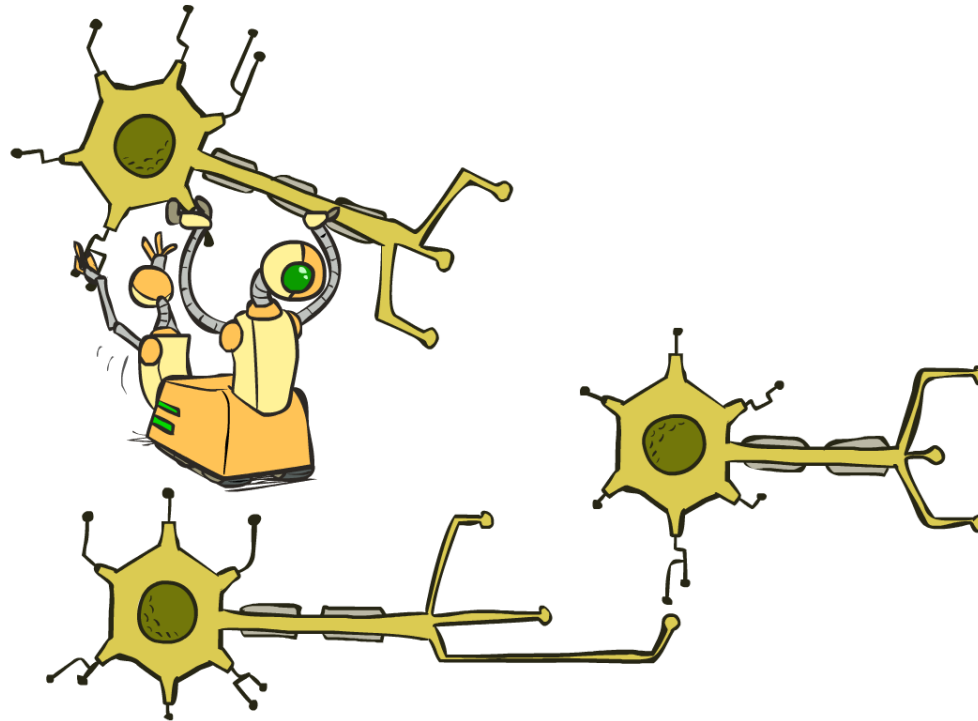# Announcements

- Gradescope
  - You must assign questions to page numbers when you submit so the TAs can easily find your answers.
  - You will get a 0 otherwise!
- HW4 due today!
- P2 due on Thursday!
- Midterm next week!

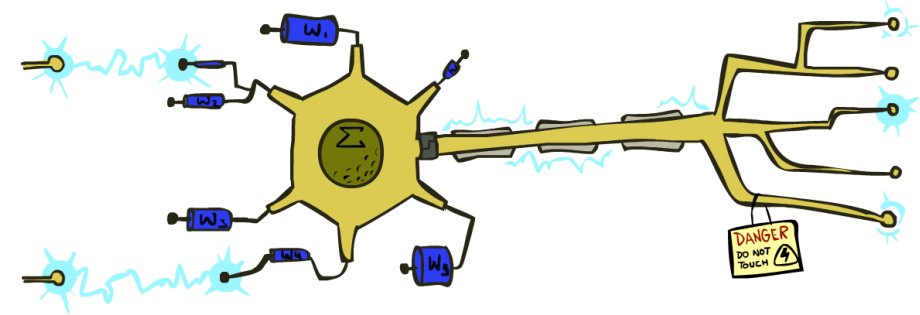# CS 188: Artificial Intelligence

## Optimization and Neural Nets

Instructor: Anca Dragan --- University of California, Berkeley

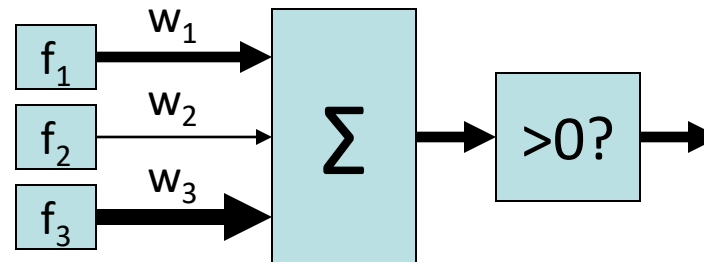# Reminder: Linear Classifiers

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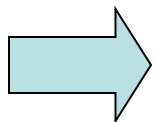  - Positive, output +1
  - Negative, output -1

# Feature Vectors

$$x \qquad f(x) \qquad y$$

```
Hello,

Do you want free printr
cartriges?  Why pay more
when you can get them
ABSOLUTELY FREE!   Just
```

➡

```
# free      : 2
YOUR_NAME   : 0
MISSPELLED  : 2
FROM_FRIEND : 0
...
```

➡

SPAM
or
Not
SPAM

➡

```
PIXEL-7,12  : 1
PIXEL-7,13  : 0
...
NUM_LOOPS   : 1
...
```

➡

"2"

# Feature Vectors
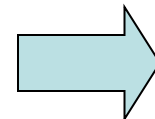
$$x \qquad\qquad f(x) \qquad\qquad y$$

```
Hello,

Do you want free printr
cartriges?  Why pay more
when you can get them
ABSOLUTELY FREE!   Just
```
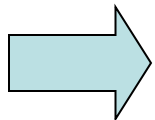
➡️

```
# free      : 2
YOUR_NAME   : 0
MISSPELLED  : 2
FROM_FRIEND : 0
...
```
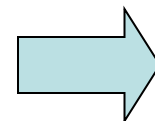
➡️

SPAM
or
Not
SPAM

➡️

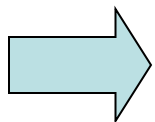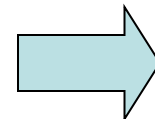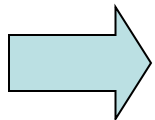```
PIXEL-7,12  : 1
PIXEL-7,13  : 0
...
NUM_LOOPS   : 1
...
```

➡️

"2"

# Learning: Binary Perceptron

- Start with weights = 0
- For each training instance:
  - Classify with current weights

$$y = \begin{cases} +1 & \text{if } w \cdot f(x) \geq 0 \\ -1 & \text{if } w \cdot f(x) < 0 \end{cases}$$

  - If correct (i.e., y=y*), no change!
  - If wrong: adjust the weight vector by adding or subtracting the feature vector. Subtract if y* is -1.

$$w = w + y^* \cdot f$$

$w$

$y^* \cdot f$

$f$

Before update $w^T f(x) > 0$ and $y^* = -1$

After update
$(w - f(x))^T f(x) = w^T f(x) - f(x)^T f(x) < w^T f(x)$

# Multiclass Decision Rule

- **If we have multiple classes:**
  - A weight vector for each class:

  $$w_y$$

  - Score (activation) of a class y:

  $$w_y \cdot f(x)$$

  - Prediction highest score wins

  $$y = \arg\max_y \ w_y \cdot f(x)$$

$w_1 \cdot f$ biggest

$w_1$

$w_3$

$w_2$

$w_2 \cdot f$ biggest

$w_3 \cdot f$ biggest

*Binary = multiclass where the negative class has weight zero*

# Learning: Multiclass Perceptron

- Start with all weights = 0
- Pick up training examples one by one
- Predict with current weights

$$y = \arg\max_y \ w_y \cdot f(x)$$

- If correct, no change!
- If wrong: lower score of wrong answer, raise score of right answer

$$w_y = w_y - f(x)$$

$$w_{y*} = w_{y*} + f(x)$$

