v} \frac{du}{dx} - \frac{u}{v^2} \frac{dv}{dx}$	$\frac{d}{dx} \sin u = \cos u \frac{du}{dx}$
$\frac{d}{dx}(u^n) = n u^{n-1} \frac{du}{dx}$	$\frac{d}{dx} \cos u = -\sin u \frac{du}{dx}$
$\frac{d}{dx}(\sqrt{u}) = \frac{1}{2\sqrt{u}} \frac{du}{dx}$	$\frac{d}{dx} \tan u = \sec^2 u \frac{du}{dx}$
$\frac{d}{dx}\left(\frac{1}{u}\right) = -\frac{1}{u^2} \frac{du}{dx}$	$\frac{d}{dx} \cot u = -\csc^2 u \frac{du}{dx}$
$\frac{d}{dx}\left(\frac{1}{u^n}\right) = -\frac{n}{u^{n+1}} \frac{du}{dx}$	$\frac{d}{dx} \sec u = \sec u \tan u \frac{du}{dx}$
$\frac{d}{dx}[f(u)] = \frac{d}{du}[f(u)] \frac{du}{dx}$	$\frac{d}{dx} \csc u = -\csc u \cot u \frac{du}{dx}$

[source: <http://hyperphysics.phy-astr.gsu.edu/hbase/Math/derfunc.html>]

How about computing all the derivatives?

But neural net f is never one of those?

No problem: CHAIN RULE:

If $f(x) = g(h(x))$

Then $f'(x) = g'(h(x))h'(x)$

→ Derivatives can be computed by following well-defined procedures

Automatic Differentiation

- Automatic differentiation software

- e.g. Theano, TensorFlow, PyTorch, Chainer
- Only need to program the function $g(x,y,w)$
- Can automatically compute all derivatives w.r.t. all entries in w
- This is typically done by caching info during forward computation pass of f , and then doing a backward pass = “backpropagation”
- Autodiff / Backpropagation can often be done at computational cost comparable to the forward pass

- Need to know this exists

- How this is done? -- outside of scope of CS188

Summary of Key Ideas

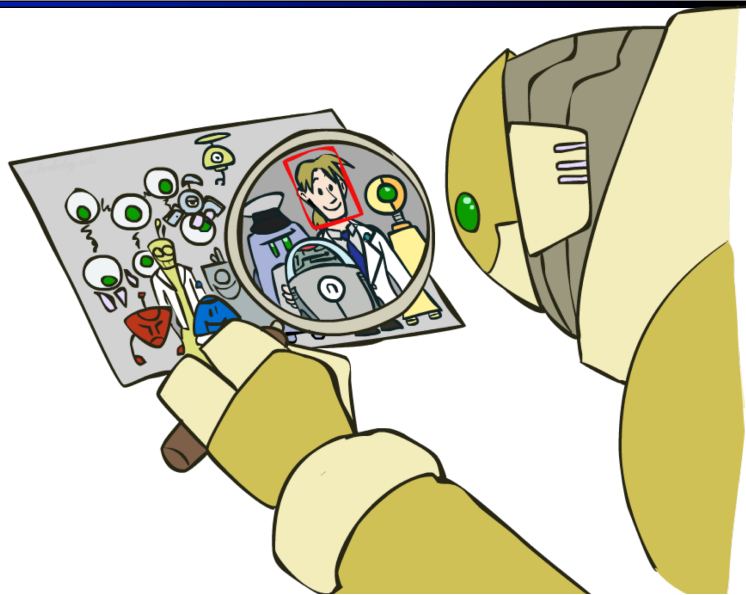
- Optimize probability of label given input $\max_w ll(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$
- Continuous optimization
 - Gradient ascent:
 - Compute steepest uphill direction = gradient (= just vector of partial derivatives)
 - Take step in the gradient direction
 - Repeat (until held-out data accuracy starts to drop = “early stopping”)
- Deep neural nets
 - Last layer = still logistic regression
 - Now also many more layers before this last layer
 - = computing the features
 - → the features are learned rather than hand-designed
 - Universal function approximation theorem
 - If neural net is large enough
 - Then neural net can represent any continuous mapping from input to output with arbitrary accuracy
 - But remember: need to avoid overfitting / memorizing the training data → early stopping!
- Automatic differentiation gives the derivatives efficiently (how? = outside of scope of 188)

How well does it work?

