

BIOS 7345: Advanced Regression for Independent Data

Andrew J. Spieker, Ph.D.

Associate Professor of Biostatistics
Vanderbilt University

Set 11: Generalized linear models (via likelihood)

Version: 10/29/2025

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

MOTIVATING EXAMPLE

Notes:

- We're about to pivot to generalized linear models (GLMs).
- You may already know:
 - ▶ Linear, logistic, and Poisson regression are all examples of GLMs.
 - ▶ Coefficients for most GLMs (including logistic and Poisson) cannot be expressed in closed form.
 - ▶ GLMs can be (but need not be) motivated in a likelihood framework.
- You may also already be aware of the `optim()` function in R for numerical optimization. The first order of business is to give you a couple of reasons why we will *not* be using the `optim()` function in our study of GLMs.
- I want you to believe fairly early on in the notes that studying GLMs in detail is a good use of your time.
 - ▶ Knowing what's going on under the hood, coding/algorithms, methodological skill building.

MOTIVATING EXAMPLE

Logistic regression: Generating data

```
1 ## Define expit function
2 expit <- function(x)
3 {
4   expitx <- exp(x)/(1 + exp(x))
5   return(expitx)
6 }
7
8 ## Set seed for reproducibility
9 set.seed(7345)
10
11 ## Generate data
12 n <- 50
13 X <- rbinom(n, 1, 0.5)
14 Y <- rbinom(n, 1, expit(-0.5 + X))
15
16 ## Display data from single replicate
17 > table(X, Y)
18    Y
19 X   0  1
20  0 11  9
21  1 13 17
```

Logistic regression: Hard-coding the likelihood

```
1 ## Negative log-likelihood for logistic regression (hard-coded)
2 negll <- function(theta, x, y)
3 {
4   beta0 <- theta[1]
5   beta1 <- theta[2]
6   lik <- dbinom(y,
7                 size = 1,
8                 prob = expit(beta0 + beta1 * x))
9   negloglik <- -1 * sum(log(lik))
10  return(negloglik)
11 }
```

Logistic regression: Numerical optimization

```
1 ## Nelder-Mead (initialized at (log(26/24), 0))
2 zz0 <- optim(c(log(26/24), 0), negll, x = X, y = Y, hessian = TRUE)
3
4 ## Report estimated coefficients
5 > zz0$par
6 [1] -0.2003658  0.4684407
7
8 ## Report estimated standard errors
9 > sqrt(diag(solve(zz0$hessian)))
10 [1] 0.4494598 0.5811689
```

MOTIVATING EXAMPLE

Logistic regression: Using the `glm()` function

```
1 ## GLM function in R
2 zz1 <- glm(Y ~ X, family = binomial(link = "logit"))
3
4 ## Report results (fairly close; not exact)
5 > summary(zz1)
6
7 Call:
8 glm(formula = Y ~ X, family = binomial(link = "logit"))
9
10 Coefficients:
11             Estimate Std. Error z value Pr(>|z|)
12 (Intercept) -0.2007      0.4495  -0.446   0.655
13 X            0.4689      0.5812   0.807   0.420
14
15 (Dispersion parameter for binomial family taken to be 1)
16
17 Null deviance: 69.235  on 49  degrees of freedom
18 Residual deviance: 68.579  on 48  degrees of freedom
19 AIC: 72.579
20
21 Number of Fisher Scoring iterations: 3
22
23 ## You may not yet know what "Fisher Scoring" means
24 ### ...but you will after this set of notes!
```

MOTIVATING EXAMPLE

Logistic regression: Testing out speed

```
1 ## Estimate coefficients by numerical approximation
2 do.optim <- function(x,y)
3 {
4   init <- c(log(sum(y)/(length(y) - sum(y))), 0)
5   zz <- optim(init, negll, x = x, y = y, hessian = TRUE)
6   return(list(COEF = as.numeric(zz$par),
7             SE = sqrt(diag(solve(zz$hessian))))))
8 }
9
10 ## Estimate coefficients by iterative approach (to discuss)
11 do.glm <- function(x,y)
12 {
13   init <- c(log(sum(y)/(length(y) - sum(y))), 0)
14   zz <- glm(Y ~ X, family = binomial(link = "logit"), start = init)
15   return(list(COEF = coef(zz),
16             SE = summary(zz)$coef[,2]))
17 }
```

MOTIVATING EXAMPLE

Logistic regression: Testing out speed ($n = 50$)

```
1 ## Number of simulations
2 nsim <- 1000
3
4 ## Sample size
5 n <- 50
6
7 ## How long does it take to do optim()?
8 > system.time(for (j in 1:nsim) {
9   X <- rbinom(n, 1, 0.5)
10  Y <- rbinom(n, 1, expit(-0.5 + X))
11  do.optim(x = X, y = Y)
12 })
13
14 >
15 user system elapsed
16 0.306 0.025 0.336
17
18 ## How long does it take to do glm()?
19 > system.time(for (j in 1:nsim) {
20   X <- rbinom(n, 1, 0.5)
21   Y <- rbinom(n, 1, expit(-0.5 + X))
22   do.glm(x = X, y = Y)
23 })
24
25 user system elapsed
26 0.445 0.019 0.470
```

MOTIVATING EXAMPLE

Logistic regression: Testing out speed ($n = 300$)

```
1 ## Number of simulations
2 nsim <- 1000
3
4 ## Sample size
5 n <- 300
6
7 ## How long does it take to do optim()?
8 > system.time(for (j in 1:nsim) {
9   X <- rbinom(n, 1, 0.5)
10  Y <- rbinom(n, 1, expit(-0.5 + X))
11  do.optim(x = X, y = Y)
12 })
13
14 >
15 user system elapsed
16 1.207 0.028 1.253
17
18 ## How long does it take to do glm()?
19 > system.time(for (j in 1:nsim) {
20   X <- rbinom(n, 1, 0.5)
21   Y <- rbinom(n, 1, expit(-0.5 + X))
22   do.glm(x = X, y = Y)
23 })
24
25 user system elapsed
26 0.596 0.010 0.611
```

MOTIVATING EXAMPLE

Logistic regression: Testing out speed ($n = 1000$)

```
1 ## Number of simulations
2 nsim <- 1000
3
4 ## Sample size
5 n <- 1000
6
7 ## How long does it take to do optim()?
8 > system.time(for (j in 1:nsim) {
9   X <- rbinom(n, 1, 0.5)
10  Y <- rbinom(n, 1, expit(-0.5 + X))
11  do.optim(x = X, y = Y)
12 })
13
14 >
15 user system elapsed
16 3.518 0.144 3.663
17
18 ## How long does it take to do glm()?
19 > system.time(for (j in 1:nsim) {
20   X <- rbinom(n, 1, 0.5)
21   Y <- rbinom(n, 1, expit(-0.5 + X))
22   do.glm(x = X, y = Y)
23 })
24
25 user system elapsed
26 1.168 0.021 1.189
```

MOTIVATING EXAMPLE

Logistic regression: Testing out speed ($n = 10,000$)

```
1 ## Number of simulations
2 nsim <- 1000
3
4 ## Sample size
5 n <- 10000
6
7 ## How long does it take to do optim()?
8 > system.time(for (j in 1:nsim) {
9   X <- rbinom(n, 1, 0.5)
10  Y <- rbinom(n, 1, expit(-0.5 + X))
11  do.optim(x = X, y = Y)
12 })
13
14 >
15   user  system elapsed
16 40.459   1.700  42.159
17
18 ## How long does it take to do glm()?
19 > system.time(for (j in 1:nsim) {
20   X <- rbinom(n, 1, 0.5)
21   Y <- rbinom(n, 1, expit(-0.5 + X))
22   do.glm(x = X, y = Y, init = c(log(sum(Y)/(n - sum(Y))), 0))
23 })
24
25   user  system elapsed
26  9.800   0.368  10.166
```

Logistic regression: Comments

- Nelder-Mead starts to wheeze a bit with large samples.
 - ▶ Unless you calculate and supply the gradient, the algorithm will approximate it at every step, which is computationally taxing.
 - ▶ May not be a “big deal” when running a single model.
 - ▶ Likely a big deal when you’re trying to implement bootstrap/imputation methods within a simulation study.
- Even if you supply the gradient, calculation of the inverse Hessian is approximate. The quality of the approximation depends in part on how well behaved the likelihood surface is.
 - ▶ Although you may not know this yet, the likelihood surface for logistic regression is quite well behaved for reasons we will discuss.

Logistic regression: Comments

- The `glm()` function utilizes an algorithm that numerically solves the exact estimating equations (though there is no closed-form expression).
- Keep in mind that our illustration is for the simple logistic regression setting; you could imagine how changing the model to be more complicated might compound these issues.
 - ▶ Informal homework: try it!
- For what it's worth, you can improve speed even further by using `glm.fit()`, but *for now*, we're going to learn all the unifying theory that underlies GLMs, including how to hard-code one from scratch.

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)**
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

Recall: Exponential families

- As in the previous set of notes, assume Y has a density function that can be written in the following very nice form:

$$f_Y(y; \theta, \phi) = \exp\left[\frac{y\theta - b(\theta)}{\phi} + c(y, \phi)\right].$$

- θ the “natural parameter” or “canonical parameter.”
- ϕ is a nuisance parameter, serving the role of a “scale” parameter.
- $E[Y] = b'(\theta)$.
- $\text{Var}[Y] = \phi b''(\theta)$.
- Some mean variance relationship, $V(\mu)$, is implied.

Pushing things forward: Regression

- Suppose we observe outcomes Y_1, \dots, Y_N , with associated covariate vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ (which you can assume for now to be fixed for ease of presentation, though we will eventually deal with the random case just as we did for OLS), with density taking the form:

$$f_Y(y|\mathbf{x}; \theta, \phi) = \exp\left[\frac{y\theta(\mathbf{x}^\top \boldsymbol{\beta}) - b(\theta(\mathbf{x}^\top \boldsymbol{\beta}))}{\phi} + c(y, \phi)\right],$$

which is to say that $\theta = \theta(\mathbf{x}^\top \boldsymbol{\beta})$ is a function of $\mathbf{x}^\top \boldsymbol{\beta}$.

