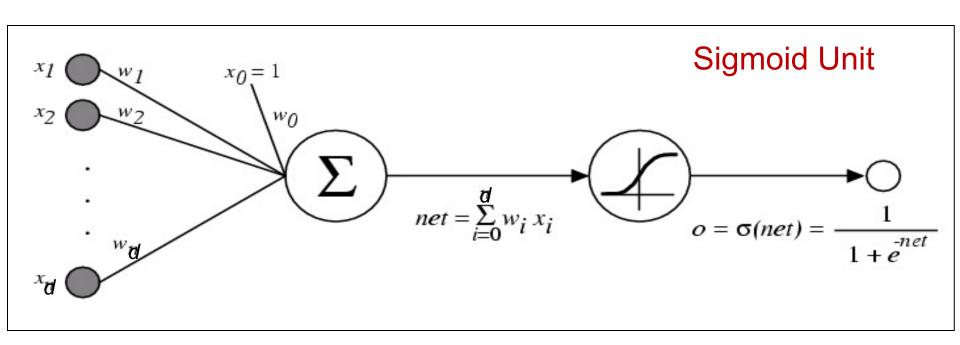
Neural Networks

Aarti Singh

Machine Learning 10-315 Sept 27, 2021

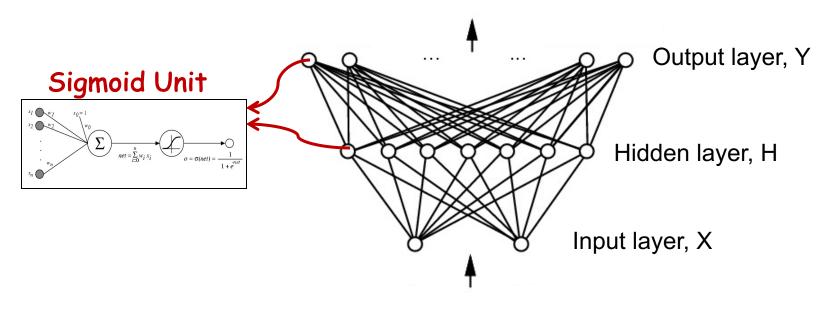
Logistic function as a Graph

Output,
$$o(\mathbf{x}) = \sigma(w_0 + \sum_i w_i X_i) = \frac{1}{1 + \exp(-(w_0 + \sum_i w_i X_i))}$$



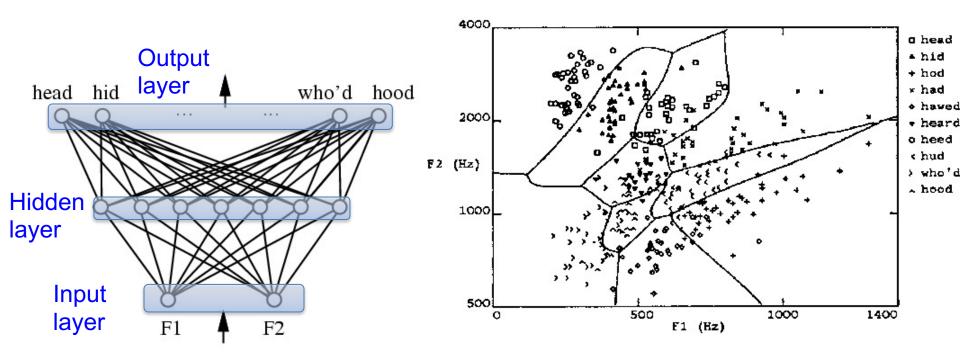
Neural Networks to learn f: X -> Y

- f can be a non-linear function
- X (vector of) continuous and/or discrete variables
- Y (vector of) continuous and/or discrete variables
- Neural networks Represent f by <u>network</u> of sigmoid (more recently ReLU – next lecture) units:



Multilayer Networks of Sigmoid Units

Neural Network trained to distinguish vowel sounds using 2 formants (features)



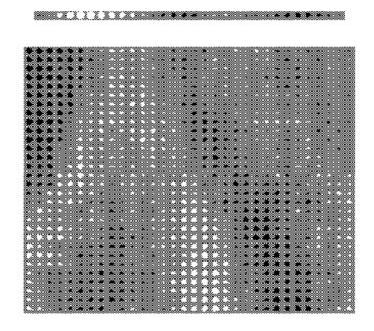
Two layers of logistic units

Highly non-linear decision surface

Neural Network trained to drive a car!