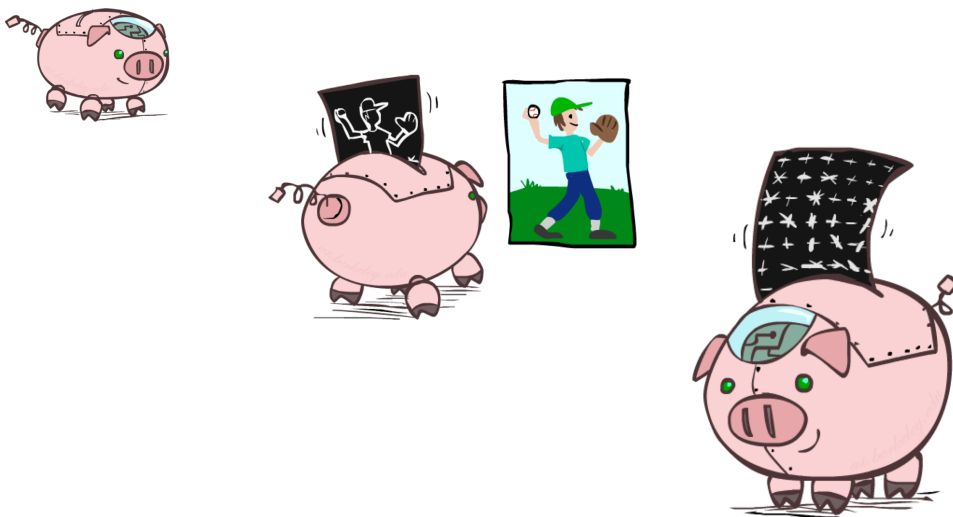
Computer Vision



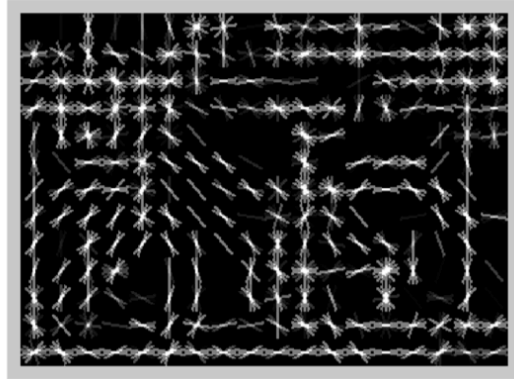
Object Detection



Manual Feature Design



Features and Generalization

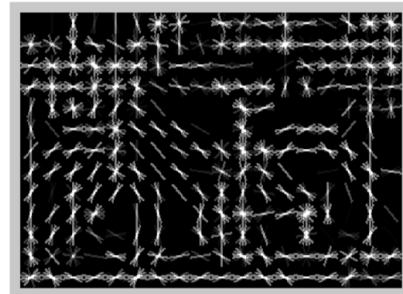


[HoG: Dalal and Triggs, 2005]

