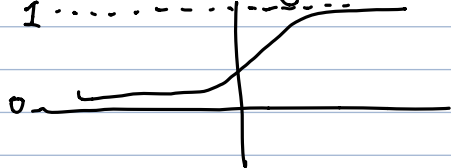


## Topics

- Perceptron
- Exponential Family
- Generalized Linear Methods
- Softmax Regression (multiclass classification)

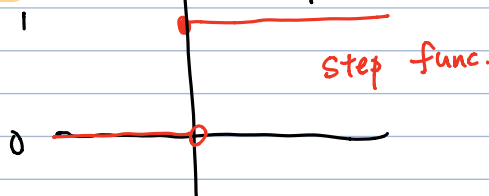
### Logistic Regression (Recap)



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

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

### Perceptron



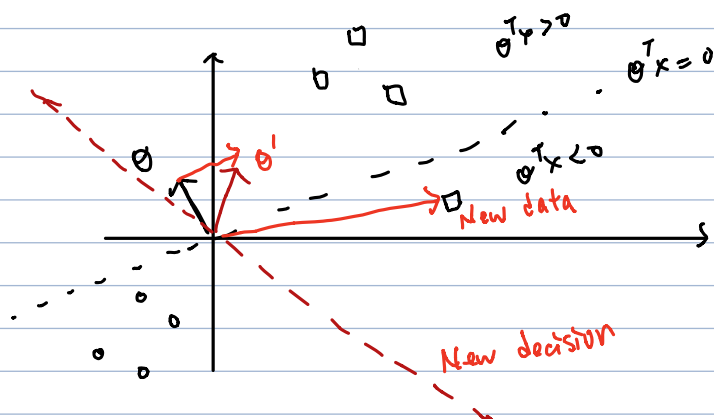
$$g(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

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

$$\theta_j := \theta_j + \alpha (y^{(i)} - \underbrace{h_{\theta}(x^{(i)})}_{\text{diff.}}) x_j^{(i)}$$

$$y^{(i)} - h_{\theta}(x^{(i)})$$

$$\begin{cases} 0: & \text{alg got it right} \\ +1: & \text{wrong } y^{(i)} = 1 \\ -1: & \text{wrong } y^{(i)} = 0. \end{cases}$$



#Note

The order of the data matters.

### Exponential Families

PDF (of  $y$  para by  $\eta$ )

$$p(y; \eta) = b(y) \exp[\eta^T T(y) - a(\eta)] = \frac{b(y) \exp[\eta^T T(y)]}{\underbrace{e^{a(\eta)}}_{\text{Normalizing func.}}}$$

$y$ : data  
 $\eta$ : natural parameter  
 $T(y)$ : sufficient statistic  
 $y$  in class  
 $b(y)$ : Base measure  
 $a(\eta)$ : log-partition function

$y$ : scalar  
 $\eta$ : vector / scalar ↗ match  
 $T(y)$ : vector / scalar ↗  
 $b(y)$ : scalar

### Examples:

• **Bernoulli**: Binary Data

$\phi$ : Probability of event

$$\begin{aligned}
 P(y; \phi) &= \phi^y (1-\phi)^{1-y} \\
 &= \exp(\log(\phi^y (1-\phi)^{1-y})) \\
 &= \exp\left[\underbrace{\log\left(\frac{\phi}{1-\phi}\right)}_{\eta} \underbrace{y}_{T(y)} + \underbrace{\log(1-\phi)}_{-a(\eta)}\right]
 \end{aligned}$$

$$b(y) = 1$$

$$T(y) = y$$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right) \Rightarrow \phi = \frac{1}{1+e^{-\eta}} \quad \text{sigmoid}$$

$$a(\eta) = -\log(1-\phi) = -\log\left(1 - \frac{1}{1+e^{-\eta}}\right) = \log(1+e^{\eta})$$

• **Gaussian** (w/ fix variance)  $\sigma^2 = 1$

$$\begin{aligned}
 P(y; \mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right) \\
 &= \underbrace{\frac{1}{\sqrt{2\pi}} e^{-y^2/2}}_{b(y)} \exp\left(\underbrace{\mu y}_{\eta T(y)} - \underbrace{\frac{1}{2}\mu^2}_{a(\eta)}\right)
 \end{aligned}$$

$$b(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right)$$

$$T(y) = y$$

$$\eta = \mu$$

$$a(\eta) = \frac{\mu^2}{2} = \frac{\eta^2}{2}$$



$\eta = g^{-1}(\eta)$  canonical link fn.

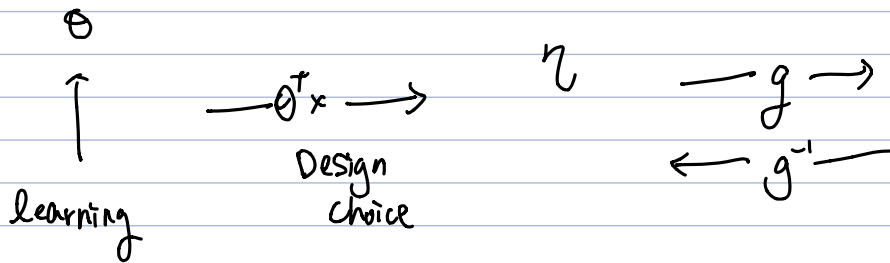
$$g(\eta) = \frac{\partial}{\partial \eta} a(\eta)$$

3 parameters

Model parameters

Natural param.

