
Exponential Family of Distributions

Prof. Nicholas Zabararas
Center for Informatics and Computational Science
<https://cics.nd.edu/>
University of Notre Dame
Notre Dame, Indiana, USA

Email: nzabararas@gmail.com
URL: <https://www.zabararas.com/>

September 6, 2018



References

- C P Robert, [The Bayesian Choice: From Decision-Theoretic Motivations to Computational Implementation](#), Springer-Verlag, NY, 2001 ([online resource](#))
- A Gelman, JB Carlin, HS Stern and DB Rubin, [Bayesian Data Analysis](#), Chapman and Hall CRC Press, 2nd Edition, 2003.
- J M Marin and C P Robert, [The Bayesian Core](#), Springer Verlag, 2007 ([online resource](#))
- D. Sivia and J Skilling, [Data Analysis: A Bayesian Tutorial](#), Oxford University Press, 2006.
- Bayesian Statistics for Engineering, [Online Course at Georgia Tech](#), B. Vidakovic.
- [Chris Bishops' PRML book](#), Chapter 2
- M. Jordan, An introduction to Probabilistic Graphical Models, Chapter 8 (pre-print)
- Kevin Murphy's, [Machine Learning: A probabilistic perspective](#), Chapters 2 and 4



Contents

- Exponential Family, Bernoulli, Beta, Gamma, Gaussian
- Conjugate Priors, Posterior Predictive
- Computing Moments
- Moment Parametrization
- MLE for the Exponential Family
- Maximum Entropy and the Exponential Family
- Generalized Linear Models



Exponential Family

- ❑ Large family of useful distributions with common properties
 - Bernoulli, beta, binomial, chi-square, Dirichlet, gamma, Gaussian, geometric, multinomial, Poisson, Weibull, . .
- ❑ Not in the family:
 - ✓ Uniform,
 - ✓ Student's T,
 - ✓ Cauchy,
 - ✓ Laplace,
 - ✓ Mixture of Gaussians,
 - ✓ . . .
- ❑ Variable can be discrete/continuous (or vectors thereof)



Exponential Family

- The exponential family of distributions over \mathbf{x} , given parameters $\boldsymbol{\eta}$, is defined to be the set of distributions of the form

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

$$p(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T u(\mathbf{x}) - A(\boldsymbol{\eta}) \}, \text{ where } A(\boldsymbol{\eta}) = -\log g(\boldsymbol{\eta})$$

\mathbf{x} is scalar/vector, discrete/continuous. **$\boldsymbol{\eta}$ are the natural parameters and $u(\mathbf{x})$ is referred to as a sufficient statistic.**

- $g(\boldsymbol{\eta})$ ensures that the distribution is normalized and satisfies

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

- The normalization factor Z and the log of it A are defined as:

$$Z(\boldsymbol{\eta}) = \frac{1}{g(\boldsymbol{\eta})}, \quad A(\boldsymbol{\eta}) = -\ln g(\boldsymbol{\eta}) = \ln Z(\boldsymbol{\eta}) = \ln \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T u(\mathbf{x}) \} d\mathbf{x}$$

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

- The space of $\boldsymbol{\eta}$ for which $\int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T u(\mathbf{x}) \} d\mathbf{x} < \infty$ is the **natural parameter space.**



Canonical or Natural Parameters

- When the parameter θ enters the exponential family as $\eta(\theta)$, we write the probability density of the exponential family as follows:

$$p(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) g(\boldsymbol{\eta}(\boldsymbol{\theta})) \exp \left\{ \boldsymbol{\eta}^T(\boldsymbol{\theta}) u(\mathbf{x}) \right\} \text{ or}$$

$$p(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) \exp \left\{ \boldsymbol{\eta}^T(\boldsymbol{\theta}) u(\mathbf{x}) - A(\boldsymbol{\eta}(\boldsymbol{\theta})) \right\},$$

$$\text{where: } A(\boldsymbol{\eta}(\boldsymbol{\theta})) = -\log g(\boldsymbol{\eta}(\boldsymbol{\theta}))$$

- $\boldsymbol{\eta}(\boldsymbol{\theta})$ are the canonical or natural parameters,
- $\boldsymbol{\theta}$ is the parameter vector of some distribution that can be written in the exponential family format

