# Exponential family models

- Definition & motivation

- Examples

- Softmax (Multiclass Classification)

Unify INFERENCE & LEARNING FOR MANY MODELS

# Exponential family

PDF. __IDEA__ "If P has special form $\Rightarrow$ some questions for free"

$$p(\underset{\nearrow \text{DATA}}{y}; \eta) = b(y) \exp[\eta^T T(y) - a(\eta)]$$

$\eta \longrightarrow$ NATURAL PARAMETERS

$T(y)$ IS called suffient STATISTIC ( we'll use $T(y) = y$ in class)

IS SAME dim AS $\eta$

$b(y)$ IS Called BASE MEASURE. Does __not__ depend on $\eta$

$a(\eta)$ IS Called les PARTITION function. Does __not__ depend on $y$

$\Rightarrow$ IT MAKES SURE $p$ is probability function

$y, a(\eta), b(y)$ ARE __SCALARS__

$\eta, T(y)$ ARE __SAME DIMENSION__

# Examples

### Bernoulli $\phi$ IS probability of an event

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

$$= \exp\left(y \log \phi + (1-y) \log (1-\phi)\right)$$

$$= \exp\left(\log \frac{\phi}{1-\phi} \cdot y + \log (1-\phi)\right)$$

CHECK fits into form:

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

$$T(y) = y \qquad \eta = \log \frac{\phi}{1-\phi} \qquad b(y) = 1$$

Claim! $\quad -a(\eta) = \log (1-\phi)$

OBSERVE: $\quad \eta = \log \frac{\phi}{1-\phi} \implies \phi = \overline{1 + e^{-\eta}}$

Here, $\quad 1-\phi = \dfrac{e^{-\eta}}{1+e^{-\eta}} = \dfrac{1}{1+e^{\eta}} \qquad$ so $-\log (1-\phi) = \log (1 + e^{-\eta})$ $\square$.

Example #2    Gaussian (w/ fixed variance) $\sigma^2 = 1$

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

$$= \underbrace{\frac{1}{\sqrt{2\pi}} e^{-y^2/2}}_{b(y)} \exp\left(\mu y - \tfrac{1}{2}\mu^2\right)$$

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

$$\eta = \mu \qquad T(y) = y \quad \text{and} \quad a(\eta) = \tfrac{1}{2}\eta^2 \quad \checkmark$$

# Why do we care about this form?

<u>Inference is "easy"</u>

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

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

<u>learning is "well defined"</u>

MLE wrt to $\eta$ is <u>CONCAVE</u>

(so negative log likelihood is convex)

# Generalized Linear Models (GLM)

Design choices $\Rightarrow$ Assumptions.

(1) $y \mid x; \theta \sim$ Exponential family

$$\text{Binary} \longrightarrow \text{Bernoulli}$$
$$\text{Real} \longrightarrow \text{Gaussian}$$
$$\text{Counts} \longrightarrow \text{Poisson}$$
$$\mathbb{R}^+ \longrightarrow \text{Gamma, Exponential}$$
$$\text{Distribution} \longrightarrow \text{Dirichlet}$$

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

(iii)   Inference @ test time

Output $\mathbb{E}[y \mid x; \theta]$   i.e. $h_\theta = \mathbb{E}[y \mid x; \theta]$

DATA     linear model

$x \longrightarrow \boxed{\theta^T x} \xrightarrow{\eta} \boxed{\begin{array}{l}\text{exp} \\ \text{model} \\ b = \\ \vdots\end{array}}$

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

$\mathbb{E}[y; \eta] = \mathbb{E}[y; \theta^T x]$

$= h_\theta(x)$ (Predict)

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

# Terminology

| Model parameter | | Natural parameter | | Canonical |
|---|---|---|---|---|

$$\Theta \quad \xrightarrow{\Theta^T x} \quad \eta \quad \xrightarrow[\xleftarrow{g^{-1}}]{g}$$

train on true

$\phi$ : Bernoulli

$\mu, \sigma^2$ : Gaussian

$\lambda$ : Poisson

$g$ is called the canonical response function
$g^{-1}$ " " the link function

$$\mu = \mathbb{E}[y ; \eta] \triangleq g(\eta)$$

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

## logistic regression (Bernoulli)

<span style="color:blue">canonical</span>   <span style="color:blue">natural</span>   <span style="color:blue">model</span>

$$h_\theta(x) = \mathbb{E}[y \mid x ; \theta] = \phi = \frac{1}{1+e^{-\eta}} = \frac{1}{1+e^{-\theta^T x}} \in [0,1)$$

### use for classification?

$$h_\theta(x) > 0.5 \implies \text{yes} \quad 1$$
$$\text{o. w.} \implies \text{No} \quad 0$$

