

# Unifying Theory of GLMs (Part II)

Rebecca C. Steorts

## Reading

These slides provide some background that will be helpful in tying together generalized linear models completely.

https:

[//statmath.wu.ac.at/courses/heather\\_turner/glmCourse\\_001.pdf](https://statmath.wu.ac.at/courses/heather_turner/glmCourse_001.pdf)

# Agenda

- ▶ Introduce alternative parameterization of exponential family that is widely used.
- ▶ Applies beyond the one-parameter setting for generalized linear models
- ▶ Using tricks and the formulation of the exponential family for GLMs, we can calculate certain summary statistics, such as the mean and variance without relying on integration

# Exponential Families

Consider the following:

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

where

- ▶  $\eta$  is the natural parameter
- ▶  $h(y)$  is the underlying measure
- ▶  $A(\eta)$  is the log normalizer (ensuring the distribution integrates to one) or **log partition function**.

Note that  $Z(\eta) = \exp\{A(\eta)\}$  or  $A(\eta) = \log Z(\eta)$ .

# Exponential Families

We can alternatively write the exponential family as follows:

$$p(y \mid \eta) = \frac{1}{Z(\eta)} h(y) \exp\{\eta^T u(y)\},$$

where

- ▶  $\eta$  are the natural parameters
- ▶  $u(y)$  are the sufficient statistics
- ▶  $Z(\eta)$  is a partition function that ensures the density is normalized.

$$Z(\eta) = \int h(y) \exp\{\eta^T u(y)\} dy.$$

# Bernouli

$$p(y | \eta) = \frac{1}{Z(\eta)} h(y) \exp\{\eta^T u(y)\}$$

Consider

$$\text{Bern}(y | \mu) = \mu^y (1 - \mu)^{1-y} \quad (1)$$

$$= (1 - \mu) \left( \frac{\mu}{1 - \mu} \right)^y \quad (2)$$

$$= (1 - \mu) \times 1 \times \exp\left\{\log\left(\frac{\mu}{1 - \mu}\right) \times y\right\} \quad (3)$$

- ▶  $Z(\eta) = \frac{1}{1 - \mu}$
- ▶  $h(y) = 1$
- ▶  $\eta = \log\left(\frac{\mu}{1 - \mu}\right)$  (canonical link or natural parameter)
- ▶  $u(y) = y$

## Bernoulli (continued)

Solving for  $\mu$ , we find that

$$\mu = \frac{e^{\eta}}{1 + e^{\eta}}.$$

Then

$$1 - \mu = \frac{1}{1 + e^{\eta}} \implies Z(\eta) = 1 + e^{\eta}.$$

Later in the lecture, we will work with the partition function to find the mean and variance of the density.

## Bernoulli continued

1.

$$\eta = \beta_0 + \sum_{i=1}^p \beta_j x_j$$

2.

$$\eta = g(\mu) \implies g^{-1}(\eta) = \mu.$$

3. Putting this together,

$$\eta = g(\mu) = \beta_0 + \sum_{i=1}^p \beta_j x_j.$$

For the Bernoulli, we thus have the following:

$$\eta = \log\left(\frac{\mu}{1-\mu}\right) = \beta_0 + \sum_{j=1}^p \beta_j x_{ji}.$$



# Summarize GLM

We can summarize GLMs in two ways:

1. We must specify the response variable as belonging to an exponential family. (Random component).
2. The systematic component, which specifies a linear predictor  $\eta_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ji}$  and a link function  $g(\mu) = \eta$  that is known, monotonic and differentiable.

# Specification

To specify this for  $y_i \sim \text{Bernoulli}(\mu_i)$ , we proceed as follows:

1.  $y_i \mid \mu_i \sim \text{Bernoulli}(\mu_i)$  (random component)
2.  $\log\left(\frac{\mu_i}{1-\mu_i}\right) = \beta_0 + \beta_1 x_i$  (systematic component)

where  $x_i$  refers to one explanatory variable.

We have already verified that the Bernoulli distribution is a member of the exponential family.

# Why are we studying the exponential family?

Many convenient properties such as:

- ▶ Sufficient statistics for maximum likelihood
- ▶ Many convenient identities for the partition function that allow for quick computation directly using the partition function (or functions of it) instead of having to work with the likelihood function.
- ▶ For instance, we can quickly compute the moments of an exponential family using the partition function.

# Maximum Likelihood and sufficient statistics

Suppose we have observations that are i.i.d  $y_1, \dots, y_n$ .

