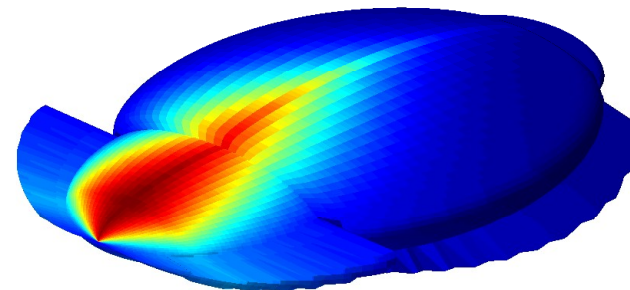
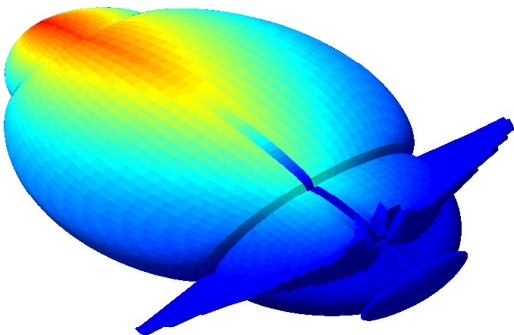


Lecture 8 Neural Classifiers



Using Backpropagation for Classification Problems

Backpropagation requires that $\frac{\partial J^q}{\partial w_{jk}}$ exists $\forall w_{ij}$

$\Rightarrow f'_i(s_i^q)$ must exist

$\Rightarrow f(u)$ (activation function) must be differentiable

$$\rightarrow \text{use } f(u) = \frac{1}{1+e^{-u}}$$
$$\text{or } f(u) = \tanh(u)$$

instead of hard threshold.

But classification requires that output neurons have 0/1 or +1/-1 output

??????????

Possible Solutions

- Use 0/1 or +1/-1 as target values y_i^q
- Set operating parameters $H, L, L < H$ such that:

$\hat{y}_i \geq H$ is considered a match for $y_i = +1$

$\hat{y}_i \leq L$ is considered a match for $y_i = 0$ or -1

- During training, use 0/1 or +1/-1 as targets but if \hat{y}_i matches the operating targets, make no weight change

Typical values for H, L are:

$H=0.75$	$L=0.25$	for 0-1 sigmoid
$H=0.75$	$L=-0.75$	for +/- sigmoid

Multi-Class Classification

If there are more than 2 classes, how should the output be represented?

Solution:

- If there are M classes, have M neurons in the output layer.
- During training, use one-hot codes as targets for each class, e.g. if $M = 3$,
Class 1 = [1 0 0] Class 2 = [0 1 0] Class 3 = [0 0 1]
- After training, for a test input x , choose the class as the output neuron with the highest output, e.g., Output = [0.1 0.6 0.2] \rightarrow [0 1 0]
- If “no choice” is allowed, require that an output must be higher than some threshold θ to be considered 1.

A better solution: softmax classifiers

Logistic Regression

Consider a case where a binary outcome (yes/no, true/false, 0/1) depends on a continuous variable (feature) x

For example:

Is patient A sick or well given their temperature?

Will borrower B pay back a loan given their income?

Will person C live until next year given their age?

The general form is:

Is F true or false given the value of x ? $F : x \rightarrow \{0,1\}$

Often, the dependence of F on x is:

- Monotonic
- Uncertain

In this situation, we can ask the question in terms of probabilities:

What is the probability that $F = 1$ given x ?

For a more detailed intro, see <http://ufldl.stanford.edu/tutorial/supervised/LogisticRegression/>

$$y \equiv P(F = 1 | x) = p(x; \mathbf{w})$$

Where

- $p(x; \mathbf{w}) \equiv$ probability estimate model
- $\mathbf{w} \equiv$ parameters of model p

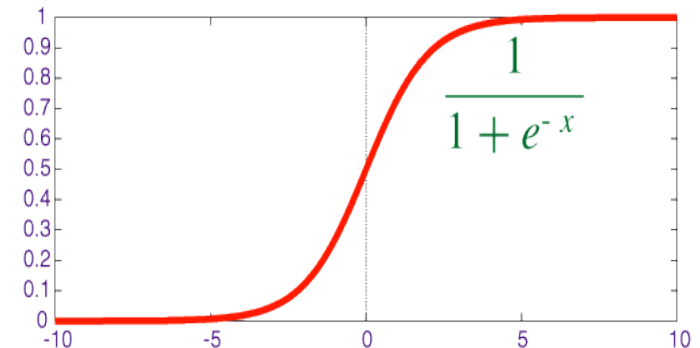
We need to tune \mathbf{w} based on data to obtain the optimal p .