- A GLM specifies two components:
 - ① The mean model: $g(E[Y]) = \mathbf{x}^\top \boldsymbol{\beta}$.
 - ② The family of distributions to which Y belongs (must be of the nice exponential family form shown above).
- Ultimately (in the next set of notes), we will argue that a GLM is uniquely determined by *less* information; but for now, we are imagining that the likelihood to be completely correct.

Now: Regression

- The log-likelihood based on observations Y_1, \dots, Y_N is given by:

$$\ell(\boldsymbol{\beta}, \phi; \mathbf{y}, \mathbf{X}) = \sum_{i=1}^N \left[\frac{y_i \theta(\mathbf{x}_i^\top \boldsymbol{\beta}) - b(\theta(\mathbf{x}_i^\top \boldsymbol{\beta}))}{\phi} + c(y_i, \phi) \right]$$

- Our first order of business is going to be to find score equations. We'll do this first for $\boldsymbol{\beta}$ (though we'll eventually have to deal with ϕ).

Notation and terminology:

- Linear predictor: $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$, or $\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta}$.
- Natural/canonical parameter: $\boldsymbol{\theta} = \boldsymbol{\theta}(\boldsymbol{\eta}) = (\theta(\mathbf{x}_1^\top \boldsymbol{\beta}), \dots, \theta(\mathbf{x}_N^\top \boldsymbol{\beta}))^\top$.
- Mean: $E[Y_i] = \mu_i$, or $E[\mathbf{y}] = \boldsymbol{\mu} = \boldsymbol{\mu}(\boldsymbol{\beta})$.
 - ▶ Recall that $E[Y_i] = \mu_i = b'(\theta_i)$, or $E[\mathbf{y}] = \boldsymbol{\mu} = \text{vec}(b'(\boldsymbol{\theta}))$.
 - ▶ I may sometimes use the notation $b'(\boldsymbol{\theta})$ for simplicity, especially on the board.
- Link function: $g(\cdot)$.
 - ▶ $g(E[Y]) = g(\boldsymbol{\mu}) = g(b'(\boldsymbol{\theta}))$.

Score equations:

- GLM mean model: $g(E[Y]) = g(\mu) = g(b'(\theta)) = \mathbf{x}^T \boldsymbol{\beta} =: \eta$.
- Note: we've already done some of the heavy lifting by developing score equations for exponential families with respect to a single real-valued parameter, θ .
- We can develop the score equations for $\boldsymbol{\beta}$ by the chain rule:

$$\left[\frac{\partial \ell(\boldsymbol{\beta}, \phi)}{\partial \boldsymbol{\beta}} \right]_{K \times 1} = \left[\frac{\partial \eta}{\partial \boldsymbol{\beta}} \right]_{K \times N} \left[\frac{\partial \boldsymbol{\mu}}{\partial \eta} \right]_{N \times N} \left[\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\mu}} \right]_{N \times N} \left[\frac{\partial \ell}{\partial \boldsymbol{\theta}} \right]_{N \times 1}$$

- Let's take each of these three components one step at a time.

Score equations: The “left two” components

- The derivative of $\boldsymbol{\mu}$ with respect to $\boldsymbol{\beta}$ is a $K \times N$ matrix:

$$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} = \frac{\partial}{\partial \boldsymbol{\beta}} \text{vec}(g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta})) = \frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\beta}} \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}}$$

- We typically denote this as $\mathbf{D}^T = \mathbf{D}^T(\boldsymbol{\beta})$.
 - ▶ Referring to the $K \times N$ derivative matrix as \mathbf{D}^T instead of simply “ \mathbf{D} ” is standard, and feels natural largely because we’re already used to expressions involving \mathbf{X}^T .
- Note that $\partial \boldsymbol{\eta} / \partial \boldsymbol{\beta} = \partial(\mathbf{X}\boldsymbol{\beta}) / \partial \boldsymbol{\beta} = \mathbf{X}^T$.
- Further note that $\partial \boldsymbol{\mu} / \partial \boldsymbol{\eta} = \text{diag}(\partial \mu_i(\boldsymbol{\beta}) / \partial \eta_i(\boldsymbol{\beta}))$.
- I will use \mathbf{D}^T notation when it is called for, but ultimately what we’ve learned is that:

$$\mathbf{D}^T(\boldsymbol{\beta}) = \mathbf{X}^T \text{diag}\left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})}\right)$$

Score equations: The “middle” component

- The derivative of $\boldsymbol{\theta}$ with respect to $\boldsymbol{\mu}$ is—in this course (independent observations)—an $N \times N$ (diagonal) matrix:

$$\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\mu}} = ???$$

- Already know: $\mu_i = b'(\theta_i)$. Therefore,

$$\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\theta}} = \text{diag}(b''(\theta_i)) = \text{diag}(V(\mu_i)) =: \mathbf{V} = \mathbf{V}(\boldsymbol{\beta}).$$

- Recall that $\text{Var}[Y_i] = \phi b''(\theta_i) \propto b''(\theta_i)$; the notation “ \mathbf{V} ” signifies the nature of the mean-variance relationship; in this set of slides $V(\mu_i)$ will be used as shorthand for $b''(\theta_i)$. Therefore,

$$\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\mu}} = \mathbf{V}^{-1} = \mathbf{V}^{-1}(\boldsymbol{\beta}) = \text{diag}(1/V(\mu_i)).$$

Score equations: The “right” component

- The derivative of ℓ with respect to $\boldsymbol{\theta}$ is an $N \times 1$ vector:

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}}.$$

- We already did the heavy lifting on deriving the score function in this nice exponential family (with respect to θ):

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}} = \frac{\mathbf{y} - \text{vec}(b'(\theta_i))}{\phi} = \frac{\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})}{\phi} = \frac{\mathbf{y} - \text{vec}(g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}))}{\phi}.$$

Score equations:

- Putting this all together,

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = \frac{\partial \eta}{\partial \boldsymbol{\beta}} \frac{\partial \boldsymbol{\mu}}{\partial \eta} \frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\mu}} \frac{\partial \ell}{\partial \boldsymbol{\theta}} = \mathbf{D}^T(\boldsymbol{\beta}) \mathbf{V}^{-1}(\boldsymbol{\beta})(\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}))/\phi.$$

- To estimate $\boldsymbol{\beta}$, we solve the following equations for $\boldsymbol{\beta}$:

$$\mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu})/\phi = \mathbf{0}.$$

- The solution, $\hat{\boldsymbol{\beta}}$, to this problem, does not depend upon ϕ , and so we can instead solve:

$$\mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}.$$

- This can be written with greater detail as:

$$\mathbf{X}^T \text{diag}((\partial \mu_i(\boldsymbol{\beta})/\partial \eta_i(\boldsymbol{\beta}))/V(\mu_i(\boldsymbol{\beta}))) (\mathbf{y} - \text{vec}(g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}.$$

Example 11.1: Normal distribution (identity link)

- Suppose our GLM is based on:
 - ▶ $Y \sim \mathcal{N}(\mu, \sigma^2)$.
 - ▶ $g(\cdot)$ given by the identity link (i.e., $g(\mu) = \mu$).
- Then,
 - 1 Since $g(\mu) = \mu = \mathbf{x}^T \boldsymbol{\beta}$, we have $\mathbf{D}^T = \partial \boldsymbol{\mu} / \partial \boldsymbol{\beta} = \mathbf{X}^T$.
 - 2 $\mathbf{V} = \text{diag}(1) = \mathbf{I}$ (constant variance).
- The likelihood-based equations below can be used to solve for $\boldsymbol{\beta}$:

$$\mathbf{D}^T \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{X}^T (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) = \mathbf{0}.$$

- We recognize these equations are the normal equations for OLS (we previously motivated the same equations without a likelihood).
- These equations have a closed-form solution (don't get used to it).

Example 11.2: Normal distribution (log link)

- Suppose our GLM is based on:
 - ▶ $Y \sim \mathcal{N}(\mu, \sigma^2)$.
 - ▶ $g(\cdot)$ given by the log link (i.e., $g(\mu) = \log(\mu)$).

• Then,

① To determine \mathbf{D}^T , note that $\boldsymbol{\mu} = \text{vec}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))$. Therefore,

$$\mathbf{D}^T = \left[\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} \right]_{K \times N} = \left[\frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\beta}} \right]_{K \times N} \left[\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right]_{N \times N} = \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta})).$$

② $\mathbf{V} = \mathbf{I}$.

- The likelihood-based equations below can be used to solve for $\boldsymbol{\beta}$:

$$\mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))(\mathbf{y} - \text{vec}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}.$$

- As another note, I may sometimes write in the style “ $\boldsymbol{\mu} = \exp(\mathbf{X}\boldsymbol{\beta})$,” which you should take to mean entries of a vector rather than, say, the matrix exponential.

Example 11.3: Bernoulli distribution (logit link)

- Suppose our GLM is based on:
 - ▶ $Y \sim \text{Bernoulli}(p)$.
 - ▶ $g(\cdot)$ given by the logit link (i.e., $g(\mu) = \text{logit}(\mu)$).
- Then,
 - ① To determine \mathbf{D}^\top , note that $\boldsymbol{\mu} = \text{vec}(\text{expit}(\mathbf{x}_i^\top \boldsymbol{\beta}))$. Therefore,

$$\begin{aligned} \mathbf{D}^\top &= \left[\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} \right]_{K \times N} = \left[\frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\beta}} \right]_{K \times N} \left[\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} \right]_{N \times N} \\ &= \mathbf{X}^\top \text{diag}(\text{expit}(\mathbf{x}_i^\top \boldsymbol{\beta})(1 - \text{expit}(\mathbf{x}_i^\top \boldsymbol{\beta}))). \end{aligned}$$

- ② $\mathbf{V} = \text{diag}(\mu_i(1 - \mu_i)) = \text{diag}(\text{expit}(\mathbf{x}_i^\top \boldsymbol{\beta})(1 - \text{expit}(\mathbf{x}_i^\top \boldsymbol{\beta})))$.
- The likelihood-based equations below can be used to solve for $\boldsymbol{\beta}$:

$$\mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{X}^\top(\mathbf{y} - \text{vec}(\text{expit}(\mathbf{x}_i^\top \boldsymbol{\beta}))) = \mathbf{0}.$$

- Notice the nice cancellation in this case (no accident)!

