

Data Graduate Program N°2 Cdiscount Academy

Lundi 20/09 – Mercredi 23/10



SOMMAIRE

- Logistic Regression



Logistic Regression



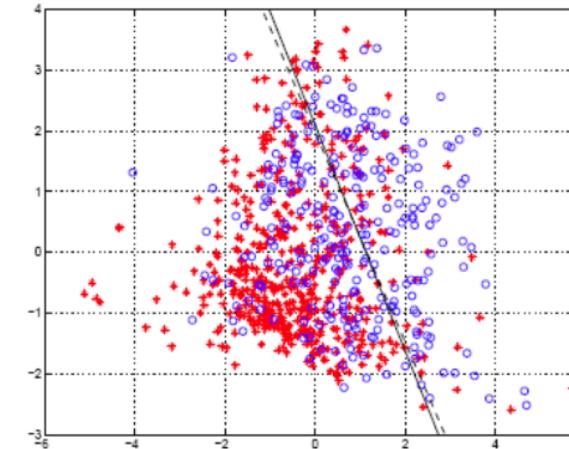
LOGISTIC REGRESSION



- Instead of just predicting the class, give the probability of the instance being that class
 - i.e., learn $p(y | x)$
- Comparison to perceptron:
 - Perceptron doesn't produce probability estimate
 - Perceptron (and other discriminative classifiers) are only interested in producing a discriminative model
- Recall that:

$$0 \leq p(\text{event}) \leq 1$$

$$p(\text{event}) + p(\neg\text{event}) = 1$$



LOGISTIC REGRESSION –



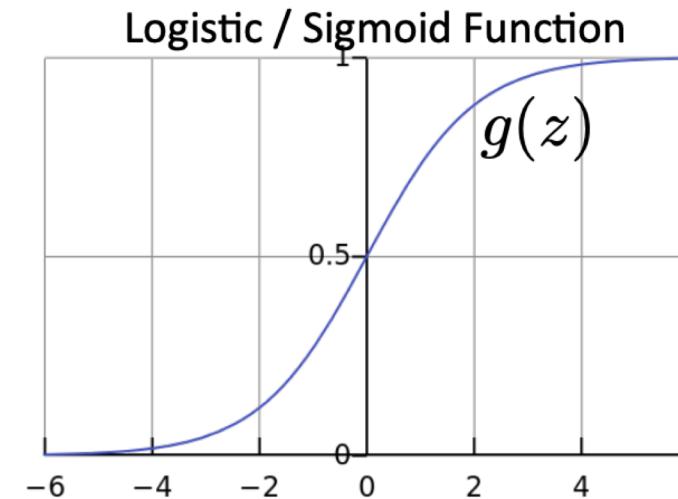
- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)
- $h_{\theta}(x)$ should give $p(y = 1 | x; \theta)$
 - Want $0 \leq h_{\theta}(x) \leq 1$
- Logistic regression model:

$$h_{\theta}(x) = g(\theta^T x)$$

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

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

Can't just use linear regression with a threshold



LOGISTIC REGRESSION – INTERPRETATION OF HYPOTHESIS OUTPUT



$$h_{\theta}(x) = \text{estimated } p(y = 1 | x; \theta)$$

Example: Cancer diagnosis from tumor size

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\theta}(\mathbf{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant



Note that: $p(y = 0 | x; \theta) + p(y = 1 | x; \theta) = 1$

Therefore, $p(y = 0 | x; \theta) = 1 - p(y = 1 | x; \theta)$

LOGISTIC REGRESSION – ANOTHER INTERPRETATION

- Equivalently, logistic regression assumes that



$$\log \frac{p(y = 1 | \mathbf{x}; \boldsymbol{\theta})}{p(y = 0 | \mathbf{x}; \boldsymbol{\theta})} = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

odds of $y = 1$

Side Note: the odds in favor of an event is the quantity $p / (1 - p)$, where p is the probability of the event

E.g., If I toss a fair dice, what are the odds that I will have a 6?