Logistic regression fits a linear model by postulating that:

$$\underbrace{\log \frac{p(x; \mathbf{w})}{1 - p(x; \mathbf{w})}}_{\text{log odds ratio}} = w_1 x + w_0$$

$$\mathbf{w} = \{w_0, w_1\}$$

$$p(x; \mathbf{w}) = \frac{e^{w_1 x + w_0}}{1 + e^{w_1 x + w_0}} = \frac{1}{1 + e^{-(w_1 x + w_0)}}$$



If there is more than one feature, $\{x_1, \dots, x_n\}$, we get:

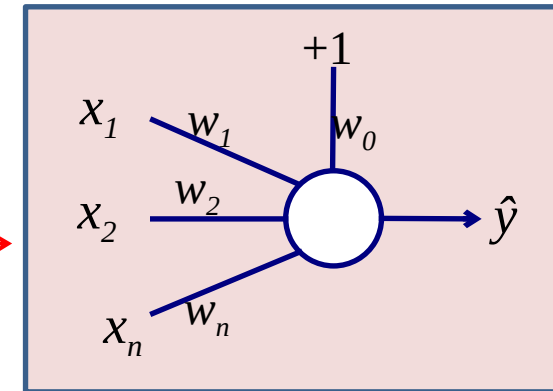
$$\log \frac{p(\mathbf{x}; \mathbf{w})}{1 - p(\mathbf{x}; \mathbf{w})} = \mathbf{w} \cdot \mathbf{x}$$

$$\mathbf{x} = \{1, x_1, \dots, x_n\}$$

$$\mathbf{w} = \{w_0, w_1, \dots, w_n\}$$

$$p(\mathbf{x}; \mathbf{w}) = \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} = \sigma(\mathbf{w} \cdot \mathbf{x}) \equiv \hat{y}$$

Probability
Model



which is the activation function of a sigmoid neuron.

Thus, the output of a sigmoid output neuron in a 2-class classification problem can be seen as a logistic estimate of the probability of Class 1 given the current input: $y = P(F = 1 | \mathbf{x})$

Think of this as the probabilistic version of linear separability.

If $P(F = 1 | \mathbf{x})$ is not a monotonic (logistic) function of \mathbf{x} , the purpose of the hidden layer is to map \mathbf{x} to a new representation, \mathbf{h} , such that $P(F = 1 | \mathbf{h})$ is a logistic function of \mathbf{h} .

In fact, logistic regression can be seen as minimizing the **cross-entropy loss function**:

$$J(\mathbf{w}) = - \left[\sum_{q=1}^M \left(y^q \log(p(\mathbf{x}^q; \mathbf{w})) + (1 - y^q) \log(1 - p(\mathbf{x}^q; \mathbf{w})) \right) \right]$$

$M \equiv$ number of data points

Which gives the gradient vector:

$$\nabla J(\mathbf{w}) = \sum_{q=1}^M (p(\mathbf{x}^q; \mathbf{w}) - y^q) \mathbf{x}^q = \sum_{q=1}^M (\hat{y}^q - y^q) \mathbf{x}^q$$

Using steepest descent results in the learning rule:

$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{w}) = \eta \sum_{q=1}^M (y^q - p(\mathbf{x}^q; \mathbf{w})) \mathbf{x}^q = \eta \sum_{q=1}^M (y^q - \hat{y}^q) \mathbf{x}^q$$

which is the perceptron learning rule, but with a sigmoid neuron.

