

CS 2750 Machine Learning

Lecture 5

Density estimation III.

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

Outline

Outline:

- **Density estimation:** ✓
 - Maximum likelihood (ML)
 - Bayesian parameter estimates
 - MAP
- **Bernoulli distribution.** ✓
- **Binomial distribution** ✓
- **Multinomial distribution** ✓
- **Normal distribution**
- **Exponential family**

Parametric density estimation

Parametric density estimation:

- A set of random variables $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$
- **A model of the distribution** over variables in X with **parameters** Θ : $\hat{p}(\mathbf{X} | \Theta)$
- **Data** $D = \{D_1, D_2, \dots, D_n\}$

Objective: find parameters Θ such that $p(\mathbf{X} | \Theta)$ describes data D the best

Parameter estimation (learning)

- Maximum likelihood (ML)

$$\Theta_{ML} = \arg \max_{\Theta} p(D | \Theta, \xi)$$

- Bayesian parameter estimation

keep the posterior density $p(\Theta | D, \xi)$

- Maximum a posteriori probability (MAP)

$$\Theta_{MAP} = \arg \max_{\Theta} p(\Theta | D, \xi)$$

- Expected value

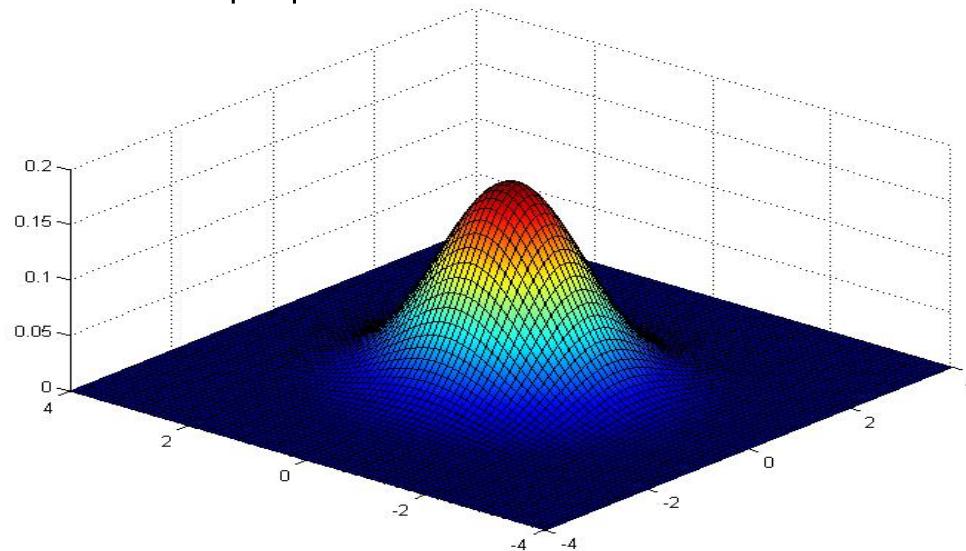
$$\Theta_{EXP} = \int_{\Theta} \Theta p(\Theta | D, \xi) d\Theta$$

Multivariate normal distribution

- **Multivariate normal:** $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$
- **Parameters:** $\boldsymbol{\mu}$ - mean
 Σ - covariance matrix
- **Density function:**

$$p(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right]$$

- **Example:**



Partitioned Gaussian Distributions

- Multivariate Gaussian:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- what are marginals and conditionals?

- Example:

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

$$\boldsymbol{\Lambda} \equiv \boldsymbol{\Sigma}^{-1}$$

$$\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

Precision matrix

Partitioned Conditionals and Marginals

- **Conditional density:**

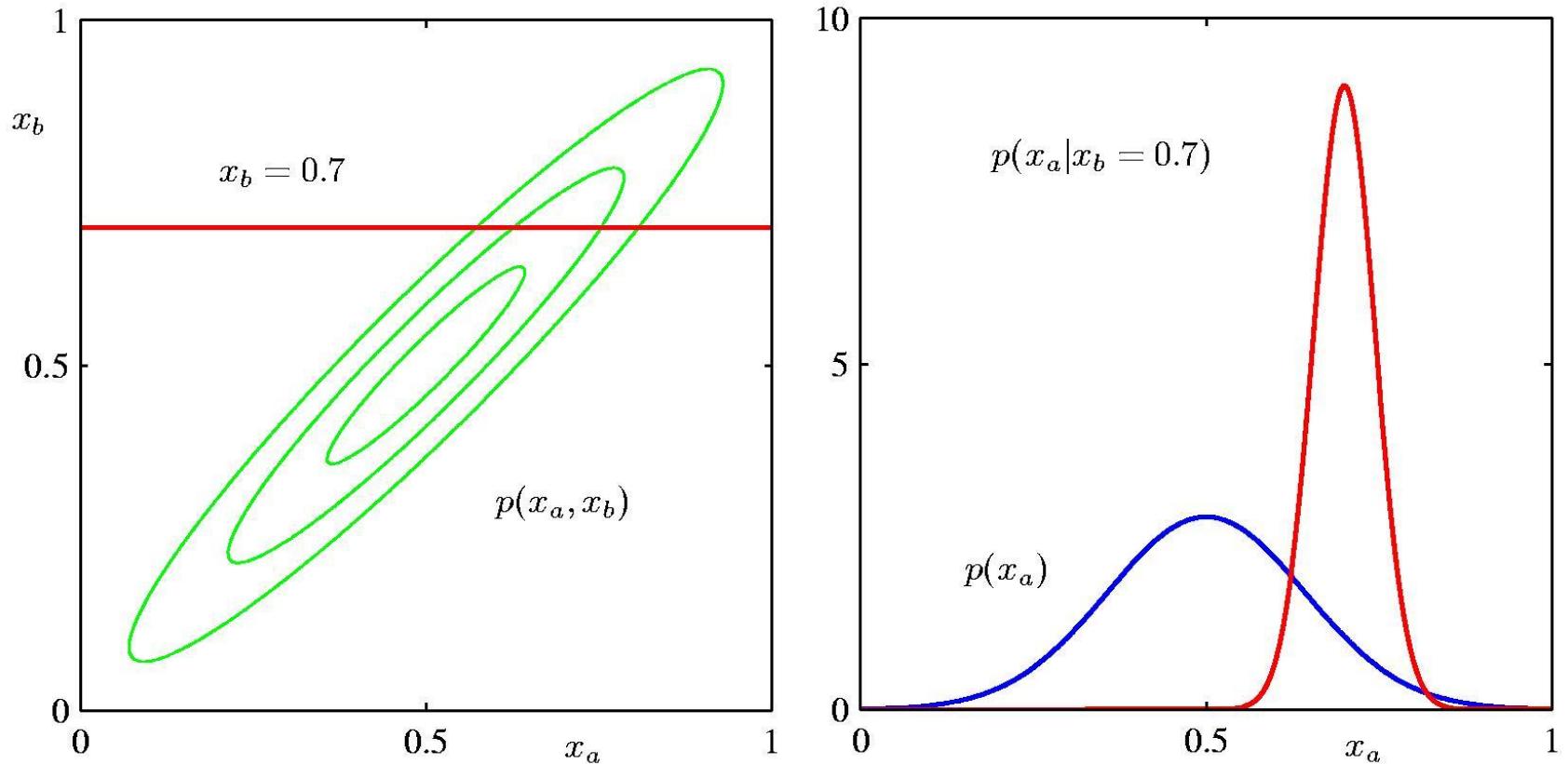
$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$$

$$\begin{aligned}\boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba} \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa}\boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)\end{aligned}$$

- **Marginal Density:**

$$\begin{aligned}p(\mathbf{x}_a) &= \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b \\ &= \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})\end{aligned}$$

Partitioned Conditionals and Marginals



Parameter estimates

- **Loglikelihood**

$$l(D, \mu, \Sigma) = \log \prod_{i=1}^n p(\mathbf{x}_i | \mu, \Sigma)$$

- **ML estimates of the mean and covariances:**

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T$$

- Covariance estimate is biased

$$E_n(\hat{\Sigma}) = E_n \left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T \right) = \frac{n-1}{n} \Sigma \neq \Sigma$$

- **Unbiased estimate:**

$$\hat{\Sigma} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu})(\mathbf{x}_i - \hat{\mu})^T$$

Posterior of a multivariate normal

- Assume a prior on the mean μ that is normally distributed:

$$p(\mu) \approx N(\mu_p, \Sigma_p)$$

- Then the posterior of μ is normally distributed

$$\begin{aligned} p(\mu | D) &\approx \left(\prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \right] \right) \\ &\quad * \frac{1}{(2\pi)^{d/2} |\Sigma_p|^{1/2}} \exp \left[-\frac{1}{2} (\mu - \mu_p)^T \Sigma_p^{-1} (\mu - \mu_p) \right] \\ &= \frac{1}{(2\pi)^{d/2} |\Sigma_n|^{1/2}} \exp \left[-\frac{1}{2} (\mu - \mu_n)^T \Sigma_n^{-1} (\mu - \mu_n) \right] \end{aligned}$$

Posterior of a multivariate normal

- Then the posterior of μ is normally distributed

$$p(\mu | D) = \frac{1}{(2\pi)^{d/2} |\Sigma_n|^{1/2}} \exp \left[-\frac{1}{2} (\mu - \mu_n)^T \Sigma_n^{-1} (\mu - \mu_n) \right]$$

$$\Sigma_n^{-1} = n\Sigma^{-1} + \Sigma_p^{-1}$$

$$\mu_n = \Sigma_p \left(\Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n x_i \right) + \frac{1}{n} \Sigma \left(\Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \mu_p$$

$$\Sigma_n = \Sigma_p \left(\Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \frac{1}{n} \Sigma$$

Sequential Bayesian parameter estimation

- **Sequential Bayesian approach**
 - Under the iid the estimates of the posterior can be computed incrementally for a sequence of data points

$$p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{\int_{\Theta} p(D | \Theta, \xi) p(\Theta | \xi) d\Theta}$$

- If we use a conjugate prior we get back the same posterior
- Assume we split the data D in the last element \mathbf{x} and the rest
- $p(D | \Theta) = P(x | \Theta)P(D_{n-1} | \Theta)$
- **Then:**

$$p(\Theta | D, \xi) = \frac{\overbrace{P(x | \Theta)P(D_{n-1} | \Theta)}^{\text{A “new” prior}} p(\Theta | \xi)}{\int_{\Theta} P(x | \Theta)P(D_{n-1} | \Theta) p(\Theta | \xi) d\Theta}$$

Exponential family

Exponential family:

- all probability mass / density functions that can be written in the exponential normal form

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- $\boldsymbol{\eta}$ a vector of **natural (or canonical) parameters**
- $t(\mathbf{x})$ a function referred to as a **sufficient statistic**
- $h(\mathbf{x})$ a function of \mathbf{x} (it is less important)
- $Z(\boldsymbol{\eta})$ a normalization constant (a **partition function**)
$$Z(\boldsymbol{\eta}) = \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}$$
- Other common form:

$$f(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x}) - A(\boldsymbol{\eta})] \quad \log Z(\boldsymbol{\eta}) = A(\boldsymbol{\eta})$$

Exponential family: examples

- Bernoulli distribution

$$\begin{aligned} p(x | \pi) &= \pi^x (1 - \pi)^{1-x} \\ &= \exp \left\{ \log \left(\frac{\pi}{1 - \pi} \right) x + \log(1 - \pi) \right\} \\ &= \exp \{ \log(1 - \pi) \} \exp \left\{ \log \left(\frac{\pi}{1 - \pi} \right) x \right\} \end{aligned}$$

- Exponential family

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp \left[\boldsymbol{\eta}^T t(\mathbf{x}) \right]$$

- Parameters

$$\boldsymbol{\eta} = ? \qquad \qquad t(\mathbf{x}) = ?$$

$$Z(\boldsymbol{\eta}) = ? \qquad \qquad h(\mathbf{x}) = ?$$

Exponential family: examples

- Bernoulli distribution

$$\begin{aligned} p(x | \pi) &= \pi^x (1 - \pi)^{1-x} \\ &= \exp \left\{ \log \left(\frac{\pi}{1 - \pi} \right) x + \log(1 - \pi) \right\} \\ &= \exp \{ \log(1 - \pi) \} \exp \left\{ \log \left(\frac{\pi}{1 - \pi} \right) x \right\} \end{aligned}$$

- Exponential family

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp \left[\boldsymbol{\eta}^T t(\mathbf{x}) \right]$$

- Parameters

$$\boldsymbol{\eta} = \log \frac{\pi}{1 - \pi} \quad (\text{note} \quad \pi = \frac{1}{1 + e^{-\eta}}) \qquad \qquad t(\mathbf{x}) = x$$

$$Z(\boldsymbol{\eta}) = \frac{1}{1 - \pi} = 1 + e^\eta \qquad \qquad h(\mathbf{x}) = 1$$

Exponential family: examples

- **Univariate Gaussian distribution**

$$\begin{aligned} p(x | \mu, \sigma) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right] \\ &= \frac{1}{2\pi} \exp\left(-\frac{\mu}{2\sigma^2} - \log \sigma\right) \exp\left\{\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2\right\} \end{aligned}$$

- **Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(x) \exp\left[\boldsymbol{\eta}^T t(x)\right]$$

- **Parameters**

$$\boldsymbol{\eta} = ?$$

$$t(\mathbf{x}) = ?$$

$$Z(\boldsymbol{\eta}) = ?$$

$$h(\mathbf{x}) = ?$$

Exponential family: examples

- **Univariate Gaussian distribution**

$$\begin{aligned} p(x | \mu, \sigma) &= \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right] \\ &= \frac{1}{2\pi} \exp\left(-\frac{\mu}{2\sigma^2} - \log \sigma\right) \exp\left\{\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2\right\} \end{aligned}$$

- **Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(x) \exp\left[\boldsymbol{\eta}^T t(x)\right]$$

- **Parameters**

$$\boldsymbol{\eta} = \begin{bmatrix} \mu / 2\sigma^2 \\ -1 / 2\sigma^2 \end{bmatrix} \quad t(\mathbf{x}) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

$$Z(\boldsymbol{\eta}) = \exp\left\{\frac{\mu}{2\sigma^2} + \log \sigma\right\} = \exp\left\{-\frac{\eta_1^2}{4\eta_2} - \frac{1}{2}\log(-2\eta_2)\right\}$$

$$h(\mathbf{x}) = 1/\sqrt{2\pi}$$

Exponential family

- For iid samples, the likelihood of data is

$$\begin{aligned} P(D \mid \boldsymbol{\eta}) &= \prod_{i=1}^n p(\mathbf{x}_i \mid \boldsymbol{\eta}) = \prod_{i=1}^n h(\mathbf{x}_i) \exp \left[\boldsymbol{\eta}^T t(\mathbf{x}_i) - A(\boldsymbol{\eta}) \right] \\ &= \left[\prod_{i=1}^n h(\mathbf{x}_i) \right] \exp \left[\sum_{i=1}^n \boldsymbol{\eta}^T t(\mathbf{x}_i) - A(\boldsymbol{\eta}) \right] \\ &= \left[\prod_{i=1}^n h(\mathbf{x}_i) \right] \exp \left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right] \end{aligned}$$

- Important:

- the dimensionality of the sufficient statistic remains the same for different sample sizes (that is, different number of examples in D)

Exponential family

- The log likelihood of data is

$$\begin{aligned} l(D, \boldsymbol{\eta}) &= \log \left[\prod_{i=1}^n h(\mathbf{x}_i) \right] \exp \left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right] \\ &= \log \left[\prod_{i=1}^n h(\mathbf{x}_i) \right] + \left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right] \end{aligned}$$

- Optimizing the loglikelihood

$$\nabla_{\boldsymbol{\eta}} l(D, \boldsymbol{\eta}) = \left(\sum_{i=1}^n t(\mathbf{x}_i) \right) - n \nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \mathbf{0}$$

- For the ML estimate it must hold

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \frac{1}{n} \left(\sum_{i=1}^n t(\mathbf{x}_i) \right)$$

Exponential family

- **Rewriting the gradient:**

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \log Z(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \log \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T t(\mathbf{x}) \} d\mathbf{x}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \frac{\int t(\mathbf{x}) h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T t(\mathbf{x}) \} d\mathbf{x}}{\int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T t(\mathbf{x}) \} d\mathbf{x}}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \int t(\mathbf{x}) h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T t(\mathbf{x}) - A(\boldsymbol{\eta}) \} d\mathbf{x}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = E(t(\mathbf{x}))$$

- **Result:**

$$E(t(\mathbf{x})) = \frac{1}{n} \left(\sum_{i=1}^n t(\mathbf{x}_i) \right)$$

- **For the ML estimate the parameters $\boldsymbol{\eta}$ should be adjusted such that the expectation of the statistic $t(\mathbf{x})$ is equal to the observed sample statistics**

Moments of the distribution

- **For the exponential family**

- The k-th moment of the statistic corresponds to the k-th derivative of $A(\eta)$
- If x is a component of $t(x)$ then we get the moments of the distribution by differentiating its corresponding natural parameter

- **Example: Bernoulli** $p(x | \pi) = \exp\left\{\log\left(\frac{\pi}{1-\pi}\right)x + \log(1-\pi)\right\}$

$$A(\eta) = \log \frac{1}{1-\pi} = \log(1+e^\eta)$$

- **Derivatives:**

$$\frac{\partial A(\eta)}{\partial \eta} = \frac{\partial}{\partial \eta} \log(1+e^\eta) = \frac{e^\eta}{(1+e^\eta)} = \frac{1}{(1+e^{-\eta})} = \pi$$

$$\frac{\partial A(\eta)}{\partial \eta^2} = \frac{\partial}{\partial \eta} \frac{1}{(1+e^{-\eta})} = \pi(1-\pi)$$

Conjugate priors

For any member of the exponential family

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T \mathbf{t}(\mathbf{x})]$$

there exists a prior:

$$p(\boldsymbol{\eta} | \boldsymbol{\chi}, \nu) = u(\boldsymbol{\chi}, \nu) g(\boldsymbol{\eta})^\nu \exp[\nu \boldsymbol{\eta}^T \boldsymbol{\chi}]$$

Such that for n examples, the posterior is

$$p(\boldsymbol{\eta} | D, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{\nu+n} \exp\left[\boldsymbol{\eta}^T \left(\left[\sum_{i=1}^n \mathbf{t}(x_i) \right] + \nu \boldsymbol{\chi} \right) \right]$$

Note that:

$$P(D | \boldsymbol{\eta}) = \left(\frac{1}{Z(\boldsymbol{\eta})} \right)^n \left[\prod_{i=1}^n h(\mathbf{x}_i) \right] \exp\left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n t(\mathbf{x}_i) \right) \right]$$

Conjugate priors

For any member of the exponential family

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T \mathbf{t}(\mathbf{x})]$$

there exists a prior:

$$p(\boldsymbol{\eta} | \chi, \nu) = u(\chi, \nu) g(\boldsymbol{\eta})^\nu \exp[\nu \boldsymbol{\eta}^T \boldsymbol{\chi}]$$

Such that for n examples, the posterior is

$$p(\boldsymbol{\eta} | D, \chi, \nu) \propto g(\boldsymbol{\eta})^{\nu+n} \exp\left[\boldsymbol{\eta}^T \left(\left[\sum_{i=1}^n \mathbf{t}(x_i) \right] + \nu \boldsymbol{\chi} \right) \right]$$

Note that:

$$P(D | \boldsymbol{\eta}) = \left(\frac{1}{Z(\boldsymbol{\eta})} \right)^n \left[\prod_{i=1}^n h(\mathbf{x}_i) \right] \exp\left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n t(\mathbf{x}_i) \right) \right]$$

Pseudo-observation