

Exponential family models

- Definition & motivation
- Examples
- Softmax (Multiclass Classification)

Unify INFERENCE &

LEARNING for many
models

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Exponential family

PDF. IDEA: "If P has special form \Rightarrow some questions for free"

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

DATA
 ↗
 ↘ NATURAL PARAMETERS

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

is same dim as γ

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

$a(n)$ is called the partition function. Does not depend on n

Assignment Project Exam Help

$\eta_{T(4)}$ Are ~~SAME~~ Dimensional
<https://powcoder.com>

Add WeChat powcoder

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:

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

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

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

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

Assignment Project Exam Help

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

<https://powcoder.com>

Add WeChat powcoder

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

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

$$= \frac{1}{\sqrt{2\pi}} \underset{b(y)}{\cancel{e^{-\frac{y^2}{2}}}} \exp\left(\cancel{ay} - \frac{1}{2}\mu^2\right)$$

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

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

Assignment Project Exam Help



<https://powcoder.com>

Add WeChat powcoder

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 want to η is concave

(so negative log likelihood is convex)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Generalized Linear Models (GLM)

Design choices \rightarrow Assumptions.

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

Binary \rightarrow Bernoulli

Real \rightarrow Gaussian

Counts \rightarrow Poisson

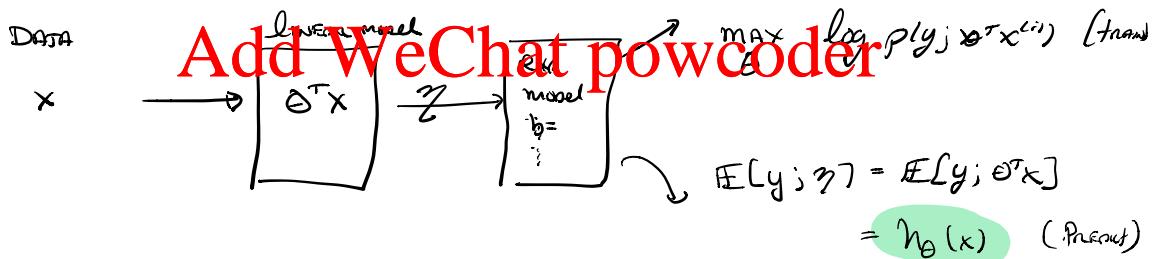
\mathbb{R}^+ \rightarrow Gamma, Exponential

Distribution \rightarrow Dirichlet

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

(iii) influence @ test time

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



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

Terminology

Model parameter

$$\theta$$

known as true

$$\theta^T x$$

Natural Parameter

$$\gamma$$

$$g \rightarrow$$

$$g^{-1}$$

Canonical

ϕ : Bernoulli

$\mu \delta^2$: Gaussian

λ : Poisson

g is called the Canonical response function

g^{-1} " the link function

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

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

logistic regression (Bernoulli)

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

use for classification?

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

Add WeChat powcoder

linear regression (Gaussian fixed variance)

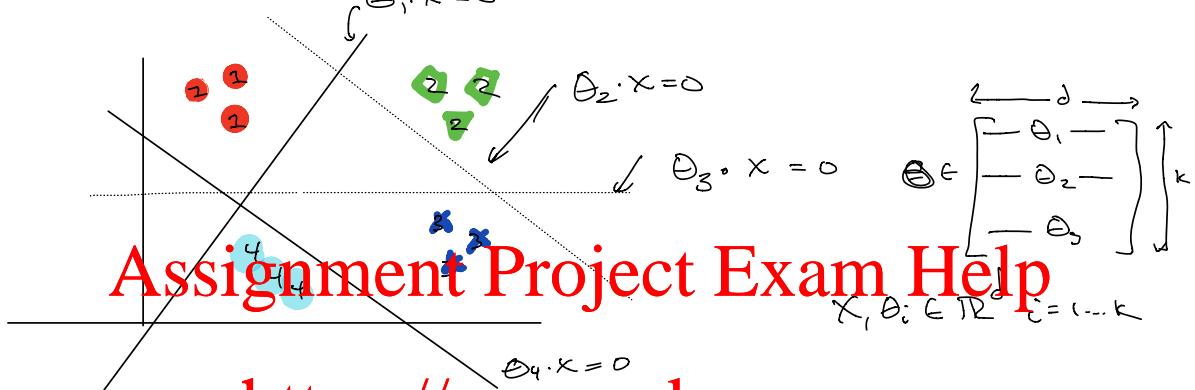
$$h_\theta(x) = \mathbb{E}[y|x; \theta] = \mu = \gamma = \theta^T x \text{ as before}$$

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



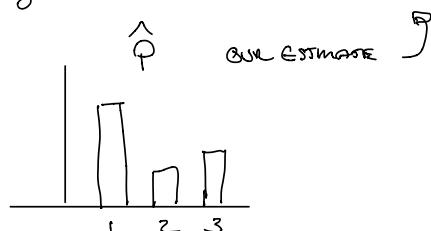
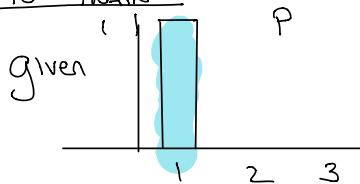
E.g. $\theta_1 \cdot x = -0.7$ Convex to global optimum $e^{-0.7} \approx 0.493$ Normalized 0.57

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

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

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

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

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

ground truth

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

Just do gradient descent

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder