

# Canonical Link Functions

Sargur N. Srihari

University at Buffalo, State University of New York  
USA

# Topics in Linear Classification using Probabilistic Discriminative Models

- Generative vs Discriminative
  1. Fixed basis functions in linear classification
  2. Logistic Regression (two-class)
  3. Iterative Reweighted Least Squares (IRLS)
  4. Multiclass Logistic Regression
  5. Probit Regression
  6. Canonical Link Functions

# Canonical Link Functions

- Objective functions in each of three problems

- linear regression (with Gaussian noise)  $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\}^2$

- logistic sigmoid (with cross-entropy)  $E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$

- softmax activation with multiclass cross-entropy

- Have derivatives wrt  $\mathbf{w}$  repeatedly of the form

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

- Used in descent  $\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - \varepsilon \nabla E(\mathbf{w})$   $\mathbf{w}^{(\text{new})} = \mathbf{w}^{(\text{old})} - H^{-1} \nabla E(\mathbf{w})$

- Contribution of data point  $n$  is

- error  $(y_n - t_n)$  times feature vector  $\phi_n$

- This is a general result if we choose:

- Conditional distribution of target is from exponential family

- Activation function is a canonical link function

# Exponential family

- General form of exponential family of distributions is:

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x}) g(\boldsymbol{\eta}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\}$$

where  $\mathbf{x}$  is a vector or scalar, discrete or continuous

$\boldsymbol{\eta}$  are natural parameters of the distribution

$\mathbf{u}(\mathbf{x})$  is some continuous function of  $\mathbf{x}$

$g(\boldsymbol{\eta})$  is a coefficient to ensure distribution is normalized and therefore satisfies

$$g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})\} d\mathbf{x} = 1$$

It is not surprising that Gaussian is exponential

We will see how discrete distributions (Bernoulli and Multinoulli) are exponential

# Bernoulli as an exponential distribution

- $p(x|\mu) = \text{Bern}(x|\mu) = \mu^x(1-\mu)^{1-x}$
- Expressing rhs as exponential of logarithm

$$\begin{aligned} p(x|\mu) &= \exp\{x \ln\mu + (1-x)\ln(1-\mu)\} \\ &= (1-\mu) \exp\{\ln(\mu/1-\mu)x\} \end{aligned}$$

Allows us to identify  $\eta = \ln(\mu/1-\mu)$

Which we can solve to get  $\mu = 1/(1 + \exp(-\eta)) = \sigma(\eta)$

- Thus Bernoulli is  $p(x|\mu) = \sigma(-\eta) \exp(\eta x)$   
i.e., exponential with  $u(x) = x$ ,  $h(x) = 1$ ,  $g(\eta) = \sigma(-\eta)$

# Multinoulli as Exponential

- Using 1-of- $M$  representation, for a single observation

$\mathbf{x} = (x_1, \dots, x_M)^T$  with parameters  $\mu = \{\mu_1, \dots, \mu_M\}$ :

$$p(\mathbf{x} | \mu) = \prod_{k=1}^M \mu_k^{x_k} = \exp \left\{ \sum_{k=1}^M x_k \ln \mu_k \right\}$$

– which can be written as  $p(\mathbf{x} | \boldsymbol{\eta}) = \exp(\boldsymbol{\eta}^T \mathbf{x})$

- where  $\eta_k = \ln \mu_k$  and we have defined  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_M)$

- Since  $\eta_k$  are not independent: They are subject to the constraint which is useful to express multinomial using softmax

$$\sum_{k=1}^M \mu_k = 1$$

$$\begin{aligned} \exp \left\{ \sum_{k=1}^M x_k \ln \mu_k \right\} &= \exp \left\{ \sum_{k=1}^{M-1} x_k \ln \mu_k + \left( 1 - \sum_{k=1}^{M-1} x_k \right) \ln \left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right\} \\ &= \exp \left\{ \sum_{k=1}^{M-1} x_k \ln \frac{\mu_k}{\left( 1 - \sum_{j=1}^{M-1} \mu_j \right)} + \ln \left( 1 - \sum_{k=1}^{M-1} \mu_k \right) \right\} \end{aligned}$$

Identifying

$$\ln \left( \frac{\mu_k}{1 - \sum_j \mu_j} \right) = \eta_k$$

Rearranging

$$\mu_k = \frac{\exp(\eta_k)}{1 + \sum_j \exp(\eta_j)}$$

This is softmax where  $\boldsymbol{\eta} = (\eta_1, \dots, \eta_{M-1}, 0)$

- Exponential form is

$$p(\mathbf{x} | \boldsymbol{\eta}) = \left( 1 + \sum_{k=1}^{M-1} \exp(\eta_k) \right)^{-1} \exp(\boldsymbol{\eta}^T \mathbf{x}) \quad \text{with } \mathbf{u}(\mathbf{x}) = \mathbf{x}, \quad h(\mathbf{x}) = 1, \quad g(\boldsymbol{\eta}) = \left( 1 + \sum_{k=1}^{M-1} \exp(\eta_k) \right)^{-1}$$

- Here we are assuming exponential distribution for target  $t$  rather than input  $\mathbf{x}$

## Exponential form also has similar update

- Restricted form of conditional distribution of target  $t$  in exponential form

$$p(t | \eta, s) = \frac{1}{s} h\left(\frac{1}{s}\right) g(\eta) \exp\left\{\frac{\eta t}{s}\right\}$$

We denote conditional mean of  $t$  by  $y$  which is related to  $\eta$  through  $\eta = \psi(y)$

- Generalized linear model has the form

$y = f(\mathbf{w}^T \phi)$  where  $f$  is the activation function and  $f^{-1}$  is known as the link function

- If we choose canonical link function  $f^{-1}(y) = \eta = \psi(y)$

– The gradient of the error function reduces to

$$\nabla \ln E(\mathbf{w}) = \frac{1}{s} \sum_{n=1}^N (y_n - t_n) \phi_n$$

- For the Gaussian  $s = \beta^{-1}$
- For the Logistic model  $s = 1$