In fact, logistic regression can be seen as minimizing the **cross-entropy loss function**:

$$J(\mathbf{w}) = - \left[\sum_{q=1}^M \left(\boxed{y^q} \log(p(\mathbf{x}^q; \mathbf{w})) + \boxed{(1 - y^q)} \log(1 - p(\mathbf{x}^q; \mathbf{w})) \right) \right]$$

$\boxed{y^q} = 1 \text{ if class of } \mathbf{x}^q = 1$
 $\boxed{(1 - y^q)} = 1 \text{ if class of } \mathbf{x}^q = 0$

Which gives the gradient vector:

$$\nabla J(\mathbf{w}) = \sum_{q=1}^M (p(\mathbf{x}^q; \mathbf{w}) - y^q) \mathbf{x}^q = \sum_{q=1}^M (\hat{y}^q - y^q) \mathbf{x}^q$$

Using steepest descent results in the learning rule:

$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{w}) = \eta \sum_{q=1}^M (y^q - p(\mathbf{x}^q; \mathbf{w})) \mathbf{x}^q = \eta \sum_{q=1}^M \boxed{(y^q - \hat{y}^q)} \mathbf{x}^q$$

$\boxed{(y^q - \hat{y}^q)}$ is a **Scalar**

which is the perceptron learning rule, but with a sigmoid neuron.

Batch mode

To derive $\nabla J(\mathbf{w}) = \sum_{q=1}^M (p(\mathbf{x}^q; \mathbf{w}) - y^q) \mathbf{x}^q$ consider

$$\begin{aligned} J^q &= -y^q \log(p(\mathbf{x}^q; \mathbf{w})) - (1 - y^q) \log(1 - p(\mathbf{x}^q; \mathbf{w})) \\ &= -y^q \log(\sigma(\mathbf{w} \cdot \mathbf{x}^q)) - (1 - y^q) \log(1 - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) \\ &= -y^q \log(\sigma(\mathbf{w} \cdot \mathbf{x}^q)) - (1 - y^q) \log(\sigma(-\mathbf{w} \cdot \mathbf{x}^q)) \\ &= \begin{cases} \log(1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)) & \text{if } y^q = 1 \\ \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)) & \text{if } y^q = 0 \end{cases} \end{aligned}$$

since

$$\log(\sigma(\mathbf{w} \cdot \mathbf{x}^q)) = \log\left(\frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}^q}}\right) = \log(1) - \log(1 + e^{-\mathbf{w} \cdot \mathbf{x}^q}) = -\log(1 + e^{-\mathbf{w} \cdot \mathbf{x}^q})$$

and

$$\log(\sigma(-\mathbf{w} \cdot \mathbf{x}^q)) = \log\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}^q}}\right) = \log(1) - \log(1 + e^{\mathbf{w} \cdot \mathbf{x}^q}) = -\log(1 + e^{\mathbf{w} \cdot \mathbf{x}^q})$$

$$J^q = \begin{cases} \log(1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)) & \text{if } y^q = 1 \\ \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)) & \text{if } y^q = 0 \end{cases}$$

$$\Rightarrow \frac{\partial J^q}{\partial \mathbf{w}_j} = \begin{cases} \frac{\partial}{\partial \mathbf{w}_j} \log(1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)) & \text{if } y^q = 1 \\ \frac{\partial}{\partial \mathbf{w}_j} \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)) & \text{if } y^q = 0 \end{cases}$$

For $y = 1$

$$\begin{aligned}
 \frac{\partial}{\partial w_j} \log(1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)) &= \left(\frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) \frac{\partial}{\partial w_j} (\exp(-\mathbf{w} \cdot \mathbf{x}^q)) \\
 &= \left(\frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) \left(\exp(-\mathbf{w} \cdot \mathbf{x}^q) \frac{\partial}{\partial w_j} (-\mathbf{w} \cdot \mathbf{x}^q) \right) \\
 &= \left(\frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) (\exp(-\mathbf{w} \cdot \mathbf{x}^q) (-x_j^q)) = - \left(\frac{\exp(-\mathbf{w} \cdot \mathbf{x}^q)}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) x_j^q \\
 &= -(1 - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q = -(y - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q
 \end{aligned}$$

Similarly, for $y = 0$

$$\begin{aligned}
 \frac{\partial}{\partial w_j} \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)) &= \left(\frac{\exp(\mathbf{w} \cdot \mathbf{x}^q)}{1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)} \right) x_j^q = (0 + \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q \\
 &= (0 - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q = -(y - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q
 \end{aligned}$$

