

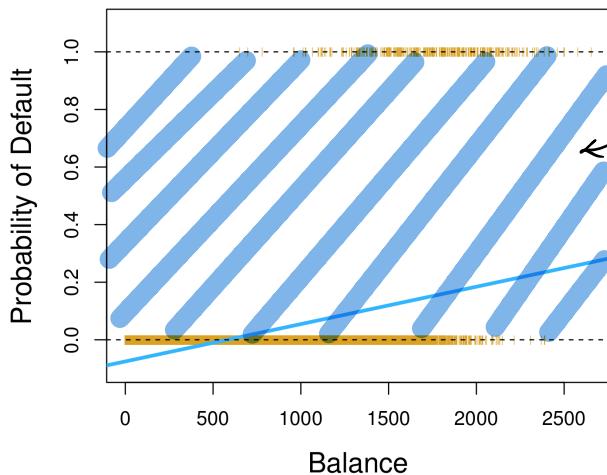
Topic 7: Logistic regression

Section 1: logistic regression

One of the issues we have encountered with MSE is the way it estimates our missclassifications, indeed if $\hat{f}_w(x) > 1$ then it contributes to the error term, but shouldn't as it just shows that we are confident in our prediction even more than for a general case $\hat{f}_w(x) \geq \frac{1}{2}$.

Also, the outcome of regression analysis is rather continuous than discrete, and the resulting MSE error penalises this, but we should allow this and maybe interpret this differently.

We need to apply a different result's digestion!



The object function should take values in $[0, 1]$.

We should also consider values of $\hat{f}_w(x)$ as a probability of belonging to one of the classes.

Transformation of linear regression

Below we consider the way of transforming the linear regression problem to adapt it to binary classification needs, but the method works for any regression problem. Let

$$f(x, w) = \langle x, w \rangle$$

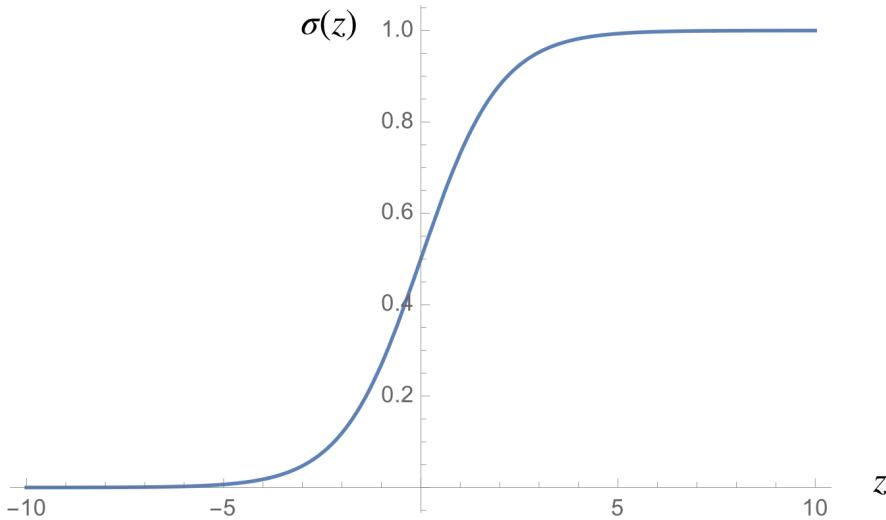
↑ weights
↑ augmented x

We will change our object/prediction function now with a composition (like in DL) of a function $\sigma: \mathbb{R} \rightarrow [0, 1]$ applied to f .