Example 11.4: Bernoulli distribution (log link)

- Suppose our GLM is based on:
 - ▶ $Y \sim \text{Bernoulli}(p)$.
 - ▶ $g(\cdot)$ given by the log link (i.e., $g(\mu) = \log(\mu)$).
- Then,
 - ① Now, note that $\boldsymbol{\mu} = \text{vec}(\exp(\mathbf{x}_i^\top \boldsymbol{\beta}))$, so $\mathbf{D}^\top = \mathbf{X}^\top \text{diag}(\exp(\mathbf{x}_i^\top \boldsymbol{\beta}))$.
 - ② $\mathbf{V} = \text{diag}(\mu_i(1 - \mu_i)) = \text{diag}(\exp(\mathbf{x}_i^\top \boldsymbol{\beta})(1 - \exp(\mathbf{x}_i^\top \boldsymbol{\beta})))$.
- The likelihood-based equations below can be used to solve for $\boldsymbol{\beta}$:

$$\mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{X}^\top \text{diag}\left(\frac{1}{1 - \exp(\mathbf{x}_i^\top \boldsymbol{\beta})}\right)(\mathbf{y} - \text{vec}(\exp(\mathbf{x}_i^\top \boldsymbol{\beta}))) = \mathbf{0}.$$

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function**
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

Choosing a link function:

- GLM mean model:

$$g(E[Y_i]) = g(\mu_i) = g(b'(\theta_i)) = \mathbf{x}_i^T \boldsymbol{\beta} =: \eta_i$$

- What seems to be a really “nice” choice of a link function?
- The choice $g^{-1}(\cdot) = b'(\cdot)$ seems to be nice in the sense that under this choice, the canonical parameter is linear in $\boldsymbol{\beta}$. Specifically, $\theta_i = \mathbf{x}_i^T \boldsymbol{\beta}$.
- In fact, this choice of $g(\cdot)$ produces some very desirable mathematical properties, and is termed the canonical link function.

THE CANONICAL LINK FUNCTION

The canonical link: $g^{-1}(\cdot) = b'(\cdot)$

- GLM mean model under the canonical link function:

$$g(E[Y_i]) = g(\mu_i) = g(b'(\theta_i)) = \theta_i = \mathbf{x}_i^\top \boldsymbol{\beta} =: \eta_i$$

- Recall the likelihood-based equations for $\boldsymbol{\beta}$:

$$\left[\frac{\partial \ell}{\partial \boldsymbol{\beta}} \right]_{K \times 1} = \left[\frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} \right]_{K \times N} \left[\frac{\partial \boldsymbol{\theta}}{\partial \boldsymbol{\mu}} \right]_{N \times N} \left[\frac{\partial \ell}{\partial \boldsymbol{\theta}} \right]_{N \times 1}$$

- Latter two terms unaffected by choice of link function beyond the fact that $\boldsymbol{\mu} = g^{-1}(\mathbf{x}^\top \boldsymbol{\beta})$. However, if $g^{-1}(\cdot) = b'(\cdot)$:

$$\mathbf{D}^\top := \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\beta}} = \frac{\partial \boldsymbol{\eta}}{\partial \boldsymbol{\beta}} \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} = \mathbf{X}^\top \frac{\partial \boldsymbol{\mu}}{\partial \boldsymbol{\eta}} = \mathbf{X}^\top \frac{\partial}{\partial \boldsymbol{\theta}} \text{vec}(b'(\theta_i)) = \mathbf{X}^\top \mathbf{V}.$$

THE CANONICAL LINK FUNCTION

The canonical link: $g^{-1}(\cdot) = b'(\cdot)$

- GLM mean model under the canonical link function:

$$g(E[Y_i]) = g(\mu_i) = g(b'(\theta_i)) = \theta_i = \mathbf{x}_i^\top \boldsymbol{\beta} =: \eta_i$$

- Under the canonical link function, the likelihood-based equations take the following form:

$$\begin{aligned} \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) &= (\mathbf{X}^\top \mathbf{V}) \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \\ &= \mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{X}^\top (\mathbf{y} - \text{vec}(g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}))) = \mathbf{0}. \end{aligned}$$

- That is, $\mathbf{D}^\top \mathbf{V}^{-1} = \mathbf{X}^\top$ under the canonical link.
- Note that while \mathbf{D}^\top and \mathbf{V}^{-1} may depend upon $\boldsymbol{\beta}$, \mathbf{X}^\top does not.

THE CANONICAL LINK FUNCTION

The canonical link: $g^{-1}(\cdot) = b'(\cdot)$

- We saw this nice reduction in two of our prior examples!
- For a normal distribution, the identity link produced a score equation of this form. We can see this from factoring this exponential family:
 - ▶ $\theta = \mu$ is the natural parameter.
 - ▶ $b(\theta) = \theta^2/2$.
 - ▶ $b'(\theta) = \theta$.
 - ▶ $g(\mu) = \mu$ is the canonical link.
- For a Bernoulli distribution, the logit link produced a score equation of this form. We can see this from factoring this exponential family:
 - ▶ $\theta = \text{logit}(p)$ is the natural parameter.
 - ▶ $b(\theta) = \log(1 + \exp(\theta))$.
 - ▶ $b'(\theta) = \text{expit}(\theta)$.
 - ▶ $g(\mu) = \text{logit}(p)$ is the canonical link.
- Notice: the canonical link is determined by how θ relates to the target parameter, or how the inverse of $b'(\cdot)$ relates to θ .

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm**
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

The Gauss-Newton approach

- Let's return to the general case (i.e., not necessarily the canonical link function). The following equations arise from the likelihood:

$$\mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}.$$

- Again, I could expand the known expression for \mathbf{D}^T , but it is convenient to do that on the back end.
- These equations need not possess a closed-form solution. In many (most) cases, we must implement this fit numerically.

The Gauss-Newton approach: Definitions/notation

- Let $\mathbf{G}(\boldsymbol{\beta}; \mathbf{x}, Y) = \mathbf{x} \frac{\partial \mu / \partial \boldsymbol{\eta}}{b''(\boldsymbol{\eta})} (Y - g^{-1}(\mathbf{x}^\top \boldsymbol{\beta}))$.
 - ▶ This is the estimating function (in this set of notes, they happen to come from the score function).
- Let $\mathbb{G}_N(\boldsymbol{\beta}; \mathbf{X}, \mathbf{y}) = \mathbf{D}^\top \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu}) = \sum_{i=1}^N \mathbf{G}_i(\boldsymbol{\beta}; \mathbf{x}_i, Y_i)$.
 - ▶ The estimating equations are given by $\mathbb{G}_N(\boldsymbol{\beta}; \mathbf{X}, \mathbf{y}) = \mathbf{0}$, which—in this set of notes—happen to arise from score equations.

The Gauss-Newton approach: Definitions/notation

- Let $\mathbf{A}(\boldsymbol{\beta}) = \mathbb{E}\left[-\frac{\partial}{\partial\boldsymbol{\theta}}\mathbf{G}(\boldsymbol{\theta}; \mathbf{x}, Y)\bigg|_{\boldsymbol{\theta}=\boldsymbol{\beta}}\right]$.
 - ▶ Proportional to the expected Fisher information for one observation.
- Let $\mathbb{A}_N(\boldsymbol{\beta}) = \mathbb{E}\left[-\frac{\partial}{\partial\boldsymbol{\theta}}\mathbb{G}_N(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y})\bigg|_{\boldsymbol{\theta}=\boldsymbol{\beta}}\right]$.
 - ▶ Proportional to the expected Fisher information for N observations.

The Gauss-Newton approach: Steps

- Our approach will be similar to Newton's method—which you may recall from Calculus I, but I am including a reminder here.
- When seeking to solve $f(x) = 0$, we begin with an initial guess, x_0 , and update as follows:

$$x_{j+1} = x_j - \frac{f(x_j)}{f'(x_j)}$$

- In the next slide, I illustrate a few iterations of Newton's method for a function of a single variable.

The Gauss-Newton approach: Steps

The Gauss-Newton approach: Steps

- We now implement a similar method, though in higher dimensions.
- Like the version you learned in single-variable calculus, we have to propose an initial value (can be an art). Call this $\boldsymbol{\beta}^{(0)}$.
 - ▶ Example: Calibrate intercept to the observed data mean (on the scale of the link function) and initialize the remaining coefficients to zero.
- The $(j + 1)$ -step is given by: $\boldsymbol{\beta}^{(j+1)} = \boldsymbol{\beta}^{(j)} + [\mathbb{A}_N(\boldsymbol{\beta}^{(j)})]^{-1} \mathbb{G}_N(\boldsymbol{\beta}^{(j)})$.
- Note that $\mathbb{A}_N(\boldsymbol{\beta})$ involves an *expectation* of the (negative) derivative, distinguishing it slightly from the version of Newton's method you're already familiar with. Because $\mathbb{A}_N(\boldsymbol{\beta})$ is proportional to the expected Fisher information, this is also called the Fisher scoring algorithm.
 - ▶ See? I promised you'd learn it!
- For the sake of giving the algorithm a bit more specificity, let's determine more explicit expressions for $\mathbb{G}_N(\boldsymbol{\beta})$ and $\mathbb{A}_N(\boldsymbol{\beta})$ for GLMs.

The Gauss-Newton approach: More explicit expressions

- We already know that $\mathbb{G}_N(\boldsymbol{\beta}) = \mathbf{D}^\top(\boldsymbol{\beta})\mathbf{V}^{-1}(\boldsymbol{\beta})(\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}))$.
- Recall: $\mathbf{D}^\top(\boldsymbol{\beta}) = \mathbf{X}^\top \text{diag}(\partial\mu_i(\boldsymbol{\beta})/\partial\eta_i(\boldsymbol{\beta}))$.
- Recall: $\mathbf{V}(\boldsymbol{\beta}) = \text{diag}(V(\mu_i(\boldsymbol{\beta})))$.
- Therefore,

$$\begin{aligned}\mathbb{G}_N(\boldsymbol{\beta}) &= \mathbf{D}^\top(\boldsymbol{\beta})\mathbf{V}^{-1}(\boldsymbol{\beta})(\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \\ &= \mathbf{X}^\top \text{diag}\left(\left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\eta_i(\boldsymbol{\beta})}\right) \left(\frac{1}{V(\mu_i(\boldsymbol{\beta}))}\right)\right)(\mathbf{y} - \text{vec}(g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}))).\end{aligned}$$

- The estimating equations for a GLM look quite a bit like those from weighted least squares, with the weights depending upon the nature of the mean model and the mean-variance relationship!