Exponential Family: The Bernoulli Distribution

- Consider the Bernoulli distribution:

$$p(x | \mu) = \mathcal{Bern}(x | \mu) = \mu^x (1 - \mu)^{1-x} = \exp \{ x \ln \mu + (1 - x) \ln(1 - \mu) \} =$$

$$= \underbrace{(1 - \mu)}_{g(\eta)} \exp \left\{ \underbrace{\ln \left(\frac{\mu}{1 - \mu} \right)}_{\eta} x \right\}$$

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\eta}) &= h(\mathbf{x}) g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T u(\mathbf{x}) \} \\ &= h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T u(\mathbf{x}) - A(\boldsymbol{\eta}) \} \end{aligned}$$

- From this we see that (note that *the relation $\mu(\eta)$ is invertible*)

$$\eta = \ln \left(\frac{\mu}{1 - \mu} \right) \Rightarrow$$

$$\mu = \sigma(\eta) = \frac{1}{1 + e^{-\eta}}$$

**Logistic sigmoid
function**

and

$$g(\eta) = 1 - \mu = 1 - \sigma(\eta) = \sigma(-\eta)$$

- Finally:

$$\begin{aligned} p(x | \eta) &= g(\eta) \exp \{ \eta x \}, u(x) = x, h(x) = 1, g(\eta) = \sigma(-\eta), \\ A(\eta) &= -\ln g(\eta) = -\log(1 - \mu) = \log(1 + e^\eta) \end{aligned}$$

Exponential Family: The Beta Distribution

□ Consider the Beta distribution

$$\text{Beta}(\mu | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \exp[(a-1)\ln \mu + (b-1)\ln(1-\mu)]$$

□ Comparing this with our exponential family:

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\eta}) &= h(\mathbf{x}) g(\boldsymbol{\eta}) \exp\{\boldsymbol{\eta}^T u(\mathbf{x})\} \\ &= h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T u(\mathbf{x}) - A(\boldsymbol{\eta})\} \end{aligned}$$

we can easily identify:

$$u(\mu) = (\ln \mu, \ln(1-\mu))^T, \boldsymbol{\eta} = (a-1, b-1)^T, h(\mu) = 1, g(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)},$$

$$A(a, b) = -\ln g(\boldsymbol{\eta}) = \ln \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$



Exponential Family: The Gaussian

□ Consider the univariate Gaussian

$$p(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}x^2 - \frac{1}{2\sigma^2}\mu^2 + \frac{\mu}{\sigma^2}x\right\}$$

□ Comparing this with our exponential family:

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

we can indentify (this is a two parameter distribution):

$$u(x) = (x, x^2)^T, \boldsymbol{\eta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right)^T, h(x) = \frac{1}{\sqrt{2\pi}}, g(\boldsymbol{\eta}) = (-2\eta_2)^{1/2} \exp \frac{\eta_1^2}{4\eta_2}$$

$$A(\boldsymbol{\eta}) = -\ln g(\boldsymbol{\eta}) = -\frac{1}{2} \ln(-2\eta_2) - \frac{\eta_1^2}{4\eta_2}$$

Conjugate Priors

- In general, for a given probability distribution $p(x|\eta)$, we can seek a prior $p(\eta)$ that is conjugate to the likelihood function, so that the posterior distribution has the same functional form as the prior.
 - For the Bernoulli, the conjugate prior is the Beta distribution
 - For the Gaussian, the conjugate prior for the mean is a Gaussian, and the conjugate prior for the precision is the Wishart distribution



Conjugate Priors

- For any member of the exponential family with likelihood

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

there exists a conjugate prior that can be written in the form

$$p(\boldsymbol{\theta} | \nu_0, \boldsymbol{\tau}_0) \propto g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_0} \exp\{\boldsymbol{\eta}^T(\boldsymbol{\theta}) \boldsymbol{\tau}_0\} = \exp\{\nu_0 \boldsymbol{\eta}^T(\boldsymbol{\theta}) \bar{\boldsymbol{\tau}}_0 - A(\boldsymbol{\eta}(\boldsymbol{\theta})) \nu_0\}, \text{ where: } \boldsymbol{\tau}_0 \equiv \nu_0 \bar{\boldsymbol{\tau}}_0$$

