

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"

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

DATA
↓
NATURAL PARAMETERS

$T(y)$ is called sufficient statistic (we'll see $T(y) = y$ in class)

is same dim as η

$b(y)$ is called base measure. Does not depend on η

$a(\eta)$ is called log partition function. Does not depend on y

\Rightarrow it makes sure P is probability functions

$y, a(\eta), b(y)$ are SCALARS

$\eta, T(y)$ are SAME DIMENSION

Examples

Bernoulli ϕ is probability of an event

$$\begin{aligned} p(y; \phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp\left(y \log \phi + (1-y) \log (1-\phi)\right) \\ &= \exp\left(\log \frac{\phi}{1-\phi} \cdot y + \log (1-\phi)\right) \end{aligned}$$

Check fits into form:

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

$$\tau(y) = y \quad \eta = \log \frac{\phi}{1-\phi} \quad b(y) = 1$$

$$\text{Claim: } -a(\eta) = \log(1-\phi)$$

$$\text{Observe: } \eta = \log \frac{\phi}{1-\phi} \Rightarrow \phi = \frac{1}{1+e^{-\eta}}$$

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

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

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

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

$$P(y|\eta) = b(y) \exp[\eta^\top \tau(y) - a(\eta)]$$

$$\eta = \mu \quad \tau(y) = y \quad \text{and} \quad a(y) = \frac{1}{2}y^2 \quad \checkmark$$

X

Why do we care about this form?

Inference is "easy"

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

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

Learning is "well defined"

MLE w.r.t to η is CONCAVE

(so negative log likelihood is convex)

Generalized Linear Models (GLM)

Design choices \rightarrow Assumptions.

$$(i) \quad y|x; \theta \sim \text{Exponential family}$$

Binary \rightarrow Bernoulli

Real \rightarrow Gaussian

Counts \rightarrow Poisson

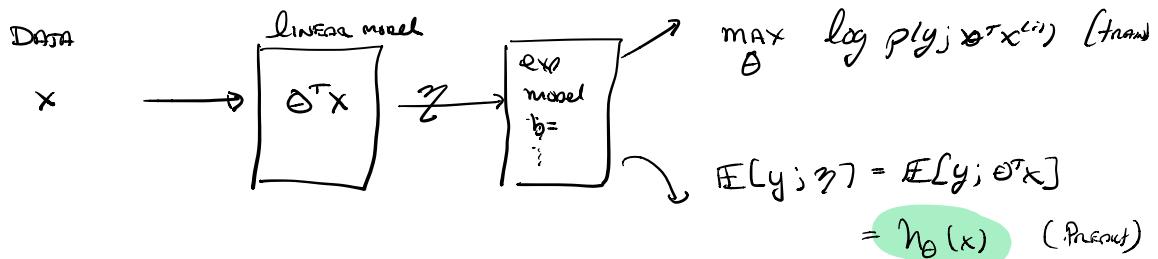
\mathbb{R}^+ \rightarrow Gamma, Exponential

Distribution \rightarrow Dirichlet

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

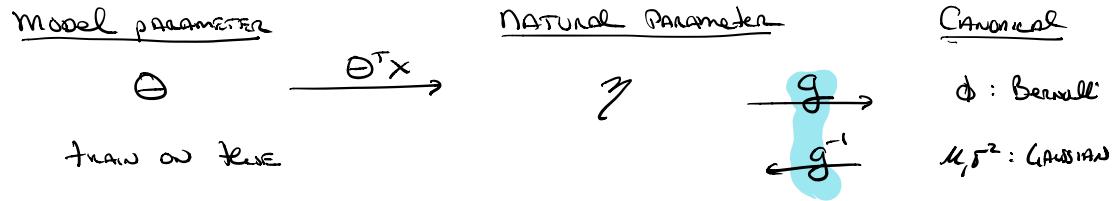
(iii) Define @ test time

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



$$\text{learning} \quad \Theta_j := \Theta_j + \alpha (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}$$

Terminology



g is called the Canonical response function

g^{-1} " the link function

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

$$\Rightarrow \frac{\partial \mu}{\partial \eta} = g'(\eta)$$

λ : poisson

ϕ : Bernoulli

logistic regression (Bernoulli)

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

use for classification?

$$h_\theta(x) > 0.5 \Rightarrow \text{yes} \quad 1$$

$$0.5 \leq h_\theta(x) < 0.5 \Rightarrow \text{no} \quad 0$$

linear regression (GAUSSIAN fixed variance)

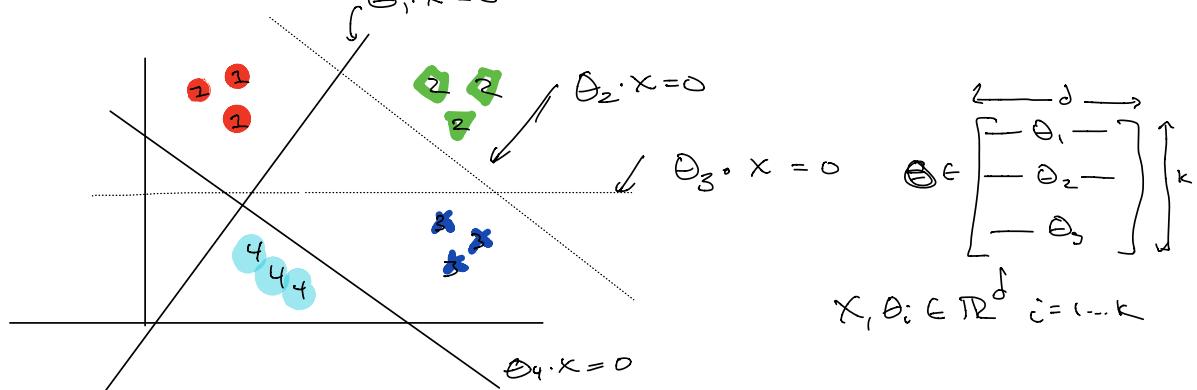
$$h_\theta(x) = \mathbb{E}[y | x; \theta] = \mu = \eta = \theta^T x \text{ as before}$$

Multiclass via SOFTMAX (Multinomial)

DISCRETE VALUES UP TO K $\{ \text{car}, \text{dog}, \text{cat}, \text{bus} \}$ $K=4$.

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 (car) $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$ is class 3 (car)



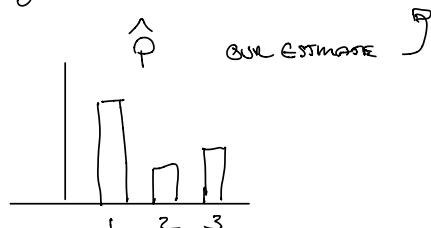
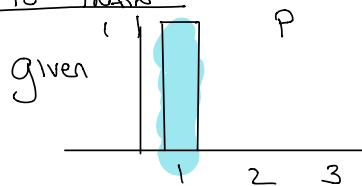
E.g. $\theta_1 \cdot x = 0.7$ Convert to probability $e^{0.7} \approx 2.013 \xrightarrow{\text{normalize}} 0.57$

$$\theta_2 \cdot x = -0.5 \implies e^{-0.5} \approx 0.606 \Rightarrow 0.17$$

$$\theta_3 \cdot x = -0.1 \implies e^{-0.1} \approx 0.904 \Rightarrow 0.256$$

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

How to train?



"the label is 1"

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

ground truth is i

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

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

Just do gradient descent