The standard way of defining σ is

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

sigmoid
function

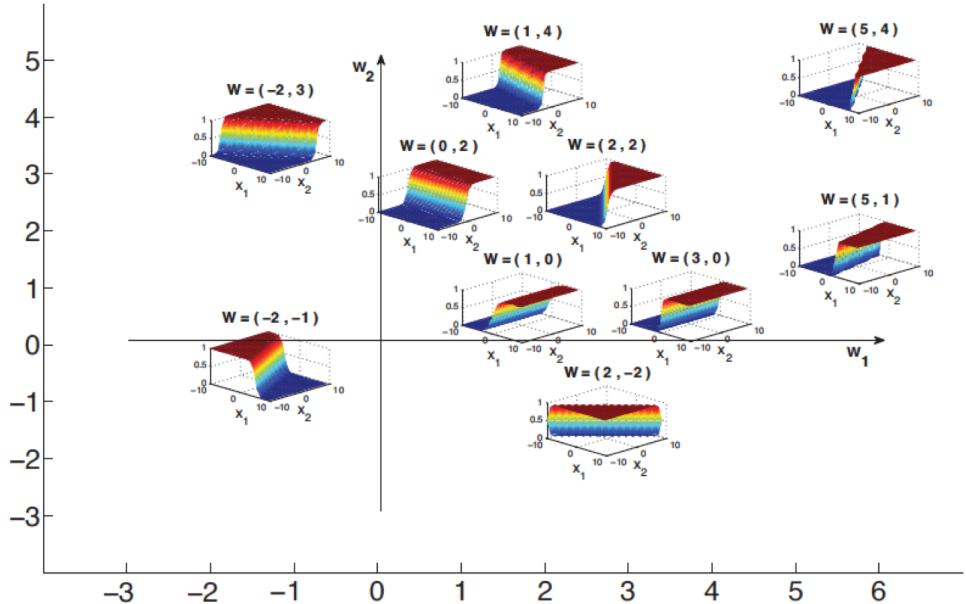


Then for any sample input \mathbf{x} we define posterior probability of having corresponding class labels as

$$p(c_1|\mathbf{x}) = \sigma(f(\mathbf{x}, \mathbf{w})) = \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$$

tendency to class 1.

$$p(c_0|\mathbf{x}) = 1 - \sigma(f(\mathbf{x}, \mathbf{w})) = 1 - \sigma(\langle \mathbf{x}, \mathbf{w} \rangle)$$



How do we obtain the optimal weights $\hat{\mathbf{w}}$?

Similar to all our previous problems, we expect optimal weights to maximise the likelihood of sampling given data sample from the distribution of class labels given by logistic function.

Optimal weights/parameters derivation

$\{x^{(i)}, y_i\}_{i=1}^s$ - i.i.d.

\uparrow
 \mathbb{R}^d

$\{0, 1\}$

$$y := \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_s \end{pmatrix}, \quad X = \begin{pmatrix} x^{(1)} \\ x^{(2)\top} \\ \vdots \\ x^{(s)\top} \end{pmatrix} = \begin{pmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(s)} & \dots & x_d^{(s)} \end{pmatrix}$$

$$g(y | X, w) = \prod_{i=1}^s g(y_i | x^{(i)}, w)$$

↑ independence ↑ augmented

$$= \prod_{i: y_i=0} g(y_i | x^{(i)}, w) \prod_{i: y_i=1} g(y_i | x^{(i)}, w)$$

Trick

$$a^\alpha b^{1-\alpha} = \begin{cases} a, & \alpha=1 \\ b, & \alpha=0 \end{cases}$$

$$= \prod_{i=1}^s \underbrace{g(x^{(i)}, w)^{y_i}}_{\text{maximise}} \underbrace{(1-g(x^{(i)}, w))^{1-y_i}}_{}$$

Let us calculate minus log-likelihood

$$-\log g(y | X, w) = -\sum_{i=1}^s y_i \log g_i + (1-y_i) \log (1-g_i)$$

$$\begin{aligned} g_i &= g(x^{(i)}, w) = \frac{e_i}{1+e_i} & \Rightarrow & y_i \log \frac{e_i}{1+e_i} + (1-y_i) \log \frac{1}{1+e_i} \\ e_i &= e^{w^T x^{(i)}} & \Rightarrow & y_i \log e_i + \log \frac{1}{1+e_i} \end{aligned}$$

So the negative log-likelihood takes the form

$$-\log p(y|X, w) = \sum_{i=1}^s \log(1 + \langle x^{(i)}, w \rangle) - y_i \langle x^{(i)}, w \rangle$$

Remark: for general regression with object function f one would obtain

$$-\log p(y|X, w) = \sum_{i=1}^s \log(1 + f(x^{(i)}, w)) - y_i f(x^{(i)}, w)$$

Optimal parameters/weights are then found as

$$\begin{aligned}\hat{w} &= \arg \min_w -\log p(y|X, w) \\ &= \arg \min_w \sum_{i=1}^s \log(1 + \langle x^{(i)}, w \rangle) - y_i \langle x^{(i)}, w \rangle\end{aligned}$$

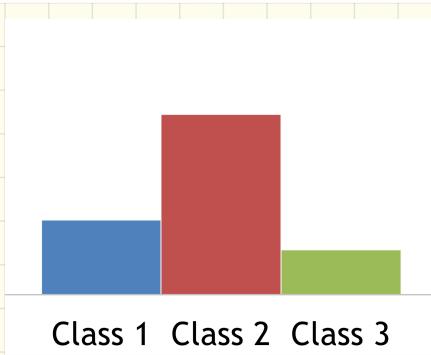
Further steps:

1. Extension of the method to multiclass classification problems
2. Ways of solving the optimisation problem

Multinomial logistic regression

Idea: instead of two probabilities $p(0|x, w)$ $p(1|x, w)$
one needs to define probabilities $p(c_j|x, w)$, $j = \overline{1, n}$
such that

$$\sum_{j=1}^n p(c_j|x, w) = 1$$



Binary case: $z \rightarrow \frac{e^z}{1+e^z}, \quad \frac{1}{1+e^z} = \frac{e^{-z}}{e^0+e^{-z}}, \quad \frac{e^0}{e^0+e^{-z}}$

Multinomial case: (z_1, z_2, \dots, z_n) n "independent" values
 $\rightarrow \text{softmax } (z_1, \dots, z_n)$

What is softmax function?

softmax ($\bar{x} = (x_1, x_2, \dots, x_n)$)

should satisfy

- softmax : $\mathbb{R}^n \rightarrow$ n-dim. probability simplex
- softmax is monotonic in all coordinates

Def.

$$g(\bar{x} = (x_1, x_2, \dots, x_n)) = \left(\frac{e^{x_1}}{\sum_{j=1}^n e^{x_j}}, \frac{e^{x_2}}{\sum_{j=1}^n e^{x_j}}, \dots, \frac{e^{x_n}}{\sum_{j=1}^n e^{x_j}} \right)$$

Remark: Component-wise the softmax function is

$$g(\bar{x})_k = \frac{e^{x_k}}{\sum_{j=1}^n e^{x_j}}$$

Remark: The softmax name comes from the fact that softmax function is a continuous approximation of an argmax function

$$\text{argmax } (x_1, x_2, \dots, x_n) = x_p$$

$$= (0 \ 0 \dots 0 \underset{p}{\underline{1}} 0 \dots 0)$$

$$\text{softmax } (x_1, \dots, x_n) = \left(\frac{e^{x_1}}{\sum_{j=1}^n e^{x_j}}, \dots, \frac{e^{x_p}}{\sum_{j=1}^n e^{x_j}}, \dots, \frac{e^{x_n}}{\sum_{j=1}^n e^{x_j}} \right)$$

$$\text{argmax } (1.5 \ 0.3 \ -3.7) = 0 = (1 \ 0 \ 0)$$

largest
and $\rightarrow 1$

if $x_p \rightarrow \infty$.