The Gauss-Newton approach: More explicit expressions

- We can determine the derivative of $\mathbb{G}_N(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$:

$$\begin{aligned} -\frac{\partial}{\partial \boldsymbol{\beta}} \mathbb{G}_N(\boldsymbol{\beta}) &= -\left[\frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) \right] \text{diag} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})} \frac{1}{V(\mu_i(\boldsymbol{\beta}))} \right) \mathbf{X} \\ &\quad - \left[\frac{\partial}{\partial \boldsymbol{\beta}} \text{vec} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})} \frac{1}{V(\mu_i(\boldsymbol{\beta}))} \right) \right] \text{diag}(y_i - \mu_i(\boldsymbol{\beta})) \mathbf{X} \\ &= \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}) \mathbf{X} - \left[\frac{\partial}{\partial \boldsymbol{\beta}} \text{vec} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})} \frac{1}{V(\mu_i(\boldsymbol{\beta}))} \right) \right] \text{diag}(y_i - \mu_i(\boldsymbol{\beta})) \mathbf{X}, \end{aligned}$$

where $\mathbf{W}(\boldsymbol{\beta}) = \text{diag} \left(\left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})} \right)^2 \frac{1}{V(\mu_i(\boldsymbol{\beta}))} \right)$.

- Now, $\mathbb{A}_N(\boldsymbol{\beta})$ is the expectation (with respect to \mathbf{y}):

$$\mathbb{A}_N(\boldsymbol{\beta}) = \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}) \mathbf{X}.$$

- Why does $\mathbb{A}_N(\boldsymbol{\beta})$ reduce down in this way?

The Gauss-Newton approach: More explicitly for GLMs

- Note: $\mathbb{A}_N(\boldsymbol{\beta}) = \mathbf{D}^\top(\boldsymbol{\beta})\mathbf{V}^{-1}(\boldsymbol{\beta})\mathbf{D}(\boldsymbol{\beta})$. The expression we derived better tells us how to calculate it:

$$\begin{aligned}\mathbb{A}_N(\boldsymbol{\beta}) &= \mathbf{D}^\top(\boldsymbol{\beta})\mathbf{V}^{-1}(\boldsymbol{\beta})\mathbf{D}(\boldsymbol{\beta}) \\ &= \mathbf{X}^\top\mathbf{W}(\boldsymbol{\beta})\mathbf{X} \\ &= \mathbf{X}^\top\text{diag}\left(\left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\eta_i(\boldsymbol{\beta})}\right)^2\frac{1}{V(\mu_i(\boldsymbol{\beta}))}\right)\mathbf{X}.\end{aligned}$$

The Gauss-Newton approach: More explicitly for GLMs

- Given these forms for $\mathbb{G}_N(\boldsymbol{\beta})$ and $\mathbb{A}_N(\boldsymbol{\beta})$, the Gauss-Newton step can be more expressly described.
- The $(j + 1)$ -step, $\boldsymbol{\beta}^{(j+1)}$, is given by:

$$\boldsymbol{\beta}^{(j)} + [\mathbb{A}_N(\boldsymbol{\beta}^{(j)})]^{-1} \mathbb{G}_N(\boldsymbol{\beta}^{(j)})$$

$$\boldsymbol{\beta}^{(j)} + [\mathbf{D}^\top(\boldsymbol{\beta}^{(j)})\mathbf{V}^{-1}(\boldsymbol{\beta}^{(j)})\mathbf{D}(\boldsymbol{\beta}^{(j)})]^{-1} \mathbf{D}^\top(\boldsymbol{\beta}^{(j)})\mathbf{V}^{-1}(\boldsymbol{\beta}^{(j)})(\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(j)}))$$

$$\boldsymbol{\beta}^{(j)} + [\mathbf{X}^\top \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{X}]^{-1} \mathbf{X}^\top \mathbf{W}(\boldsymbol{\beta}^{(j)}) \left(\frac{\partial \eta_i(\boldsymbol{\beta}^{(j)})}{\partial \mu_i(\boldsymbol{\beta}^{(j)})} \right) (\mathbf{y} - \text{vec}(g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}^{(j)}))),$$

where $\mathbf{W}(\boldsymbol{\beta}^{(j)}) = \text{diag}((\partial \mu_i(\boldsymbol{\beta}^{(j)}) / \partial \eta_i(\boldsymbol{\beta}^{(j)}))^2 (1/V(\mu_i(\boldsymbol{\beta}^{(j)}))))$.

- Don't worry—there will be examples of how to code this later in the notes—it's not as bad as it may appear.

The Gauss-Newton approach: GLMs with canonical link

- Under the canonical link, $\mathbf{D}^T(\boldsymbol{\beta})\mathbf{V}^{-1}(\boldsymbol{\beta}) = \mathbf{X}^T$ and $\mathbf{D}(\boldsymbol{\beta}) = \mathbf{V}(\boldsymbol{\beta})\mathbf{X}$.
- Therefore, the estimating equations are given by:

$$\mathbb{G}_N(\boldsymbol{\beta}) = \mathbf{X}^T(\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta})) = \mathbf{0},$$

and $\mathbb{A}_N(\boldsymbol{\beta})$ is given by:

$$\mathbb{A}_N(\boldsymbol{\beta}) = \mathbf{D}^T(\boldsymbol{\beta})\mathbf{V}^{-1}(\boldsymbol{\beta})\mathbf{D}(\boldsymbol{\beta}) = \mathbf{X}^T\mathbf{V}(\boldsymbol{\beta})\mathbf{X}.$$

- With this in mind, the $(j + 1)$ -step, $\boldsymbol{\beta}^{(j+1)}$, is given by:

$$\boldsymbol{\beta}^{(j)} + \left[\mathbb{A}_N(\boldsymbol{\beta}^{(j)})\right]^{-1} \mathbb{G}_N(\boldsymbol{\beta}^{(j)})$$

$$\boldsymbol{\beta}^{(j)} + \left[\mathbf{D}^T(\boldsymbol{\beta}^{(j)})\mathbf{V}^{-1}(\boldsymbol{\beta}^{(j)})\mathbf{D}(\boldsymbol{\beta}^{(j)})\right]^{-1} \mathbf{D}^T(\boldsymbol{\beta}^{(j)})\mathbf{V}^{-1}(\boldsymbol{\beta}^{(j)})(\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(j)}))$$

$$\boldsymbol{\beta}^{(j)} + \left[\mathbf{X}^T\mathbf{V}(\boldsymbol{\beta}^{(j)})\mathbf{X}\right]^{-1} \mathbf{X}^T(\mathbf{y} - \text{vec}(g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)}))).$$

The Newton-Raphson approach:

- You can base the algorithm on the *observed* derivative, $\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta})$, rather than its expectation (referred to as the Newton-Raphson algorithm), though it can be somewhat messier to do so:

$$\begin{aligned}\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta}) &= -\frac{\partial}{\partial \boldsymbol{\beta}} \mathbb{G}_N(\boldsymbol{\beta}) \\ &= \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}) \mathbf{X} - \left[\frac{\partial}{\partial \boldsymbol{\beta}} \text{vec} \left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})} \frac{1}{V(\mu_i(\boldsymbol{\beta}))} \right) \right] \text{diag}(y_i - \mu_i(\boldsymbol{\beta})) \mathbf{X},\end{aligned}$$

where $\mathbf{W}(\boldsymbol{\beta}^{(j)}) = \text{diag}((\partial \mu_i(\boldsymbol{\beta}^{(j)}) / \partial \eta_i(\boldsymbol{\beta}^{(j)}))^2 (1/V(\mu_i(\boldsymbol{\beta}))))$.

- Caveat: if the canonical link is used, $\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta})$ and $\mathbb{A}_N(\boldsymbol{\beta})$ are the same! Why?

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance**
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

Score equations:

- If forming estimates of the variance based on likelihood theory, we need to revisit the exact score equations. We did this for $\boldsymbol{\beta}$ already:

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = \mathbf{D}^T \mathbf{V}^{-1}(\mathbf{y} - \boldsymbol{\mu})/\phi.$$

- If $\phi = 1$, there is no need to worry about estimating ϕ , with caveats that we will discuss when we cover quasi-likelihood.
- Now, if there is a nuisance parameter, we do need to give it some attention; the score equations for ϕ take the following form:

$$\begin{aligned} \frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^N \frac{\partial}{\partial \phi} \left[\frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi) \right] \\ &= \sum_{i=1}^N \left[-\frac{1}{\phi^2} (y_i \theta_i - b(\theta_i)) + c'(y_i, \phi) \right] \end{aligned}$$

Information:

- The (expected) information matrix takes the following form:

$$\mathcal{I}(\boldsymbol{\beta}, \phi) = -E \begin{bmatrix} \frac{\partial^2 \ell(\boldsymbol{\beta}, \phi)}{\partial \boldsymbol{\beta}^2} & \frac{\partial^2 \ell(\boldsymbol{\beta}, \phi)}{\partial \boldsymbol{\beta} \partial \phi} \\ \frac{\partial^2 \ell(\boldsymbol{\beta}, \phi)}{\partial \phi \partial \boldsymbol{\beta}} & \frac{\partial^2 \ell(\boldsymbol{\beta}, \phi)}{\partial \phi^2} \end{bmatrix}$$

- Note the equivalence of the off-diagonal blocks.
- Further, the lower-right block is not of primary interest and we won't compute it explicitly.

Information:

- Tackling the upper-left block:

$$-E \left[\frac{\partial^2 \ell(\boldsymbol{\beta}, \phi)}{\partial \boldsymbol{\beta}^2} \right] = \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \phi$$

- See presentation of the Gauss-Newton algorithm for this derivation (we're just including the term ϕ this time).