Features and Generalization



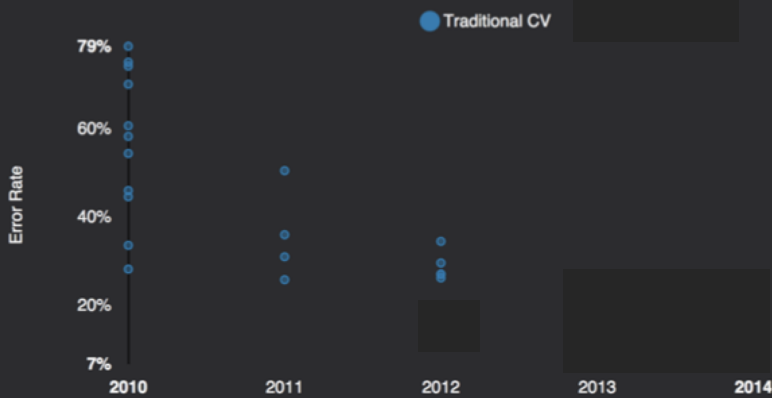
Image



HoG

Performance

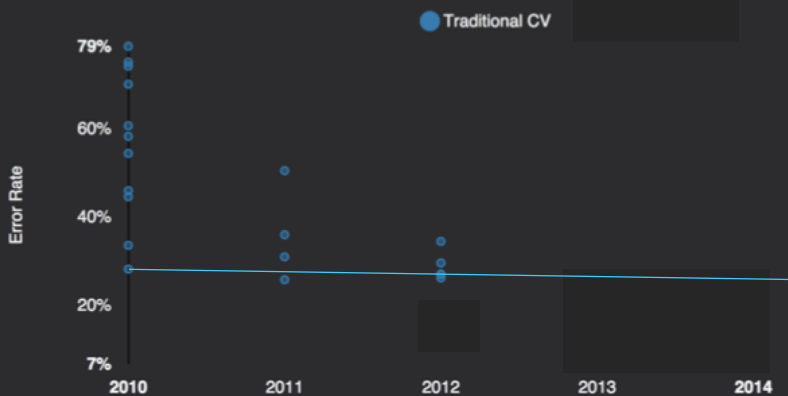
ImageNet Error Rate 2010-2014



graph credit Matt
Zeiler, Clarifai

Performance

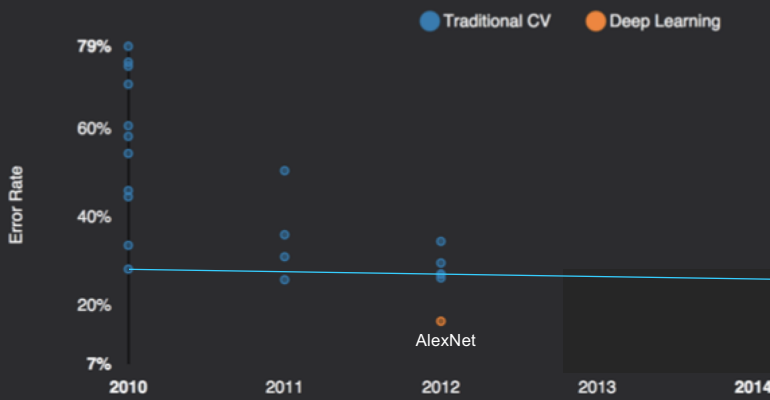
ImageNet Error Rate 2010-2014



graph credit Matt
Zeiler, Clarifai

Performance

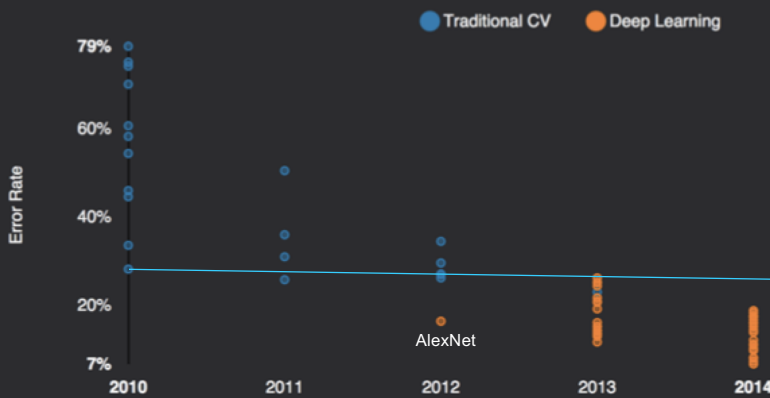
ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

Performance

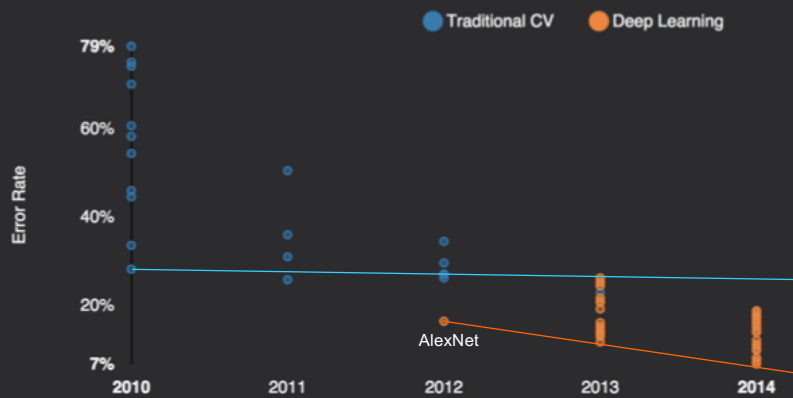
ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

Performance

ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

MS COCO Image Captioning Challenge



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



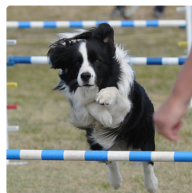
"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."