# Problems with the Perceptron

- Noise: if the data isn't separable, weights might thrash
  - Averaging weight vectors over time can help (averaged perceptron)

- Mediocre generalization: finds a "barely" separating solution

- Overtraining: test / held-out accuracy usually rises, then falls
  - Overtraining is a kind of overfitting

# Improving the Perceptron

# Non-Separable Case: Probabilistic Decision

# How to get probabilistic decisions?

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

- Sigmoid function

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



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

# A 1D Example



$$P(\text{red}|x) = \frac{e^{w_{\text{red}} \cdot x}}{e^{w_{\text{red}} \cdot x} + e^{w_{\text{blue}} \cdot x}}$$

probability increases exponentially as we move away from boundary

normalizer

# The *Soft* Max



$P(\text{red}|x)$

$$\dfrac{e^{5w_{\text{red}}\cdot x}}{e^{5w_{\text{red}}\cdot x}+e^{5w_{\text{blue}}\cdot x}}$$

$$\dfrac{e^{100w_{\text{red}}\cdot x}}{e^{100w_{\text{red}}\cdot x}+e^{100w_{\text{blue}}\cdot x}}$$

$$\dfrac{e^{w_{\text{red}}\cdot x}}{e^{w_{\text{red}}\cdot x}+e^{w_{\text{blue}}\cdot x}}$$

looks like $\max_y w_y \cdot x$

$x$

$$P(\text{red}|x) = \dfrac{e^{w_{\text{red}}\cdot x}}{e^{w_{\text{red}}\cdot x}+e^{w_{\text{blue}}\cdot x}}$$

# Multiclass Logistic Regression

- Recall Perceptron:

  - A weight vector for each class: $w_y$

  - Score (activation) of a class y: $w_y \cdot f(x)$

