### Exponential Family and Generalized Linear Models

## Piyush Rai IIT Kanpur

#### Probabilistic Machine Learning (CS772A)

Jan 20, 2016

Probabilistic Machine Learning (CS772A)

イロト イポト イヨト イヨト

#### **Generalized Linear Models**

- Models we have seen so far..
  - (Probabilistic) Linear regression: when y is real-valued

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^{\top}\mathbf{x}, \beta^{-1})$$

• Logistic regression: when y is binary (0/1)

$$p(y|x, w) = \text{Bernoulli}(\sigma(w^{\top}x)) = [\sigma(w^{\top}x)]^{y}[1 - \sigma(w^{\top}x)]^{1-y}$$

where 
$$\sigma(\mathbf{w}^{\top}\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^{\top}\mathbf{x})} = \frac{\exp(\mathbf{w}^{\top}\mathbf{x})}{1 + \exp(\mathbf{w}^{\top}\mathbf{x})}$$

- In both, the model depends on the inputs x linearly via  $w^{\top}x$
- Both are special cases of a more general class: Generalized Linear Models

$$p(y|\eta) = h(y) \exp(\eta y - A(\eta))$$

- .. a special type of exponential family distribution
- GLM can be used to also model responses that aren't reals/binary (can be any exponential family distribution in general)

### **Exponential Family Distributions**

• An exponential family distribution is of the form

$$p(y|\eta) = h(y) \exp(\eta^{\top} T(y) - A(\eta))$$

- $\eta$  is called the natural parameter
- h(y) is usually a constant w.r.t.  $\eta$
- T(y) is the sufficient statistics:  $p(y|\eta)$  depends on y only through T(y)
- $A(\eta)$ : log partition function or cumulant function

$$A(\eta) = \log \int h(y) \exp(\eta^{\top} T(y)) dy$$

.. can also be seen as the log of a normalization factor

イロン イヨン イヨン イヨン

• Bernoulli in the usual form:

$$\mathsf{Bernoulli}(y|p) = p^{y}(1-p)^{1-y} = \exp\left(y\log\left(\frac{p}{1-p}\right) + \log(1-p)\right)$$

• Comparing it as  $p(y|\eta) = h(y) \exp(\eta^{\top} T(y) - A(\eta))$ , we have

• h(y) = 1•  $\eta = \log\left(\frac{p}{1-p}\right)$ • T(y) = y•  $A(\eta) = -\log(1-p)$ 

イロン 不良 とうほう 不良 とうほ

• Gaussian in the usual form:

$$\mathcal{N}(y|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{\mu}{\sigma^2}y - \frac{1}{2\sigma^2}y^2 - \frac{\mu^2}{2\sigma^2} - \log\sigma\right)$$

• Comparing it as  $p(y|\eta) = h(y) \exp(\eta^{\top} T(y) - A(\eta))$ , we have

• 
$$h(y) = \frac{1}{\sqrt{2\pi}}$$
  
•  $\eta = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)^T$   
•  $T(y) = (y, y^2)^T$   
•  $A(\eta) = \frac{\mu^2}{2\sigma^2} + \log \sigma$ 

イロト 不得下 イヨト イヨト 二日

- The log partition function  $A(\eta)$  has several useful properties
- First derivative of  $A(\eta)$  w.r.t.  $\eta$  is the expectation of the sufficient statistics

$$rac{dA(\eta)}{d\eta} = \mathbb{E}[T(y)]$$
 (proof done on board)

• Second derivative of  $A(\eta)$  w.r.t.  $\eta$  is the variance of sufficient statistics

$$\frac{d^2 A(\eta)}{d\eta^2} = \operatorname{var}[T(y)]$$

• Note:  $A(\eta)$  is also convex (because second derivative is non-negative)

イロン イロン イヨン イヨン 三日

#### **MLE for Exponential Family Distributions**

• The log-likelihood is given by

$$L(\eta) = \log p(Y|\eta) = \log \prod_{n=1}^{N} p(y_n|\eta) = \log \prod_{n=1}^{N} h(y_n) \exp(\eta^{\top} T(y_n) - A(\eta))$$
$$= \log \prod_{n=1}^{N} h(y_n) + \eta^{\top} (\sum_{n=1}^{N} T(y_n)) - NA(\eta)$$