Information:

- Tackling the off-diagonal blocks:

$$\begin{aligned}
 -\mathbb{E} \left[\frac{\partial^2 \ell(\boldsymbol{\beta}, \phi)}{\partial \phi \partial \boldsymbol{\beta}} \right] &= \mathbb{E} \left[\sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} \left[\frac{1}{\phi^2} (y_i \theta_i - b(\theta_i)) + c'(y_i, \phi) \right] \right] \\
 &= \frac{1}{\phi^2} \mathbb{E} \left[\sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\beta}} (y_i \theta_i - b(\theta_i)) \right] \\
 &= \frac{1}{\phi^2} \mathbb{E} [\mathbf{D}^T \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})] = \mathbf{0}
 \end{aligned}$$

- See derivation of the score equations for $\boldsymbol{\beta}$ for this result, and recognize the appearance of the unbiased estimating equations.

Information:

- Here's the bottom line:

$$\mathcal{I}(\boldsymbol{\beta}, \phi) = \begin{bmatrix} \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D} / \phi & \mathbf{0} \\ \mathbf{0}^T & \dots \end{bmatrix}$$

- Because of the block-diagonal structure of the Fisher information, we need not propagate uncertainty in estimation of ϕ in estimating the variance of $\hat{\boldsymbol{\beta}}$.
 - ▶ Why is this the case? This fact is worth understanding and not one to simply take at face value.

Asymptotic distribution:

- Under standard likelihood theory,

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}(\boldsymbol{\beta}, \phi(\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D})^{-1}).$$

- To estimate $\text{Cov}[\hat{\boldsymbol{\beta}}]$, we could be in one of two cases:
 - 1 $\phi = 1$, in which case $\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = (\mathbf{D}^T(\hat{\boldsymbol{\beta}})\mathbf{V}^{-1}(\hat{\boldsymbol{\beta}})\mathbf{D}(\hat{\boldsymbol{\beta}}))^{-1}$.
 - 2 Otherwise, we require a consistent estimate of ϕ . This one will do:

$$\hat{\phi} = \frac{1}{N - K} \sum_{i=1}^N \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)},$$

and we have that $\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = \hat{\phi}(\mathbf{D}^T(\hat{\boldsymbol{\beta}})\mathbf{V}^{-1}(\hat{\boldsymbol{\beta}})\mathbf{D}(\hat{\boldsymbol{\beta}}))^{-1}$, where again, you can use the more detailed formulas for \mathbf{D}^T .

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs**
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

EXAMPLE: LOGISTIC REGRESSION

Example 11.5: Bernoulli distribution (logit link)

- Mean model: $\text{logit}(E[Y]) = \mathbf{x}^T \boldsymbol{\beta}$ ($\mu_i = \text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}) = \text{expit}(\eta_i)$).
- Mean-variance relationship: $\mathbf{V} = \text{diag}(\mu_i(1 - \mu_i))$.
- Solving the following equations (previously derived) gives us $\hat{\boldsymbol{\beta}}$:

$$\mathbf{X}^T(\mathbf{y} - \text{vec}(\text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}.$$

- Because we are in the setting of the canonical link,

$$\mathbb{A}_N(\boldsymbol{\beta}) = \mathbf{X}^T \mathbf{V}(\boldsymbol{\beta}) \mathbf{X} = \mathbf{X}^T \text{diag}(\text{expit}(\mathbf{x}_i^T \boldsymbol{\beta})(1 - \text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}))) \mathbf{X}.$$

- Each Gauss-Newton step, $\boldsymbol{\beta}^{(j+1)}$, is given by:

$$\boldsymbol{\beta}^{(j)} + (\mathbf{X}^T \text{diag}(\text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)})(1 - \text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)}))) \mathbf{X})^{-1} \mathbf{X}^T(\mathbf{y} - \text{vec}(\text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)}))).$$

- Further, the variance can be estimated as:

$$\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = (\mathbf{X}^T \text{diag}(\text{expit}(\mathbf{x}_i^T \hat{\boldsymbol{\beta}})(1 - \text{expit}(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))) \mathbf{X})^{-1}.$$

EXAMPLE: LOGISTIC REGRESSION

Example 11.5: Simulation (Setup)

```
1 ## Define expit function
2 expit <- function(x)
3 {
4   expitx <- exp(x)/(1 + exp(x))
5   return(expitx)
6 }
7
8 ## Set seed for reproducibility
9 set.seed(7345)
10
11 ## Set sample size
12 n <- 100
13
14 ## Set true value of beta
15 beta <- matrix(c(0.7, 1.2), nrow = 2)
```

Example 11.5: Simulation (Data generation/set tolerance)

```
1 ## Generate data
2 X <- matrix(cbind(1, rnorm(n, 0, 1)), ncol = 2)
3 y <- rbinom(n, 1, expit(X %*% beta))
4
5 ## Initialize beta
6 betaj <- c(log(mean(y)/(1 - mean(y))), 0)
7
8 ## Set tolerance and count iterations
9 tol <- 1
10 iter <- 1
```

EXAMPLE: LOGISTIC REGRESSION

Example 11.5: Simulation (Gauss-Newton iterations)

```
1 while(tol > 1e-12 & iter < 50)
2 {
3   ## Store prior iteration
4   betaj.prior <- betaj
5
6   ## Linear predictor
7   etaj <- c(X %*% betaj)
8
9   ## Estimating function
10  Gn <- t(X) %*% (y - expit(etaj))
11
12  ## Mean-variance relationship
13  Vn <- as.numeric(expit(etaj)*(1 - expit(etaj)))
14
15  ## An
16  An <- t(X * Vn) %*% X
17
18  ## Update to next iteration
19  betaj <- betaj + solve(An) %*% Gn
20
21  ## Determine if conditions are met
22  tol <- sum((betaj - betaj.prior)^2)
23  iter <- iter + 1
24 }
```

Example 11.5: Simulation (Report results)

```
1 ## Report number of iterations
2 > print(iter)
3 [1] 7
4
5 ## Report tolerance
6 > print(tol)
7 [1] 1.218104e-23
```

EXAMPLE: LOGISTIC REGRESSION

Example 11.5: Simulation (Report estimates/standard errors)

```
1 ## Report estimate
2 > betaj
3           [,1]
4 [1,] 0.7115884
5 [2,] 1.2642435
6
7 ## Obtain and report standard errors (re-update)
8 Vn <- as.numeric(expit(c(X %*% betaj))*(1 - expit(c(X %*% betaj))))
9 An <- t(X * Vn) %*% X
10
11 > sqrt(diag(solve(An)))
12 [1] 0.2394554 0.3358070
13
14 ## Compare to R's glm capabilities
15 zz <- glm(y ~ X[,2], family = binomial(link = "logit"))
16
17 > summary(zz)$coef
18           Estimate Std. Error z value Pr(>|z|)
19 (Intercept) 0.711588  0.239450 2.97176 0.002960970
20 X[, 2]      1.264244  0.335792 3.76497 0.000166572
21
22 ## More stringent convergence threshold (perfect agreement)
23 zz2 <- glm(y ~ X[,2], family = binomial(link = "logit"), control = list(epsilon = 1e-12))
24
25 > summary(zz2)$coef
26 Coefficients:
27           Estimate Std. Error z value Pr(>|z|)
28 (Intercept) 0.7115884  0.2394554 2.97169 0.00296161 **
29 X[, 2]      1.2642435  0.3358070 3.76479 0.00016669 ***
```

Advantages of hard-coding:

- If you create a function to use the `glm()` function to run logistic regression 10,000 times, it takes about four seconds (on my computer).
 - ▶ Running the same number of iterations of the hard-coded version takes about one second.
- When we talk about different variance estimation procedures, you'll have more flexibility to modify them as appropriate without being at the mercy of others' pre-written code.
- Know what's going on under the hood!

A note about computation:

- Keep in mind that I employed a computational efficiency trick in hard-coding the GLM.
- We often write various estimating equations, algorithms, and variance estimators using the notation $\text{diag}(\cdot)$.
- A large diagonal matrix has a lot of zeros; multiplying diagonal matrices the “literal” way is computationally wasteful.
- In R, the following two commands produce the same answer:

```
1 out1 <- t(X) %*% diag(z)
2 out2 <- t(X * z)
```

- ... but the difference in computation time is *not* subtle when your data set is of moderate or larger size. Let's see a speed test in action.

A note about computation: Speed test (small data)

```
1 ## Set seed for reproducibility (computing time varies)
2 set.seed(7345)
3
4 ## Sample size
5 N <- 100
6
7 ## Number of parameters
8 K <- 2
9
10 ## Matrix
11 X <- matrix(rnorm(N*K), nrow = N)
12
13 ## Vector
14 z <- rnorm(N)
15
16 ## Using our symbolic representation
17 > system.time(for (j in 1:50) {t(X) %*% diag(z)})
18   user  system elapsed
19  0.004   0.001   0.004
20
21 ## Computational trick
22 > system.time(for (j in 1:50){t(X * z)})
23   user  system elapsed
24  0.001   0.000   0.001
25
26 ## No big deal so far...
```

A note about computation: Speed test (less-small data)

```
1 ## Set seed for reproducibility (computing time varies)
2 set.seed(7345)
3
4 ## Sample size
5 N <- 500
6
7 ## Number of parameters
8 K <- 3
9
10 ## Matrix
11 X <- matrix(rnorm(N*K), nrow = N)
12
13 ## Vector
14 z <- rnorm(N)
15
16 ## Using our symbolic representation
17 > system.time(for (j in 1:50) {t(X) %*% diag(z)})
18   user  system elapsed
19  0.061   0.006   0.069
20
21 ## Computational trick
22 > system.time(for (j in 1:50){t(X * z)})
23   user  system elapsed
24  0.001   0.000   0.001
25
26 ## Still not convinced...
```

A note about computation: Speed test (healthy-sized data)

```
1 ## Set seed for reproducibility (computing time varies)
2 set.seed(7345)
3
4 ## Sample size
5 N <- 1000
6
7 ## Number of parameters
8 K <- 4
9
10 ## Matrix
11 X <- matrix(rnorm(N*K), nrow = N)
12
13 ## Vector
14 z <- rnorm(N)
15
16 ## Using our symbolic representation
17 > system.time(for (j in 1:50) {t(X) %*% diag(z)})
18   user  system elapsed
19  0.263   0.025   0.289
20
21 ## Computational trick
22 > system.time(for (j in 1:50){t(X * z)})
23   user  system elapsed
24  0.002   0.000   0.002
25
26 ## I might be getting closer to being convinced...
```