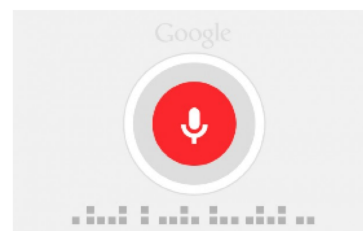
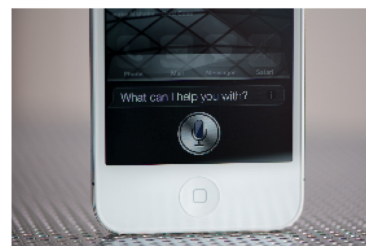
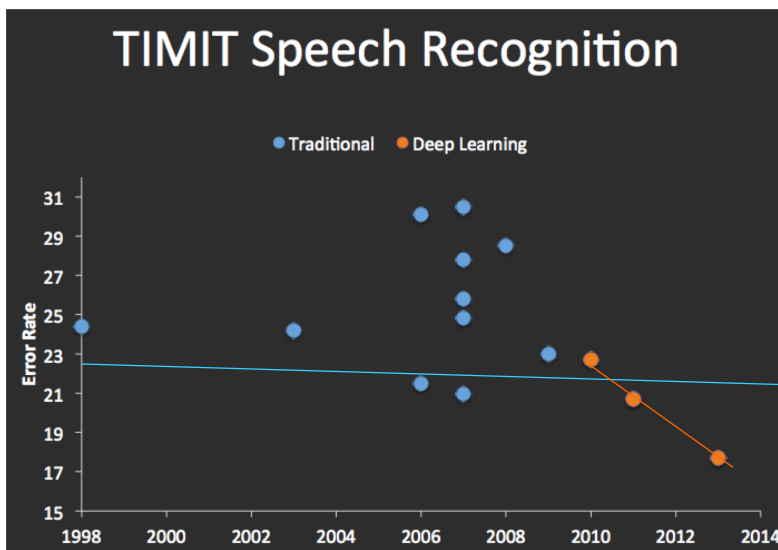
Karpathy & Fei-Fei, 2015; Donahue et al., 2015; Xu et al, 2015; many more

Visual QA Challenge

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh

 <p>What vegetable is on the plate? Neural Net: broccoli Ground Truth: broccoli</p>	 <p>What color are the shoes on the person's feet ? Neural Net: brown Ground Truth: brown</p>	 <p>How many school busses are there? Neural Net: 2 Ground Truth: 2</p>	 <p>What sport is this? Neural Net: baseball Ground Truth: baseball</p>
 <p>What is on top of the refrigerator? Neural Net: magnets Ground Truth: cereal</p>	 <p>What uniform is she wearing? Neural Net: shorts Ground Truth: girl scout</p>	 <p>What is the table number? Neural Net: 4 Ground Truth: 40</p>	 <p>What are people sitting under in the back? Neural Net: bench Ground Truth: tent</p>

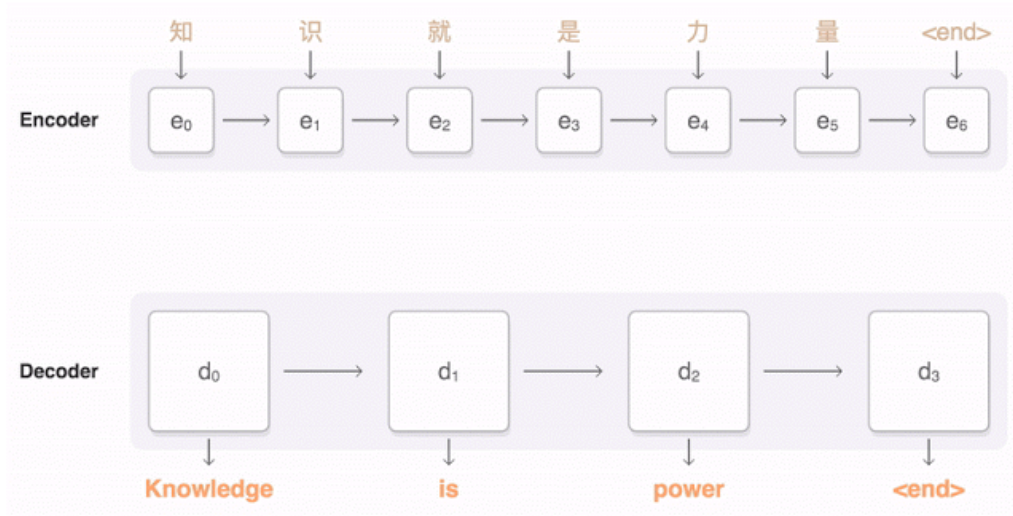
Speech Recognition



graph credit Matt Zeiler, Clarifai

Machine Translation

Google Neural Machine Translation (in production)



Next: More Neural Net Applications!