S harp Left Straight Ahead Shaip Right 30 Output Units 4 Hidden Units 30x32 Sensor Input Retina Weights to output units from one hidden unit



Weights of each pixel for one hidden unit

Connectionist Models

Consider humans:

- Neuron switching time ~ .001 second
- Number of neurons ~ 10¹⁰
- Connections per neuron $\sim 10^{4-5}$
- Scene recognition time ~ .1 second
- 100 inference steps doesn't seem like enough
- \rightarrow much parallel computation

Properties of artificial neural nets (ANN's):

- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed process

Prediction using Neural Networks

Prediction – Given neural network (hidden units and weights), use it to predict the label of a test point

Forward Propagation -

Start from input layer

For each subsequent layer, compute output of sigmoid unit

Sigmoid unit:

$$o(\mathbf{x}) = \sigma(w_0 + \sum_i w_i x_i)$$

$$o(\mathbf{x}) = \sigma \left(w_0 + \sum_h w_h \sigma(w_0^h + \sum_i w_i^h x_i) \right)$$

Training Neural Networks – 12 loss

$$W \leftarrow \arg\min_{W} E[W]$$

$$W \leftarrow \arg\min_{W} \sum_{l} (y^{l} - \hat{f}(x^{l}))^{2}$$

Learned neural network

Where $\widehat{f}(x^l) = o(x^l)$, output of neural network for training point \mathbf{x}^{l}

Train weights of all units to minimize sum of squared errors of predicted network outputs

Minimize using Gradient Descent

For Neural Networks, *E[w]* no longer convex in w

Gradient

$$\nabla E[\vec{w}] \equiv \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$$

Training rule:

$$\Delta \vec{w} = -\eta \nabla E[\vec{w}]$$

i.e.,

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

Incremental (Stochastic) Gradient Descent

Batch mode Gradient Descent:

Do until satisfied

- 1. Compute the gradient $\nabla E_D[\vec{w}]$ Using all training data D
- $2. \vec{w} \leftarrow \vec{w} \eta \nabla E_D[\vec{w}]$ $E_D[\vec{w}] \equiv \frac{1}{2} \sum_{\mathbf{l} \in D} (\mathbf{y}^{\mathbf{l}} o^{\mathbf{l}})^2$

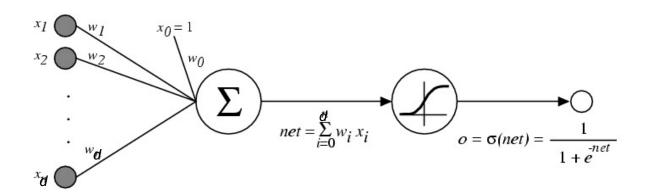
Incremental mode Gradient Descent:

Do until satisfied

- \bullet For each training example | in D
 - 1. Compute the gradient $\nabla E_{\vec{i}}[\vec{w}]$
 - $2. \ \vec{w} \leftarrow \vec{w} \eta \nabla E_{\parallel}[\vec{w}]$ $E_{\parallel}[\vec{w}] \equiv \frac{1}{2} (\mathbf{y}^{\parallel} o^{\parallel})^2$

Incremental Gradient Descent can approximate Batch Gradient Descent arbitrarily closely if η made small enough

Training Neural Networks



 $\sigma(x)$ is the sigmoid function

$$\frac{1}{1 + e^{-x}}$$

Nice property:
$$\frac{d\sigma(x)}{dx} =$$

Differentiable

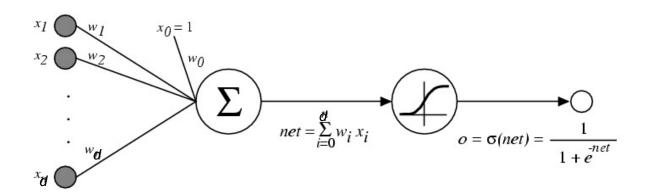
A.
$$\sigma(x)(1-\sigma(x))$$

C.
$$-\sigma(x)$$

B.
$$\sigma(x) \sigma(-x)$$

D.
$$\sigma(x)^2$$

Training Neural Networks



 $\sigma(x)$ is the sigmoid function

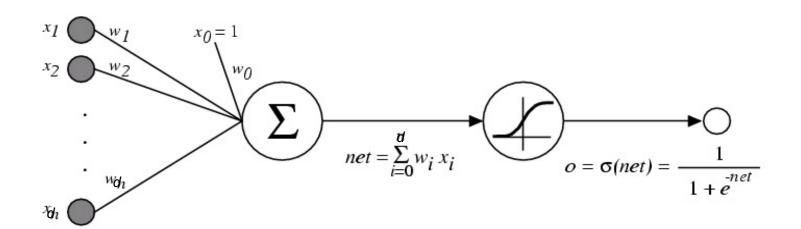
$$\frac{1}{1+e^{-x}}$$

Nice property:
$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$$
 Differentiable

We can derive gradient decent rules to train

- One sigmoid unit
- $Multilayer\ networks$ of sigmoid units \rightarrow Backpropagation

Gradient Descent for 1 sigmoid unit



$$\frac{\partial E}{\partial w_i} \, = \, \frac{\partial}{\partial w_i} \, \frac{1}{2} \sum_{\mathsf{I} \in D} (\mathsf{y}^{\mathsf{I}} - o^{\mathsf{I}})^2 \ = \ \sum_{\mathsf{I}} (\mathsf{y}^{\mathsf{I}} - o^{\mathsf{I}}) \left(-\frac{\partial o^{\mathsf{I}}}{\partial w_i} \right)^2$$

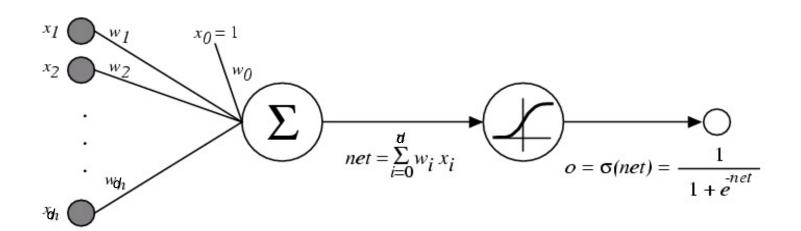
Gradient of the sigmoid function output wrt its input

$$\frac{\partial \sigma(net)}{\partial net} = \sigma(net)(1 - \sigma(net)) = o(1 - o)$$

Gradient of the sigmoid unit output wrt input weights

$$\frac{\partial o}{\partial w_i} =$$

Gradient Descent for 1 sigmoid unit



$$\frac{\partial E}{\partial w_i} \, = \, \frac{\partial}{\partial w_i} \, \frac{1}{2} \sum_{\mathsf{I} \in D} (\mathsf{y}^{\mathsf{I}} - o^{\mathsf{I}})^2 \ = \ \sum_{\mathsf{I}} (\mathsf{y}^{\mathsf{I}} - o^{\mathsf{I}}) \left(-\frac{\partial o^{\mathsf{I}}}{\partial w_i} \right)^2$$

Gradient of the sigmoid function output wrt its input

$$\frac{\partial \sigma(net)}{\partial net} = \sigma(net)(1 - \sigma(net)) = o(1 - o)$$

Gradient of the sigmoid unit output wrt input weights

$$\frac{\partial o}{\partial w_i} = \frac{\partial o}{\partial net} \cdot \frac{\partial net}{\partial w_i} = o(1 - o)x_i$$

Gradient descent for training NNs

Gradient descent via Chain rule for computing gradients

= Back-propagation algorithm for training NNs

Backpropagation Algorithm (MLE) using Stochastic gradient descent

1 final output unit

Initialize all weights to small random numbers. Until satisfied, Do

- For each training example, Do
 - 1. Input the training example to the network and compute the network outputs

2.

$$\delta \leftarrow o(1-o)(y-o)$$

3.

4. Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where

$$\Delta w_{i,j} = \eta \delta_j o_i$$

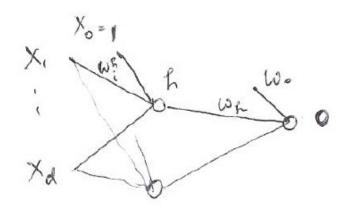
Using Forward propagation

y = label of current training example

 w_{ij} = wt from i to j

Note: if i is input variable, $o_i = x_i$

Gradient Descent for 1 hidden layer 1 output NN



$$\omega_{k}$$
 ω_{k}
 ω_{k

$$\frac{\partial E}{\partial w_i} \, = \, \frac{\partial}{\partial w_i} \, \frac{1}{2} \sum_{\mathbf{l} \in D} (\mathbf{y}^{\mathbf{l}} - o^{\mathbf{l}})^2 \ = \ \sum_{\mathbf{l}} (\mathbf{y}^{\mathbf{l}} - o^{\mathbf{l}}) \left(-\frac{\partial o^{\mathbf{l}}}{\partial w_i} \right)^2$$

Gradient of the output with respect to w_h

$$\frac{\partial o}{\partial w_h} = o(1 - o)o_h$$

Gradient Descent for 1 hidden layer 1 output NN

$$0 = \sigma(\omega_0 + Z\omega_R O_R) = \sigma(Z\omega_R O_R)$$

$$0_R = \sigma(\omega_0^R + Z\omega_I^R X_i) = \sigma(Z\omega_I^R X_i)$$

$$\frac{\partial E}{\partial w_i} \, = \, \frac{\partial}{\partial w_i} \, \frac{1}{2} \sum_{\mathbf{l} \in D} (\mathbf{y}^{\mathbf{l}} - o^{\mathbf{l}})^2 \ = \ \sum_{\mathbf{l}} (\mathbf{y}^{\mathbf{l}} - o^{\mathbf{l}}) \left(-\frac{\partial o^{\mathbf{l}}}{\partial w_i} \right)^2$$

Gradient of the output with respect to input weights whi

$$\frac{\partial o}{\partial w_i^h} = \frac{\partial o}{\partial o_h} \cdot \frac{\partial o_h}{\partial w_i^h}$$
$$= o(1 - o)o_h(1 - o_h)w_h x_i$$

1 final output unit

Initialize all weights to small random numbers. Until satisfied, Do

- For each training example, Do
 - 1. Input the training example to the network and compute the network outputs

2.

$$\delta \leftarrow o(1-o)(y-o)$$

3. For each hidden unit h

$$\delta_h \leftarrow o_h (1 - o_h) w_h \delta$$

4. Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where

$$\Delta w_{i,j} = \eta \delta_j o_i$$

y = label of current

training example

Using Forward propagation

o_(h) = unit output (obtained by forward propagation)

 w_{ij} = wt from i to j

Note: if i is input variable, $o_i = x_i$

Backpropagation Algorithm (MLE) using Stochastic gradient descent

Initialize all weights to small random numbers. Until satisfied, Do

- For each training example, Do
 - 1. Input the training example to the network and compute the network outputs
 - 2. For each output unit k

$$\delta_k \leftarrow o_k(1-o_k)(y_k-o_k)$$

3. For each hidden unit h

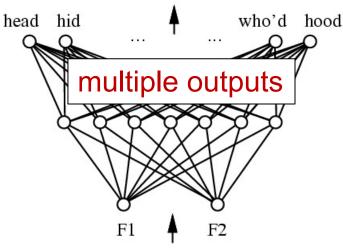
$$\delta_h \leftarrow o_h(1 - o_h) \sum_{k \in outputs} w_{h,k} \delta_k$$

4. Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j}$$

where

$$\Delta w_{i,j} = \eta \delta_j o_i$$



Using Forward propagation

y_k = label of current training example for output unit k

 $o_{k(h)}$ = unit output (obtained by forward propagation)

 w_{ij} = wt from i to j

Note: if i is input variable, $o_i = x_i$

HW2

Classification – cross-entropy error metric

Can implement backpropagation with matrix-vector products – uses matrix-vector calculus heavily

Note on numerator vs denominator layout

More on Backpropagation

- Gradient descent over entire *network* weight vector
- Easily generalized to arbitrary directed graphs
- Will find a local, not necessarily global error minimum
 - In practice, often works well (can run multiple times)

- Minimizes error over *training* examples
 - Will it generalize well to subsequent examples?
- Training can take thousands of iterations → slow!
- Using network after training is very fast

Objective/Error no longer convex in weights