

# Unit 4: Generalized linear models: General theory

Eugene Katsevich

October 27, 2021

Units 1-3 focused on the most common class of models used in applications: linear models. Despite their versatility, linear models do not apply in all situations. In particular, they are not designed to deal with binary or count responses. In Unit 4, we introduced *generalized linear models* (GLMs), a generalization of linear models that encompasses a wide variety of incredibly useful models including logistic regression and Poisson regression.

We'll start Unit 4 by introducing exponential family models (Section 1), a generalization of the Gaussian distribution that serves as the backbone of GLMs. Then we formally define a GLM, demonstrating logistic regression and Poisson regression as special cases (Section 2). Next we discuss maximum likelihood inference in GLMs (Section 3). Finally, we discuss how to carry out statistical inference in GLMs (Section 4).

## 1 Exponential family distributions

**Definition and examples.** Let's start with the Gaussian distribution, taking variance  $\sigma^2 = 1$  for simplicity. If  $y \sim N(\mu, 1)$ , then it has density

$$f(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right) = \exp\left(\mu y - \frac{1}{2}\mu^2\right) \cdot \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right). \quad (1)$$

Here is a way of generalizing this density:

$$f_\eta(y) = \exp(\eta y - \psi(\eta))h(y). \quad (2)$$

Here  $\eta$  is called the *natural parameter*,  $\psi$  is called the *log-partition function*, and  $h$  is called the *base measure*. The distribution with density  $f_\eta$  is called a *one-parameter natural exponential family*. Therefore,  $y \sim N(\mu, 1)$  is in the exponential family with

$$\eta = \mu, \quad \psi(\eta) = -\frac{1}{2}\eta^2, \quad h(y) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right). \quad (3)$$

Several other well-known distributions are in the exponential family as well. For example, consider  $y \sim \text{Ber}(\pi)$ . Then, we have

$$f(y) = \pi^y(1 - \pi)^{1-y} = \exp\left(y \log \frac{\pi}{1 - \pi} + \log(1 - \pi)\right). \quad (4)$$

Therefore, we have  $\eta = \log \frac{\pi}{1 - \pi}$ , so that  $\log(1 - \pi) = -\log(1 + e^\eta)$ . It follows that

$$\eta = \log \frac{\pi}{1 - \pi}, \quad \psi(\eta) = \log(1 + e^\eta), \quad h(y) = 1. \quad (5)$$

As another example, consider the Poisson distribution  $y \sim \text{Poi}(\mu)$ . We have

$$f(y) = e^{-\mu} \frac{\mu^y}{y!} = \exp(y \log \mu - \mu) \frac{1}{y!}. \quad (6)$$

Therefore, we have  $\eta = \log \mu$ , so that  $\mu = e^\eta$ . It follows that

$$\eta = \log \mu, \quad \psi(\eta) = e^\eta, \quad h(y) = \frac{1}{y!}. \quad (7)$$

**Moments of exponential family distributions.** It turns out that the derivatives of the log-partition function  $\psi$  give the moments of  $y$ . Indeed, let's start with the relationship

$$\int f_\eta(y) dy = \int \exp(\eta y - \psi(\eta)) h(y) dy = 1. \quad (8)$$

Differentiating in  $\eta$  and interchanging the derivative and the integral, we obtain

$$0 = \frac{d}{d\eta} \int f_\eta(y) dy = \int (y - \dot{\psi}(\eta)) f_\eta(y) dy, \quad (9)$$

from which it follows that

$$\dot{\psi}(\eta) = \int \dot{\psi}(\eta) f_\eta(y) dy = \int y f_\eta(y) dy = \mathbb{E}_\eta[y] \equiv \mu_\eta. \quad (10)$$

Thus, the first derivative of the log partition function is the mean of  $y$ . Differentiating again, we get

$$\ddot{\psi}(\eta) = \int y(y - \dot{\psi}(\eta)) f_\eta(y) dy = \int y(y - \mu_\eta) f_\eta(y) dy = \int (y - \mu_\eta)^2 f_\eta(y) dy = \text{Var}_\eta[y]. \quad (11)$$

Thus, the second derivative of the log-partition function is the variance of  $y$ .

**Relationship between mean and natural parameter.** The log-partition function  $\psi$  induces a connection (10) between the natural parameter  $\eta$  and the mean  $\mu$ . Because

$$\frac{d\mu}{d\eta} = \frac{d}{d\eta} \dot{\psi}(\eta) = \ddot{\psi}(\eta) = \text{Var}_\eta[y] > 0, \quad (12)$$

it follows that  $\mu$  is a strictly increasing function of  $\eta$ , so in particular the mapping between  $\mu$  and  $\eta$  is bijective. Therefore, we can think of equivalently parameterizing the distribution via  $\mu$  or  $\eta$ . In the context of GLMs (see Section 2), the mean-variance relationship is quantified in terms of the *canonical link function*  $g$ , which maps the mean to the natural parameter:

$$\eta = \dot{\psi}^{-1}(\mu) \equiv g(\mu). \quad (13)$$

**Relationship between mean and variance.** Note that the mean of an exponential family distribution determines its variance (since it determines the natural parameter  $\eta$ ). For example, a Poisson random variable with mean  $\mu$  has variance  $\mu$  and a Bernoulli random variable with mean  $\mu$  has variance  $\mu(1 - \mu)$ . The mean-variance relationship turns out to characterize the exponential family distribution, i.e. an exponential family distribution with mean equal to its variance is the Poisson distribution.

## 2 Generalized linear models and examples

In this class, the focus is on building models that tie a vector of predictors ( $\mathbf{x}_{i*}$ ) to a response  $y_i$ . For linear regression, the mean of  $y$  was modeled as a linear combination of the predictors  $\mathbf{x}_{i*}^T \boldsymbol{\beta}$ :  $\mu = \mathbf{x}_{i*}^T \boldsymbol{\beta}$ . Typically, the “right” thing to do is to model the response linearly on the scale of the natural parameter  $\eta$  rather than on the scale of the mean parameter  $\mu$ . It just happens for linear models (where the underlying distribution is Gaussian) that these two parameters coincide.

**Definition.** We define  $\{(y_i, \mathbf{x}_{i*})\}_{i=1}^n$  as following a generalized linear model based on the exponential family  $f_\eta$  if

$$y_i \stackrel{\text{ind}}{\sim} f_{\eta_i}, \quad \eta_i = \mathbf{x}_{i*}^T \boldsymbol{\beta}. \quad (14)$$

GLMs are often written in terms of their link functions  $g$ , which relate the mean of  $y$  to the linear predictor  $\mathbf{x}_{i*}^T \boldsymbol{\beta}$ . When modeling the natural parameter as a linear function in the predictors, as in the definition (14), we get a GLM with *canonical link function*  $g = \dot{\psi}^{-1}$ :

$$g(\mathbb{E}[y_i]) = \dot{\psi}^{-1}(\mathbb{E}[y_i]) = \mathbf{x}_{i*}^T \boldsymbol{\beta}. \quad (15)$$

**Examples.** For example, *logistic regression* is the GLM based on the Bernoulli distribution:

$$y_i \stackrel{\text{ind}}{\sim} \text{Ber}(\pi_i); \quad \eta_i = \log \frac{\pi_i}{1 - \pi_i} = \mathbf{x}_{i*}^T \boldsymbol{\beta}. \quad (16)$$

Thus the canonical link function for logistic regression is the *logistic link function*  $g(\mu) = \log \frac{\mu}{1-\mu}$ . As another example, *Poisson regression* is the GLM based on the Poisson distribution:

$$y_i \stackrel{\text{ind}}{\sim} \text{Poi}(\mu_i); \quad \eta_i = \log \mu_i = \mathbf{x}_{i*}^T \boldsymbol{\beta}. \quad (17)$$

Thus the canonical link function for Poisson regression is the *log link function*  $g(\mu) = \log \mu$ .

## 3 Maximum likelihood estimation in GLMs

**GLM normal equations.** Recall that the least squares estimate  $\hat{\boldsymbol{\beta}}$  is also the maximum likelihood estimate. For general GLMs, we also estimate  $\boldsymbol{\beta}$  via maximum likelihood. To derive this estimates, let's write down the GLM likelihood and then take a derivative. The GLM likelihood is

$$\mathcal{L}(\boldsymbol{\beta}) = \prod_{i=1}^n f_{\eta_i}(y_i) = \prod_{i=1}^n \exp(\eta_i y_i - \psi(\eta_i)) h(y_i). \quad (18)$$

Taking a logarithm, we have

$$\log \mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^n (\eta_i y_i - \psi(\eta_i)) + \sum_{i=1}^n \log h(y_i) = \sum_{i=1}^n (\mathbf{x}_{i*}^T \boldsymbol{\beta} y_i - \psi(\mathbf{x}_{i*}^T \boldsymbol{\beta})) + \sum_{i=1}^n \log h(y_i). \quad (19)$$

Taking a gradient in  $\boldsymbol{\beta}$ , we get

$$\nabla_{\boldsymbol{\beta}} \log \mathcal{L}(\boldsymbol{\beta}) = \sum_{i=1}^n (\mathbf{x}_{i*} y_i - \mathbf{x}_{i*} \dot{\psi}(\mathbf{x}_{i*}^T \boldsymbol{\beta})) = \mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})). \quad (20)$$

Setting this expression to zero, we get the normal equations:

$$\mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}})) = 0. \quad (21)$$

Recall that, for least squares, we got the same equation, with  $\boldsymbol{\mu}(\hat{\boldsymbol{\beta}}) = \mathbf{X}\hat{\boldsymbol{\beta}}$ . We can interpret the normal equations as stating that  $\boldsymbol{\mu}(\hat{\boldsymbol{\beta}})$  is a projection of  $\mathbf{y}$  onto the model “space”

$$C_{\mu}(\mathbf{X}) \equiv \{\boldsymbol{\mu} = \dot{\psi}(\boldsymbol{\eta}) = \dot{\psi}(\mathbf{X}\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\}. \quad (22)$$

parallel to the columns of  $\mathbf{X}$ . Note that the subscript  $\mu$  on  $C_{\mu}(\mathbf{X})$  indicates that we are considering the “space” (actually, *set*) of possible  $\boldsymbol{\mu}$  as opposed to the space  $C_{\eta}(\mathbf{X})$  of possible  $\boldsymbol{\eta}$ , which we denoted in Unit 1 simply as  $C(\mathbf{X})$ . For linear models, it is the case that  $C_{\mu}(\mathbf{X}) = C_{\eta}(\mathbf{X})$ , but in general, these two are different. Note that  $C_{\mu}(\mathbf{X})$  in general is a manifold as opposed to a linear subspace of  $\mathbb{R}^n$ , while  $C_{\eta}(\mathbf{X})$  is always a linear subspace.

**Log-concavity of GLM likelihood.** Unlike linear regression, in general GLMs the function  $\boldsymbol{\mu}(\boldsymbol{\beta})$  is nonlinear. Therefore, there is in general no closed-form solution to the GLM normal equations (21). We must instead iteratively compute the maximum likelihood estimate  $\hat{\boldsymbol{\beta}}$ . Before talking about the computation of the MLE  $\hat{\boldsymbol{\beta}}$ , we state the important fact that  $\log \mathcal{L}(\boldsymbol{\beta})$  is a concave function of  $\boldsymbol{\beta}$ , which implies that this function is “easy to optimize”, i.e. has no local maxima.

**Proposition 3.1.** *The function  $\log \mathcal{L}(\boldsymbol{\beta})$  defined in (19) is concave in  $\boldsymbol{\beta}$ .*

*Proof.* We claim it suffices to show that  $\psi$  is a convex function. Indeed, then  $\log \mathcal{L}(\boldsymbol{\beta})$  would be the sum of a linear function of  $\boldsymbol{\beta}$  and the composition of a concave function with a linear function. To verify that  $\psi$  is convex, it suffices to recall that  $\ddot{\psi}(\boldsymbol{\eta}) = \text{Var}_{\boldsymbol{\eta}}[y] \geq 0$ .  $\square$

Proposition (3.1) gives us confidence that an iterative algorithm will converge to the global maximum of the likelihood. We present such an iterative algorithm next.

