

CIS603 - Artificial Intelligence

Logistic regression

Vasileios Megalooikonomou

(some material adopted from notes by M. Hauskrecht)

CIS603 - AI

Supervised learning

Data: $D = \{d_1, d_2, \dots, d_n\}$ a set of n examples

$$d_i = \langle \mathbf{x}_i, y_i \rangle$$

\mathbf{x}_i is input vector, and y is desired output (given by a teacher)

Objective: learn the mapping $f : X \rightarrow Y$

$$\text{s.t. } y_i \approx f(x_i) \quad \text{for all } i = 1, \dots, n$$

Two types of problems:

- **Regression:** Y is **continuous**
Example: earnings, product orders \rightarrow company stock price
- **Classification:** Y is **discrete**
Example: temperature, heart rate \rightarrow disease

Now: **BINARY classification problems**

CIS603 - AI

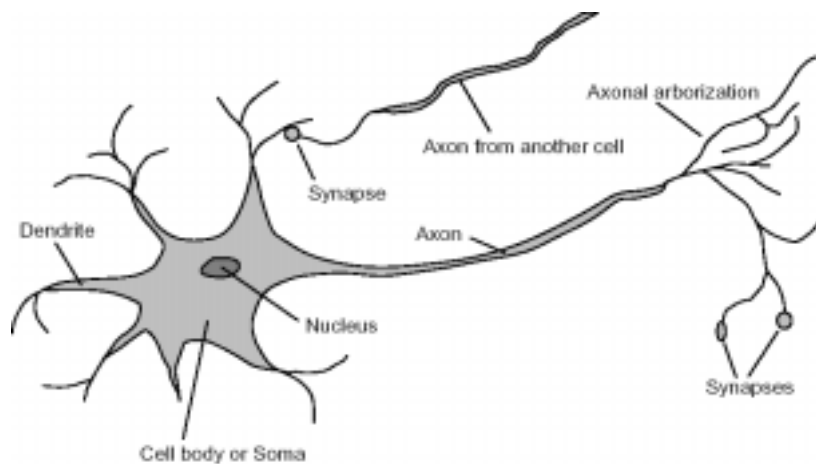
Binary classification

- Two classes $Y = \{0,1\}$
- Our goal is to learn to classify correctly two types of examples
 - Class 0 – labeled as 0,
 - Class 1 – labeled as 1
- We would like to learn $f : X \rightarrow \{0,1\}$
- **First step:** we need to devise a model of the function f
- **Inspiration:** *neuron (nerve cells)*

CIS603 - AI

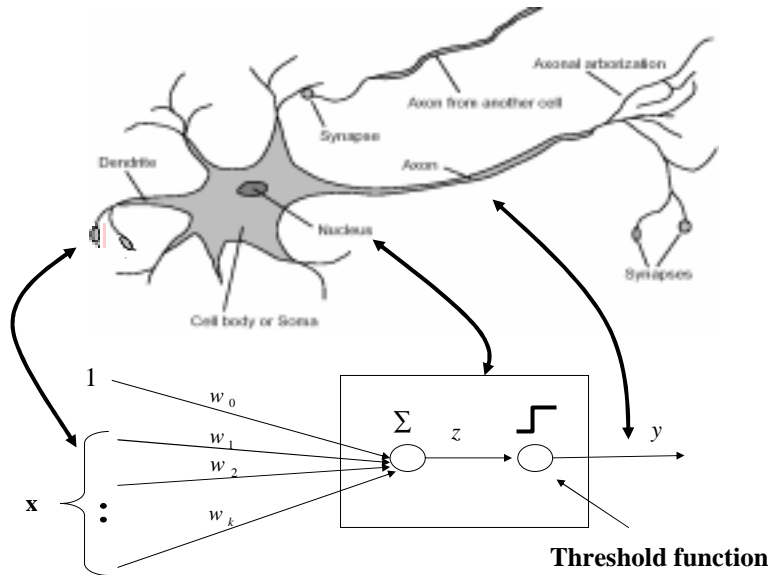
Neuron

- **neuron (nerve cell) and its activities**



CIS603 - AI

Neuron-based binary classification model



CIS603 - AI

Binary classification

- Instead of learning the mapping to discrete values 0,1
 $f : X \rightarrow \{0,1\}$
- It is easier to learn a probabilistic function
 $f' : X \rightarrow [0,1]$
 - where f' describes the probability of a class 1 given x
 $p(y = 1 | \mathbf{x})$
- Transformation to discrete class values:

If $p(y = 1 | \mathbf{x}) \geq 1/2$ then choose **1**
Else choose **0**

- **Logistic regression model** uses a probabilistic function

CIS603 - AI

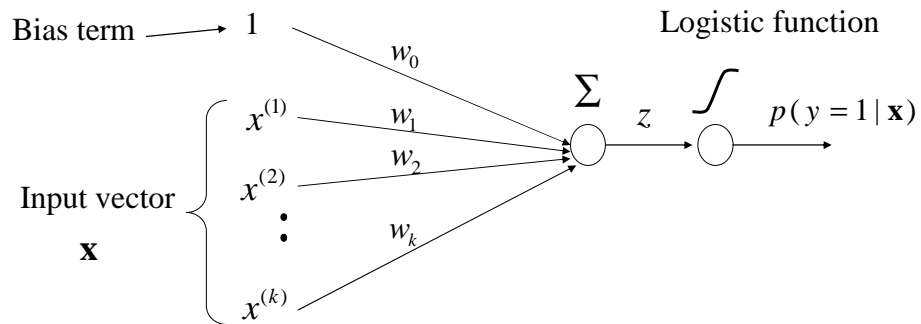
Logistic regression

- **Logistic regression:**

$$p(y = 1 | \mathbf{x}, \mathbf{w}) = g(z) = g(w_0 + w_1 x^{(1)} + \dots + w_k x^{(k)})$$

where \mathbf{w} are parameters of the models

and $g(z)$ is a **logistic function** $g(z) = 1/(1 + e^{-z})$

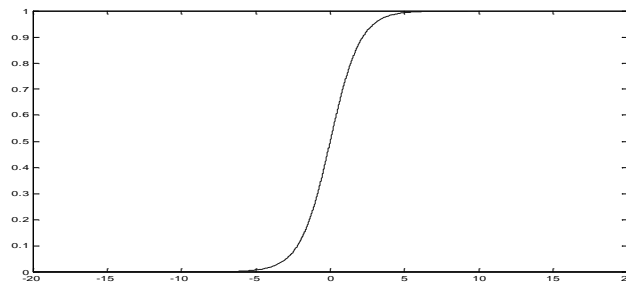


CIS603 - AI

Logistic function

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

- also referred to as sigmoid function
- replaces threshold function with smooth switching
- takes a real number and outputs the number in the interval $[0,1]$

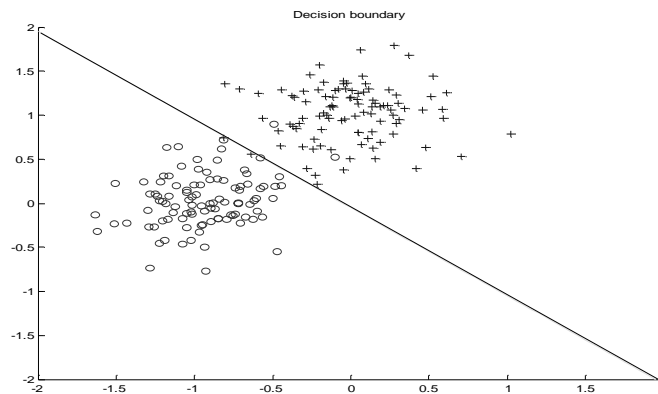


CIS603 - AI

Logistic regression - Decision boundary

Logistic regression model defines a linear decision boundary

- Example: 2 classes (crosses and circles)



CIS603 - AI

Binary classification - Error

- Two classes $Y = \{0,1\}$
- Our goal is to classify correctly as many examples as possible
- Zero-one error function

$$Error(x_i, y_i) = \begin{cases} 1 & f(\mathbf{x}_i, \mathbf{w}) \neq y_i \\ 0 & f(\mathbf{x}_i, \mathbf{w}) = y_i \end{cases}$$

- Error we would like to minimize: $E_{(x,y)}(Error(x, y))$
- The error is minimized if we choose:
 $y = 1$ if $p(y = 1 | \mathbf{x}, \mathbf{w}) > p(y = 0 | \mathbf{x}, \mathbf{w})$
 $y = 0$ otherwise
- We construct a probabilistic version of the error function based on the **likelihood of the data** $L(D, \mathbf{w}) = P(D | \mathbf{w})$

Inverse optimization problem $Error(D, \mathbf{w}) = -L(D, \mathbf{w})$

CIS603 - AI

Logistic regression: parameter learning

- **Likelihood of data** $L(D, \mathbf{w}) = P(D | \mathbf{w}) = \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{1-y_i}$

where $d_i = \langle \mathbf{x}_i, y_i \rangle$

$$\mu_i = p(y_i = 1 | \mathbf{x}_i, \mathbf{w}) = g(z_i) = g(w_0 + \sum_{j=1}^k w_j x_i^{(j)})$$

We want weights \mathbf{w} that maximize the likelihood of data

- **Trick: maximize the log-likelihood of data instead**

$$\begin{aligned} l(D, \mathbf{w}) &= \log \prod_{i=1}^n \mu_i^{y_i} (1 - \mu_i)^{1-y_i} = \sum_{i=1}^n \log \mu_i^{y_i} (1 - \mu_i)^{1-y_i} = \\ &= \sum_{i=1}^n y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i) = \sum_{i=1}^n -J_{\text{online}}(d_i, \mathbf{w}) \end{aligned}$$

- Rational: The optimal weights are the same for both the likelihood and the log-likelihood

CIS603 - AI

Logistic regression: parameter estimation

- **Log likelihood**

$$l(D, \mathbf{w}) = \sum_{i=1}^n -J_{\text{online}}(d_i, \mathbf{w}) = \sum_{i=1}^n y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)$$

- **On-line component of the log-likelihood**

$$-J_{\text{online}}(d_i, \mathbf{w}) = y_i \log \mu_i + (1 - y_i) \log(1 - \mu_i)$$

- **Derivatives of the online error component (in terms of weights)**

$$\frac{\partial}{\partial w_0} J_{\text{online}}(d_i, \mathbf{w}) = -(y_i - g(z_i))$$

⋮

$$\frac{\partial}{\partial w_j} J_{\text{online}}(d_i, \mathbf{w}) = -x_i^{(j)} (y_i - g(z_i))$$

CIS603 - AI

Logistic regression. Online gradient.

- We want to find the set of parameters optimizing the log-likelihood of data or minimizing the error
- **On-line learning update for weight w** $J_{online}(d_i, \mathbf{w})$

$$w \leftarrow w^* - \alpha \frac{\partial}{\partial w} [J_{online}(d_i, w) |_{w^*}]$$

- **(i+1)th update for the logistic regression** and $d = \langle \mathbf{x}, y \rangle$

$$w_0^{(i+1)} \leftarrow w_0^{(i)} + \alpha(i+1)(y - g(w_0 + \sum_{u=1}^k w_u x^{(u)}))$$

•

$$w_j^{(i+1)} \leftarrow w_j^{(i)} + \alpha(i+1)(y - g(w_0 + \sum_{u=1}^k w_u x^{(u)}))x^{(j)}$$

α - annealed learning rate (depends on the number of updates)

The same, easy update rule as used in the linear regression !!!