For $y = 1$

$$\begin{aligned}
 \frac{\partial}{\partial w_j} \log(1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)) &= \left(\frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) \frac{\partial}{\partial w_j} (\exp(-\mathbf{w} \cdot \mathbf{x}^q)) \\
 &= \left(\frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) \left(\exp(-\mathbf{w} \cdot \mathbf{x}^q) \frac{\partial}{\partial w_j} (-\mathbf{w} \cdot \mathbf{x}^q) \right) \\
 &= \left(\frac{1}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) (\exp(-\mathbf{w} \cdot \mathbf{x}^q) (-x_j^q)) = - \left(\frac{\exp(-\mathbf{w} \cdot \mathbf{x}^q)}{1 + \exp(-\mathbf{w} \cdot \mathbf{x}^q)} \right) x_j^q \\
 &= -(1 - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q = \boxed{-(y - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q}
 \end{aligned}$$

Similarly, for $y = 0$

$$\begin{aligned}
 \frac{\partial}{\partial w_j} \log(1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)) &= \left(\frac{\exp(\mathbf{w} \cdot \mathbf{x}^q)}{1 + \exp(\mathbf{w} \cdot \mathbf{x}^q)} \right) x_j^q = (0 + \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q \\
 &= (0 - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q = \boxed{-(y - \sigma(\mathbf{w} \cdot \mathbf{x}^q)) x_j^q}
 \end{aligned}$$

same

→ $\nabla J^q(\mathbf{w}) = \left[\frac{\partial J^q}{\partial w_1} \frac{\partial J^q}{\partial w_2} \cdots \frac{\partial J^q}{\partial w_n} \right]^T = (\sigma(\mathbf{w} \cdot \mathbf{x}^q) - y) \mathbf{x}^q$ Note $(\sigma(\mathbf{w} \cdot \mathbf{x}^q) - y)$ is scalar

$$\nabla J^q = \left[\frac{\partial J^q}{\partial w_1} \frac{\partial J^q}{\partial w_2} \dots \frac{\partial J^q}{\partial w_n} \right]^T = (\sigma(\mathbf{w} \cdot \mathbf{x}^q) - y) \mathbf{x}^q$$

$$\hookrightarrow \nabla J(\mathbf{w}) = \sum_{q=1}^M \nabla J^q(\mathbf{w}) = \sum_{q=1}^M (\hat{y}^q - y^q) \mathbf{x}^q$$

$$\text{Since } p(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w} \cdot \mathbf{x}) \equiv \hat{y}$$

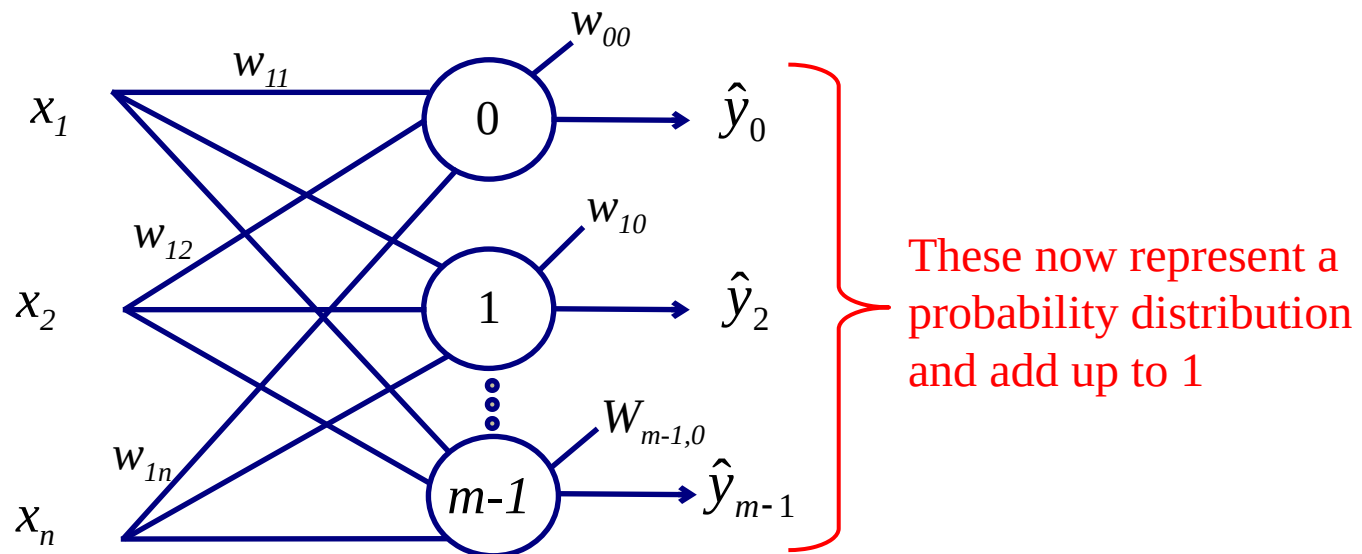
Softmax Classifier

The softmax classifier is the multi-class version of the logistic classifier.