- In normalized form, we write:

$$p(\boldsymbol{\theta} | \nu_0, \boldsymbol{\tau}_0) = \frac{1}{Z(\nu_0, \boldsymbol{\tau}_0)} g(\boldsymbol{\eta}(\boldsymbol{\theta}))^{\nu_0} \exp\{\boldsymbol{\eta}^T(\boldsymbol{\theta}) \boldsymbol{\tau}_0\} = \frac{1}{Z(\nu_0, \boldsymbol{\tau}_0)} \exp\{\nu_0 \boldsymbol{\eta}^T(\boldsymbol{\theta}) \bar{\boldsymbol{\tau}}_0 - A(\boldsymbol{\eta}(\boldsymbol{\theta})) \nu_0\}$$
$$\text{where: } Z(\nu_0, \boldsymbol{\tau}_0) = \int \exp\{\nu_0 \boldsymbol{\eta}^T(\boldsymbol{\theta}) \bar{\boldsymbol{\tau}}_0 - A(\boldsymbol{\eta}(\boldsymbol{\theta})) \nu_0\} d\boldsymbol{\theta}$$

Conjugate Priors

$$p(X | \theta) = \prod_{n=1}^N \left(h(\mathbf{x}_n) g(\boldsymbol{\eta}(\theta)) \exp\{\boldsymbol{\eta}^T(\theta) u(\mathbf{x}_n)\} \right) = \prod_{n=1}^N \left(h(\mathbf{x}_n) g(\boldsymbol{\eta}(\theta)) \right)^N \exp\left\{ \boldsymbol{\eta}^T(\theta) \sum_{n=1}^N u(\mathbf{x}_n) \right\}$$

$$p(\theta | v_0, \boldsymbol{\tau}_0) = \frac{1}{Z(v_0, \boldsymbol{\tau}_0)} g(\boldsymbol{\eta}(\theta))^{v_0} \exp\{\boldsymbol{\eta}^T(\theta) \boldsymbol{\tau}_0\} = \frac{1}{Z(v_0, \boldsymbol{\tau}_0)} \exp\{v_0 \boldsymbol{\eta}^T(\theta) \bar{\boldsymbol{\tau}}_0 - A(\boldsymbol{\eta}(\theta)) v_0\}$$

□ Using $\bar{\mathbf{u}} = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$, the posterior becomes (this form justifies $\bar{\boldsymbol{\tau}}_0$):

$$p(\theta | X, \boldsymbol{\chi}, v) \propto g(\boldsymbol{\eta}(\theta))^{v_0+N} \exp\left\{ \boldsymbol{\eta}^T(\theta) \left(\sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) + v_0 \bar{\boldsymbol{\tau}}_0 \right) \right\} = g(\boldsymbol{\eta}(\theta))^{v_0+N} \exp\{\boldsymbol{\eta}^T(\theta) (N\bar{\mathbf{u}} + v_0 \bar{\boldsymbol{\tau}}_0)\}$$

□ The parameter v_0 can be interpreted as *effective number of fictitious observations* in the prior each of which has a value for the sufficient statistic equal to $\bar{\boldsymbol{\tau}}_0$.

$$p(\theta | X, v_N, \boldsymbol{\tau}_N) = \frac{1}{Z(v_N, \boldsymbol{\tau}_N)} g(\boldsymbol{\eta}(\theta))^{v_N} \exp\left\{ (N + v_0) \boldsymbol{\eta}^T(\theta) \frac{N\bar{\mathbf{u}} + v_0 \bar{\boldsymbol{\tau}}_0}{N + v_0} \right\} = \frac{1}{Z(v_N, \boldsymbol{\tau}_N)} g(\boldsymbol{\eta}(\theta))^{v_N} \exp\{v_N \boldsymbol{\eta}^T(\theta) \bar{\boldsymbol{\tau}}_N\},$$

$$\text{where } v_N = v_0 + N, \bar{\boldsymbol{\tau}}_N = \frac{N\bar{\mathbf{u}} + v_0 \bar{\boldsymbol{\tau}}_0}{N + v_0}, \boldsymbol{\tau}_N = v_N \bar{\boldsymbol{\tau}}_N = N\bar{\mathbf{u}} + v_0 \bar{\boldsymbol{\tau}}_0 = \sum_{i=1}^N \mathbf{u}(\mathbf{x}_i) + \boldsymbol{\tau}_0$$



Posterior Predictive

□ Let $u(X) = \sum_{i=1}^N u(x_i)$, $u(X') = \sum_{i=1}^{N'} u(x'_i)$, the posterior predictive is then:

$$\begin{aligned} p(X' | X) &= \int p(X' | \theta) p(\theta | X) d\theta \\ &= \prod_{i=1}^{N'} h(x'_i) \int g(\eta)^{N'} \exp\{\eta^T(\theta) u(X')\} \frac{1}{Z(\nu_0 + N, u(X) + \tau_0)} g(\eta(\theta))^{\nu_N} \exp\{\eta^T(\theta)(u(X) + \tau_0)\} d\theta \end{aligned}$$

□ This is simplified as follows:

$$\begin{aligned} p(X' | X) &= \prod_{i=1}^{N'} h(x'_i) \frac{1}{Z(\nu_0 + N, u(X) + \tau_0)} \int g(\eta(\theta))^{N' + \nu_N} \exp\{\eta^T(\theta)(u(X') + u(X) + \tau_0)\} d\theta \\ &= \prod_{i=1}^{N'} h(x'_i) \frac{Z(\nu_0 + N + N', u(X') + u(X) + \tau_0)}{Z(\nu_0 + N, u(X) + \tau_0)} \end{aligned}$$

□ If $N = 0$, this becomes the marginal likelihood of X' , which reduces to the normalizer of the posterior divided by the normalizer of the prior multiplied by a constant.



Beta/Bernoulli: Posterior Predictive

- Consider a Bernoulli likelihood with a Beta prior. The likelihood takes the familiar exponential distribution form:

$$p(\mathcal{D} | \theta) = \theta^{\sum_i x_i} (1-\theta)^{N-\sum_i x_i} = (1-\theta)^N \exp \left(\log \frac{\theta}{1-\theta} \sum_i x_i \right)$$

- The conjugate prior is a Beta: $p(\theta | \nu_0, \tau_0) = \theta^{\tau_0} (1-\theta)^{\nu_0-\tau_0} \propto (1-\theta)^{\nu_0} \exp \left(\log \left(\frac{\theta}{1-\theta} \right) \tau_0 \right)$

$$p(\theta | \nu_0, \tau_0) = \text{Beta}(\alpha, \beta), \alpha = \tau_0 + 1, \beta = \nu_0 - \tau_0 + 1,$$

- Thus the posterior becomes: $p(\theta | \mathcal{D}) \propto \theta^{\tau_0+s} (1-\theta)^{\nu_0-\tau_0+N-s} \Rightarrow$

$$p(\theta | \mathcal{D}) = \text{Beta}(\alpha_N, \beta_N), \alpha_N = \alpha + s, \beta_N = \beta + (N - s), s = \sum_i \mathbb{I}(x_i = 1)$$

- Let s the number of heads in the past data. The probability of $s' = \sum_{i=1}^m \mathbb{I}(x'_i = 1)$ future heads in m trials is then:

$$p(s' | \mathcal{D}, m) = \int \theta^{s'} (1-\theta)^{m-s'} \frac{\Gamma(\alpha_N + \beta_N)}{\Gamma(\alpha_N) \Gamma(\beta_N)} \theta^{\alpha_N-1} (1-\theta)^{\beta_N-1} d\theta = \frac{\Gamma(\alpha_N + \beta_N)}{\Gamma(\alpha_N) \Gamma(\beta_N)} \frac{\Gamma(\alpha_{N+m}) \Gamma(\beta_{N+m})}{\Gamma(\alpha_{N+m} + \beta_{N+m})}$$

$$\alpha_{N+m} = \alpha_N + s', \beta_{N+m} = \beta_N + (m - s')$$



Computing Moments of Sufficient Statistics $u(x)$

- Differentiate wrt η the $\int p(x | \eta) dx = 1$ for the exponential family:

$$\int p(x | \eta) dx = \int h(x) g(\eta) \exp\{\eta^T u(x)\} dx = 1$$

$$\nabla g(\eta) \int h(x) \exp\{\eta^T u(x)\} dx + g(\eta) \int h(x) \exp\{\eta^T u(x)\} u(x) dx = 0 \Rightarrow$$

$$-\frac{\nabla g(\eta)}{g(\eta)} = \int h(x) \exp\{\eta^T u(x)\} u(x) dx = \int p(x | \eta) u(x) dx = \mathbb{E}[u(x)]$$

- The above equation can be further simplified if written in terms of the partition function $Z = 1/g(\eta)$ or $A = \log Z = -\log g(\eta)$:

$$\nabla A(\eta) = \mathbb{E}[u(x)]$$

- Let us re-write explicitly the above equation as:

$$\nabla A(\eta) = \int h(x) \exp\{\eta^T u(x)\} u(x) dx$$

- We can compute the variance of $u(x)$ by differentiating the Eq. above with respect to η .



Computing Moments of Sufficient Statistics $u(x)$

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

$$\nabla^2 A(\boldsymbol{\eta}) = \underbrace{\nabla g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T u(\mathbf{x})\} u(\mathbf{x}) d\mathbf{x}}_{-\mathbb{E}[u(\mathbf{x})]\mathbb{E}[u(\mathbf{x})^T]} + \underbrace{g(\boldsymbol{\eta}) \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T u(\mathbf{x})\} u(\mathbf{x}) u(\mathbf{x})^T d\mathbf{x}}_{\mathbb{E}[u(\mathbf{x}) u(\mathbf{x})^T]}$$

- Thus the covariance of $u(\mathbf{x})$ can be expressed in terms of the 2nd derivatives of $A(\boldsymbol{\eta})$ and similarly for higher order moments.

$$\nabla^2 A(\boldsymbol{\eta}) = \text{cov}[u(\mathbf{x})] \text{ so } A(\boldsymbol{\eta}) \text{ is convex}$$

- Provided we can normalize a distribution from the exponential family, we can always find its moments by simple differentiation.

Computing Moments of Sufficient Statistics $u(x)$

$$\nabla A(\boldsymbol{\eta}) = \mathbb{E}[u(x)]$$

$$\nabla^2 A(\boldsymbol{\eta}) = \text{cov}[u(x)] \text{ so } A(\boldsymbol{\eta}) \text{ is convex}$$

□ Let us check these relations for the [Univariate Gaussian](#):

$$g(\eta) = (-2\eta_2)^{1/2} \exp \frac{\eta_1^2}{4\eta_2}$$

$$A(\boldsymbol{\eta}) = -\frac{1}{2} \ln(-2\eta_2) - \frac{\eta_1^2}{4\eta_2}, \boldsymbol{\eta} = \left(\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right)^T, u(x) = (x, x^2)^T$$

$$\frac{\partial A(\boldsymbol{\eta})}{\partial \eta_1} = -\frac{\eta_1}{2\eta_2} = \mu = \mathbb{E}[X], \frac{\partial A(\boldsymbol{\eta})}{\partial \eta_2} = -\frac{1}{2\eta_2} + \frac{\eta_1^2}{4\eta_2^2} = \mu^2 + \sigma^2 = \mathbb{E}[X^2]$$

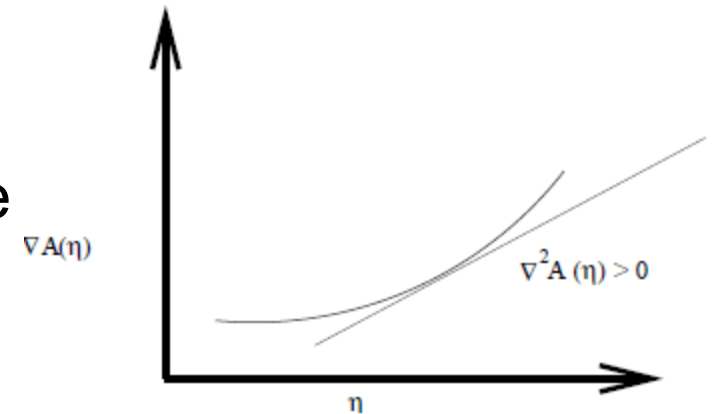
$$\frac{\partial^2 A(\boldsymbol{\eta})}{\partial \eta_1^2} = -\frac{1}{2\eta_2} = \sigma^2 = \text{var}[X], \text{ etc.}$$

Moment Parametrization

- We have shown that we can compute the mean of the distribution $\mu = E[u(\mathbf{x})]$ in terms of the canonical parameter η :

$$\mu = \mathbb{E}[u(\mathbf{x})] = \nabla A(\eta)$$

- We have also shown that $A(\eta)$ is a convex function. Since for a convex function there is one-to-one relation between the argument of the function and its derivative, the mapping $\mu(\eta)$ is invertible.



- Thus the exponential family of distributions can also be parameterized in terms of μ (*moment parametrization*) exactly as we started this course.

MLE for the Exponential Family

- The joint density for a data $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is itself an exp. distribution with sufficient statistics $\sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$

$$p(\mathbf{X} | \boldsymbol{\eta}) = \prod_{n=1}^N \left(h(\mathbf{x}_n) g(\boldsymbol{\eta}) \exp \left\{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}_n) \right\} \right) = \prod_{n=1}^N \left(h(\mathbf{x}_n) \right) g(\boldsymbol{\eta})^N \exp \left\{ \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) \right\} \Rightarrow$$

$$\ln p(\mathbf{X} | \boldsymbol{\eta}) = \sum_{n=1}^N h(\mathbf{x}_n) + N \ln g(\boldsymbol{\eta}) + \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) = \sum_{n=1}^N h(\mathbf{x}_n) - N A(\boldsymbol{\eta}) + \boldsymbol{\eta}^T \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$$

- The exponential family is the only family of distributions **with finite sufficient statistics** (size independent of the data set size).
- The log likelihood is concave (A convex) and has a unique maximum.
- Maximizing wrt $\boldsymbol{\eta}$ gives: $\nabla A(\boldsymbol{\eta}_{ML}) = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n) \Rightarrow \mathbb{E}[\mathbf{u}(\mathbf{x})] = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$
- At the MLE, the empirical average of the sufficient statistic is equal the model's theoretical expected sufficient statistics (moment matching).
- Thus to find the expected value of the sufficient statistics, one can use directly the data without having to estimate $\boldsymbol{\eta}$. When $\mathbf{u}(\mathbf{x}) = \mathbf{x}$, the above allows us to compute the expectation of \mathbf{x} directly from the data.



MLE for the Exponential Family

$$\nabla A(\boldsymbol{\eta}_{ML}) = \mathbb{E}[\mathbf{u}(\mathbf{x})] = \frac{1}{N} \sum_{n=1}^N \mathbf{u}(\mathbf{x}_n)$$

- Using the sufficient statistic, one can in principle invert the above equ. to compute $\boldsymbol{\eta}_{MLE}$. For example, for the Bernoulli distribution,

$$p(x | \eta) = g(\eta) \exp\{\eta x\}, u(x) = x, h(x) = 1,$$

$$\mu = \frac{1}{1 + e^{-\eta}}, g(\eta) = \frac{1}{1 + e^{\eta}}, \eta = \ln\left(\frac{\mu}{1 - \mu}\right)$$

and thus:

$$\mathbb{E}[X] = p(X = 1) = \bar{\mu} \equiv \mu_{MLE} = \frac{1}{N} \sum_{n=1}^N \mathbb{I}(x_n = 1)$$

and

$$\eta_{MLE} = \ln\left(\frac{\bar{\mu}}{1 - \bar{\mu}}\right)$$

MLE and Kullback-Leibler Distance

- A useful property for the MLE (and not just a property for the exponential family of distributions) is the following:
- Minimizing the KL distance to the empirical distribution is equivalent to maximizing the likelihood.
- Indeed, let us consider the model $\log p(x|\theta)$ and the empirical distribution:

$$p_{emp}(x) = \frac{1}{N} \sum_{n=1}^N \delta(x, x_n)$$

- We can then derive the following:

$$\sum_x p_{emp}(x) \log p(x|\theta) = \frac{1}{N} \sum_{n=1}^N \sum_x \delta(x, x_n) \log p(x|\theta) = \frac{1}{N} \sum_{n=1}^N \log p(x_n|\theta) = \frac{1}{N} \ell(\theta | \mathcal{D})$$

and from this:

$$\begin{aligned} KL(p_{emp}(x), p(x|\theta)) &= \sum_x p_{emp}(x) \log \frac{p_{emp}(x)}{p(x|\theta)} = \sum_x p_{emp}(x) \log p_{emp}(x) - \sum_x p_{emp}(x) \log p(x|\theta) \\ &= \sum_x p_{emp}(x) \log p_{emp}(x) - \frac{1}{N} \ell(\theta | \mathcal{D}) \end{aligned}$$

- Since the 1st term is independent of θ , the assertion is proved.



Maximum Entropy and Exponential Family

- If nothing is known about a distribution except that it belongs to a certain class, then the distribution with the largest entropy should be chosen as the default.^a

- The entropy is defined as

➤ discrete case $\mathbb{H}(\pi) = -\sum_k \pi(\theta_k) \log(\pi(\theta_k))$

- When some statistics (moments) of the distribution are known,

$$\mathbb{E}_\pi [g_k(\theta)] = w_k, k = 1, \dots, K$$

the maximum entropy distribution is of the form (λ 's are the Lagrange multipliers enforcing the constraints):

$$\pi(\theta_i) = \frac{\exp\left(-\sum_{k=1}^K \lambda_k g_k(\theta_i)\right)}{\sum_j \exp\left(-\sum_{k=1}^K \lambda_k g_k(\theta_j)\right)}, \lambda_k = \text{Lagrange multipliers}$$

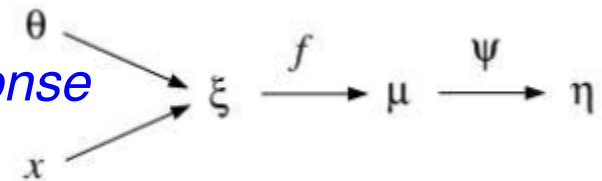
- Thus the MaxEnt distribution has the form of the exponential family.

^a C. P. Robert, [The Bayesian Choice](#), Springer, 2nd edition, [chapter 3](#) (full text available)



Generalized Linear Models

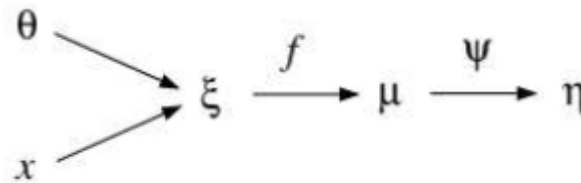
- ❑ We now study regression given data X and Y and a GLIM.
- ❑ We choose a particular conditional expectation of Y . We denote the modeled value of conditional expectation as $m = f(\theta^T x)$, $\xi = \theta^T x$.
- ❑ For linear regression, *GLIM extends these ideas beyond the Gaussian, Bernoulli and multinomial setting to the more general exponential family.*
- ❑ X enters linearly as $\theta^T x$ and f is called a response function. Ψ is a one-to-one map of μ to η .
- ❑ To specify a GLIM we need (a) a choice of exponential family distribution, and (b) a choice of the response function $f(\cdot)$.
- ❑ Choosing the exponential family distribution is strongly constrained by the nature of the data.
- ❑ Note that $f(\cdot)$ needs to be both monotonic and differentiable. However, *there is a particular response function (canonical response function) that is uniquely associated with a given exponential family distribution.*



Canonical Response Function

- Canonical response function:

$$f(\cdot) = \Psi^{-1}(\cdot)$$
$$\xi = \eta$$



- *If we decide to use the canonical response function, the choice of the exponential family density completely determines the GLIM.*

$$\xi = f^{-1}(\mu) = \Psi(\mu) = \eta$$

MLE & Canonical Response Function

- Consider a regression problem with data $\mathcal{D} = \{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$. The log likelihood for a GLIM is as:

$$\ell(\theta, \mathcal{D}) = \sum_{n=1}^N \log h(y_n) + \sum_{n=1}^N \left(\eta_n y_n - A(\eta_n) \right)_{\psi(\mu_n)}, \text{ where: } \mu_n = f(\xi_n) \text{ with } \xi_n = \theta^T \mathbf{x}_n$$

- For a canonical response, $\eta = \xi = \theta^T \mathbf{x}$, and this is simplified as:

$$\ell(\theta, \mathcal{D}) = \sum_{n=1}^N \log h(y_n) + \theta^T \underbrace{\sum_{n=1}^N \mathbf{x}_n y_n}_{\text{Sufficient statistic for } \theta} - \sum_{n=1}^N A(\eta_n) \quad \text{expectation}$$

- Regardless of N , the size of the sufficient statistic is fixed: the dimension of \mathbf{x}_n - important reason for using canonical response.

$$\nabla_{\theta} \ell(\theta, \mathcal{D}) = \sum_{n=1}^N (y_n - A'(\eta_n)) \nabla_{\theta} \eta_n = \sum_{n=1}^N (y_n - \mu_n) \nabla_{\theta} \eta_n = \sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n \text{ or } \nabla_{\theta} \ell(\theta, \mathcal{D}) = \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu})$$

- This is a general expression for GLM with exponential family distributions and the canonical response function.



Iterative Reweighted Least Squares (IRLS)

- The Hessian can now be computed from

$$\nabla_{\theta} \ell(\theta, \mathcal{D}) = \sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n \text{ or } \nabla_{\theta} \ell(\theta, \mathcal{D}) = \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu})$$

as:

$$\mathbf{H} = \nabla_{\theta}^2 \ell(\theta, \mathcal{D}) = -\sum_{n=1}^N \frac{d\mu_n}{d\eta_n} \mathbf{x}_n \mathbf{x}_n^T \text{ or } \nabla_{\theta}^2 \ell(\theta, \mathcal{D}) = -\mathbf{X}^T \mathbf{W} \mathbf{X}, \text{ where } \mathbf{W} = \left\{ \frac{d\mu_1}{d\eta_1}, \dots, \frac{d\mu_n}{d\eta_n} \right\}$$

- To estimate parameters in the canonical response function choice, one can use the [iteratively reweighted least squares \(IRLS\) algorithm](#)
- The batch [Newton algorithm](#) now takes the familiar IRLS form:

$$\begin{aligned} \theta^{t+1} &= \theta^t + \left(\mathbf{X}^T \mathbf{W}^t \mathbf{X} \right)^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}^t) = \left(\mathbf{X}^T \mathbf{W}^t \mathbf{X} \right)^{-1} \left(\mathbf{X}^T \mathbf{W}^t \mathbf{X} \theta^t + \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}^t) \right) \\ &= \left(\mathbf{X}^T \mathbf{W}^t \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}^t \left(\underset{\eta^t}{\mathbf{X} \theta^t + \mathbf{W}^{t-1} (\mathbf{y} - \boldsymbol{\mu}^t)} \right) = \left(\mathbf{X}^T \mathbf{W}^t \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{W}^t \left(\boldsymbol{\eta}^t + \mathbf{W}^{t-1} (\mathbf{y} - \boldsymbol{\mu}^t) \right) \end{aligned}$$

- For non-canonical response functions, the Hessian has an extra term that contains the factor $(\mathbf{y} - \boldsymbol{\mu})$. When we take expectations this term vanishes! So using the expected Hessian in the Newton method the algorithm looks essentially the same (Fisher Scoring algorithm).



Sequential Estimation - LMS

- An on-line estimation algorithm can be obtained by following the stochastic gradient of the log likelihood function.

$$\boldsymbol{\theta}^{t+1} = \boldsymbol{\theta}^t + \rho \left(y_n - \mu_n^t \right) x_n, \mu_n^t = f \left(\boldsymbol{\theta}^{t^T} x_n \right)$$

- If we do not use the canonical response function, then the gradient also includes the derivatives of $f(\cdot)$ and $\Psi(\cdot)$. These can be viewed as scaling coefficients that alter the step size, but otherwise leave the general LMS form intact.
- *The LMS algorithm is the generic stochastic gradient algorithm for models throughout the GLIM family.*