**Newton-Raphson.** We can solve the equation (21) using the Newton Raphson algorithm, which involves the gradient and Hessian of the function we’d like to maximize. We already computed the gradient in equation (20). To compute the Hessian, we take another gradient in  $\boldsymbol{\beta}$ . We have

$$\begin{aligned} \nabla_{\boldsymbol{\beta}}^2 \log \mathcal{L}(\boldsymbol{\beta}) &= \nabla_{\boldsymbol{\beta}}(\mathbf{X}^T(\mathbf{y} - \dot{\psi}(\mathbf{X}\boldsymbol{\beta}))) = -\nabla_{\boldsymbol{\beta}} \mathbf{X}^T \dot{\psi}(\mathbf{X}\boldsymbol{\beta}) \\ &= -\mathbf{X}^T \text{diag}(\ddot{\psi}(\mathbf{X}\boldsymbol{\beta})) \mathbf{X} \equiv -\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}) \mathbf{X}. \end{aligned} \quad (23)$$

Here,  $\dot{\psi}$  and  $\ddot{\psi}$  applied to vectors are interpreted element-wise and  $\mathbf{W}(\boldsymbol{\beta}) \in \mathbb{R}^{n \times n}$  is the diagonal matrix such that

$$W_{ii}(\boldsymbol{\beta}) = \text{Var}_{\boldsymbol{\beta}}[y_i]. \quad (24)$$

The Newton-Raphson iteration is therefore

$$\hat{\boldsymbol{\beta}}^{(t+1)} = \hat{\boldsymbol{\beta}}^{(t)} - (\nabla_{\boldsymbol{\beta}}^2 \log \mathcal{L}(\hat{\boldsymbol{\beta}}^{(t)}))^{-1} \nabla_{\boldsymbol{\beta}} \log \mathcal{L}(\hat{\boldsymbol{\beta}}^{(t)}) = \hat{\boldsymbol{\beta}}^{(t)} + (\mathbf{X}^T \mathbf{W}(\hat{\boldsymbol{\beta}}^{(t)}) \mathbf{X})^{-1} \mathbf{X}^T(\mathbf{y} - \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}^{(t)})). \quad (25)$$

**Iteratively reweighted least squares (IRLS).** A nice interpretation of the Newton-Raphson algorithm is as a sequence of weighted least squares fits, known as the iteratively reweighted least squares (IRLS) algorithm. Suppose that we have a current estimate  $\hat{\boldsymbol{\beta}}^{(t)}$ , and suppose we are looking for a vector  $\boldsymbol{\beta}$  near  $\hat{\boldsymbol{\beta}}^{(t)}$  that fits the model even better. We have

$$\mathbb{E}_{\boldsymbol{\beta}}[\mathbf{y}] = \dot{\psi}(\mathbf{X}\boldsymbol{\beta}) \approx \dot{\psi}(\mathbf{X}\hat{\boldsymbol{\beta}}^{(t)}) + \text{diag}(\ddot{\psi}(\mathbf{X}\hat{\boldsymbol{\beta}}^{(t)}))(\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}}^{(t)}) = \boldsymbol{\mu}(\hat{\boldsymbol{\beta}}^{(t)}) + \mathbf{W}(\hat{\boldsymbol{\beta}}^{(t)})(\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\hat{\boldsymbol{\beta}}^{(t)}).$$

and

$$\text{Var}_{\boldsymbol{\beta}}[\mathbf{y}] \approx \mathbf{W}(\hat{\boldsymbol{\beta}}^{(t)}).$$

Thus, up to the first two moments, near  $\beta = \hat{\beta}^{(t)}$  the distribution of  $\mathbf{y}$  is approximately

$$\mathbf{y} = \mu(\hat{\beta}^{(t)}) + \mathbf{W}(\hat{\beta}^{(t)})(\mathbf{X}\beta - \mathbf{X}\hat{\beta}^{(t)}) + \epsilon, \quad \epsilon \sim N(\mathbf{0}, \mathbf{W}(\hat{\beta}^{(t)})), \quad (26)$$

or, equivalently,

$$\mathbf{z}^{(t)} \equiv \mathbf{W}(\hat{\beta}^{(t)})^{-1}(\mathbf{y} - \mu(\hat{\beta}^{(t)})) + \mathbf{X}\hat{\beta}^{(t)} = \mathbf{X}\beta + \epsilon', \quad \epsilon' \sim N(\mathbf{0}, \mathbf{W}(\hat{\beta}^{(t)})^{-1}). \quad (27)$$