We have $F : x \rightarrow \{0, 1, \dots, m - 1\}$

Since F can now take m different values (instead of two), we want to estimate:
 $P(F = 0 \mid \mathbf{x}), P(F = 1 \mid \mathbf{x}), \dots, P(F = m-1 \mid \mathbf{x})$

This can be done using an output layer of m neurons with $\hat{y}_i = P(F = i \mid \mathbf{x})$



Hypothesize that:

$$P(F = i | \mathbf{x}; \mathbf{w}) = \frac{e^{\mathbf{w}_i \cdot \mathbf{x}}}{\sum_{r=0}^{m-1} e^{\mathbf{w}_r \cdot \mathbf{x}}} \quad \leftarrow \text{Softmax function}$$

where

- \mathbf{w}_i is the weight vector of the output neuron for class i
- \mathbf{w} is the weight matrix of the output layer

Then

$$p(\mathbf{x}; \mathbf{w}) = \begin{bmatrix} P(F = 0 | \mathbf{x}; \mathbf{w}) \\ P(F = 1 | \mathbf{x}; \mathbf{w}) \\ \dots \\ P(F = m-1 | \mathbf{x}; \mathbf{w}) \end{bmatrix} = \frac{1}{\sum_{r=0}^{m-1} e^{\mathbf{w}_r \cdot \mathbf{x}}} \begin{bmatrix} e^{\mathbf{w}_0 \cdot \mathbf{x}} \\ e^{\mathbf{w}_1 \cdot \mathbf{x}} \\ \dots \\ e^{\mathbf{w}_{m-1} \cdot \mathbf{x}} \end{bmatrix}$$

The weights can be learned by optimizing the cross-entropy loss function:

$$J(\mathbf{w}) = - \left[\sum_{q=1}^M \sum_{i=0}^{m-1} 1\{y^q = i\} \log \frac{e^{\mathbf{w}_i \cdot \mathbf{x}^q}}{\sum_{r=0}^{m-1} e^{\mathbf{w}_r \cdot \mathbf{x}^q}} \right]$$

\mathbf{x}^q = q th input vector
 y^q = class label of \mathbf{x}^q
 $1\{u\} = 1$ if u is true, else 0

The gradient vector for the weight vector \mathbf{w}_i is:

$$\nabla J_i^q(\mathbf{w}) = - \sum_{q=1}^M \left[\left(1\{y^q = i\} - \frac{e^{\mathbf{w}_i \cdot \mathbf{x}^q}}{\sum_{r=0}^{m-1} e^{\mathbf{w}_r \cdot \mathbf{x}^q}} \right) \mathbf{x}^q \right]$$

The j th component of this can be used to adjust weight w_{ij}

The weights can be learned by optimizing the cross-entropy loss function:

$$J(\mathbf{w}) = - \left[\sum_{q=1}^M \sum_{i=0}^{m-1} 1\{y^q = i\} \log \frac{e^{\mathbf{w}_i \cdot \mathbf{x}^q}}{\sum_{r=0}^{m-1} e^{\mathbf{w}_r \cdot \mathbf{x}^q}} \right]$$

\mathbf{x}^q = q th input vector
 y^q = class label of \mathbf{x}^q
 $1\{u\} = 1$ if u is true, else 0

The gradient vector for the weight vector \mathbf{w}_i is:

$$\nabla J_i^q(\mathbf{w}) = - \sum_{q=1}^M \left[\left(1\{y^q = i\} - \frac{e^{\mathbf{w}_i \cdot \mathbf{x}^q}}{\sum_{r=0}^{m-1} e^{\mathbf{w}_r \cdot \mathbf{x}^q}} \right) \mathbf{x}^q \right]$$

scalar

The j th component of this can be used to adjust weight w_{ij}