A note about computation: Speed test (moderately-sized data)

```
1 ## Set seed for reproducibility (computing time varies)
2 set.seed(7345)
3
4 ## Sample size
5 N <- 5000
6
7 ## Number of parameters
8 K <- 4
9
10 ## Matrix
11 X <- matrix(rnorm(N*K), nrow = N)
12
13 ## Vector
14 z <- rnorm(N)
15
16 ## Using our symbolic representation
17 > system.time(for (j in 1:50) {t(X) %*% diag(z)})
18   user  system elapsed
19  9.931   0.827  10.772
20
21 ## Computational trick
22 > system.time(for (j in 1:50){t(X * z)})
23   user  system elapsed
24  0.004   0.001   0.005
25
26 ## I think I'm convinced...
```

A note about computation: Speed test (hefty data)

```
1 ## Set seed for reproducibility (computing time varies)
2 set.seed(7345)
3
4 ## Sample size
5 N <- 10000
6
7 ## Number of parameters
8 K <- 4
9
10 ## Matrix
11 X <- matrix(rnorm(N*K), nrow = N)
12
13 ## Vector
14 z <- rnorm(N)
15
16 ## Using our symbolic representation
17 > system.time(for (j in 1:50) {t(X) %*% diag(z)})
18   user   system elapsed
19 31.672   3.296  35.307
20
21 ## Computational trick
22 > system.time(for (j in 1:50){t(X * z)})
23   user   system elapsed
24  0.012   0.004   0.017
25
26 ## Getting out of hand!!
```

A note about computation:

- The bottom line is this: we'll continue to use $\text{diag}(\cdot)$ notation as our typical symbolic representation of estimating equations, algorithms, variances, etc.
- For coding, however, you will save yourself time by taking advantage of these simple computational tricks.
- Let's try another example of a GLM we've probably not encountered before.

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Example 11.6: Normal regression (log link)

- Mean model: $\log(E[Y]) = \mathbf{x}^T \boldsymbol{\beta}$ ($\mu_i = \exp(\mathbf{x}_i^T \boldsymbol{\beta}) = \exp(\eta_i)$).
- Mean-variance relationship: $\mathbf{V} = \mathbf{I}$.
- $\theta_i = \mu_i$ is the natural parameter; $b(\theta_i) = \mu_i^2/2$, and $b''(\theta_i) = 1$ (see prior notes in which we factored density into exponential form).
- Solving the following equations (previously derived) gives us $\hat{\boldsymbol{\beta}}$:

$$\mathbf{D}^T \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta})) (\mathbf{y} - \text{vec}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}.$$

- Because this is a non-canonical link,

$$\begin{aligned} \mathbb{A}_N(\boldsymbol{\beta}) &= \mathbf{D}^T \mathbf{V}^{-1} \mathbf{D}^{-1} = \mathbf{X}^T \text{diag} \left(\left(\frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \eta_i(\boldsymbol{\beta})} \right)^2 \frac{1}{V(\mu_i(\boldsymbol{\beta}))} \right) \mathbf{X} \\ &= \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))^2 \mathbf{X}. \end{aligned}$$

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Example 11.6: Normal regression (log link)

- Each Gauss-Newton step, $\boldsymbol{\beta}^{(j+1)}$, is given by:

$$\boldsymbol{\beta}^{(j)} + (\mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)}))^2 \mathbf{X})^{-1} \mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)})) (\mathbf{y} - \text{vec}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)}))).$$

- Here, we have a nuisance parameter that must be estimated:

$$\hat{\phi} = \frac{1}{N - K} \sum_{i=1}^N \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)} = \frac{1}{N - K} \sum_{i=1}^N (y_i - \exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))^2.$$

- Therefore, this gives rise to the following estimator for $\text{Cov}[\hat{\boldsymbol{\beta}}]$:

$$\widehat{\text{Cov}}[\hat{\boldsymbol{\beta}}] = \left(\frac{1}{N - K} \sum_{i=1}^N (y_i - \exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))^2 \right) \left(\mathbf{X}^T \text{diag}(\exp(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}))^2 \mathbf{X} \right)^{-1}.$$

Example 11.6: Simulation (Setup)

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Set sample size
5 n <- 50
6
7 ## Set true value of beta
8 beta <- matrix(c(-0.5, 0.6), nrow = 2)
```

Example 11.6: Simulation (Data generation/set tolerance)

```
1 ## Generate data
2 X <- matrix(cbind(1, rnorm(n, 0, 1)), ncol = 2)
3 y <- rnorm(n, exp(X %*% beta))
4
5 ## Initialize beta
6 beta_j <- c(log(mean(y)), 0)
7
8 ## Set tolerance and count iterations
9 tol <- 1
10 iter <- 1
```

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Example 11.6: Simulation (Gauss-Newton iterations)

```
1 while(tol > 1e-15 & iter < 50)
2 {
3   ## Store prior iteration
4   betaj.prior <- betaj
5
6   ## Linear predictor
7   etaj <- c(X %*% betaj)
8
9   ## Estimating function
10  Gn <- t(X * exp(etaj)) %*% (y - exp(etaj))
11
12  ## An
13  An <- t(X * exp(etaj)^2) %*% X
14
15  ## Update to next iteration
16  betaj <- betaj + solve(An) %*% Gn
17
18  ## Nuisance
19  phi <- sum((y - exp(etaj))^2)/(n - 2)
20
21  ## Determine if conditions are met
22  tol <- sum((betaj - betaj.prior)^2)
23  iter <- iter + 1
24 }
```

Example 11.6: Simulation (Report results)

```
1 ## Report number of iterations
2 > print(iter)
3 [1] 6
4
5 ## Report tolerance
6 > print(tol)
7 [1] 4.620906e-17
```

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Example 11.6: Simulation (Report estimates/standard errors)

```
1 ## Report estimate
2 > beta_j
3           [,1]
4 [1,] -0.2832756
5 [2,]  0.4282233
6
7 ## Obtain and report standard errors (re-update)
8 An <- t(X * exp(c(X %*% beta_j))^2) %*% X
9
10 > sqrt(diag(phi * solve(An)))
11 [1] 0.1804808 0.1373220
12
13 ## Compare to R's glm capabilities
14 zz <- glm(y ~ X[,2], start = c(log(mean(y)), 0), family = gaussian(link = "log"))
15
16 > summary(zz)$coef
17           Estimate Std. Error z value Pr(>|z|)
18 (Intercept) -0.2833      0.1805  -1.570  0.12309
19 X[, 2]       0.4282      0.1373   3.118  0.00307 **
```

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Example 11.6: Normal regression (log link)

- Let's try re-doing this based on the Newton-Raphson algorithm.
- Recall: $(\partial/\partial\boldsymbol{\beta})\text{vec}(f(\mathbf{x}_i^T\boldsymbol{\beta})) = \mathbf{X}^T\text{diag}(f'(\mathbf{x}_i^T\boldsymbol{\beta}))$.
- More detail: $\mu_i(\boldsymbol{\beta}) = \exp(\mathbf{x}_i^T\boldsymbol{\beta})$ and $b''(\theta_i(\boldsymbol{\beta})) = 1$ in this example, so:

$$\begin{aligned}\frac{\partial}{\partial\boldsymbol{\beta}}\left[\text{vec}\left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\eta_i(\boldsymbol{\beta})}\frac{1}{V(\mu_i(\boldsymbol{\beta}))}\right)\right] &= \frac{\partial}{\partial\boldsymbol{\beta}}\left[\text{vec}\left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\eta_i(\boldsymbol{\beta})}\right)\right] = \frac{\partial}{\partial\boldsymbol{\beta}}[\text{vec}(\exp(\mathbf{x}_i^T\boldsymbol{\beta}))] \\ &= \begin{bmatrix} \exp(\mathbf{x}_1^T\boldsymbol{\beta}) & \exp(\mathbf{x}_2^T\boldsymbol{\beta}) & \cdots & \exp(\mathbf{x}_N^T\boldsymbol{\beta}) \\ x_{11}\exp(\mathbf{x}_1^T\boldsymbol{\beta}) & x_{21}\exp(\mathbf{x}_2^T\boldsymbol{\beta}) & \cdots & x_{N1}\exp(\mathbf{x}_N^T\boldsymbol{\beta}) \\ \vdots & \vdots & \ddots & \vdots \\ x_{1K}\exp(\mathbf{x}_1^T\boldsymbol{\beta}) & x_{2K}\exp(\mathbf{x}_2^T\boldsymbol{\beta}) & \cdots & x_{NK}\exp(\mathbf{x}_N^T\boldsymbol{\beta}) \end{bmatrix} \\ &= \mathbf{X}^T\text{diag}(\exp(\mathbf{x}_i^T\boldsymbol{\beta})).\end{aligned}$$

- Therefore,

$$\begin{aligned}\mathbb{A}_N^{\text{obs}}(\boldsymbol{\beta}) &= \mathbb{A}_N(\boldsymbol{\beta}) - \frac{\partial}{\partial\boldsymbol{\beta}}\left[\text{vec}\left(\frac{\partial\mu_i(\boldsymbol{\beta})}{\partial\eta_i(\boldsymbol{\beta})}\frac{1}{V(\mu_i(\boldsymbol{\beta}))}\right)\right]\text{diag}(y_i - \mu_i(\boldsymbol{\beta}))\mathbf{X} \\ &= \mathbf{X}^T\text{diag}(\exp(\mathbf{x}_i^T\boldsymbol{\beta})^2 - \exp(\mathbf{x}_i^T\boldsymbol{\beta})(y_i - \exp(\mathbf{x}_i^T\boldsymbol{\beta})))\mathbf{X}.\end{aligned}$$

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Example 11.6: Simulation (Newton-Raphson iterations)