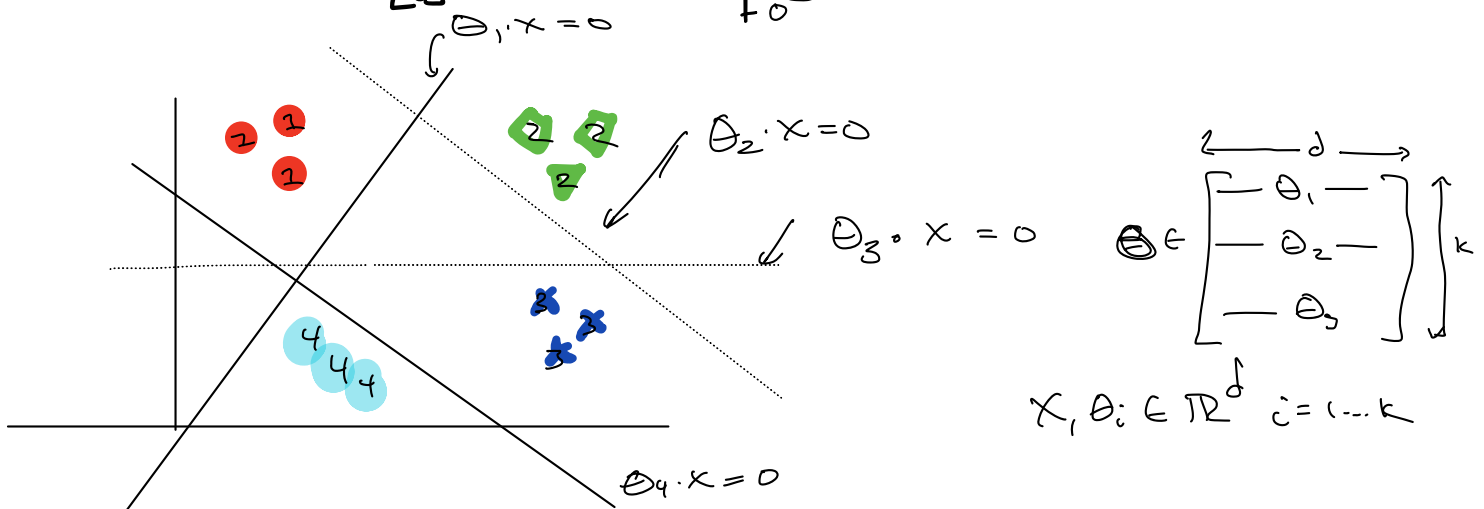## linear regression (Gaussian fixed variance)

$$h_\theta(x) = \mathbb{E}[y \mid x ; \theta] = \mu = \eta = \theta^T x \quad \text{as before!}$$

# Multiclass via Softmax (Multinomial)

① Discrete values up to $k$   {CAT, DOG, CAR, BUS}  $k=4$.

Encoded as one-hot vector $\Rightarrow$ $y \in \{0,1\}^k$

E.g. $k=3$ $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ is class 1 (CAT) $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \leftarrow 0$ is class 3 (CAR)
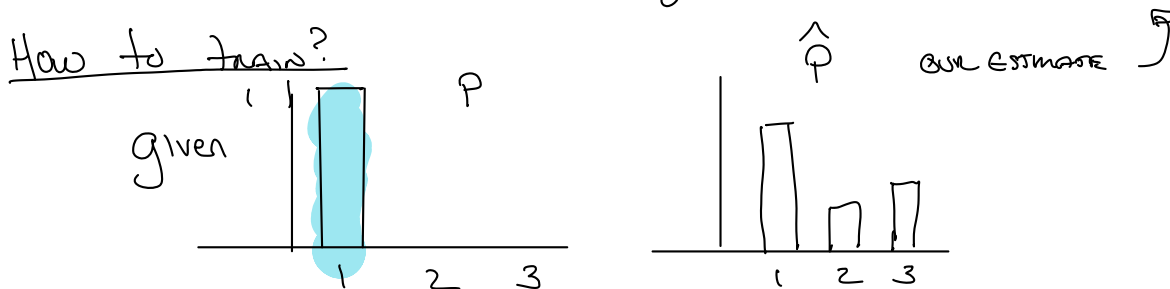
$\theta_1 \cdot x = 0$

$\theta_2 \cdot x = 0$

$\theta_3 \cdot x = 0$

$\theta_4 \cdot x = 0$

$\theta \in \begin{bmatrix} - \theta_1 - \\ - \theta_2 - \\ - \theta_3 - \end{bmatrix} \Big\} k$

$\overset{\longleftarrow \ d \ \longrightarrow}{}$

$x, \theta_i \in \mathbb{R}^d \quad i = 1 \ldots k$

e.g.

$\theta_1 \cdot x = 0.7$    Convert to p. dist.

$\theta_2 \cdot x = -0.5$ $\Longrightarrow exp \Rightarrow$

$\theta_3 \cdot x = -0.1$

$e^{0.7} \approx 2.013$     normalize  $0.57$

$e^{-0.5} \approx 0.606 \Rightarrow 0.17$

$e^{-0.1} \approx 0.904 \Rightarrow .256$

$$P(y=k \mid x; \theta) = \frac{\exp(\theta_k \cdot x)}{\sum_{j=1 \ldots k} \exp(\theta_j \cdot x)}$$

## How to train?

given $P$

"the label is 1"

$\hat{q}$    our estimate $\hat{q}$

1  2  3

1  2  3

$$\min \quad \text{Cross Entropy}(p, \hat{p}) = - \sum_{y=1}^{k} p(y) \log(\hat{p}(y))$$

ground truth is $i$

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

ground truth

$$= -\log \frac{\exp(\theta_i \cdot x)}{\sum_{j=1} \exp(\theta_j \cdot x)}$$

Just Do gradient descent