

Topic 7:

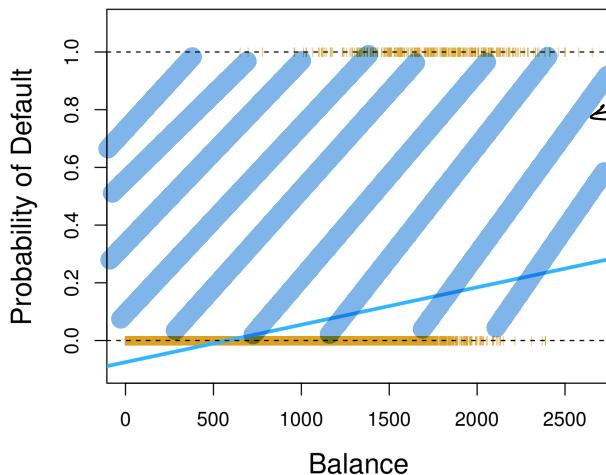
Logistic regression

1 Motivation

Idea of solving binary classification problem:

- 1.) Build a regression $f_w(x)$
- 2.) Say that $y(x) = \begin{cases} 0, & f_w(x) < \frac{1}{2} \\ 1, & f_w(x) \geq \frac{1}{2} \end{cases}$

Problems to solve:



- Move values of f_w into interval $[0, 1]$
- Think about interpretation of continuous values into $\{0, 1\}$.

understand f_w as probability

values into $\{0, 1\}$.

1.1 Transformation of $f_w(x)$.

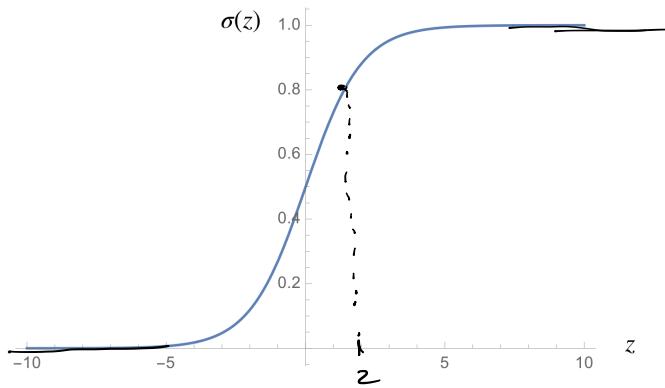
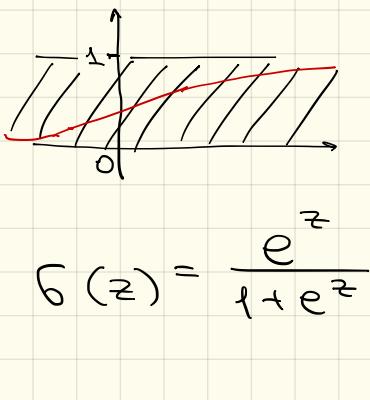
Assumption: we need to solve a binary classification problem "using" linear regression

$$f_w(x) = \langle x, w \rangle \quad \begin{matrix} \rightarrow \in \mathbb{R}^{d+1}, \text{weights} \\ \text{augmented} \in \mathbb{R}^{d+1} \end{matrix}$$

input $\in \mathbb{R}^d$

- ① Change object/prediction function to $\sigma(f_w(x))$ such that

$$\sigma: \mathbb{R} \rightarrow [0, 1]$$



Using the above function we would define such probabilities:

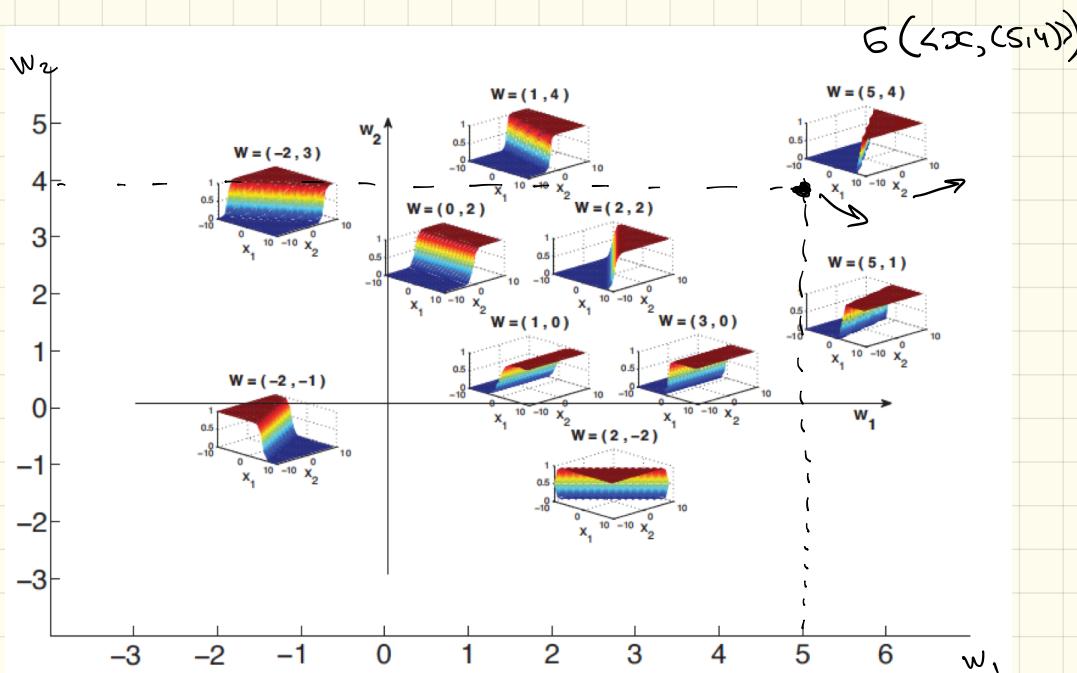
$$g(1 | x, w) = \sigma(f_w(x))$$

↑ class label ↑ given ↓

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

$$g(0 | x, w) = 1 - \sigma(f_w(x))$$

$$= 1 - \sigma(x, w).$$



How to define \hat{w} ?

1.2 Optimisation of weights.

Similar to classical regression problems we want to maximise the likelihood of getting y_i values for given $x^{(i)}$ and w over varying w .

$$\left\{ \begin{matrix} \mathbf{x}^{(i)} \\ \mathbb{R}^d \end{matrix}, y_i \right\}_{i=1}^n - \text{i.i.d.}$$

\uparrow
 $\{0,1\}$

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

split into
2 groups:
 $\{y_i = 0\}$
 $\{y_i = 1\}$

$$\Rightarrow = \prod_{i: y_i=0} \underbrace{g(y_i | x^{(i)}, w)}_{n-\sigma}$$

$$\cdot \prod_{i: y_i=1} \underbrace{g(y_i | x^{(i)}, w)}_{\sigma}$$

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

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

split into
2 groups:
 $\begin{cases} y_i = 0 \\ y_i = 1 \end{cases}$

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

Trick: $\frac{a^\alpha b^{1-\alpha}}{\Gamma(\alpha)} = \begin{cases} a, \alpha=1 \\ b, \alpha=0 \end{cases}$

$$= \prod_{i=1}^s \sigma_i^{y_i} (1 - \sigma_i)^{1-y_i}$$

$$L(w) = -\log g(y | X, w)$$

$$\sigma_i = \sigma(x^{(i)}, w) = -\log \prod_{i=1}^s \sigma_i^{y_i} (1 - \sigma_i)^{1-y_i}$$

$$e_i = e^{-\sum_{i=1}^s y_i \log \sigma_i + (1-y_i) \log (1-\sigma_i)}$$

$$= -\sum_{i=1}^s y_i \log \frac{e_i}{1+e_i} + (1-y_i) \log \left(1 - \frac{e_i}{1+e_i}\right)$$

$$= -\sum_{i=1}^s \underbrace{y_i \log e_i}_{\text{constant}} - \underbrace{y_i \log (1+e_i)}_{\text{linear}} - \underbrace{(1-y_i) \log (1+e_i)}_{\text{quadratic}}$$

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

$$L(w) = \sum_{i=1}^S \log(1 + e^{\langle x^{(i)}, w \rangle}) - y_i \langle x^{(i)}, w \rangle$$

$$\hat{w} = \arg \min_w L(w)$$

Questions to answer:

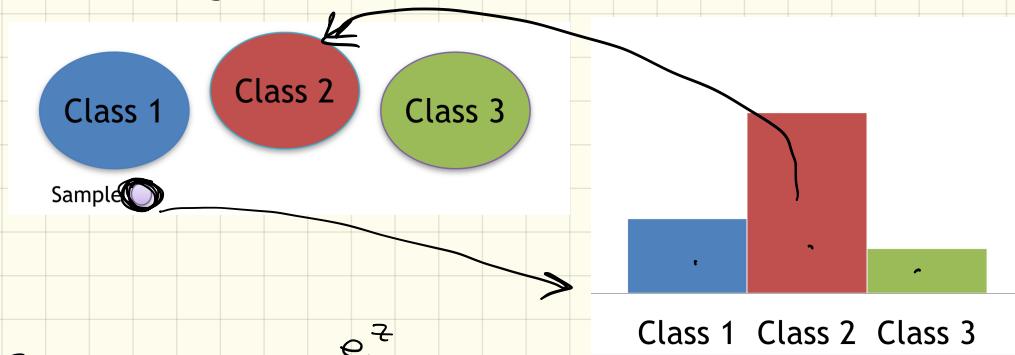
1. How to extend the method to multinomial regression?
2. How to solve corresponding optimisation problems?

2. Multinomial logistic regression

If one wants to classify inputs into n classes then one would be interested in defining $p(j | \mathbf{x}, \mathbf{w})$ for $j = \overline{1, n}$.

Restrictions:

- $p(j | \mathbf{x}, \mathbf{w}) \in [0, 1]$
- $\sum_{j=1}^n p(j | \mathbf{x}, \mathbf{w}) = 1$



$$\text{Binary: } z \rightarrow \frac{e^z}{1+e^z}$$

$$\text{Multinomial case: } (z_1, z_2, \dots, z_n) \xrightarrow{\text{softmax}} (p_1, \dots, p_n)$$

Def: The softmax function

$$g: \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \text{is}$$

$$\text{softmax} = \sigma(z = (z_1, \dots, z_n))$$

$$= \begin{pmatrix} e^{z_1} \\ \sum_{j=1}^n e^{z_j} & e^{z_2} \\ \sum_{j=1}^n e^{z_j} & \dots & e^{z_n} \\ \sum_{j=1}^n e^{z_j} \end{pmatrix}$$

Remark: binary multinomial

$$\frac{e^z}{1+e^z} = \sigma(z) \iff \sigma((0,z))$$

$$\bar{x} \sim \langle x, w \rangle = z \xrightarrow{\text{plug}} (0, z)$$

e^w $\left(\frac{1}{1+e^z}, \frac{e^z}{1+e^z} \right)$

Origin of softmax name?

$$\arg \max \left(x_1, x_2, \dots, x_n \right) = p$$

\downarrow
 $= (0 \ 0 \ \dots \ 0 \ \underbrace{1}_{p} \ 0 \ \dots \ 0)$

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

largest, $\rightarrow 1, x_p \rightarrow \infty$
 ↓
 p

2.2 Setting up a multinomial logistic regression problem

How to get n dimensional vector from d-dimensional input?

$$f_w(\bar{x}) = \begin{pmatrix} \langle \bar{x}, w^{(1)} \rangle & \langle \bar{x}, w^{(2)} \rangle & \dots & \langle \bar{x}, w^{(n)} \rangle \end{pmatrix}$$

$$\sigma(\bar{x}, W) = \begin{pmatrix} e^{\langle \bar{x}, w^{(1)} \rangle} \\ \sum_{j=1}^n e^{\langle \bar{x}, w^{(j)} \rangle} \\ \dots \\ \sum_{j=1}^n e^{\langle \bar{x}, w^{(j)} \rangle} \end{pmatrix}$$

$$g(y_i = k \mid \bar{x}^{(i)}, W) = (\sigma(\bar{x}^{(i)}, W))_k$$

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

2.3

Solving logistic regression.