The regression of the *adjusted response variable*  $\mathbf{z}^{(t)}$  on  $\mathbf{X}$  leaves us with a weighted linear regression, whose maximum likelihood estimate is

$$\hat{\beta}^{(t+1)} = (\mathbf{X}^T \mathbf{W}(\hat{\beta}^{(t)}) \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}(\hat{\beta}^{(t)}) \mathbf{z}^{(t)}, \quad (28)$$

which we define as our next iterate. It's easy to verify that the IRLS iteration (28) is equivalent to the Newton-Raphson iteration (25).

## 4 Inference in GLMs

**Asymptotic normality and GLM standard errors.** Inference in GLMs is based on asymptotic likelihood theory. Using the Hessian computation (23), we can compute the Fisher information matrix

$$\mathbf{I}(\beta) = -\mathbb{E}_\beta[\nabla_\beta^2 \log \mathcal{L}(\beta)] = \mathbf{X}^T \mathbf{W}(\beta) \mathbf{X}, \quad (29)$$

recalling the definition of  $\mathbf{W}$  in equation (24). Therefore, likelihood theory tells us that, as the sample size  $n$  grows, we have

$$\hat{\beta} \sim N(\beta, (\mathbf{X}^T \mathbf{W}(\beta) \mathbf{X})^{-1}). \quad (30)$$

Using the plug-in variance estimate, we typically construct GLM standard errors based on

$$\widehat{\text{Var}}[\hat{\beta}] \equiv (\mathbf{X}^T \mathbf{W}(\hat{\beta}) \mathbf{X})^{-1}. \quad (31)$$

**Confidence intervals.** A confidence interval for each coordinate  $\beta_j$  can be obtained via

$$\text{CI}(\hat{\beta}_j) \equiv \hat{\beta}_j \pm 2 \cdot \text{SE}(\hat{\beta}_j), \quad \text{where} \quad \text{SE}(\beta_j) \equiv \sqrt{(\mathbf{X}^T \mathbf{W}(\hat{\beta}) \mathbf{X})_{jj}^{-1}}. \quad (32)$$

A confidence interval for  $\eta_i = \mathbf{x}_{i*}^T \beta$  can be obtained via

$$\text{CI}(\hat{\eta}_i) \equiv \mathbf{x}_{i*}^T \hat{\beta} \pm 2 \cdot \text{SE}(\eta_i), \quad \text{where} \quad \text{SE}(\hat{\eta}_i) \equiv \sqrt{\mathbf{x}_{i*}^T (\mathbf{X}^T \mathbf{W}(\hat{\beta}) \mathbf{X})^{-1} \mathbf{x}_{i*}}. \quad (33)$$

A confidence interval for  $\mu_i \equiv \mathbb{E}_\beta[y_i] = \dot{\psi}(\eta_i)$  can be obtained by applying the strictly increasing function  $\dot{\psi}$  to the endpoints of the confidence interval for  $\eta_i$ . Note that the resulting confidence interval may be asymmetric.

**Testing a single coordinate.** We can invert the confidence interval (32) to get a test of the hypothesis  $H_0 : \beta_j = \beta_j^0$  for any  $\beta_j^0 \in \mathbb{R}$ :

$$\phi(\mathbf{X}, \mathbf{y}) = \mathbb{1}(|z(\mathbf{X}, \mathbf{y})| > z_{1-\alpha/2}), \quad \text{where} \quad z(\mathbf{X}, \mathbf{y}) \equiv \frac{\hat{\beta}_j}{\text{SE}(\hat{\beta}_j)}. \quad (34)$$

This is the analog of the  $t$ -test for a linear regression.

**Testing multiple coordinates.** Likelihood ratio test, deviances, goodness of fit testing. (TBD)

## **5 Bells and whistles**

**Offsets.**

**Exponential dispersion families.**

**Non-canonical links.**