```
1 while(tol > 1e-15 & iter < 50)
2 {
3   ## Store prior iteration
4   betaj.prior <- betaj
5
6   ## Linear predictor
7   etaj <- c(X %*% betaj)
8
9   ## Estimating function
10  Gn <- t(X * exp(etaj)) %*% (y - exp(etaj))
11
12  ## An
13  An <- t(X * exp(etaj)^2) %*% X
14
15  ## AnObs
16  AnObs <- An - t(X * c(exp(etaj) * (y - exp(etaj)))) %*% X
17
18  ## Update to next iteration
19  betaj <- betaj + solve(AnObs) %*% Gn
20
21  ## Nuisance
22  phi <- sum((y - exp(etaj))^2)/(n - 2)
23
24  ## Determine if conditions are met
25  tol <- sum((betaj - betaj.prior)^2)
26  iter <- iter + 1
27 }
```

Example 11.6: Simulation (Report results)

```
1 ## Report number of iterations
2 > print(iter)
3 [1] 8
4
5 ## Report tolerance
6 > print(tol)
7 [1] 2.501604e-23
8
9 ## Report estimate
10 >      betaj
11          [,1]
12 [1,] -0.2832756
13 [2,]  0.4282233
```

EXAMPLE: NORMAL REGRESSION WITH A LOG LINK

Notes:

- It is typical for the Gauss-Newton algorithm to converge in fewer iterations. The algorithms converge to the **same** quantity.
- Note that the likelihood-based variance estimator is still based on $\mathbb{A}_N(\hat{\boldsymbol{\beta}})$. We don't change the variance estimator based on the algorithm (although we will discuss in the next set of notes other uses of $\mathbb{A}_N^{\text{obs}}(\hat{\boldsymbol{\beta}})$ for variance estimation).

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Inverse-Gaussian

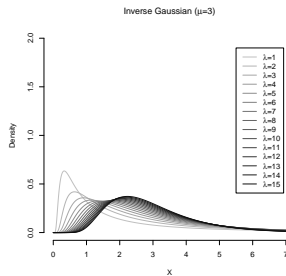
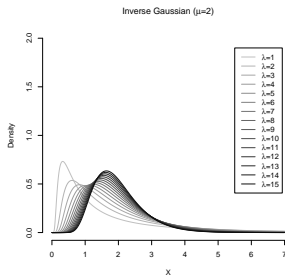
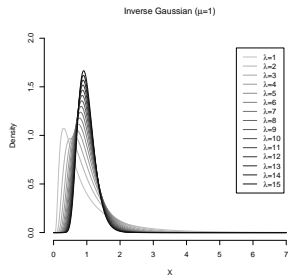
- Let's go through an example that we haven't worked with as closely. Consider modeling the outcome as an inverse Gaussian distribution, which has the following density function:

$$f_X(x; \mu, \lambda) = \sqrt{\frac{\lambda}{2\pi x^3}} \exp\left(-\frac{\lambda(x - \mu)^2}{2\mu^2 x}\right) \mathbf{1}(x > 0).$$

- Parameters: $\mu > 0$ (mean), and $\lambda > 0$ (shape) or $1/\lambda$ (dispersion).
- Goals for this example:
 - Write this exponential family in the canonical form.
 - Determine $E[Y]$ and $\text{Var}[Y]$ based on this canonical form.
 - Determine the canonical link based on the corresponding GLM.
 - Determine the score equations for $\boldsymbol{\beta}$ based on the canonical link.
 - Determine the variance of $\hat{\boldsymbol{\beta}}$ based on the above model.
 - Use R to iterate the Gauss-Newton algorithm on a simulated example.

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Density function



EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Density function in canonical form

- Time for some algebra:

$$\begin{aligned}f_y(y; \mu, \lambda) &= \sqrt{\frac{\lambda}{2\pi y^3}} \exp\left(-\frac{\lambda(y - \mu)^2}{2\mu^2 y}\right) \\&= \exp\left(-\frac{\lambda(y^2 - 2\mu y + \mu^2)}{2\mu^2 y} + \frac{1}{2}(\log(\lambda) - \log(2\pi y^3))\right) \\&= \exp\left(-\frac{y/2\mu^2 - 1/\mu + 1/2y}{1/\lambda} + \frac{1}{2}(\log(\lambda) - \log(2\pi y^3))\right) \\&= \exp\left(\frac{y(-1/(2\mu^2)) - (-1/\mu)}{1/\lambda} + \frac{1}{2}\left(\frac{\lambda}{y} + \log(\lambda) - \log(2\pi y^3)\right)\right).\end{aligned}$$

- Natural parameter: $\theta = -1/(2\mu^2)$; nuisance parameter: $\phi = 1/\lambda$.
- $b(\theta) = -(-2\theta)^{1/2}$.
- $E[Y] = b'(\theta) = (-2\theta)^{-1/2} = \mu$.
- $b''(\theta) = (-2\theta)^{-3/2} \Rightarrow V(\mu) = \mu^3$.
- $\text{Var}[Y] = \phi b''(\theta) = \dots = \mu^3/\lambda$.

Example 11.7: Some points of caution

- What's to stop us from choosing $\theta = 1/\mu^2$ and $\phi = -2/\lambda$?
 - ▶ Well... nothing. In fact, a prior iteration of these slides had exactly that parameterization and everything worked out fine...
- However, there is some benefit to choosing the nuisance parameter such that $b''(\theta)$ is a non-scaled transformation of $\mu = b'(\theta)$.
- In particular, having a standard—even if arbitrary—procedure helps us avoid getting tripped up by ambiguity.
 - ▶ Defaulting to $V(\mu) = \mu^3$ (equality rather than proportionality) helps with what you might think of as “GLM bookkeeping.”

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Canonical link function

- The canonical link function is based on $g^{-1}(\cdot) = b'(\cdot)$.
- More explicitly for an inverse Gaussian GLM:

$$\begin{aligned}g^{-1}(\eta_i) &= b'(\theta) \\ &= 1/\sqrt{-2\theta} \\ \Rightarrow g(\mu_i) &= -1/(2\mu_i^2)\end{aligned}$$

- The above formulation of the canonical link is *technically correct*, but in practice we often choose $g(\mu_i) = 1/\mu_i^2$ (this does not change the mechanics of the GLM procedures).
- Under this “near-canonical” link:

$$\begin{aligned}1/\mu^2 = 1/(E[Y])^2 &= \eta = \mathbf{x}^T \boldsymbol{\beta}, \text{ or} \\ \mu = E[Y] &= 1/\sqrt{\eta} = 1/\sqrt{\mathbf{x}^T \boldsymbol{\beta}}.\end{aligned}$$

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Score equations

- Under this “near-canonical” link, the score function for $\boldsymbol{\beta}$ is given by:

$$\begin{aligned}\mathbb{G}_N(\boldsymbol{\beta}) &= \mathbf{X}^\top \text{diag}((\partial \mu_i / \partial \eta_i)) / V_i (\mathbf{y} - \text{vec}(\boldsymbol{\mu}_i)) / \phi \\ &= \mathbf{X}^\top \text{diag} \left(-\frac{1}{2} (\mathbf{x}_i^\top \boldsymbol{\beta})^{-3/2} \frac{1}{(\mathbf{x}_i^\top \boldsymbol{\beta})^{-3/2}} \right) (\mathbf{y} - \text{vec}((\mathbf{x}_i^\top \boldsymbol{\beta})^{-1/2})) / \phi \\ &= (-1/2) \mathbf{X}^\top (\mathbf{y} - \text{vec}((\mathbf{x}_i^\top \boldsymbol{\beta})^{-1/2})) / \phi.\end{aligned}$$

- While you can ignore the nuisance when solving the estimating equations, I believe you’ll ultimately be happier leaving the “ $-1/2$ ” in so that you never have to worry if you threw off a constant in the $V(\cdot)$ function.
- Any such constant will be absorbed into estimation of ϕ , but...
 - ▶ What if $\phi = 1$ is known?
 - ▶ Even if ϕ is unknown, you might be wondering why you’re not estimating it consistently (say, in a simulation).

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Gauss-Newton implementation

- We can determine $\mathbb{A}_N(\boldsymbol{\beta})$ as follows:

$$\begin{aligned}\mathbb{A}_N(\boldsymbol{\beta}) &= \mathbf{X}^\top \text{diag} \left(\left(-\frac{1}{2}(\mathbf{x}_i^\top \boldsymbol{\beta})^{-3/2} \right)^2 \frac{1}{(\mathbf{x}_i^\top \boldsymbol{\beta})^{-3/2}} \right) \mathbf{X} \\ &= \frac{1}{4} \mathbf{X}^\top \text{diag} \left((\mathbf{x}_i^\top \boldsymbol{\beta})^{-3/2} \right) \mathbf{X}.\end{aligned}$$

- Having a “near-canonical” link just means that the usual mechanics of a canonical link are off by a constant.
- Each Gauss-Newton step, $\boldsymbol{\beta}^{(j+1)}$, is given by:

$$\boldsymbol{\beta}^{(j+1)} = \boldsymbol{\beta}^{(j)} - 2(\mathbf{X}^\top \text{diag}((\mathbf{x}_i^\top \boldsymbol{\beta}^{(j)})^{-3/2})\mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{y} - \text{vec}((\mathbf{x}_i^\top \boldsymbol{\beta}^{(j)})^{-1/2})).$$

- What do you think the Newton-Raphson algorithm will look like?

Example 11.7: Variance

- We can also determine the covariance matrix for $\hat{\boldsymbol{\beta}}$:

$$\text{Cov}[\hat{\boldsymbol{\beta}}] = 4\phi(\mathbf{X}^T \text{diag}((\mathbf{x}_i^T \boldsymbol{\beta})^{-3/2}))\mathbf{X})^{-1}.$$

- We can estimate the nuisance parameter in the usual way:

$$\hat{\phi} = \frac{1}{N-K} \sum_{i=1}^N \frac{(y_i - \hat{\mu}_i)^2}{V(\hat{\mu}_i)},$$

where—by construction—we chose the parameterization in which $V(\mu) = \mu^3$, such that $\hat{\phi}$ will be consistent for ϕ .

Example 11.7: Simulation (Setup)

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Important library for inverse Gaussian
5 library("statmod")
6
7 ## Set sample size
8 n <- 5000
9
10 ## Set true value of beta
11 beta <- matrix(c(0.5, 0.2), nrow = 2)
```