CIS603 - AI

Online logistic regression algorithm

Online-logistic-regression (D , number of iterations)

initialize weights $w_0, w_1, w_2 \dots w_k$

for $i=1:1$: number of iterations

do **select** a data point $d = \langle \mathbf{x}, y \rangle$ from D

set $\alpha = 1/i$

update weights (in parallel)

$$w_0 = w_0 + \alpha[y - p(y=1 | \mathbf{x}, \mathbf{w})]$$

•

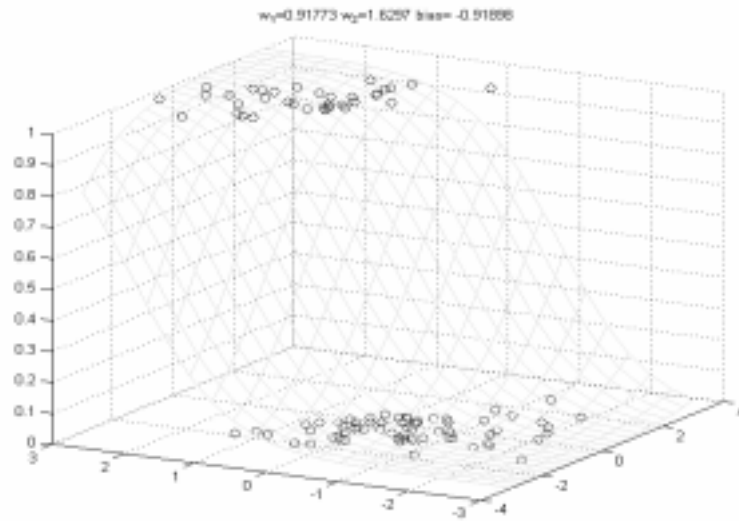
$$w_j = w_j + \alpha[y - p(y=1 | \mathbf{x}, \mathbf{w})]x^{(j)}$$

end for

return weights

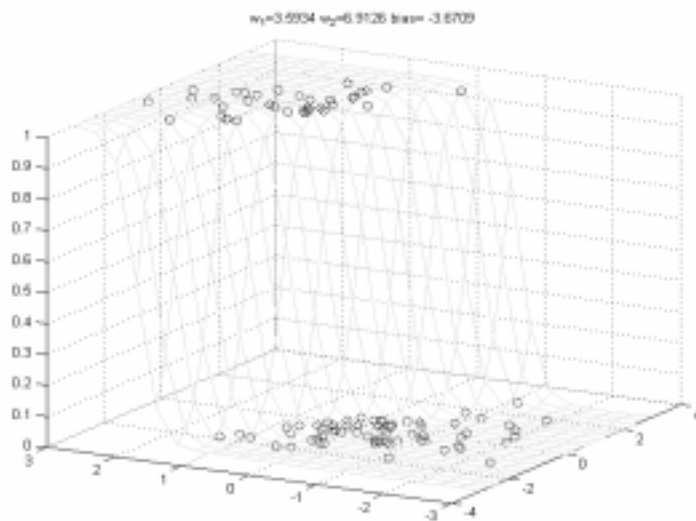
CIS603 - AI

Online algorithm. Example.



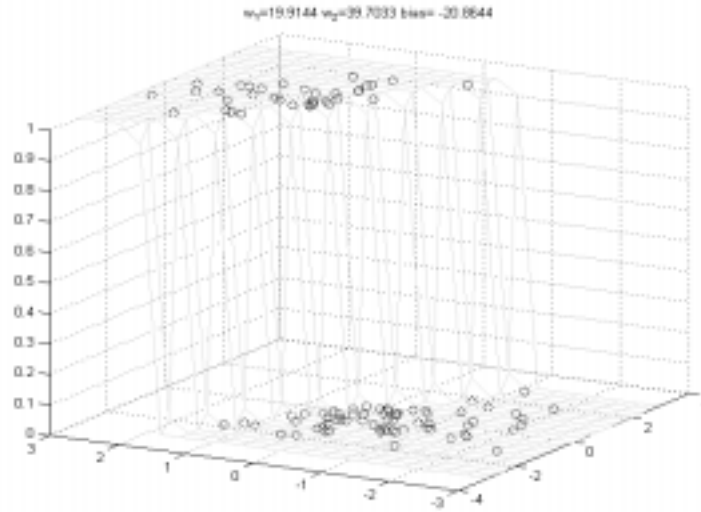
CIS603 - AI

Online algorithm. Example.



CIS603 - AI

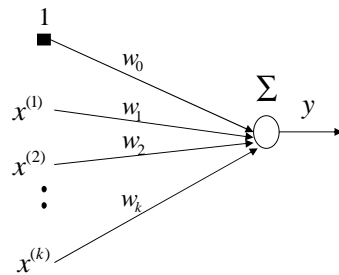
Online algorithm. Example.



CIS603 - AI

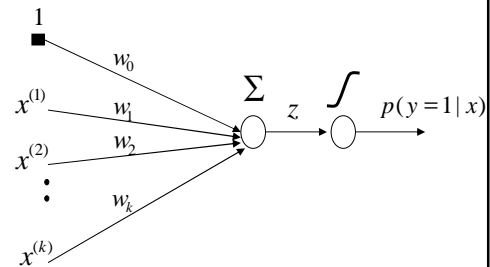
Limitations of basic linear units

Linear regression



Function linear in inputs

Logistic regression



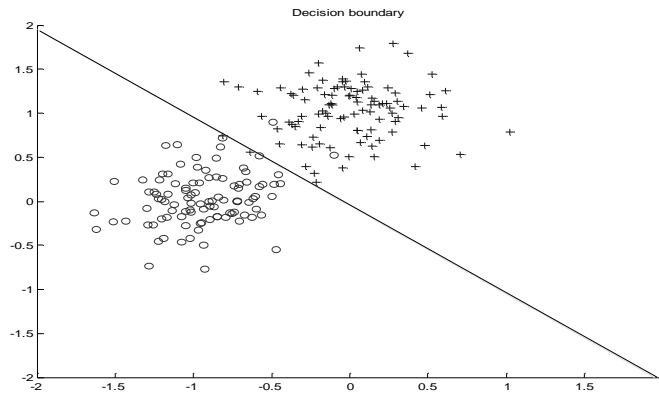
Linear decision boundary

CIS603 - AI

Logistic regression - Decision boundary

Logistic regression model defines a linear decision boundary

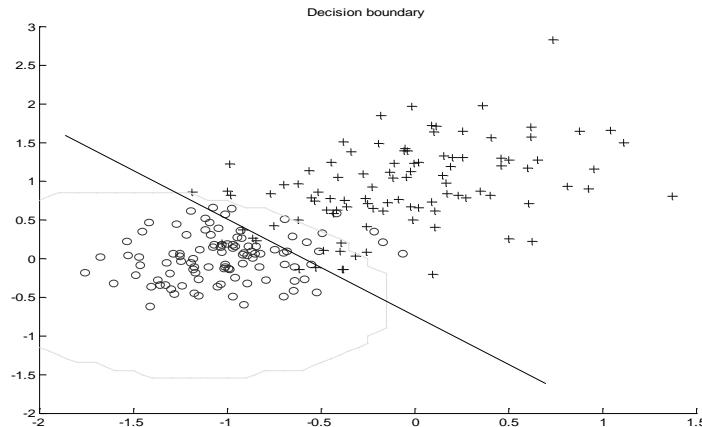
- Example: 2 classes (crosses and circles)



CIS603 - AI

Linear decision boundary

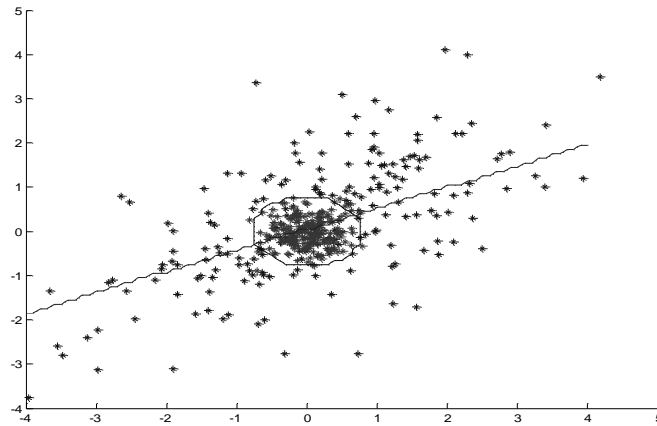
- Example when logistic regression model is not optimal, but not that bad



CIS603 - AI

When logistic regression fails?

