

Logistic Regression

- Consider a data set $\{(x_n, t_n)\}$, $t_n \in \{0, 1\}$, $n = 1, \dots, N$.

Example: Test (GRE) scores and admission to a graduate school

x	272	331	295	287	315	266	303	294	317	309
t	0	1	1	0	1	0	0	0	1	1

$t = 1$ accepted; $t = 0$ rejected

If $x = 299$, what is the **probability** to be accepted?

- This problem is called **logistic regression**.

- Idea: Transform x_n into probability y_n of admission so that $t_n = 1$ with probability y_n .

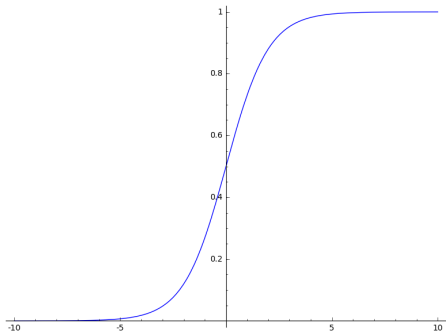
In other words, y_n is the probability of success for score x_n .

- We need a function from $(-\infty, \infty)$ to $(0, 1)$.

Use the logistic sigmoid function

$$\sigma(x) = \frac{e^x}{e^x + 1} = \frac{1}{1 + e^{-x}}$$

Graph of sigmoid



- First, set $\mathbf{x}_n := [x_n, 1]$ and

$$a_n := w_1 x_n + w_2 = \mathbf{x}_n \mathbf{w}$$

for some (unknown) $\mathbf{w} := [w_1, w_2]^\top$.

This is exactly a linear regression.

- Second,

$$y_n := p(t_n = 1 | x_n) = \sigma(a_n).$$

- Third, each result is determined by a Bernoulli trial.

$$x_n \rightsquigarrow y_n \rightsquigarrow t_n = 0, 1$$

Let T be a **Bernoulli** random variable.

$$\Pr(T = 1) = y \quad \text{and} \quad \Pr(T = 0) = 1 - y$$

We consider y as the probability of success.

- The probability mass function $\text{Ber}(\cdot|y)$ is given by

$$\text{Ber}(1|y) = y \quad \text{and} \quad \text{Ber}(0|y) = 1 - y.$$

- We write

$$T \sim \text{Ber}(t|y) = y^t(1 - y)^{1-t}, \quad t = 0, 1.$$

- $E(T) = y$ and $\text{Var}(T) = y(1 - y)$

- data set: $\{(x_n, t_n)\}$, $t_n \in \{0, 1\}$, $n = 1, \dots, N$

(score) \rightsquigarrow (probability of success) \rightsquigarrow (success or failure)

$$x_n \rightsquigarrow y_n = \sigma(w_1 x_n + w_0) \rightsquigarrow t_n$$

- Assume that the test scores are independent.

Each score brings about a Bernoulli random variable.

The likelihood function can be written

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N \text{Ber}(t_n|y_n) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n},$$

where $\mathbf{t} = (t_1, \dots, t_N)$.

- Task: Determine $\mathbf{w} = [w_1, w_2]^T$ to obtain **maximum likelihood**.

- We want to maximize

$$\ln p(\mathbf{t}|\mathbf{w}) = \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\},$$

or equivalently, we want to minimize

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}.$$

It is called the **cross-entropy** error function.

- How can we minimize this error function?

We can use **Gradient Descent** or **Newton's Method**.

k features

- Data set: $\{(x_{n,1}, x_{n,2}, \dots, x_{n,k}; t_n)\}$, $t_n = 0, 1$, $n = 1, 2, \dots, N$

Example

$x_{n,1}$ = GRE score, $x_{n,2}$ = GPA, ...

- Set $\mathbf{x}_n := [x_{n,1}, x_{n,2}, \dots, x_{n,k}, 1]$ and

$$a_n := w_1 x_{n,1} + w_2 x_{n,2} + \dots + w_k x_{n,k} + x_{k+1} = \mathbf{x}_n \mathbf{w}$$

for some (unknown) $\mathbf{w} := [w_1, w_2, \dots, w_{k+1}]^\top$.

- Set $y_n = \sigma(a_n)$ and

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}.$$

- We have

$$\nabla E(\mathbf{w}) = \left[\sum_{n=1}^N (y_n - t_n) x_{nj} \right] = \mathbf{X}^\top (\mathbf{y} - \mathbf{t}),$$

where

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1k} & 1 \\ x_{21} & \cdots & x_{2k} & 1 \\ \vdots & & \vdots & \vdots \\ x_{N1} & \cdots & x_{Nk} & 1 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \quad \text{and } \mathbf{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{bmatrix}.$$

- Gradient Descent

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \eta \mathbf{X}^\top (\mathbf{y} - \mathbf{t}).$$



$$\mathbf{HE} = \left[\sum_{n=1}^N y_n(1 - y_n) \mathbf{x}_{ni} \mathbf{x}_{nj} \right] = \mathbf{X}^\top \mathbf{R} \mathbf{X},$$

where $\mathbf{R} = \text{diag}(y_n(1 - y_n))$.

- Newton's Method

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \eta (\mathbf{X}^\top \mathbf{R} \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{t}),$$

where \mathbf{R} and \mathbf{y} are determined by \mathbf{w}_i in each step.

Example (continued): GRE scores and admission

- Choose $k = 1$. Write $a_n = w_1 x_n + w_2$.

x	272	331	295	287	315	266	303	294	317	309
t	0	1	1	0	1	0	0	0	1	1

- Choose $\mathbf{w}_0 = [0, 0]^\top$. Then \mathbf{w}_i converges to

$$[0.1910, -57.2937]^\top.$$

- If $x = 299$, then the probability of admission is

$$y = \sigma(0.1910 * 299 - 57.2937) \approx 0.454.$$