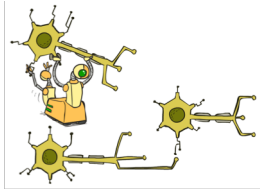


CS 687

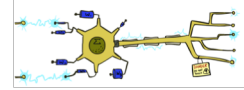
Optimization and Neural Nets



[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/>.]
Modifications and additions J. Kozsecka GMLU

Reminder: Linear Classifiers

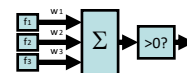
- Inputs are **feature values**
- Each feature has a **weight**
- Sum is the **activation**



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

- If the activation is:

- Positive, output +1
- Negative, output -1

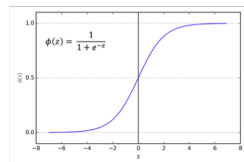


How to get probabilistic decisions?

- 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

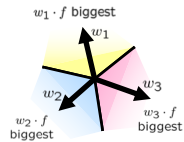
Multiclass Logistic Regression

- 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:

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

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

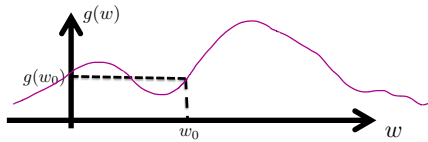
- Start wherever
- Repeat: move to the best neighboring state
- If no neighbors better than current, quit



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

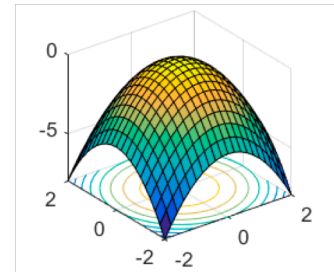
- Optimization over a continuous space
 - Infinitely many neighbors!
 - How to do this efficiently?

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



Source: offconvex.org

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

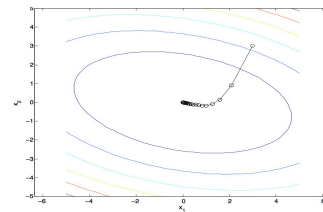


Figure source: Mathworks

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

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

Optimization Procedure: Gradient Ascent

```

▪ init w
▪ for iter = 1, 2, ...
    w ← w + α * ∇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 iter = 1, 2, ...
    w ← w + α * ∑_i ∇ 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

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

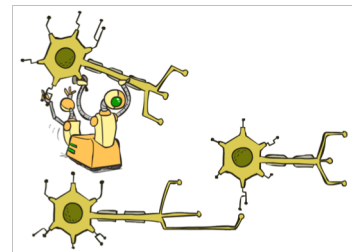
Observation: gradient over small set of training examples (=mini-batch) can be computed in parallel, might as well do that instead of a single one

```
▪ 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?

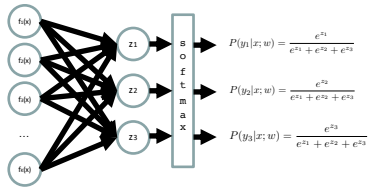
- We'll talk about that once we covered neural networks, which are a generalization of logistic regression

Neural Networks

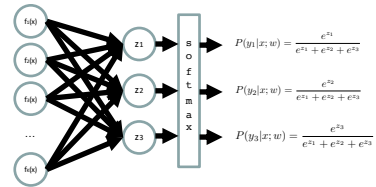


Multi-class Logistic Regression

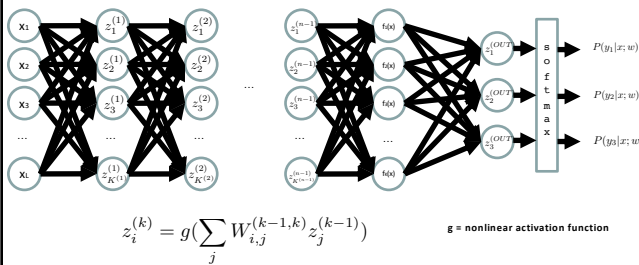
- = special case of neural network



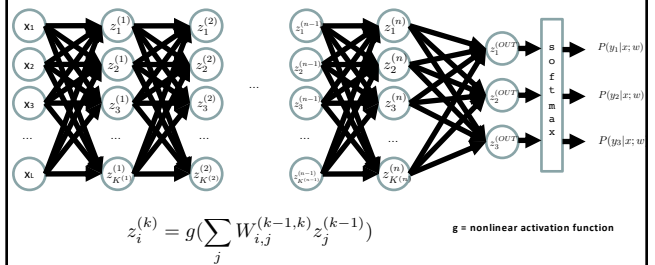
Deep Neural Network = Also learn the features!



