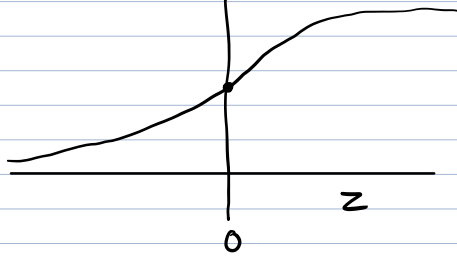


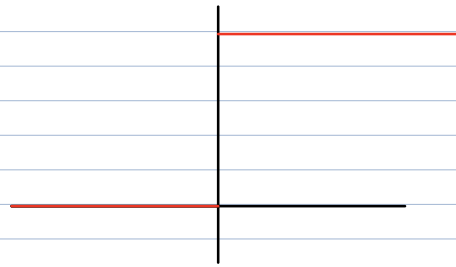
- Perceptron
- Exponential Family
- Generalized Linear Models
- Softmax Regression (Multiclass Classification)

## Logistic Regression



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

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



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

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

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

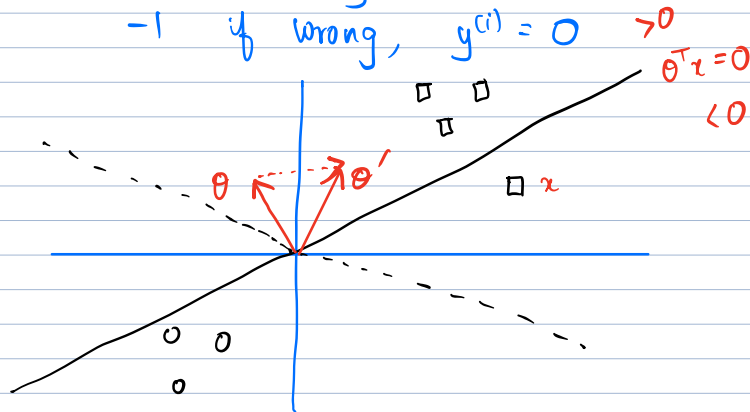
different

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

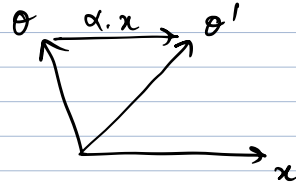
0 : algorithm got it right

+1/-1 : +1 if wrong,  $y^{(i)} = 1$

-1 if wrong,  $y^{(i)} = 0$



$$\begin{aligned}\theta' &= \theta + \alpha x \\ \theta'^T x &= (\theta + \alpha x)^T x \\ &= \theta^T x + \underbrace{\alpha \cdot x^T x}_{>0}\end{aligned}$$



## Exponential Families

PDF

$$p(y; \eta) = b(y) \exp[\eta^T T(y) - a(\eta)]$$

$y$ : data

$\eta$ : natural parameter

$T(y)$ : sufficient statistic  
:  $y$  today

$b(y)$ : Base measure

$a(\eta)$ : log-partition function

$$p(y; \eta) = \frac{b(y) \exp(\eta^T T(y))}{e^{a(\eta)}}$$

$y$ : scalar

$\eta$ : vector / scalar

$T(y)$ : ———

$b(y)$ : scalar

## Bernoulli (Binary Data)

$\phi$ : probability of event

$$p(y; \phi) = \phi^y (1-\phi)^{1-y}$$

$$= \exp(\log(\phi^y (1-\phi)^{1-y}))$$

$$= \exp \left[ \underbrace{\log\left(\frac{\phi}{1-\phi}\right)}_{\eta} \underbrace{y}_{T(y)} + \underbrace{\log(1-\phi)}_{a(\eta)} \right]$$

$$b(y) = 1$$

$$\tau(y) = y$$

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

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

Gaussian (w. fixed variance) Assume  $\sigma^2 = 1$

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

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

$$\tau(y) = y$$

$$\eta = \mu$$

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

Properties ↙ natural parameter

(a) MLE w.r.t.  $\eta$  is concave

negative Log Likelihood (NLL) is convex

(b)  $E[y; \eta] = \frac{\partial}{\partial \eta} a(\eta)$

(c)  $\text{Var}[y; \eta] = \frac{\partial^2}{\partial \eta^2} a(\eta)$

# GLM

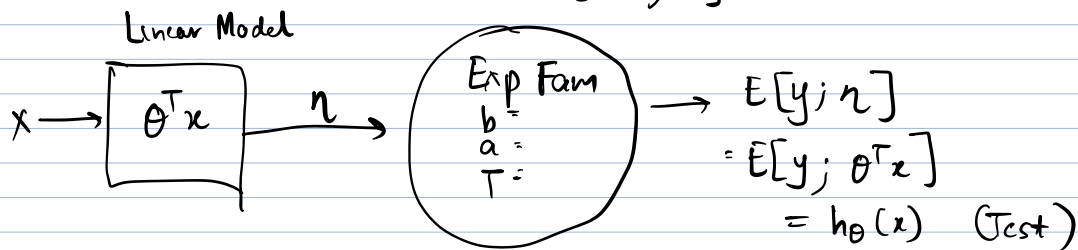
## Assumptions / Design Choices

(i)  $y|x; \theta \sim$  Exponential family

Real — Gaussian  
 Binary — Bernoulli  
 Count — Poisson  
 $\mathbb{R}^+$  — Gamma, Exponential  
 Dist<sup>n</sup> — Beta, Dirichlet

(ii)  $\eta = \theta^T x$        $x = \mathbb{R}^d$   
                                   $\theta = \mathbb{R}^d$

(iii) Test time: Output  $E[y|x; \theta]$   
 $h_\theta(x) = E[y|x; \theta]$



$$\max_{\theta} \log p(y^{(i)}; \theta^T x^{(i)}) \quad (\text{Train})$$

Learning Update Rule

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

↑  
 plug in appropriate  $h_\theta(x)$

Terminology:  $\eta$ : natural parameter

$$\mu = E[y; \eta] = g(\eta) \rightarrow \text{canonical response fn}$$

$$\eta = g^{-1}(\mu) \rightarrow \text{link fn}$$

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

### 3 parametrizations

Model parameters

Natural param

Canonical param

$\theta$   
↑

learning

$-\theta^T x \rightarrow \eta$

Design  
choice

$\rightarrow g \rightarrow$   
 $\leftarrow g^{-1} \leftarrow$

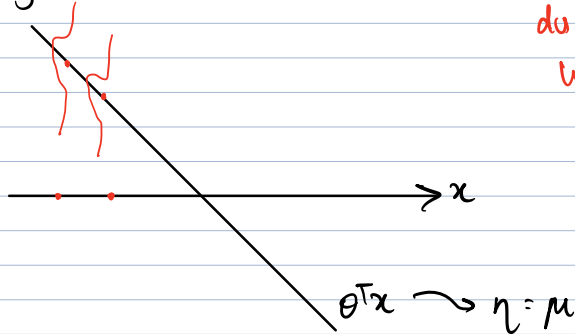
$\phi$ : Bernoulli  
 $\mu, \sigma^2$ : Gaussian  
 $\lambda$ : Poisson

### Logistic Regression

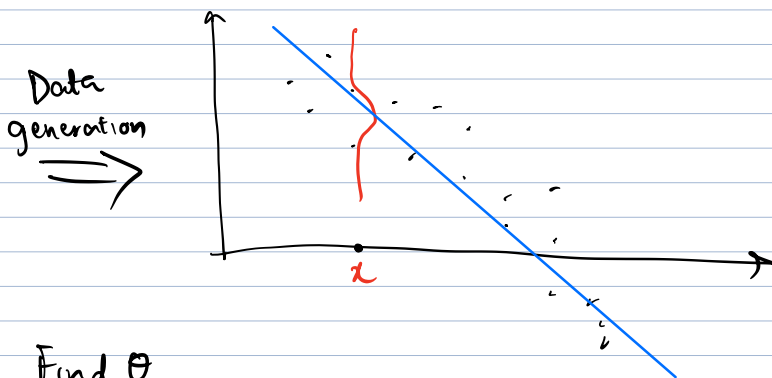
$$h_{\theta}(x) = E[y|x; \theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}}$$

### Assumptions

Regression



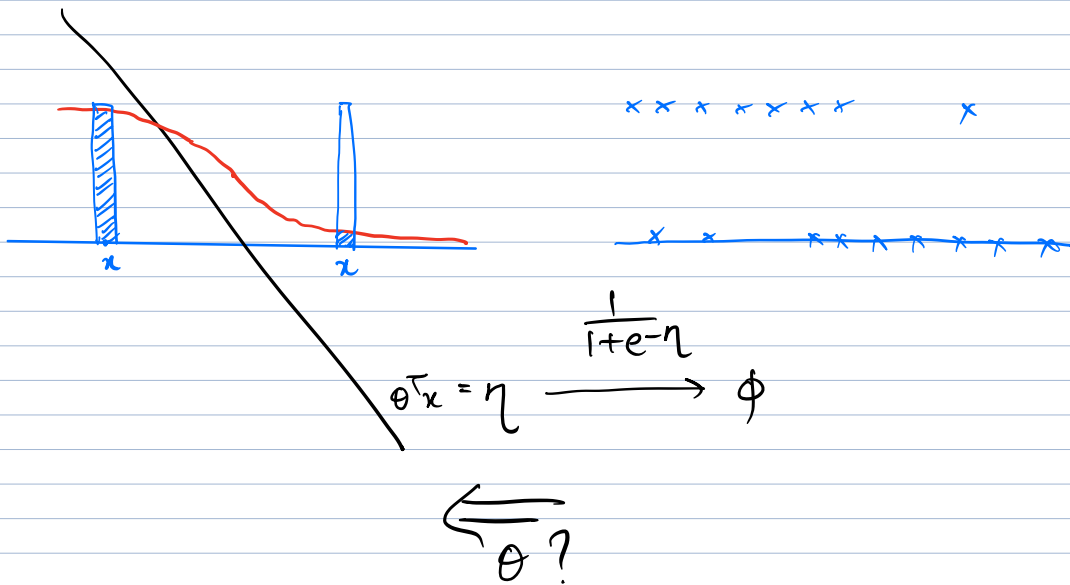
Assumption:  $y$  for any  $x$   
distributed as gaussian  
w. mean  $\mu = \theta^T x$