Canonical param.

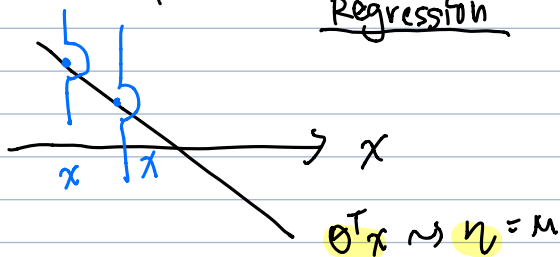


$\phi$ : Bernoulli  
 $\mu, \sigma^2$ : Gaussian  
 $\lambda$ : Poisson

Logistic Regression

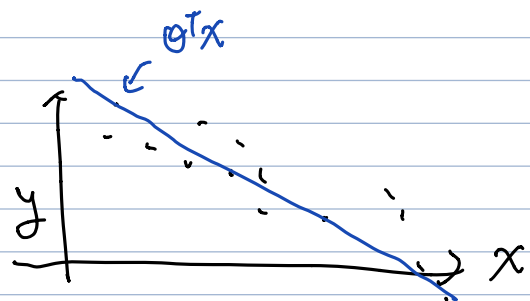
$$h_\theta(x) = E[y|x; \theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}}$$

Assumptions: Linear Regression



corresponds to

Find  $\theta$

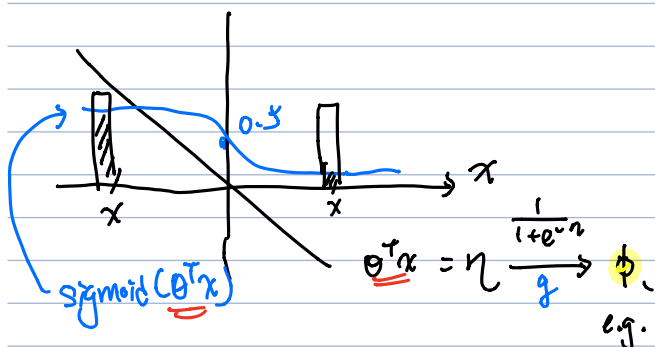


$\forall x$ , corresponding  $y$  is a gaussian with mean  $\mu$ .

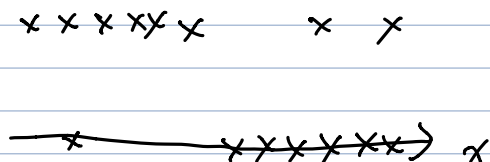
Choose a line, that the means of the distributions lies on the line

Classification

Data



$\theta?$

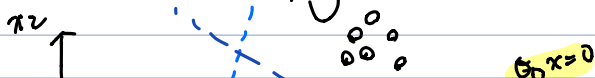


e.g. bernoulli(p).  $p = \phi$

Softmax Regression

(Multiclass regression)

Cross Entropy minimization



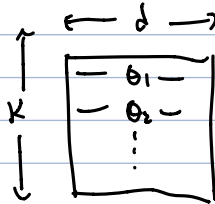


$$x^{(i)} \in \mathbb{R}^d$$

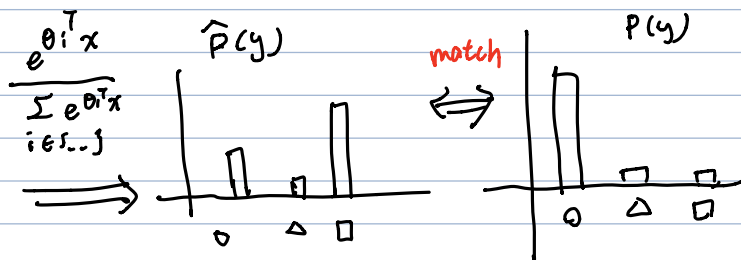
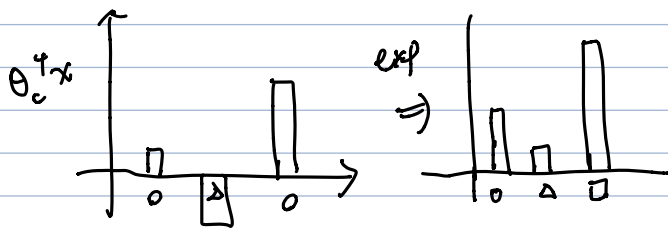
$$y^{(i)} \in \{0, 1\}^x \text{ "one-hot" vector } [0, 0, 1, 0]$$

$$\theta_{\text{class}} \in \mathbb{R}^d$$

$$\text{class} \in \{\Delta, \square, \circ\}$$



$$\theta_{\text{class}}^T x$$



goal: min distance between 2 distributions

$$\min \text{Cross Entropy } (p, \hat{p}) : - \sum_{y \in \{\Delta, \square, \circ\}} p(y) \log(\hat{p}(y))$$

$$= - \log \hat{p}(y_0)$$

$$= - \log \frac{e^{\theta_0^T x}}{\sum_{c \in \{\Delta, \square, \circ\}} e^{\theta_c^T x}}$$

Find  $\theta_0, \theta_\Delta, \theta_\square$  through gradient descent.