$$\text{softmax } (1.5 \ 0.3 \ -3.7) = (0.765 \ 0.231 \ 0.004)$$

Setting up the multinomial regression

$\left\{ \underset{\substack{\in \mathbb{R}^d \\ \{1, 2, \dots, K\}}}{x^{(i)}}, y_i \right\}_{i=1}^s$ - data samples

Linear regression assumption is not essential but makes the notations clearer.

$$f(\bar{x}, W) = \langle \bar{x}, w^{(1)} \rangle \quad \langle \bar{x}, w^{(2)} \rangle \quad \dots \quad \langle \bar{x}, w^{(k)} \rangle$$

$$\sigma(\bar{x}, W) = \left(\frac{e^{\langle \bar{x}, w^{(1)} \rangle}}{\sum_{j=1}^k e^{\langle \bar{x}, w^{(j)} \rangle}}, \dots, \frac{e^{\langle \bar{x}, w^{(k)} \rangle}}{\sum_{j=1}^k e^{\langle \bar{x}, w^{(j)} \rangle}} \right)$$

$$g(y_i=k | x^{(i)}, W) = \sigma(\bar{x}, W)_k = \frac{e^{\langle \bar{x}, w^{(k)} \rangle}}{\sum_{j=1}^k e^{\langle \bar{x}, w^{(j)} \rangle}}$$

The likelihood maximisation

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_s \end{pmatrix} \quad X = \begin{pmatrix} x^{(1)\top} \\ x^{(2)\top} \\ \vdots \\ x^{(s)\top} \end{pmatrix} \quad x^{(i)} = (1 \ x^{(i)}_1 \dots x^{(i)}_d)$$

$$W = \begin{pmatrix} w^{(1)} & w^{(2)} & \dots & w^{(k)} \end{pmatrix} \quad w^{(i)} = \begin{pmatrix} w_0^{(i)} & w_1^{(i)} & \dots & w_d^{(i)} \end{pmatrix}$$

$$g(y | X, W) = \prod_{i=1}^s g(y_i | x^{(i)}, W)$$

split into groups

$$= \prod_{i: y_i=1} g(1 | x^{(i)}, W) \cdot \prod_{i: y_i=2} g(2 | x^{(i)}, W) \cdot \dots \cdot \prod_{i: y_i=k} g(k | x^{(i)}, W)$$

$$\text{Trick: } \prod_{p=1}^K d_p^{1_{\{p=j\}}} = d_j$$

$$p(y|X, W) = \prod_{i=1}^s \prod_{k=1}^K p(k|x^{(i)}, W)^{1_{y_i=k}}$$

$$\hat{W} = \arg \min_w -\log p(y|X, W)$$

$$= \arg \min_w - \sum_{i=1}^s \sum_{k=1}^K \log p(k|x^{(i)}, W)^{1_{y_i=k}}$$

$$= \arg \min_w - \sum_{i=1}^s \sum_{k=1}^K 1_{y_i=k} \cdot \log p(k|x^{(i)}, W)$$

$$= \arg \min_w - \sum_{i=1}^s \sum_{k=1}^K 1_{y_i=k} \cdot \log \frac{e^{<x^{(i)}, w^{(k)}>}}{\sum_{j=1}^K e^{<x^{(i)}, w^{(j)}>}}$$

$$= \arg \min_w - \sum_{i=1}^s \sum_{k=1}^K 1_{y_i=k} \left[<x^{(i)}, w^{(k)}> - \log \left(\sum_{j=1}^K e^{<x^{(i)}, w^{(j)}>} \right) \right]$$

$$= \arg \min_w \sum_{i=1}^s \sum_{k=1}^K 1_{y_i=k} \cdot \log \left(\sum_{j=1}^K e^{<x^{(i)}, w^{(j)}>} \right) - \sum_{i=1}^s \sum_{k=1}^K 1_{y_i=k} <x^{(i)}, w^{(k)}>$$

$$\hat{W} = \arg \min_w \sum_{i=1}^s \log \left(\sum_{j=1}^K e^{<x^{(i)}, w^{(j)}>} \right) - \sum_{i=1}^s \sum_{k=1}^K 1_{y_i=k} <x^{(i)}, w^{(k)}>$$

Section 2: Solving Logistic regression

Previously we have derived the optimisation/minimisation problems that one needs to solve to classify the data using logistic regression.