Binary classification:

$$\hat{w} = \arg \min_w \left\{ \sum_{i=1}^s \log \left(1 + e^{-\langle x^{(i)}, w \rangle} \right) - \sum_{i=1}^s y_i \cdot \underbrace{\langle x^{(i)}, w \rangle}_{\text{green bracket}} \right\}$$

$w^{(0)} = (1, \dots)$
 $w^{(1)} = w$

Multinomial classification

$$\hat{W} = \arg \min_W \left\{ \sum_{i=1}^s \log \left(\sum_{j=1}^k e^{\langle x^{(i)}, w^{(j)} \rangle} \right) - \sum_{i=1}^s \sum_{j=1}^k \mathbb{1}_{y_i=j} \cdot \underbrace{\langle x^{(i)}, w^{(j)} \rangle}_{\text{green bracket}} \right\}$$

2 classes {0, 1}

$$- \sum_{i=1}^s \underbrace{\mathbb{1}_{y_i=0} \cdot \langle x^{(i)}, w_0 \rangle}_{j=0} + \underbrace{\mathbb{1}_{y_i=1} \cdot \langle x^{(i)}, w_1 \rangle}_{j=1}$$

$$y_i = \begin{cases} 0, & \mathbb{1}_{y_i=0} \\ 1, & \mathbb{1}_{y_i=1} \end{cases}$$

Gradient descent?

2.3.1 Gradient descent for logistic regression.

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

$$L_n(w) = \sum_{i=1}^s \left\{ \log \left(\sum_{j=1}^k e^{<x^{(i)}, w^{(j)}>} \right) - \sum_{j=1}^k y_{i=j} <x^{(i)}, w^{(j)}> \right\}$$

$$W^{(k+1)} = W^{(k)} - \tau \nabla L_m(W^{(k)})$$

\vec{v} vector if $m=2$

matrix if $m > 2$.

Binary-

$\nabla L_2(w) - ? \Rightarrow$ need to find $\frac{\partial L_2(w)}{\partial w_p}$

- What is the derivative $\frac{\partial}{\partial w_p} f(x^{(i)}, w)$

$$\frac{\partial}{\partial w_p} f \left(\underbrace{x_0^{(i)} w_0 + x_1^{(i)} w_1 + \dots + x_d^{(i)} w_d}_{\text{only one has } w_p} \right) = f' \left(\langle x^{(i)}, w \rangle \right) \cdot x_p^{(i)}$$

$$\frac{\partial}{\partial w_p} L_2(w) = \sum_{i=1}^S \left\{ \frac{\partial}{\partial w_p} \log(1 + e^{x^{(i)} \cdot w}) - \frac{\partial}{\partial w_p} y_i \cdot \Delta x^{(i)} \cdot w \right\}$$

$$= \sum_{i=1}^s \left\{ \delta \left(\underbrace{\langle \alpha^{(i)}, w \rangle}_{\text{f}(z) = y_i \cdot z} \right) \cdot x_P^{(i)} - y_i \cdot x_P^{(i)} \right\}$$

$f(z) = y_i \cdot z$
 $f'(z) = y_i$

$$X = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ \vdots & & & & \\ 1 & x_1^{(s)} & x_2^{(s)} & \dots & x_d^{(s)} \end{pmatrix} = X_{ip}$$

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

$$f(z) = \log(1 + e^z)$$

$$f'(z) = \sigma(z)$$

$$= \sum_{i=1}^s \left\{ \sigma(\underbrace{(x_w)_i}_{<x^{(i)}, w>}) \cdot x_p^{(i)} - y_i \cdot x_p^{(i)} \right\}$$

$$f(z) = y_i \cdot z$$

$$f'(z) = y_i$$

$$X = \begin{pmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ \vdots & & & & \\ 1 & x_1^{(s)} & x_2^{(s)} & \dots & x_d^{(s)} \end{pmatrix}$$

$$X_w = \begin{pmatrix} x^{(1)} \\ x^{(2)} \end{pmatrix} \begin{pmatrix} w \end{pmatrix} = \begin{pmatrix} \overleftarrow{x^{(1)}, w} \\ \overrightarrow{x^{(2)}, w} \end{pmatrix} \Rightarrow \langle x^{(i)}, w \rangle = (X_w)_i$$

$$= \sum_{i=1}^s X_{ip} \left(\sigma((X_w)_i) - y_i \right) = \boxed{X^T \sigma(X_w) - y}$$

$\sigma(X_w)$ - ? coordinate-wise application of σ .

$$\nabla L_2$$