# 2. Model Calibration

$h_\theta(x) \longrightarrow$ Why we might treat the output as a probability.

$\Rightarrow$ When the probabilities outputted by a model match empirical observation, the model is said to be <mark>well-calibrated</mark>.

    $\hookrightarrow$ Logistic regression tends to output well calibrated probabilities.

       $\left(\begin{array}{l}\text{This is often not true with other classifiers}\\\text{such as Naive Bayes, or SVM}\end{array}\right)$

$\Rightarrow$ Suppose we have a training set $\{(x^{(i)}, y^{(i)}) \mid x^{(i)} \in R^{n+1} \& y^{(i)} \in \{0,1\}\}$

$\Rightarrow$ Let $\Theta \in R^{n+1}$ be the maximum likelihood parameters learned after training a logistic regression model.

$\Rightarrow$ In order for the model to be considered well-calibrated, given any range of probabilities $(a,b)$ such that $0 \leqslant a \leqslant b \leqslant 1$, and training examples $x^{(i)}$ where the model outputs $h_\theta(x^{(i)})$ fall in the range $(a,b)$, the fraction of positive in that set of example should be equal to the average of the model output for those examples.

$$\frac{\sum\limits_{i \in I_{a,b}} P(y^{(i)}=1 \mid x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} = \frac{\sum\limits_{i \in I_{a,b}} \mathbb{1}\{y^{(i)}=1\}}{|\{i \in I_{a,b}\}|}$$

$$I_{a,b} = \{ i \mid i \in \{1, \cdots m\}, h_\theta(x^{(i)}) \in (a,b) \}$$

$|S| =$ Size of Set $S$

(a) Show that the property holds true for

$$\left( \frac{\text{logistic regression}}{\text{model}} \right) \longrightarrow (a,b) = (0,1)$$

over the range

$\Rightarrow$ For Logistic regression:

$$\theta = \underset{\theta}{\text{argmax}} \; \boxed{P(y^{(1)} \cdots y^{(m)} \mid x^{(1)} \cdots x^{(m)}; \theta)}$$

$$\left\{ h_\theta(x^{(i)}) = \frac{1}{1+e^{-\theta^T x^{(i)}}} \right\}$$

$$\downarrow$$

$$\prod_{i=1}^{m} P(y^{(i)} \mid x^{(i)}; \theta)$$

$$\downarrow$$

$$\prod_{i=1}^{m} h_\theta(x^{(i)})^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{(1-y^{(i)})}$$

$$\theta = \underset{\theta}{\text{argmax}} \; \boxed{\sum_{i=1}^{m} y^{(i)} \log h_\theta(x^{(i)}) + (1-y^{(i)}) \log (1 - h_\theta(x^{(i)}))}$$

$$\longrightarrow \ell(\theta)$$

$\Rightarrow$ To prove:

$$\frac{1}{m} \sum_{i=1}^{m} h_\theta(x^{(i)}) = \frac{1}{m} \sum_{i=1}^{m} \mathbb{I}\{y^{(i)} = 1\}$$

$$\sum_{i=1}^{m} h_\theta(x^{(i)}) = \sum_{i=1}^{m} \mathbb{I}\{y^{(i)} = 1\} = \sum_{i=1}^{m} y^{(i)}$$

$$\nabla l(\theta) = \sum_{i=1}^{m} (y^{(i)} - h_\theta(x^{(i)}))\theta$$

$\Rightarrow$ For the best $\theta$    $\nabla l(\theta) = 0$

$$\sum_{i=1}^{m} y^{(i)} = \sum_{i=1}^{m} h_\theta(x^{(i)})$$

(b) Both the answers are no.

(c) When a regularization $\lambda\|\theta\|$ is added, the equation becomes

$$\sum_{i=1}^{m} y^{(i)} = \sum_{i=1}^{m} h_\theta(x^{(i)}) + 2\lambda\theta_0$$

Where $\theta_0$ is the parameter for the Intercept. In general, we will not penalize this term, and in this case regularization will have no effect.