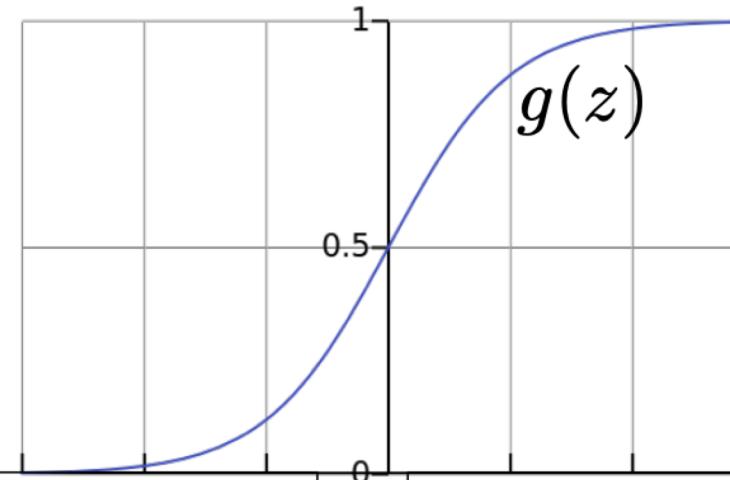
- In other words, logistic regression assumes that the log odds is a linear function of \mathbf{x}

LOGISTIC REGRESSION



$$h_{\theta}(x) = g(\theta^T x)$$

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

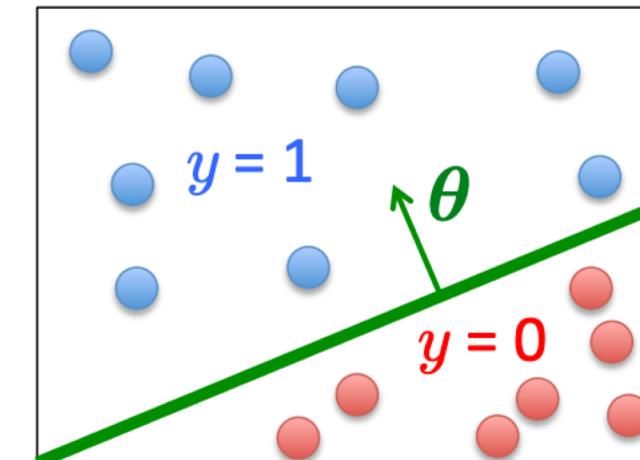


$\theta^T x$ should be large negative values for negative instances

$\theta^T x$ should be large positive values for positive instances



- Assume a threshold and...
 - Predict $y = 1$ if $h_{\theta}(x) \geq 0.5$
 - Predict $y = 0$ if $h_{\theta}(x) < 0.5$

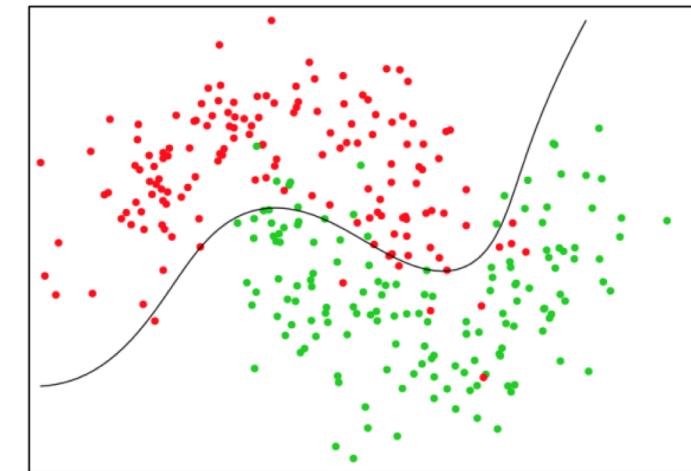


LOGISTIC REGRESSION – NON-LINEAR DECISION BOUNDARY

- Can apply basis function expansion to features, same as with linear regression



$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ \vdots \end{bmatrix}$$



LOGISTIC REGRESSION



- Given $\left\{ \left(\mathbf{x}^{(1)}, y^{(1)} \right), \left(\mathbf{x}^{(2)}, y^{(2)} \right), \dots, \left(\mathbf{x}^{(n)}, y^{(n)} \right) \right\}$
where $\mathbf{x}^{(i)} \in \mathbb{R}^d$, $y^{(i)} \in \{0, 1\}$
- Model: $h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^\top \mathbf{x})$

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



$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = [1 \ x_1 \ \dots \ x_d]$$

LOGISTIC REGRESSION – OBJECTIVE FUNCTION



- Can't just use squared loss as in linear regression:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

- Using the logistic regression model

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}}}$$

results in a non-convex optimization



LOGISTIC REGRESSION – DERIVING COST FUNCTION USING MLE



- Likelihood of data is given by: $l(\theta) = \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; \theta)$
- So, looking for the θ that maximizes the likelihood

$$\theta_{\text{MLE}} = \arg \max_{\theta} l(\theta) = \arg \max_{\theta} \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; \theta)$$