  - Prediction highest score wins $y = \arg\max_y \; w_y \cdot f(x)$

$w_1 \cdot f$ biggest

$w_1$

$w_3$

$w_2$

$w_2 \cdot f$ biggest

$w_3 \cdot f$ biggest

- How to make the scores into probabilities?

$$z_1, z_2, z_3 \rightarrow \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}, \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

original activations

softmax activations

# Best w?

- Maximum likelihood estimation:

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

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**

# How do we learn in this setting?

- Optimization

  - i.e., how do we solve:

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

# Hill Climbing

- **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?

# 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:

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

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

# Gradient Ascent

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

# Gradient in n dimensions

$$\nabla g = \begin{bmatrix} \dfrac{\partial g}{\partial w_1} \\ \dfrac{\partial g}{\partial w_2} \\ \cdots \\ \dfrac{\partial g}{\partial w_n} \end{bmatrix}$$

# Optimization Procedure: Gradient Ascent

- `Init` $w$
- `for iter = 1, 2, …`

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

- $\alpha$: learning rate --- tweaking parameter that needs to be chosen carefully
- How? Try multiple choices

# 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, …

$$w \leftarrow w + \alpha * \sum_i \nabla \log P(y^{(i)}|x^{(i)};w)$$

# 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

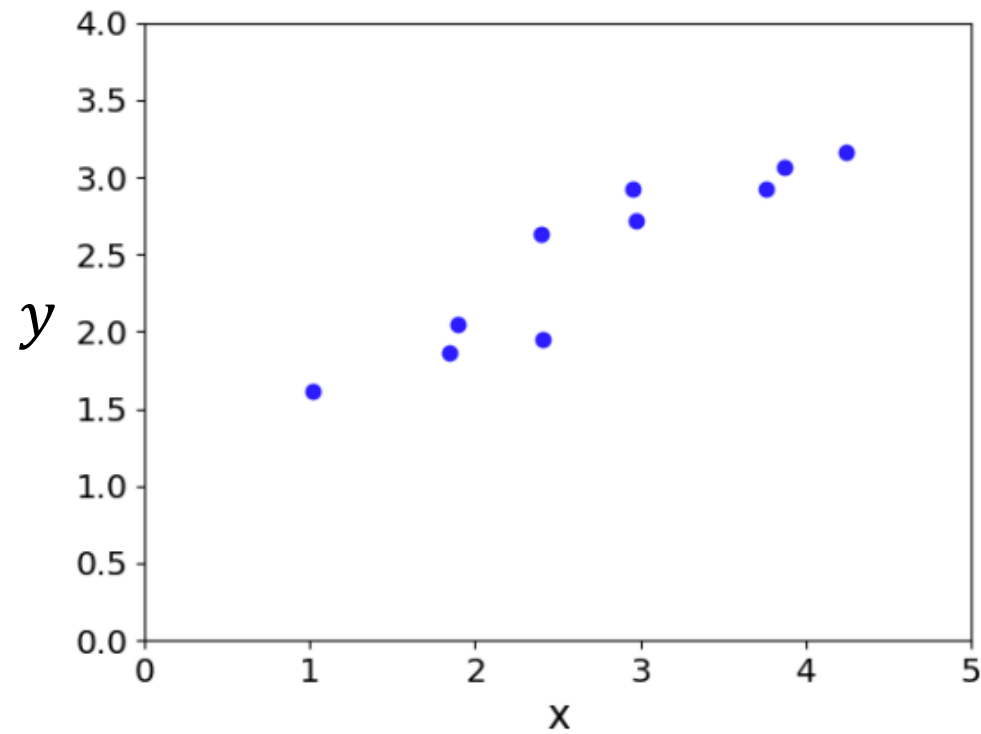$$\max_{w} \; ll(w) = \max_{w} \; \sum_{i} \log P(y^{(i)}|x^{(i)}; w)$$

**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`

  $$w \leftarrow w + \alpha * \sum_{j \in J} \nabla \log P(y^{(j)}|x^{(j)}; w)$$

# Regression

# Linear Regression

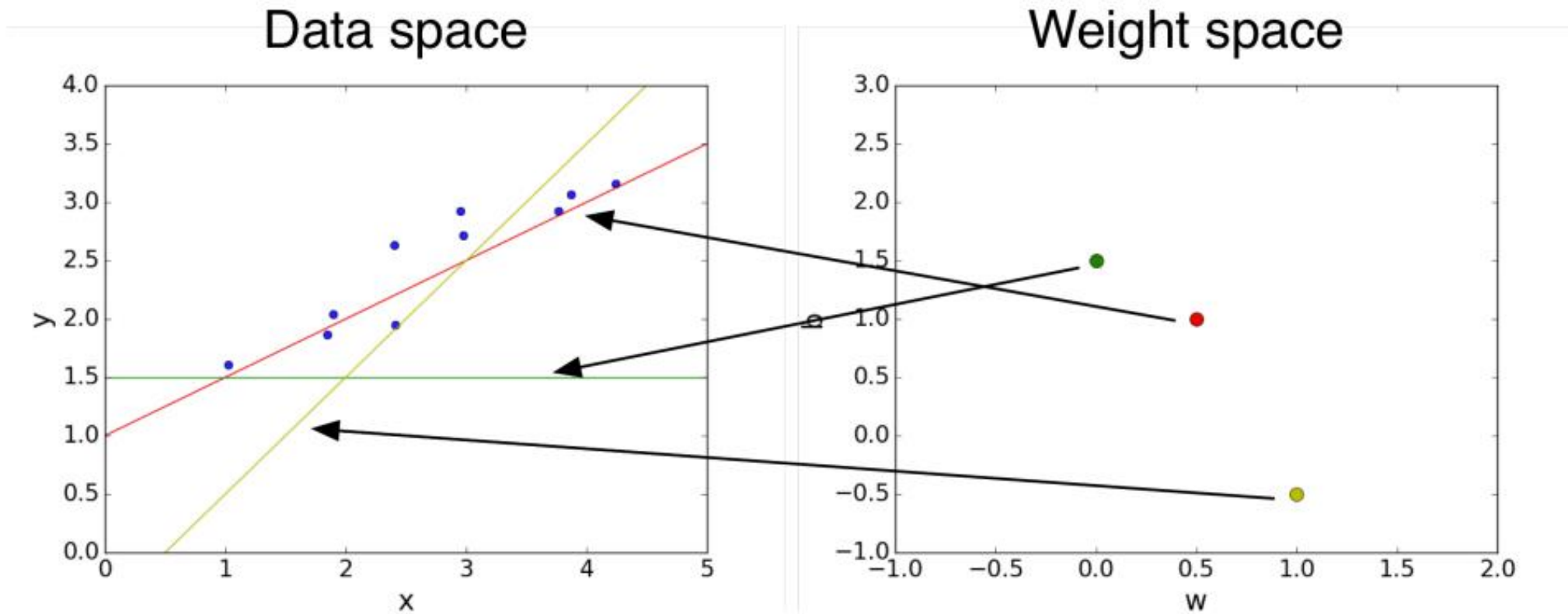- Learn to map from inputs $x$ to outputs $y \in \mathbb{R}$
- $\hat{y} = w^T f(x) = w \circ f(x) = \sum_{i=1}^{d} w_i \cdot x_i$
  - Where $f(x)$ maps from the raw input x into a set of features

- How can we measure how good a set of weights $w$ are?
  - Loss Function
    - Squared Error

$$\mathcal{L}(w) = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}\left(y - w^T f(x)\right)^2$$

# Linear Regression

- We want to find w that minimizes this loss function

$$\mathcal{L}(w) = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}\left(y - w^T f(x)\right)^2$$

- Calculus to the rescue!

$$\frac{\partial \mathcal{L}}{\partial w_i} = -\left(y - w^T f(x)\right)f_i(x) = \left(w^T f(x) - y\right)f_i(x)$$

# Linear Regression

- How can we measure how good a set of weights $w$ are?
  - Mean Squared Error

$$\mathcal{L}_{MSE}(w) = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2}(y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{2}(y_i - w^T f(x))^2$$

$$\frac{\partial \mathcal{L}_{MSE}}{\partial w_i} = \frac{1}{N} \sum_{i=1}^{N} (w^T f(x) - y) f_i(x)$$
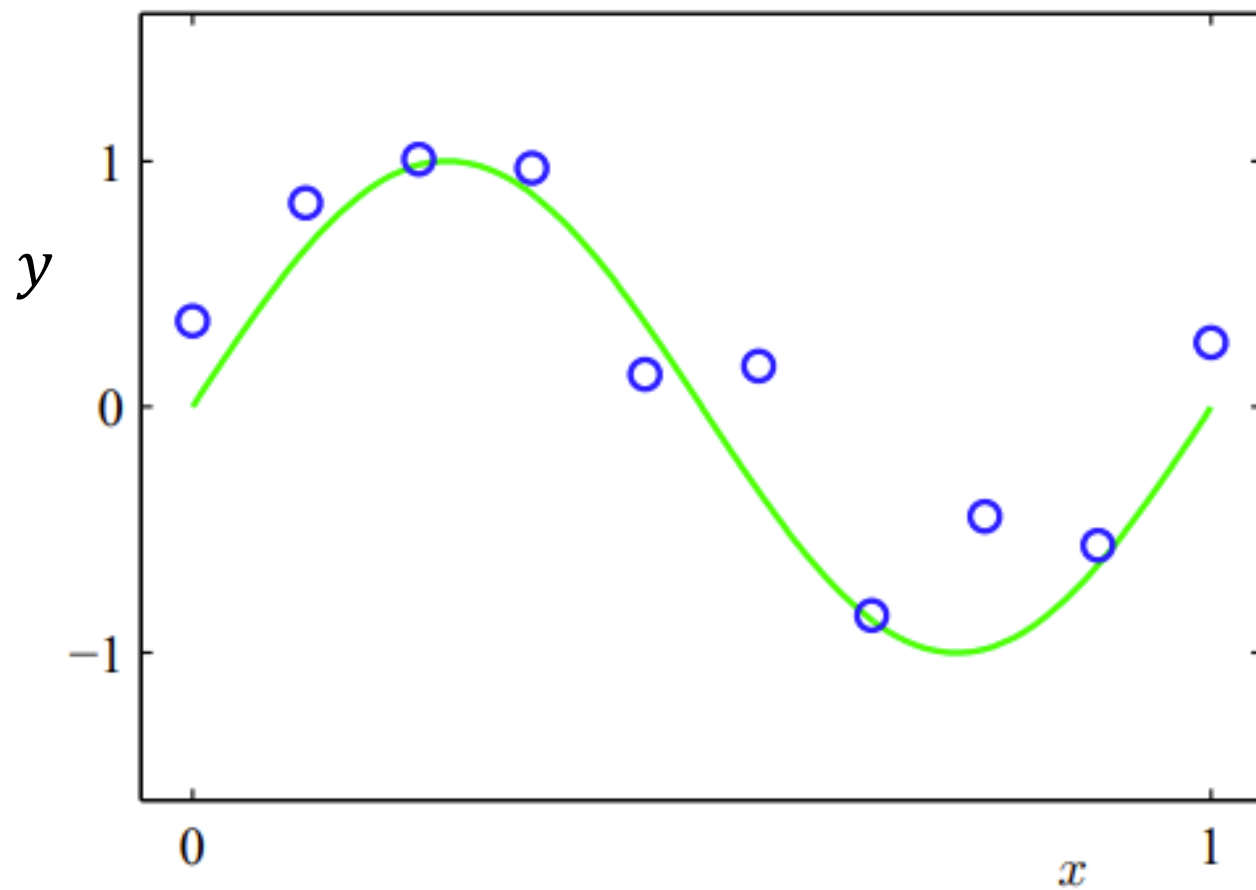
# How to optimize?

- Gradient Descent

  - Init $w$
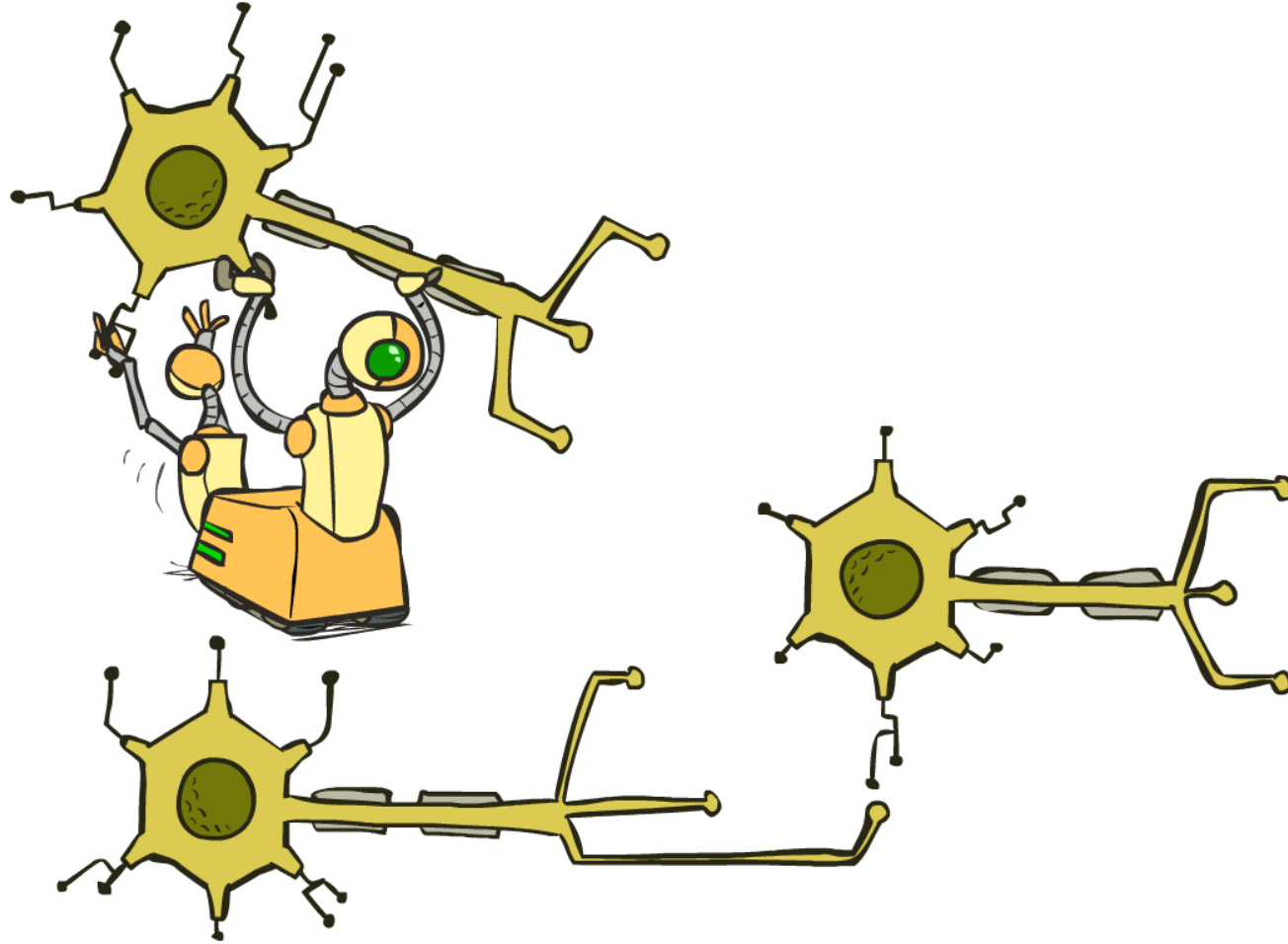  - for iter = 1, 2, …

    $$w \leftarrow w - \alpha \cdot \nabla \mathcal{L}(w)$$

- Usually we will use stochastic gradient descent using mini-batches, just like we talked about for classification.
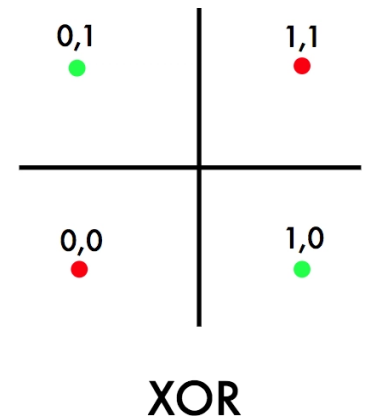
# Non-Linear Regression?

# Neural Networks

# History Lesson

- ## 1943: Artificial Neuron
  - McCulloch and Pitts showed simple threshold logic

- ## 1957: Perceptron
  - Rosenblatt introduced algorithm for single-layer neural network
  - Explored "Multi-Layer Perceptrons" but lacked good learning algorithms

- ## 1969: AI Winter
  - Minsky and Papert publish book called "Perceptrons"

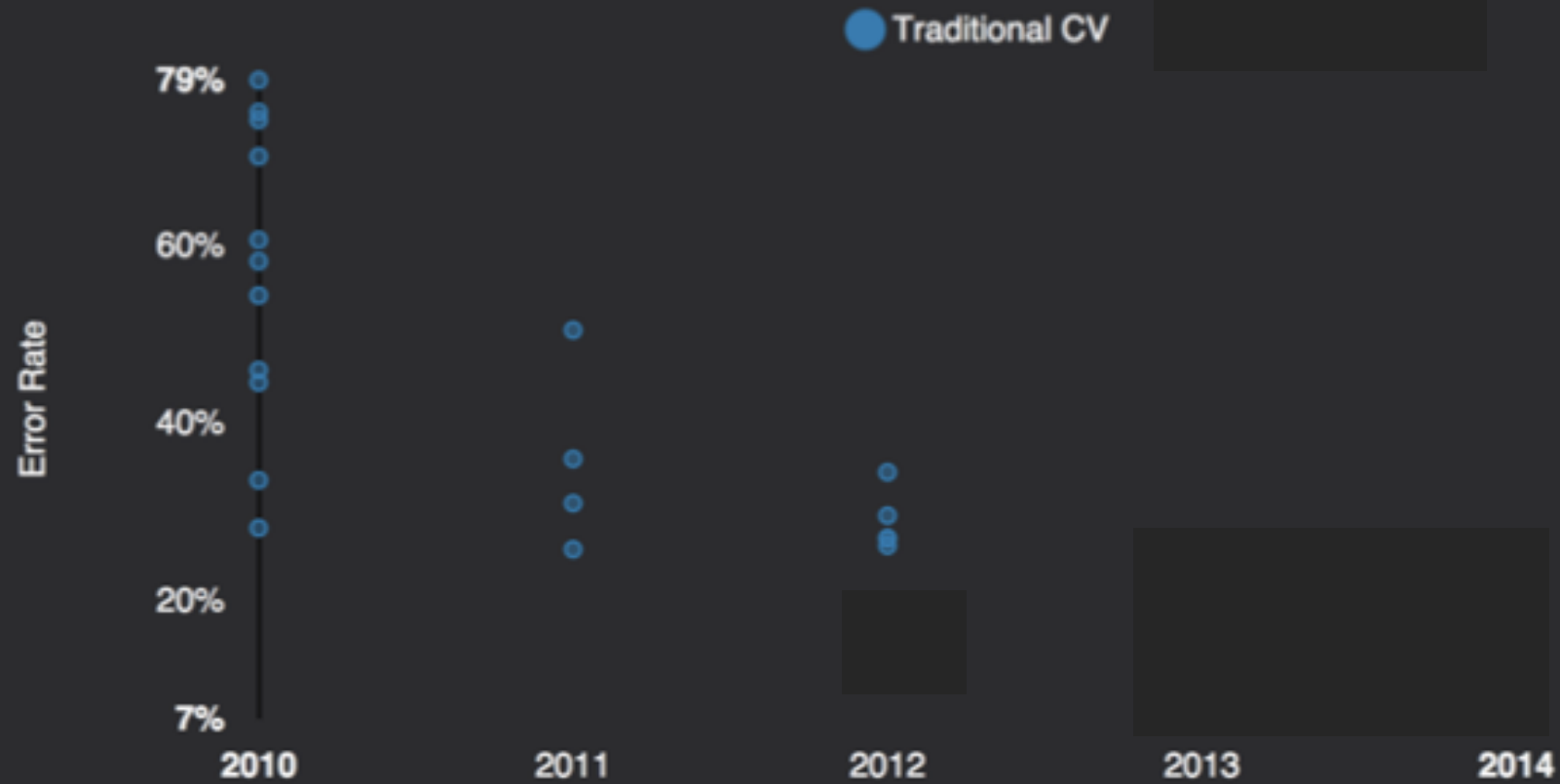0,1          1,1

0,0          1,0

XOR

# History Lesson

- 1986: Backpropagation

  - Rumelhart, Hinton, and Williams show power of multi-layer perceptrons trained via backpropagation

- Early 2010s: Hardware and algorithms converge with big data

- 2012: AlexNet and the image recognition breakthrough

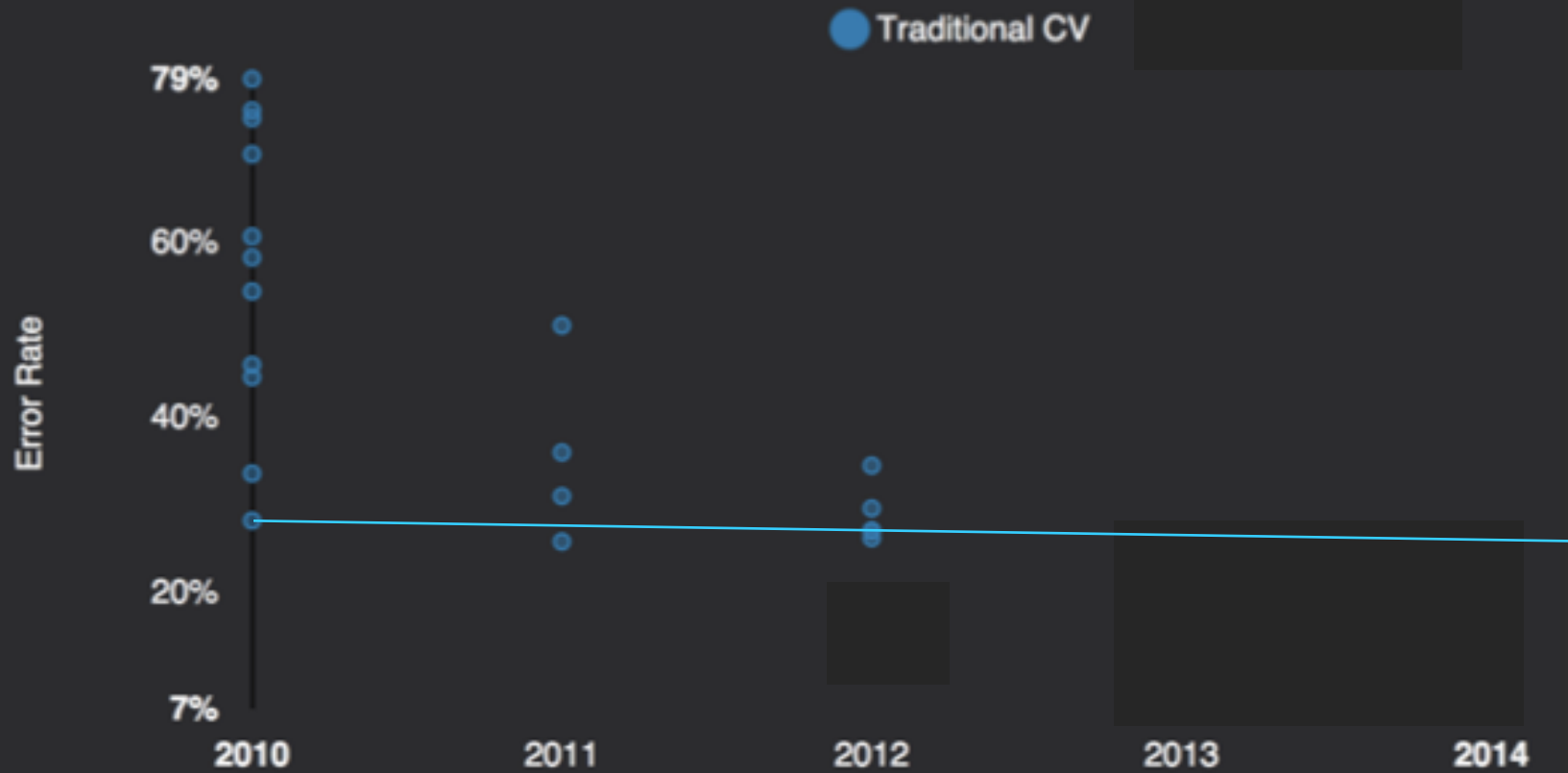- 2012–present: The deep learning era

# Performance



ImageNet Error Rate 2010-2014

*graph credit Matt Zeiler, Clarifai*

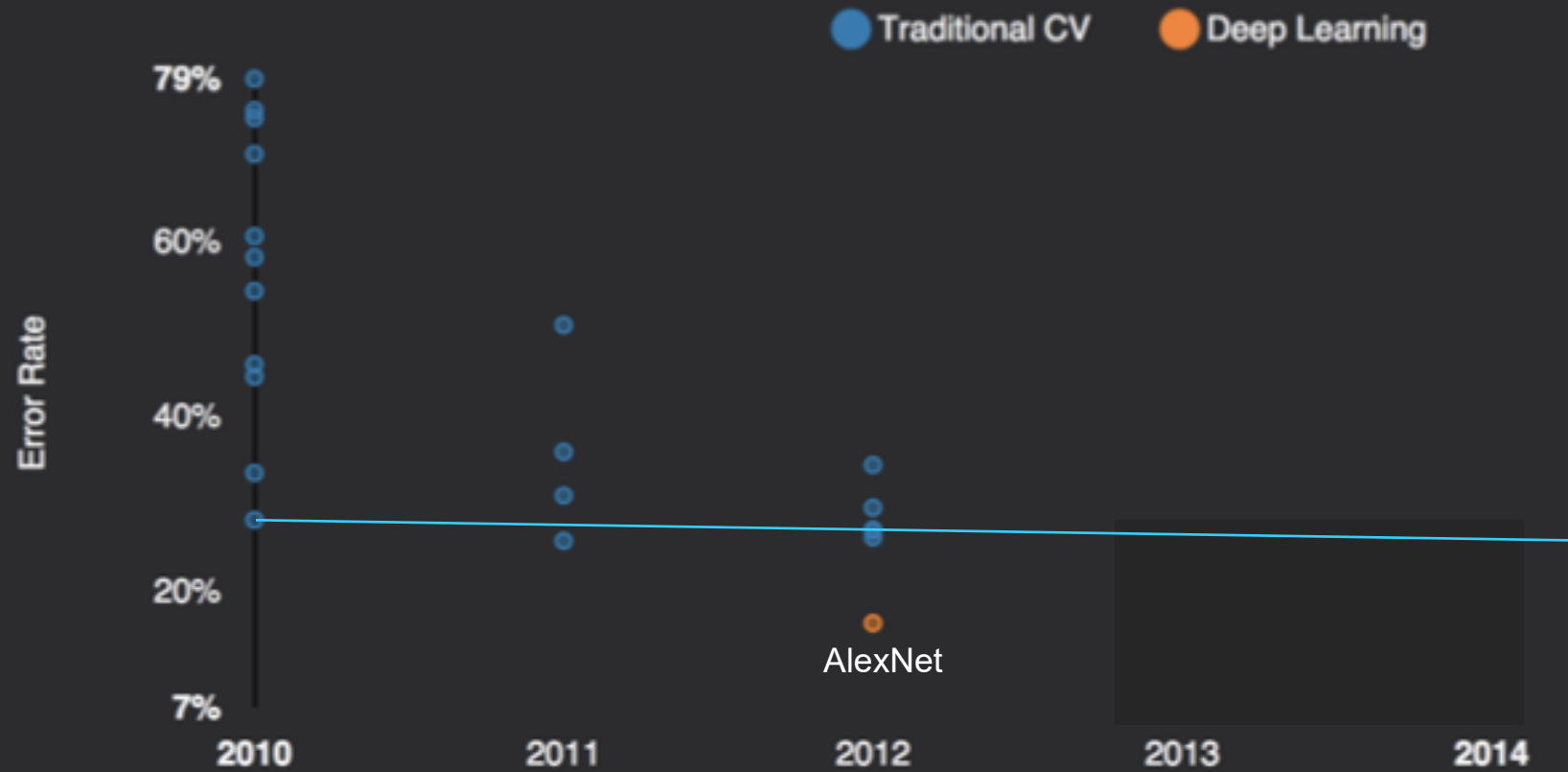# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance



*graph credit Matt Zeiler, Clarifai*

# Performance
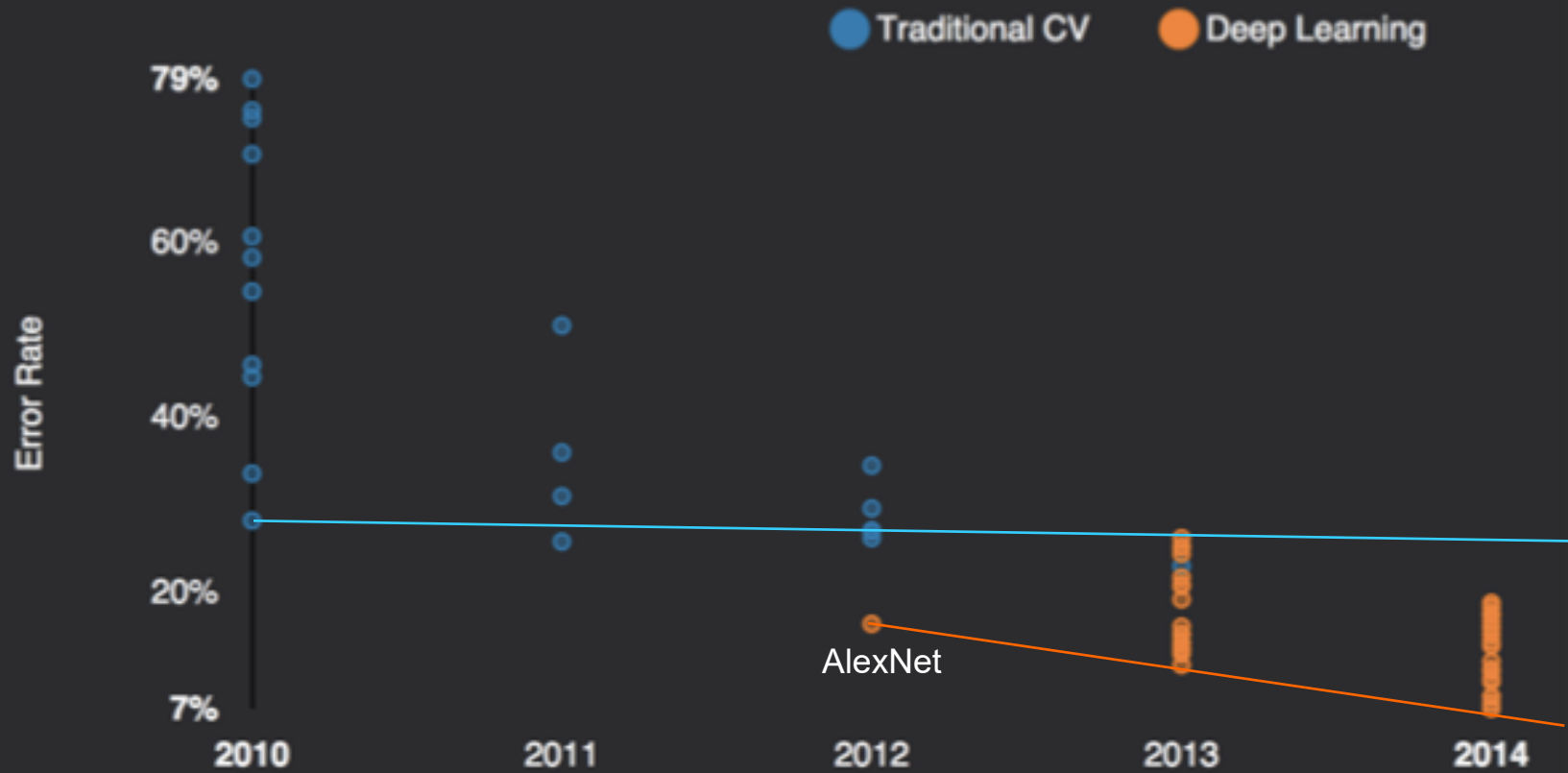


ImageNet Error Rate 2010-2014

*graph credit Matt Zeiler, Clarifai*

# Performance



*graph credit Matt Zeiler, Clarifai*

# Multi-class Logistic Regression

- = special case of neural network



$$P(y_1|x;w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x;w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x;w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Deep Neural Network = Also learn the features!



$$P(y_1|x; w) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_2|x; w) = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$P(y_3|x; w) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# 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!



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

**g = nonlinear activation function**

# Common Activation Functions

## Sigmoid Function



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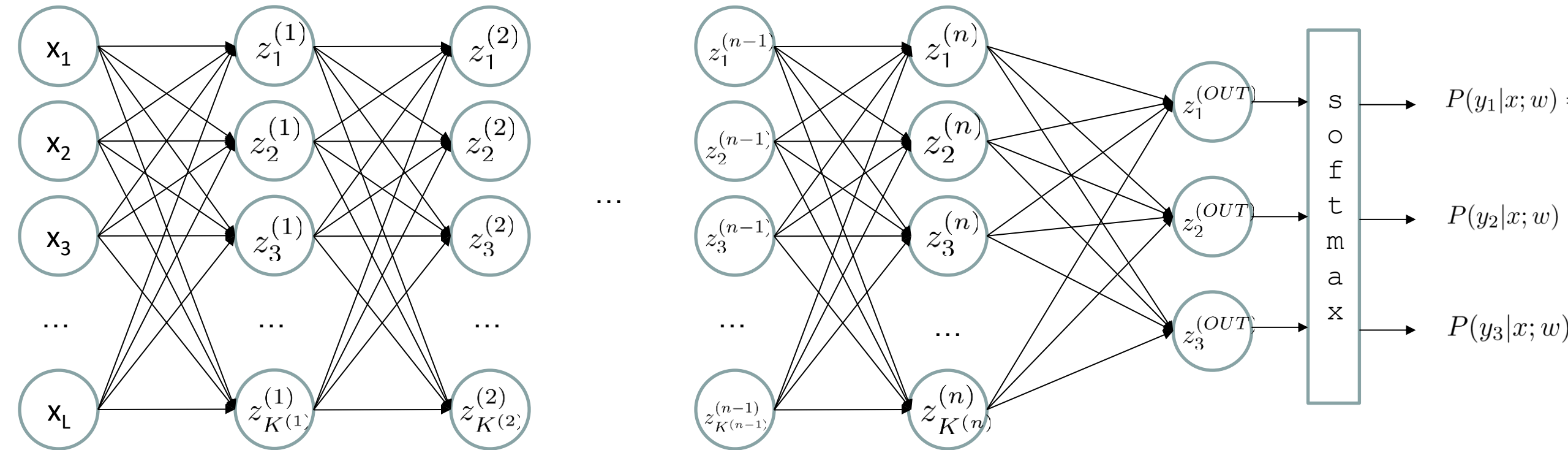
$$g'(z) = g(z)(1 - g(z))$$

## Hyperbolic Tangent



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

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

## Rectified Linear Unit (ReLU)



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

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

# 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

# Deep Neural Networks for Regression

# Deep Neural Networks for Regression

# 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.

# Fun Neural Net Demo Site

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

# How about computing all the derivatives?

- Derivatives tables:

$$\frac{d}{dx}(a) = 0$$

$$\frac{d}{dx}(x) = 1$$

$$\frac{d}{dx}(au) = a\frac{du}{dx}$$

$$\frac{d}{dx}(u+v-w) = \frac{du}{dx} + \frac{dv}{dx} - \frac{dw}{dx}$$

$$\frac{d}{dx}(uv) = u\frac{dv}{dx} + v\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{u}{v}\right) = \frac{1}{v}\frac{du}{dx} - \frac{u}{v^2}\frac{dv}{dx}$$

$$\frac{d}{dx}(u^n) = nu^{n-1}\frac{du}{dx}$$

$$\frac{d}{dx}(\sqrt{u}) = \frac{1}{2\sqrt{u}}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u}\right) = -\frac{1}{u^2}\frac{du}{dx}$$

$$\frac{d}{dx}\left(\frac{1}{u^n}\right) = -\frac{n}{u^{n+1}}\frac{du}{dx}$$

$$\frac{d}{dx}[f(u)] = \frac{d}{du}[f(u)]\frac{du}{dx}$$

$$\frac{d}{dx}[\ln u] = \frac{d}{dx}[\log_e u] = \frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}[\log_a u] = \log_a e \frac{1}{u}\frac{du}{dx}$$

$$\frac{d}{dx}e^u = e^u\frac{du}{dx}$$

$$\frac{d}{dx}a^u = a^u \ln a \frac{du}{dx}$$

$$\frac{d}{dx}(u^v) = vu^{v-1}\frac{du}{dx} + \ln u \; u^v \frac{dv}{dx}$$

$$\frac{d}{dx}\sin u = \cos u \frac{du}{dx}$$

$$\frac{d}{dx}\cos u = -\sin u \frac{du}{dx}$$

$$\frac{d}{dx}\tan u = \sec^2 u \frac{du}{dx}$$

$$\frac{d}{dx}\cot u = -\csc^2 u \frac{du}{dx}$$

$$\frac{d}{dx}\sec u = \sec u \tan u \frac{du}{dx}$$

$$\frac{d}{dx}\csc u = -\csc u \cot u \frac{du}{dx}$$

[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. PyTorch, TensorFlow, Jax
  - 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 our class