

A3. i) ~~$P \in [0,1]$~~ $P \in (0,1)$ as
 $g(z) = 1/(1 + \exp(-z))$ which is
bound between $(0,1)$ as
it is parameterized by w .

We can change w_0, w_1 as needed
to produce $e^{-w_0'}$. $e^{-w_1'x} = e^{-k}$
where $P(y=1 | x; k) \in (0,1)$

ii) logit function $l(x) = \log(x/(1-x))$

1. $l(x)$ is continuous in domain $[0,1]$

Proof: for any $c \in [0,1]$

i) $\log(c/(1-c))$ is defined and

$$\lim_{c \rightarrow 0} \log(c/(1-c)) = -\infty \text{ and}$$

$$\lim_{c \rightarrow 1} \log(c/(1-c)) = \infty \quad \text{--- (1)}$$

ii) $\lim_{c \rightarrow 1} \log(c/(1-c)) = l(c)$ and

$\lim_{c \rightarrow 0} l(c)$ are defined right and

left limit.

Hence $l(x)$ is continuous
in domain $[0,1]$

from (1) we conclude
that $\text{range}(f(x)) \in (-\infty, \infty)$

Q4.

$$a) J(\theta) = \frac{1}{n} \sum_{i=1}^n |h_{\theta}(x^{(i)}) - y^{(i)}|$$

$$\text{Gradient update} = \theta_i \leftarrow \theta_i - \alpha \cdot \frac{\partial J(\theta)}{\partial \theta_i}$$

where,

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{n} \sum_i (-1) \cdot x_j^{(i)}, \text{ if } h_{\theta}(x^{(i)}) - y^{(i)} < 0$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{n} \sum_i (-1) \cdot x_j^{(i)}, \text{ if } h_{\theta}(x^{(i)}) - y^{(i)} < 0$$

$$= 0, \text{ if } h_{\theta}(x^{(i)}) - y^{(i)} = 0$$

$$= \frac{1}{n} \sum_i (-1) \cdot x_j^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{n} \sum \text{sgn}(h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

where $\text{sgn}(\cdot)$ is the signum function.

- This is difficult to optimize because of the fixed size parameter ~~update~~

update. As it is not looking at the

cost incurred, the fixed ~~update~~ update will keep overshooting the optimal unless α is small.

In case of MSE, step size decreased w/ decrease in cost, so step sizes decreased as the optimum was approached.

b) Used for datasets w/ high variance or many outliers, since the ~~weight update~~ ~~parameter update~~ update gives the same weight to all costs of the same sign.

Thus noisy ~~ex~~ examples are ignored.

penalize

c) We can specifically ~~reduce~~ over-estimation or under-estimation more than the other.

This also provides a range or interval of predictions as compared to exact values. This becomes relevant for models data with non-constant variance, which cannot be solved properly with MAE ~~loss~~ like losses, which will give bad estimates.