• Taking derivative w.r.t.  $\eta$  and setting it to zero

$$N\frac{dA(\eta)}{d\eta} = \sum_{n=1}^{N} T(y_n)$$

• Defining  $\mu = \mathbb{E}[T(y)] = \frac{dA(\eta)}{d\eta}$ , we get

 $\hat{\mu}_{MLE} = \frac{1}{N} \sum_{n=1}^{N} T(y_n)$  (can be used for parameter estimation via moment-matching)

• Note that the estimate only depends on data via the sufficient statistics T(y)

(ロ) (四) (主) (主) (三)

#### **Generalized Linear Models**

- An exp. fam. model for  $x \rightarrow y$  is a Generalized Linear Model if:
  - **(**) Observed inputs  $x_n$  enter the model via linear combination  $w^{\top}x_n$
  - **2** Conditional mean of response  $y_n$  depends on  $\boldsymbol{w}^\top \boldsymbol{x}_n$  via a response function f

$$\mu_n = \mathbb{E}[y_n] = f(\mathbf{w}^\top \mathbf{x}_n)$$

- for linear regression  $\mu_n = f(\mathbf{w}^\top \mathbf{x}_n) = \mathbf{w}^\top \mathbf{x}_n$ ,
- for logistic regression  $\mu_n = f(\mathbf{w}^\top \mathbf{x}_n) = \exp(\mathbf{w}^\top \mathbf{x}_n)/(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))$

$$T(y) = y$$

Form of a GLM

$$p(y|\eta) = h(y) \exp(\eta y - A(\eta))$$

where natural parameter  $\eta = \psi(\mu)$ ,  $\mu$ : conditional mean,  $\psi$ : link function



• Note: Some GLM can be represented as  $p(y|\eta, \phi) = h(y, \phi) \exp(\frac{\eta y - A(\eta)}{\phi})$ where  $\phi$  is a dispersion parameter (Gaussian/gamma GLMs use this rep.)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─ のへの

#### **GLM** with Canonical Response Function

- A GLM has a canonical response function f if  $f = \psi^{-1}$
- For such a GLM,  $\eta_n = \psi(\mu_n) = \psi(f(\boldsymbol{w}^\top \boldsymbol{x}_n)) = \boldsymbol{w}^\top \boldsymbol{x}_n$
- E.g., for logistic regression  $\eta_n = \log \frac{\mu_n}{1-\mu_n} = \mathbf{w}^\top \mathbf{x}_n$  (exercise: verify by recalling the exponential family representation of Bernoulli distribution)
- Thus, for Canonical GLMs

$$p(y|\eta) = h(y) \exp(\eta y - A(\eta))$$
  
=  $h(y) \exp(y \mathbf{w}^{\top} \mathbf{x} - A(\eta))$ 

• Such design choices in the canonical GLM make parameter estimation easy

イロン イロン イヨン イヨン 三日

#### **MLE for Generalized Linear Models**

Log likelihood

$$L(\eta) = \log p(Y|\eta) = \log \prod_{n=1}^{N} h(y_n) \exp(y_n \mathbf{w}^\top \mathbf{x}_n - A(\eta_n)) = \sum_{n=1}^{N} \log h(y_n) + \mathbf{w}^\top \sum_{n=1}^{N} y_n \mathbf{x}_n - \sum_{n=1}^{N} A(\eta_n)$$

• Convexity of  $A(\eta)$  guarantees a global optima. Taking derivative w.r.t. **w** 

$$\sum_{n=1}^{N} \left( y_n \mathbf{x}_n - A'(\eta_n) \frac{d\eta_n}{d\mathbf{w}} \right) = \sum_{n=1}^{N} (y_n \mathbf{x}_n - \mu_n \mathbf{x}_n) = \sum_{n=1}^{N} (y_n - \mu_n) \mathbf{x}_n$$

where  $\mu_n = f(\mathbf{w}^{\top} \mathbf{x}_n)$  and 'f'  $(= \psi^{-1})$  depends on type of response y, e.g.,

- Real-valued y (linear regression): f is identity, i.e.,  $\mu_n = \mathbf{w}^\top \mathbf{x}_n$
- Binary y (logistic regression): f is logistic function, i.e.,  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$
- Count-valued y (Poisson regression):  $\mu_n = \exp(\mathbf{w}^\top \mathbf{x}_n)$
- Positive reals y (gamma regression):  $\mu_n = -(\mathbf{w}^\top \mathbf{x}_n)^{-1}$
- To estimate *w*, either set the derivative to zero or use iterative methods (e.g., gradient descent, iteratively reweighted least squares, etc.)

# Next class: Clustering via Gaussian Mixture Models

< 🗇 🕨