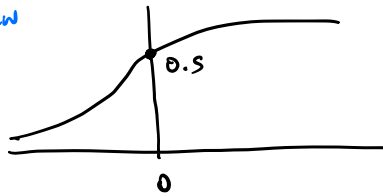


## Perceptron

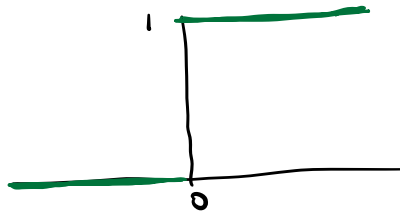
Log key review



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

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

Perceptron function



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

Algorithm:

Repeatedly:

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

0: algo got it right

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

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

Based on this: update decision boundary.

## Exponential families.

Exponential Prob. Density fn (PDF): standard form

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

$y$ : data

$\eta$ : natural parameter

$T(y)$ : sufficient statistic

$b(y)$ : Base measure

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

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

### Bernoulli (for binary data)

$\phi$ : probability of event

$$\begin{aligned} p(y; \phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp(\log(\phi^y (1-\phi)^{1-y})) \\ &= \exp(\underbrace{\log(\phi)}_{\eta} y + \underbrace{\log(1-\phi)}_{a(\eta)}) \end{aligned}$$

$$b(y) = 1; T(y) = y$$

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

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

### Gaussian (w. fixed variance)

Assume  $\sigma^2 = 1$

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

$$b(y) = \frac{1}{\sqrt{2\pi}} e^{-y^2/2}$$

$$T(y) = y$$

$$\eta = \mu$$

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

## Properties

1. MLE wrt.  $\eta$  is concave  
negative log likelihood (NLL) is convex
2.  $E[y; \eta] = \frac{d}{d\eta} a(\eta)$
3.  $\text{Var}[y; \eta] = \frac{d^2}{d\eta^2} a(\eta)$

## Generalized Linear Models

### Assumptions / Design choices

- 1)  $y|x; \theta \sim \text{Exponential family}$

Real - Gaussian

Binary - Bernoulli

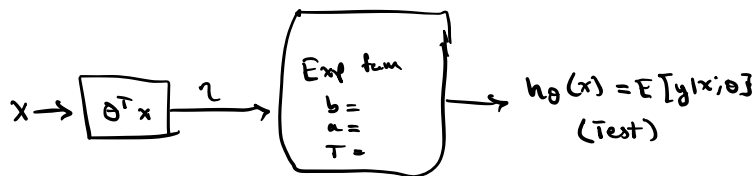
Count - Poisson

$\mathbb{R}^+$  - Gamma, Exponential

$\text{Dist}^\wedge$  - Beta, Dirichlet

$$2) \eta = \theta^T x \quad \begin{array}{l} x \in \mathbb{R}^d \\ \theta \in \mathbb{R}^d \end{array}$$

- 3) Test time: output  $h_\theta(x) = E[y|x; \theta]$



Goal: train  $\theta$  to predict parameter of exp fam  
whose mean is hypothesis of parameter to be output

At train time:

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

Learning update rule:

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

Terminology:  $\eta$  natural parameter

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

$$\eta = g^{-1}(\mu)$$

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

3 parametrizations

Model parameters

$\theta$

$\uparrow$

learning

$\longleftrightarrow$

$\theta^T x$

Design choice

Natural param

$\eta$

$\xrightarrow{\eta}$

$\xleftarrow{g^{-1}}$

Canonical param

$\Phi$ : Bernoulli

$\mu, \sigma^2$ : Gaussian

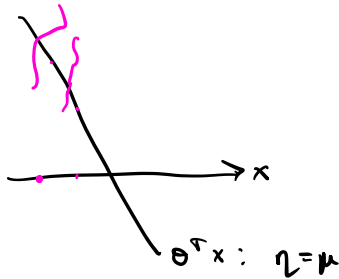
$\lambda$ : poisson

Logistic regression in GLM framework

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

Assumptions

Regression

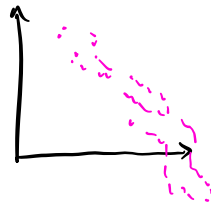


Data generation

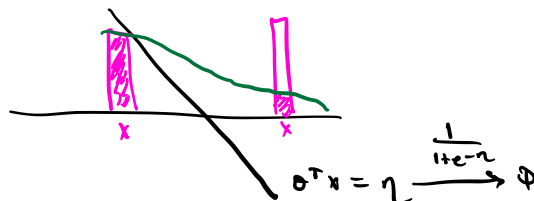
$\Rightarrow$

find  $\theta$

$\Leftarrow$



Classification

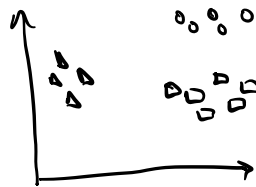


# Softmax Regression

Member of GLM family

Cross Entropy Minimization

Multiclass Classification



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

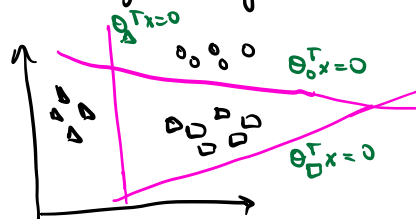
Label  $y \in \{0, 1\}^k$  i.e.  $[0, 0, 1, 0]$  "one-hot vector"

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

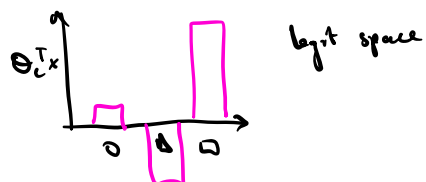
class  $\in \{0, 1, 2, \dots, k-1\}$   
 $\leftarrow d \rightarrow$

$$\begin{matrix} \uparrow k \\ \downarrow \end{matrix} \begin{bmatrix} \leftarrow \Theta_0 \rightarrow \\ \leftarrow \Theta_1 \rightarrow \\ \vdots \end{bmatrix}$$

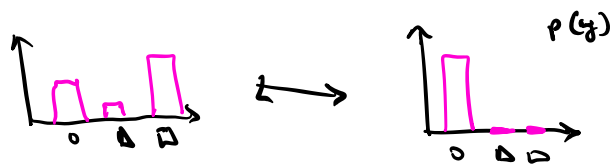
Softmax regression: generalization of logistic regression



Given  $x$



once exp'd, normalized



Goal: minimize dist. between distributions  
 by minimizing cross entropy

## Cross Entropy

$$\begin{aligned}\text{Cross Entropy}(y, \hat{p}) &= - \sum_{y \in \{0,1,D\}} p(y) \log(\hat{p}(y)) \\ &= - \log \hat{p}(y_0) \\ &= - \log \frac{e^{\theta_0^T x}}{\sum_{i \in \{0,1,D\}} e^{\theta_i^T x}}\end{aligned}$$

Minimize via gradient descent.