Find  $\eta$  that maximizes  $p(Y \mid \eta)$ .

$$p(Y \mid \eta) = \prod_{i=1}^n \frac{1}{Z(\eta)} \left( h(y_i) \exp\{\eta^T u(y_i)\} \right) \quad (4)$$

$$= \left( \frac{1}{Z(\eta)} \right)^n \prod_{i=1}^n h(y_i) \exp\left\{ \eta^T \sum_{i=1}^n u(y_i) \right\} \quad (5)$$

## Maximum Likelihood and sufficient statistics

It follows that

$$\log p(Y | \eta) = -n \log(Z(\eta)) + \eta^T \sum_{i=1}^n u(y_i) + \log\left(\prod_{i=1}^n h(y_i)\right)$$

$$\nabla_{\eta} \log p(Y | \eta) = -n \nabla_{\eta} \log Z(\eta) + \sum_{i=1}^n u(y_i) := 0$$

This implies that

$$\nabla_{\eta} \log Z(\eta) = \frac{1}{n} \sum_{i=1}^n u(y_i).$$

The maximum likelihood solution only depends on  $\sum_{i=1}^n u(y_i)$ .

Thus,  $u(y)$  is called the **sufficient statistic**.

## Maximum Likelihood and sufficient statistics

We showed that the maximum likelihood estimator (that belongs to an exponential family) only depends on the sufficient statistic.

# Sufficiency

The notion of a sufficient statistic is a fundamental one in statistical theory and its applications.

Fisher (1922) introduced this into the literature in the attempt to formalize the notion of no loss of information.

A sufficient statistic should solely contain all of the information about the unknown parameters of the underlying distribution that the entire sample could have provided.

In that sense, there is nothing to lose by restricting attention to just a sufficient statistic.

## Bernoulli (continued)

Bernoulli:

- ▶ Recall  $u(y) = y$ .
- ▶ Only need to store  $\sum_i y_i$ .



## Distributions not in the Exponential family

Can you think of some distributions that would not be in the exponential family and why?

# Identities

1.

$$\nabla_{\eta} \log Z(\eta) = E_{y \sim p(y|\eta)}[u(y)] =: \xi \quad (\text{moments})$$

2.

$$\nabla_{\eta} \log p(Y | \eta) = u(y) - E_{y \sim p(y|\eta)}[u(y)].$$

3.

$$\nabla_{\eta}^2 \log Z(\eta) = \text{Cov}(u(x)) = -\nabla_{\eta}^2 \log p(y | \eta).$$

These show that to calculate the mean and the variance for a distribution in the exponential family, we can work directly with the partition function.

## Bernoulli (continued)

Find the mean and variance using the derivation of the Bernoulli being in the exponential family.

Specifically, show that the mean is  $\mu$  and the variance is  $\mu(1 - \mu)$  using the identities above.

## Bernoulli (continued)

Recall that  $Z(\eta) = 1 + e^\eta$ .

Consider

$$\nabla_\eta \log Z(\eta) = \nabla_\eta (1 + e^\eta) = \frac{e^\eta}{1 + e^\eta} = \mu.$$

Consider

$$\nabla_\eta^2 \log Z(\eta) = \nabla_\eta \left( \frac{e^\eta}{1 + e^\eta} \right) = \frac{e^{\eta^2}}{(1 + e^\eta)^2} = \mu(1 - \mu).$$

## Lemma 1

$$Z(\eta) = \int h(y) \exp\{\eta^T u(y)\} dy.$$

$$\nabla_{\eta} \log Z(\eta) = E_{y \sim p(y|\eta)}[u(y)] =: \xi \quad (\text{moments})$$

## Proof of Lemma 1

$$Z(\eta) = \int h(y) \exp\{\eta^T u(y)\} dy.$$

$$\nabla_{\eta} \log Z(\eta) = E_{y \sim p(y|\eta)}[u(y)] =: \xi$$

By the chain rule,

$$\nabla_{\eta} \log Z(\eta) = \frac{\nabla_{\eta} Z(\eta)}{Z(\eta)}.$$

## Proof of Lemma 1 (continued)

It then follows

$$\nabla_{\eta} \log Z(\eta) = \frac{1}{Z(\eta)} \nabla_{\eta} \int h(y) \exp\{\eta^T u(y)\} dy \quad (6)$$

$$= \frac{1}{Z(\eta)} \int h(y) \exp\{\eta^T u(y)\} u(y) dy \quad (7)$$

$$= E_{y \sim p(y|\eta)}[u(y)] =: \xi \quad (8)$$

1. Observe there is a one-to-one mapping between  $\eta$  and  $\xi$ .
2.  $\xi$  provides an alternate parameterization for the exponential family.