Binary $\hat{w} = \arg \min_w \left\{ \sum_{i=1}^s \log(1 + e^{-\langle x^{(i)}, w \rangle}) - y_i \langle x^{(i)}, w \rangle \right\}$

Multinomial

$$\hat{w} = \arg \min_w \left\{ \sum_{i=1}^s \log \left(\sum_{j=1}^k e^{\langle x^{(i)}, w^{(j)} \rangle} \right) - \sum_{i=1}^s \sum_{j=1}^k 1_{y_i=j} \langle x^{(i)}, w^{(j)} \rangle \right\}$$

Gradient descent?

$$L_2(w) = \sum_{i=1}^s \left\{ \log (1 + e^{-\langle x^{(i)}, w \rangle}) - y_i \langle x^{(i)}, w \rangle \right\}$$

$$L_n(w) = \sum_{i=1}^s \log \left(\sum_{j=1}^k e^{\langle x^{(i)}, w^{(j)} \rangle} \right) - \sum_{i=1}^s \sum_{j=1}^k 1_{y_i=j} \langle x^{(i)}, w^{(j)} \rangle$$

$$w^{(k+1)} = w^{(k)} - \tau \nabla L(w^{(k)})$$

↑
matrix or vector

We solely concentrate on a binary logistic regression

$$L_2(w) = \sum_{i=1}^s \log \left(1 + e^{<\mathbf{x}^{(i)}, w>} \right) - y_i <\mathbf{x}^{(i)}, w>$$

To calculate the gradient we will need to calculate two sorts of derivatives:

- $\ell(z) = \log \left(1 + e^z \right)$

$$\ell'(z) = \frac{e^z}{1+e^z} = \sigma(z)$$

- $\frac{\partial}{\partial w_p} <\mathbf{x}^{(i)}, w> = \frac{\partial}{\partial w_p} \sum_{j=0}^d x_j^{(i)} w_j$

$$= x_p^{(i)}$$

$$\begin{aligned} \frac{\partial}{\partial w_p} L_2(w) &= \sum_{i=1}^s \frac{\partial}{\partial w_p} \ell(<\mathbf{x}^{(i)}, w>) - y_i <\mathbf{x}^{(i)}, w> \\ &= \sum_{i=1}^s \frac{\partial}{\partial w_p} <\mathbf{x}^{(i)}, w> \cdot \frac{d}{dz_i} (\ell(z_i) - y_i \cdot z_i) \end{aligned}$$

where $z_i = <\mathbf{x}^{(i)}, w>$

$$= \sum_{i=1}^s x_p^{(i)} (\sigma(<\mathbf{x}^{(i)}, w>) - y_i)$$

$$\nabla L_2(w) = X^T (g(Xw) - y),$$

where $g(u) = (g(u_1) \ g(u_2) \dots g(u_s))$

Remark: The above is very similar to the classical regression gradient descent result.

Will the gradient descent converge to a global minimiser?

There are two parts:

1. Would the result be a minimiser?
2. would the method converge?

Convexity of L

L is a sum of identical functions applied to vectors $x^{(i)}$. So let us check every single piece individually.

L can be considered as a composition of a function applied to linear function, thus we need to check convexity of an external one only.

$$l(z) = \log(1 + e^z) \Rightarrow l'(z) = g(z)$$

$$g'(z) = g(z)(1 - g(z)) > 0$$



Remark: We can also consider regularisation of logistic regression

$$\hat{w} = \arg \min_w \{ L_2(w) + \lambda R(w) \}$$

If R is differentiable, then one can use the gradient descent directly. If R is not differentiable, then one can use the proximal gradient descent.