Find  $\theta$

$\Leftarrow$

## Classification

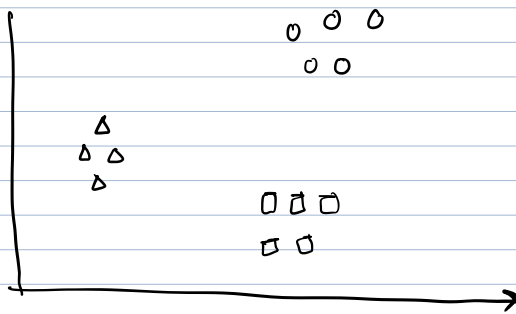


## Softmax Regression

Member of GLM family

Cross Entropy Minimization

Multiclass Classification



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

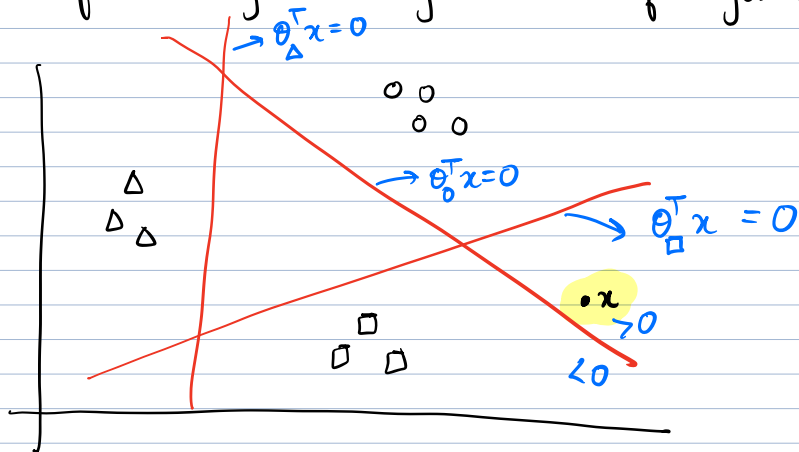
Label  $y \in \{0, 1\}^K$  eg.  $[0, 0, 1, 0]$  "one hot vector"

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

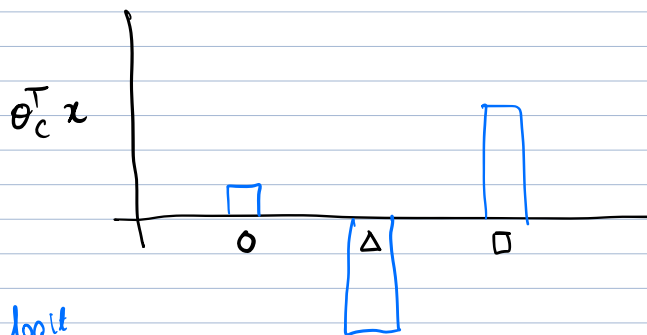
$$\text{class} \in \{\Delta, \circ, \square, \dots\}$$

$$\begin{matrix} \uparrow \\ k \\ \downarrow \end{matrix} \begin{bmatrix} \text{---} \theta_1 \text{---} \\ \text{---} \theta_2 \text{---} \end{bmatrix}^d$$

Softmax regression generalization of logistic regression

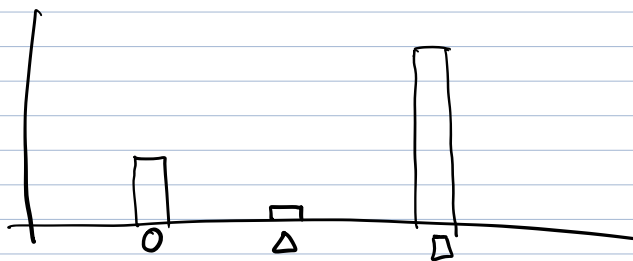


Given  $x$



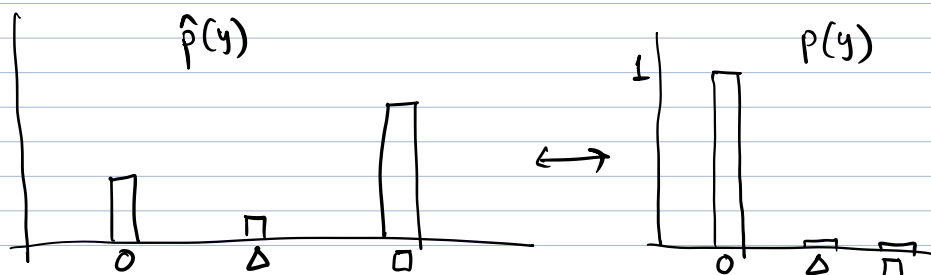
logit space

$\Rightarrow$  exp



normalize

$$\Rightarrow \frac{e^{\theta_c^T x}}{\sum e^{\theta_i^T x}}$$



goal: min distance between distrib<sup>n</sup>

min cross entropy

$$\text{Cross Entropy}(p, \hat{p}) = - \sum_{y \in \{0, \Delta, 0\}} p(y) \log \hat{p}(y)$$

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

$$= - \log \frac{e^{\theta_0^T x}}{\sum_{i \in \{\Delta, 0, 0\}} e^{\theta_i^T x}}$$



Gradient descent