- Can take the log without changing the solution:



$$\theta_{\text{MLE}} = \arg \max_{\theta} \log \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^n \log p(y^{(i)} | \mathbf{x}^{(i)}; \theta)$$

LOGISTIC REGRESSION – DERIVING COST FUNCTION USING MLE

- Expand as follows:



$$\begin{aligned}\theta_{\text{MLE}} &= \arg \max_{\theta} \sum_{i=1}^n \log p(y^{(i)} | \mathbf{x}^{(i)}; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^n \left[y^{(i)} \log p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \theta) + (1 - y^{(i)}) \log (1 - p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \theta)) \right]\end{aligned}$$

- Substitute in model, and take negative to yield

Logistic regression objective:



$$\min_{\theta} J(\theta)$$

$$J(\theta) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(\mathbf{x}^{(i)})) \right]$$

LOGISTIC REGRESSION – INTUITION BEHIND THE OBJECTIVE



$$J(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

- Cost of a single instance:

$$\text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

- Can re-write objective function as



$$J(\boldsymbol{\theta}) = \sum_{i=1}^n \text{cost}\left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)}\right)$$

Compare to linear regression: $J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$

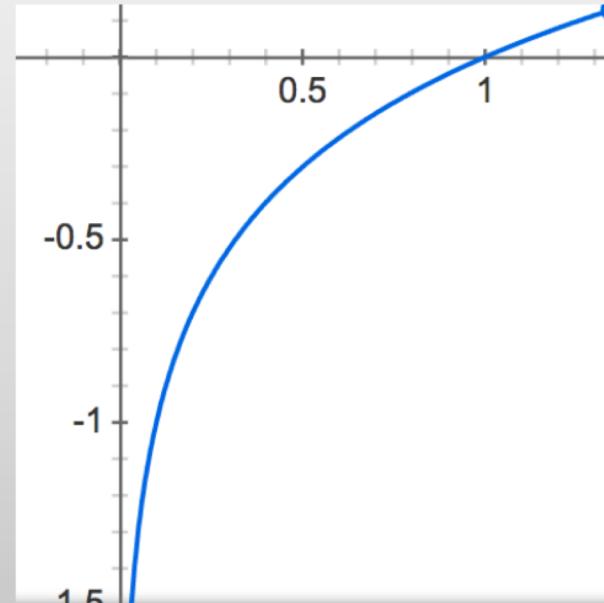
LOGISTIC REGRESSION – INTUITION BEHIND THE OBJECTIVE



$$\text{cost } (h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$



Aside: Recall the plot of $\log(z)$



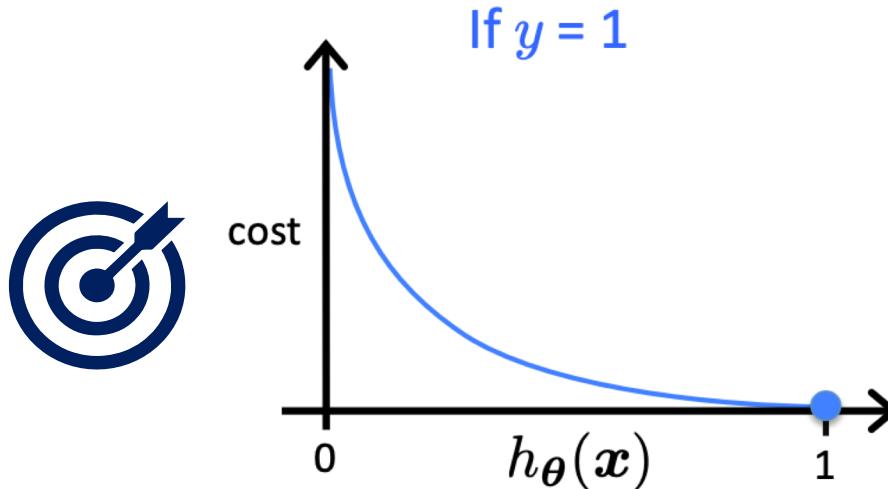
LOGISTIC REGRESSION – INTUITION BEHIND THE OBJECTIVE



$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$

If $y = 1$

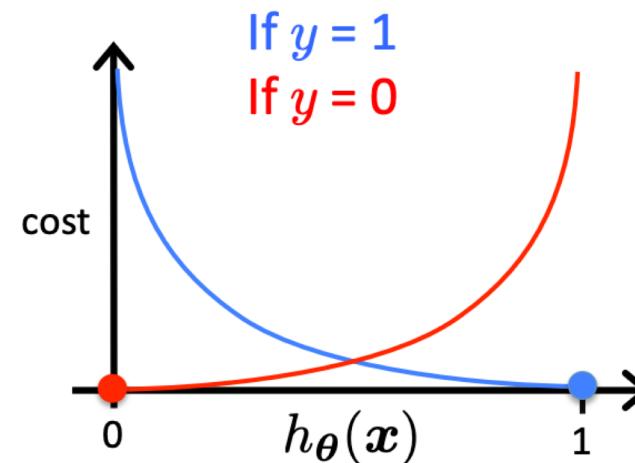
- Cost = 0 if prediction is correct
- As $h_{\theta}(x) \rightarrow 0$, cost $\rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties
 - e.g., predict $h_{\theta}(x) = 0$, but $y = 1$



LOGISTIC REGRESSION – INTUITION BEHIND THE OBJECTIVE



$$\text{cost } (h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$



If $y = 0$

- Cost = 0 if prediction is correct
- As $(1 - h_{\theta}(x)) \rightarrow 0$, cost $\rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties

LOGISTIC REGRESSION – REGULARIZED LOGISTIC REGRESSION



$$J(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

- We can regularize logistic regression exactly as before:

$$\begin{aligned} J_{\text{regularized}}(\boldsymbol{\theta}) &= J(\boldsymbol{\theta}) + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2 \\ &= J(\boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2 \end{aligned}$$



LOGISTIC REGRESSION – GRADIENT DESCENT



$$J_{\text{reg}}(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update
for $j = 0 \dots d$



Use the natural logarithm ($\ln = \log_e$) to cancel with the $\exp()$ in $h_{\boldsymbol{\theta}}(\mathbf{x})$

LOGISTIC REGRESSION – GRADIENT DESCENT



$$J_{\text{reg}}(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^n \left(h_{\boldsymbol{\theta}} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[\sum_{i=1}^n \left(h_{\boldsymbol{\theta}} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} + \lambda \theta_j \right]$$



LOGISTIC REGRESSION – GRADIENT DESCENT



- Initialize θ
- Repeat until convergence (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[\sum_{i=1}^n \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} + \lambda \theta_j \right]$$



This looks IDENTICAL to linear regression!!!

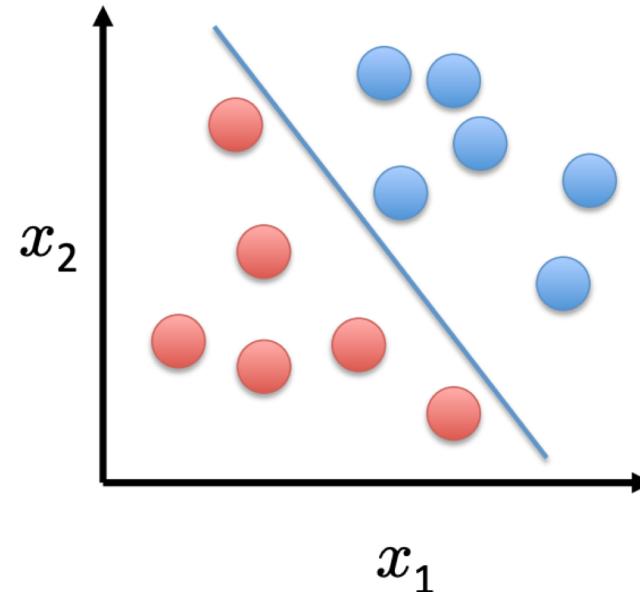
- Ignoring the $1/n$ constant
- However, the form of the model is very different:

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

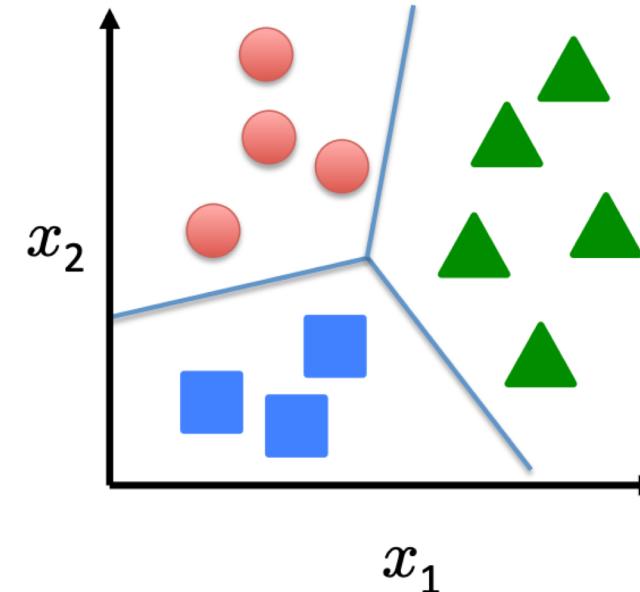
LOGISTIC REGRESSION – MULTI CLASS CLASSIFICATION



Binary classification:



Multi-class classification:



Disease diagnosis: healthy / cold / flu / pneumonia

Object classification: desk / chair / monitor / bookcase

LOGISTIC REGRESSION – MULTI CLASS LOGISTIC REGRESSION

- For 2 classes:



$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)} = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}$$

weight assigned
to $y = 0$ weight assigned
to $y = 1$

- For C classes $\{1, \dots, C\}$:



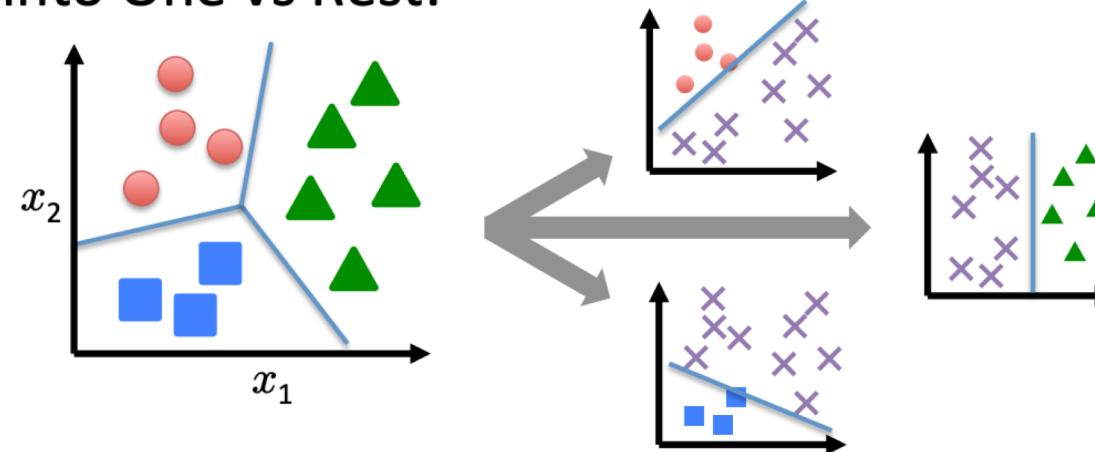
$$p(y = c \mid x; \theta_1, \dots, \theta_C) = \frac{\exp(\theta_c^T x)}{\sum_{c=1}^C \exp(\theta_c^T x)}$$

– Called the **softmax** function

LOGISTIC REGRESSION – MULTI CLASS LOGISTIC REGRESSION



Split into One vs Rest:



- Train a logistic regression classifier for each class i to predict the probability that $y = i$ with

$$h_c(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\boldsymbol{\theta}_c^\top \mathbf{x})}$$

LOGISTIC REGRESSION – MULTI CLASS LOGISTIC REGRESSION



- Use $h_c(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\boldsymbol{\theta}_c^\top \mathbf{x})}$ as the model for class c
- Gradient descent simultaneously updates all parameters for all models
 - Same derivative as before, just with the above $h_c(\mathbf{x})$



- Predict class label as the most probable label

$$\max_c h_c(\mathbf{x})$$



Logistic Regression from scratch



LOGISTIC REGRESSION



Logistic regression is an extension on linear regression (both are generalized linear methods). We will still learn to model a line (plane) that models y given X . Except now we are dealing with classification problems as opposed to regression problems so we'll be predicting probability distributions as opposed to discrete values. We'll be using the softmax operation to normalize our logits (XW) to derive probabilities.

Our goal is to learn a logistic model \hat{y} that models y given X .

$$\hat{y} = \frac{e^{XW_y}}{\sum_j e^{XW}}$$

Variable	Description
N	total numbers of samples
\hat{y}	predictions $\in \mathbb{R}^{NX1}$
X	inputs $\in \mathbb{R}^{NxD}$
W	weights $\in \mathbb{R}^{DX1}$



LOGISTIC REGRESSION



LOGISTIC REGRESSION – M





www.keyrus.com

Shriman TIWARI

Tech Lead/Manager Data Science

Mobile: +33 (0)6 49 71 80 68
shriman.tiwari@keyrus.com

KEYRUS