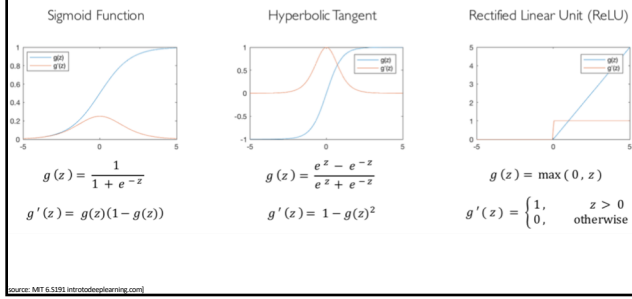
Deep Neural Network = Also learn the features!



Deep Neural Network = Also learn the features!



Common Activation Functions



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 \mathbb{R}^k such that $\int_{\mathbb{R}^k} |f(x)|^p d\mu(x) < \infty$) 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 $X \subseteq \mathbb{R}^k$, 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 $f(x)$.

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"

[illegible]

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

- Derivatives tables:

$$\begin{array}{ll}
 \frac{d}{dx}(0) = 0 & \frac{d}{dx}[\log u] = \frac{d}{dx}[\log v] = \frac{1}{u} \frac{du}{dx} \\
 \frac{d}{dx}(x) = 1 & \frac{d}{dx}[\log v] = \frac{1}{v} \frac{dv}{dx} \\
 \frac{d}{dx}(c) = \frac{dc}{dx} & \frac{d}{dx} e^u = e^u \frac{du}{dx} \\
 \frac{d}{dx}(u+v) = \frac{du}{dx} + \frac{dv}{dx} & \frac{d}{dx} u^a = a u^{a-1} \frac{du}{dx} \\
 \frac{d}{dx}(uv) = u \frac{dv}{dx} + v \frac{du}{dx} & \frac{d}{dx} (u^a v^b) = u^a v^b \left(a \frac{du}{u} + b \frac{dv}{v} \right) \\
 \frac{d}{dx} \left(\frac{u}{v} \right) = \frac{1}{v^2} u \frac{dv}{dx} - \frac{v}{v^2} \frac{du}{dx} & \frac{d}{dx} \sin u = \cos u \frac{du}{dx} \\
 \frac{d}{dx}(u^a) = au^{a-1} \frac{du}{dx} & \frac{d}{dx} \cos u = -\sin u \frac{du}{dx} \\
 \frac{d}{dx} \left(\frac{1}{u} \right) = -\frac{1}{u^2} \frac{du}{dx} & \frac{d}{dx} \tan u = \sec^2 u \frac{du}{dx} \\
 \frac{d}{dx} \left(\frac{1}{v} \right) = -\frac{1}{v^2} \frac{dv}{dx} & \frac{d}{dx} \cot u = -\csc^2 u \frac{du}{dx} \\
 \frac{d}{dx} \left(\frac{1}{u} \right) = -\frac{1}{u^2} \frac{du}{dx} & \frac{d}{dx} \operatorname{cosec} u = -\operatorname{cosec} u \cot u \frac{du}{dx} \\
 \frac{d}{dx} f(u) = \frac{du}{dx} f'(u) & \frac{d}{dx} \operatorname{cosec} u = -\operatorname{cosec} u \cot u \frac{du}{dx}
 \end{array}$$

- 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



Backpropagation: applications

- Perhaps the most successful and widely used learning algorithm for NNs;
- Used in a variety of domains:
 - clinical diagnosis,
 - predicting protein structure,
 - character recognition,
 - fingerprint recognition,
 - modeling residual chlorine decay in water,
 - weather forecast,
 - waveform recognition,
 - backgammon, etc.

38