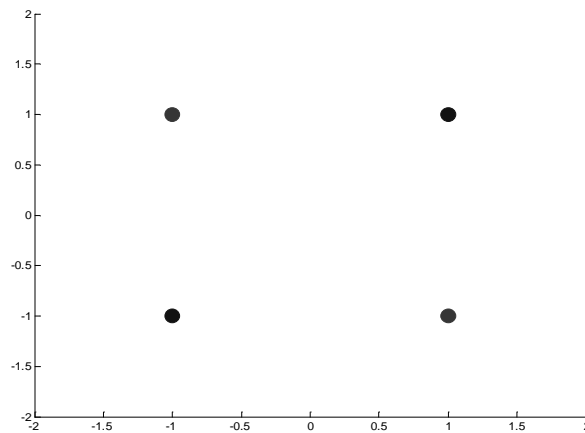
- Example in which the logistic regression model fails



CIS603 - AI

Limitations of logistic regression.

- **parity function** - no linear decision boundary



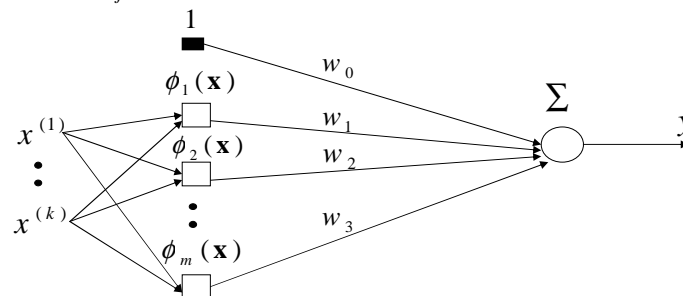
CIS603 - AI

Extensions of simple linear units

Replace inputs to linear units with **feature (basis) functions** to model **nonlinearities**

$$f(\mathbf{x}) = w_0 + \sum_{j=1}^m w_j \phi_j(\mathbf{x})$$

$\phi_j(\mathbf{x})$ - an arbitrary function of \mathbf{x}



The same trick can be done for the logistic regression

CIS603 - AI

Extension of simple linear units

- **Example: Fitting of a polynomial of degree m**

- **Data points:** pairs of $\langle x, y \rangle$

- **Feature functions:**

$$\phi_i(x) = x^i$$

- **Function to learn:**

$$f(x, \mathbf{w}) = w_0 + \sum_{i=1}^m w_i x^i$$

- **On line update** for $\langle x, y \rangle$ pair

$$w_0 = w_0 + \alpha(y - f(\mathbf{x}, \mathbf{w}))$$

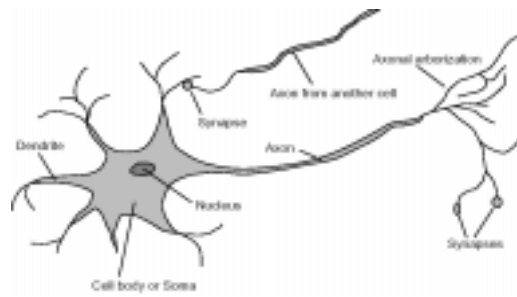
⋮

$$w_i = w_i + \alpha(y - f(\mathbf{x}, \mathbf{w}))x^i$$

CIS603 - AI

Multi-layered neural networks

- Alternative way to introduce nonlinearities to regression/classification models
- **Idea:** Cascade several simple neural models (based on logistic regression). Much like neuron connections.



CIS603 - AI