## Lemma 2

$$p(y \mid \eta) = \frac{1}{Z(\eta)} h(y) \exp\{\eta^T u(y)\}$$

$$\nabla_{\eta} \log p(Y \mid \eta) = u(y) - E_{y \sim p(y|\eta)}[u(y)].$$



## Proof of Lemma 2

$$p(y \mid \eta) = \frac{1}{Z(\eta)} h(y) \exp\{\eta^T u(y)\}$$

$$\nabla_{\eta} \log p(Y \mid \eta) = \nabla_{\eta} \left[ -\log Z(\eta) + \log h(y) + \eta^T u(y) \right] \quad (9)$$

$$= -\nabla_{\eta} \log Z(\eta) + \nabla_{\eta} [\eta^T u(y)] \quad (10)$$

$$= -E_{y \sim p(y \mid \eta)} [u(y)] + u(y) \quad (11)$$

## Connection with Maximum likelihood and Moments

Recall in maximum likelihood, we have

$$\nabla_{\eta} \sum_{i=1}^n \log p(y_i | \eta) = n \left[ \underbrace{\frac{1}{n} \sum_{i=1}^n u(y_i)}_{\text{empirical moments } \hat{\xi}} - \underbrace{\nabla_{\eta} \log Z(\eta)}_{\text{moments } \xi} \right] \quad (12)$$

$$= n(\hat{\xi} - \xi). \quad (13)$$

This shows that maximum likelihood leads to moment matching.

# Moment Matching

The moment matching method (MME) is a widely used method of estimation of parameters.

The idea is to find values of the unknown parameters that result in a match between the theoretical (or population) and sample moments evaluated from data.

## Lemma 3

$$\nabla_{\eta}^2 \log Z(\eta) = \text{Cov}(u(x)) = -\nabla_{\eta}^2 \log p(y \mid \eta).$$

## Proof of Lemma 3

$$p(y \mid \eta) = \frac{1}{Z(\eta)} h(y) \exp\{\eta^T u(y)\}$$

Show

$$\nabla_{\eta}^2 \log p(y \mid \eta) = -\nabla_{\eta}^2 \log Z(\eta)$$

Recall that

$$\nabla_{\eta} \log p(Y \mid \eta) = \nabla_{\eta} \left[ -\log Z(\eta) + \log h(y) + \eta^T u(y) \right] \quad (14)$$

$$= -\nabla_{\eta} \log Z(\eta) + \nabla_{\eta} [\eta^T u(y)] \quad (15)$$

$$= -\nabla_{\eta} \log Z(\eta) + u(y), \implies \quad (16)$$

$$\nabla_{\eta}^2 \log p(y \mid \eta) = -\nabla_{\eta}^2 \log Z(\eta).$$

## Proof of Lemma 3

Show

$$\nabla_{\eta}^2 \log Z(\eta) = \text{Cov}(u(y)).$$

## Other relationships

$$\nabla_{\eta}^2 \log Z(\eta) = F_{\eta} := E[-\nabla_{\eta}^2 \log p(y \mid \eta)]$$

(Fisher information)

$$F_{\eta} = \nabla_{\eta} \xi = J_{\xi, \eta}.$$

(Jacobian mapping  $\eta \rightarrow \xi$ ).

## Exercise: Normal ( $\mu$ random)

$$p(y \mid \eta) = \frac{1}{Z(\eta)} h(y) \exp\{\eta^T u(y)\}$$

$$p(y \mid \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{\frac{-1}{2\sigma^2}(y - \mu)^2\right\} \quad (17)$$

$$= \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{\frac{-1}{2\sigma^2}y^2 + \frac{\mu}{\sigma^2}y - \frac{\mu^2}{2\sigma^2}\right\} \quad (18)$$

$$= \underbrace{\frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{\frac{-1}{2\sigma^2}y^2\right\}}_{h(y)} \underbrace{\exp\left\{-\frac{\mu^2}{2\sigma^2}\right\}}_{-Z(\eta)} \exp\left\{\underbrace{\frac{\mu}{\sigma^2}}_{\eta^T} \underbrace{y}_{u(y)}\right\} \quad (19)$$



## Exercise (continued)

The natural parameter is  $\eta = \frac{\mu}{\sigma^2} \implies \mu = \eta\sigma^2$  and the partition function is

$$Z(\eta) = \exp\left\{\frac{\mu^2}{2\sigma^2}\right\} = \exp\left\{\frac{\eta^2\sigma^4}{2\sigma^2}\right\} \exp\left\{\frac{\eta^2\sigma^2}{2}\right\}$$

This implies that

$$\log Z(\eta) = \frac{\eta^2\sigma^2}{2}.$$

Using our identities from earlier, it follows that

$$\nabla_{\eta} \log Z(\eta) = \frac{2\eta\sigma^2}{2} = \frac{2\mu\sigma^2}{2\sigma^2} = \mu$$

and

$$\nabla_{\eta}^2 \log Z(\eta) = \nabla_{\eta} \eta\sigma^2 = \sigma^2.$$

## Exercise: Normal ( $\mu$ and $\sigma^2$ )

$$p(y \mid \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{\frac{-1}{2\sigma^2}(y - \mu)^2\right\} \quad (20)$$

$$= \frac{1}{(2\pi)^{1/2}} \sigma^{-1} \exp\left\{\frac{-1}{2\sigma^2}y^2 + \frac{\mu}{\sigma^2}y - \frac{\mu^2}{2\sigma^2}\right\} \quad (21)$$

$$= \underbrace{\frac{1}{(2\pi)^{1/2}}}_{h(y)} \underbrace{\sigma^{-1} \exp\left\{-\frac{\mu^2}{2\sigma^2}\right\}}_{-Z(\eta)} \exp\left\{\underbrace{\begin{bmatrix} \frac{\mu}{\sigma^2} \\ -1 \\ \frac{1}{2\sigma^2} \end{bmatrix}}_{\eta^T}^T \underbrace{\begin{bmatrix} y \\ y^2 \end{bmatrix}}_{u(y)}\right\} \quad (22)$$

$$(23)$$

## Normal (continued)

The partition function can be written as:

$$Z(\mu, \sigma^2) = \exp\left\{\frac{\mu^2}{2\sigma^2} + \log(\sigma)\right\} \implies \log Z(\mu, \sigma^2) = \frac{\mu^2}{2\sigma^2} + \log(\sigma).$$

From the natural parameters  $\eta_1 = \frac{\mu}{\sigma^2}$  and  $\eta_2 = -\frac{1}{2\sigma^2}$ , we solve for  $\mu$  and  $\sigma^2$  in terms of  $\eta_1$  and  $\eta_2$ :

$$\mu = \eta_1 \sigma^2, \quad \sigma^2 = -\frac{1}{2\eta_2}$$

## Normal (continued)

Substituting these into the log-partition function, we compute the terms:

$$\begin{aligned}\frac{\mu^2}{2\sigma^2} &= \frac{(\eta_1\sigma^2)^2}{2\sigma^2} = \frac{\eta_1^2\sigma^2}{2} \\ \frac{\eta_1^2\sigma^2}{2} &= \frac{\eta_1^2}{2} \cdot \left(-\frac{1}{2\eta_2}\right) = -\frac{\eta_1^2}{4\eta_2}\end{aligned}$$

Thus, the term  $\log(\sigma)$  simplifies

$$\begin{aligned}\log(\sigma) &= \log\left(\sqrt{-\frac{1}{2\eta_2}}\right) = \frac{1}{2}\log\left(-\frac{1}{2\eta_2}\right) \\ \log(\sigma) &= -\frac{1}{2}\log(2) - \frac{1}{2}\log(\eta_2)\end{aligned}$$

## Normal (continued)

Combining these terms, the log-partition function in terms of  $\eta_1$  and  $\eta_2$  is:

$$\log Z(\eta_1, \eta_2) = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(2) - \frac{1}{2} \log(\eta_2)$$

## Normal (continued)

$$\log Z(\eta_1, \eta_2) = -\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(2) - \frac{1}{2} \log(\eta_2)$$

Consider

$$\nabla_{\eta_1} \log Z(\eta_1, \eta_2) = -\frac{\eta_1}{2\eta_2}$$

and

$$\nabla_{\eta_1}^2 \log Z(\eta_1, \eta_2) = -\frac{1}{2\eta_2}$$

## Normal (continued)

Recall  $\eta_1 = \frac{\mu}{\sigma^2}$  and  $\eta_2 = -\frac{1}{2\sigma^2}$ .

It follows that

$$\nabla_{\eta_1} \log Z(\eta_1, \eta_2) = -\frac{\mu}{2\sigma^2}(-2\sigma^2) = \mu.$$

and

$$\nabla_{\eta_1}^2 \log Z(\eta_1, \eta_2) = -\frac{1}{2\eta_2} = \sigma^2.$$

Think about why this derivation gives the mean and the variance of  $y$ .