Example 11.7: Simulation (Data generation/set tolerance)

```
1 ## Generate data
2 X <- matrix(cbind(1, runif(n, 0, 5)), ncol = 2)
3 y <- rinvgauss(n, mean = 1/sqrt((X %*% beta)), dispersion = 1.2)
4
5 ## Initialize beta
6 betaj <- c(1/mean(y)^2, 0)
7
8 ## Set tolerance and count iterations
9 tol <- 1
10 iter <- 1
```

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Simulation (Gauss-Newton iterations)

```
1 while(tol > 1e-12 & iter < 50)
2 {
3   ## Store prior iteration
4   betaj.prior <- betaj
5
6   ## Linear predictor
7   etaj <- c(X %*% betaj)
8
9   ## Estimating function
10  Gn <- -(1/2)*t(X) %*% (y - 1/sqrt(etaj))
11
12  ## An
13  An <- t(X * c((1/4)*(etaj^(-3/2)))) %*% X
14
15  ## Update to next iteration
16  betaj <- betaj + solve(An) %*% Gn
17
18  ## Nuisance
19  num <- (y - 1/(X %*% betaj))^(1/2)
20  denom <- (X %*% betaj)^(-3/2)
21  phi <- (1/(n - 2)) * sum((num^2/denom))
22
23  ## Determine if conditions are met
24  tol <- sum((betaj - betaj.prior)^2)
25  iter <- iter + 1
26 }
```

Example 11.7: Simulation (Report results)

```
1 ## Report number of iterations
2 > print(iter)
3 [1] 6
4
5 ## Report tolerance
6 > print(tol)
7 [1] 2.407682e-17
```

EXAMPLE: INVERSE GAUSSIAN DISTRIBUTION

Example 11.7: Simulation (Report estimates/standard errors)

```
1 ## Report estimates
2 > betaj
3           [,1]
4 [1,] 0.4680579
5 [2,] 0.2093309
6
7 ## Estimate of phi
8 > phi
9 [1] 1.118703
10
11 ## Obtain and report standard errors (re-update)
12 etaj <- c(X %*% betaj)
13 An <- t(X * c((1/4)*(etaj^(-3/2)))) %*% X
14
15 > sqrt(diag(phi*solve(An)))
16 [1] 0.04374994 0.01953502
17
18 ## Compare to R's glm capabilities
19 zz <- glm(y ~ X[,2], start = c(1/mean(y)^2, 0), family = inverse.gaussian(),
20          control = list(epsilon = 1e-12))
21
22 > summary(zz)$coef
23           Estimate Std. Error t value Pr(>|t|)
24 (Intercept)  0.46806    0.04375   10.70 <2e-16 ***
25 X[, 2]       0.20933    0.01954   10.72 <2e-16 ***
26
27 (Dispersion parameter for inverse.gaussian family taken to be 1.118703)
```

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS**
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

Rewriting the Gauss-Newton algorithm:

- We can expand the $(j + 1)$ step of the Gauss-Newton algorithm:

$$\begin{aligned}
 \boldsymbol{\beta}^{(j+1)} &= \boldsymbol{\beta}^{(j)} + [\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{X}]^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \left[\text{diag} \left(\left. \frac{\partial \boldsymbol{\eta}(\boldsymbol{\beta})}{\partial \boldsymbol{\mu}(\boldsymbol{\beta})} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(j)}} \right) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(j)})) \right] \\
 &= [\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{X}]^{-1} \left[\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{X} \boldsymbol{\beta}^{(j)} + \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \left[\text{diag} \left(\left. \frac{\partial \boldsymbol{\eta}(\boldsymbol{\beta})}{\partial \boldsymbol{\mu}(\boldsymbol{\beta})} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(j)}} \right) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(j)})) \right] \right] \\
 &= [\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{X}]^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \left[\mathbf{X} \boldsymbol{\beta}^{(j)} + \left[\text{diag} \left(\left. \frac{\partial \boldsymbol{\eta}(\boldsymbol{\beta})}{\partial \boldsymbol{\mu}(\boldsymbol{\beta})} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(j)}} \right) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(j)})) \right] \right] \\
 &= [\mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{X}]^{-1} \mathbf{X}^T \mathbf{W}(\boldsymbol{\beta}^{(j)}) \mathbf{z}(\boldsymbol{\beta}^{(j)}),
 \end{aligned}$$

where $\mathbf{z}(\boldsymbol{\beta}^{(j)}) = \mathbf{X} \boldsymbol{\beta}^{(j)} + \left[\text{diag} \left(\left. \frac{\partial \boldsymbol{\eta}(\boldsymbol{\beta})}{\partial \boldsymbol{\mu}(\boldsymbol{\beta})} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(j)}} \right) (\mathbf{y} - \boldsymbol{\mu}(\boldsymbol{\beta}^{(j)})) \right]$

- $r_{ij}^W = \left. \frac{\partial \eta_i(\boldsymbol{\beta})}{\partial \mu_i(\boldsymbol{\beta})} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(j)}} (y_i - \mu_i(\boldsymbol{\beta}^{(j)}))$ is known as the *working residual*.

IRLS:

- The Gauss-Newton algorithm and IRLS are equivalent!
- As an informal exercise, try to hard-code one of the previous examples using IRLS and verify that you're getting the same answer.

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)**
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

EXAMPLE: BINOMIAL REGRESSION

Example 11.8: Binomial experiments in trial/success format

- Consider the following data set emerging from a set of binomial experiments within each of seven levels of X (treated continuously):

x	n_x	y_x
0	50	22
1	25	7
2	60	30
3	70	50
4	20	14
5	50	39
6	30	28

- n_x represents number of trials; y_x represents number of successes.
- $N = \sum_x n_x = 305$. Could represent the data as a 305×2 array with y_x converted to binary format and run logistic regression, but it is faster to use the compact data representation.

EXAMPLE: BINOMIAL REGRESSION

Example 11.8: Binomial experiments in trial/success format

- Note: $E[y_x|X = x] = n_x p_x$.
- Note: $\text{Var}[y_x|X = x] = n_x p_x (1 - p_x)$.
- Estimating equations (canonical link):

$$\mathbf{X}^T(\mathbf{y} - \text{diag}(n_X)\text{vec}(\text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}.$$

- Gauss-Newton step:

$$\boldsymbol{\beta}^{(j+1)} = \boldsymbol{\beta}^{(j)} + (\mathbf{X}^T \mathbf{V}(\boldsymbol{\beta}^{(j)}) \mathbf{X})^{-1} \mathbf{X}^T(\mathbf{y} - \text{diag}(n_X)\text{vec}(\text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)}))),$$

where $\mathbf{V}(\boldsymbol{\beta}^{(j)}) = \text{diag}(n_x \text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)})(1 - \text{expit}(\mathbf{x}_i^T \boldsymbol{\beta}^{(j)})))$.

- Using `glm()` in R:

```
1 glm(cbind(y, n-y) ~ x, family = binomial(link = "logit"))
```

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms**
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices

EXAMPLE: POISSON REGRESSION

Example 11.9: Offsets

- Poisson regression is commonly used to analyze count data ($Y \in \mathbb{N}$).
- Sometimes, counts measured over variable ranges of space-time.
 - ▶ Number of asthma exacerbations over a length of time that may vary from patient to patient.
 - ▶ Number of SARS-CoV-2 cases emerging in counties that have different population sizes.
- Including an *offset term* standardizes the region of comparison.
- Model (log link): $\log(E[Y]/W) = \mathbf{x}^T \boldsymbol{\beta} \Rightarrow \log(E[Y]) = \log(W) + \mathbf{x}^T \boldsymbol{\beta}$.
- Estimating equations (convince yourself that this is correct):

$$\mathbf{X}^T(\mathbf{y} - \text{diag}(w_i)\text{vec}(\exp(\mathbf{x}_i^T \boldsymbol{\beta}))) = \mathbf{0}$$

- Using `glm()` in R:

```
1 glm(y ~ x, offset = log(w), family = poisson(link = "log"))
```

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation**
- 11 Quick note: Fixed vs. random design matrices

Linear models: Discrete strata

- Suppose you have N observations, $(X_1, Y_1), \dots, (X_N, Y_N)$.
 - ▶ Consider X binary (0/1).
- Consider fitting the following regression model:

$$E[Y|X = x] = \beta_0 + \beta_1 x.$$

- You may already know:
 - ▶ Each of the two group gets its own mean.
 - ▶ No borrowing of information.
 - ▶ Essentially a two-sample t -test.
- It is possible to build saturated models when the strata defined by \mathbf{x} are discrete.

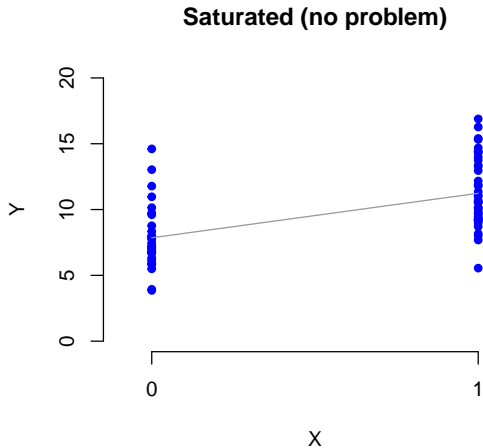
SATURATION AND SUPER-SATURATION

Linear models: Discrete strata

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Generate data
5 n <- 70
6 x <- rbinom(n, 1, 0.5)
7 y <- 7 + 4*x + rnorm(n, 0, 3)
8
9 ## Simple linear regression (t-test, effectively)
10 zz <- lm(y ~ x)
11
12 > summary(zz)
13 Call:
14 lm(formula = y ~ x)
15
16 Residuals:
17     Min       1Q   Median       3Q      Max
18 -5.686 -1.877 -0.666  1.851  6.753
19
20 Coefficients:
21             Estimate Std. Error t value Pr(>|t|)
22 (Intercept)   7.854     0.485   16.2   <2e-16 ***
23 x             3.382     0.626    5.4    9e-07 ***
24 ---
25 Residual standard error: 2.57 on 68 degrees of freedom
26 Multiple R-squared:  0.3, Adjusted R-squared:  0.29
27 F-statistic: 29.2 on 1 and 68 DF, p-value: 9.02e-07
28
29 ## No concerns!
```

SATURATION AND SUPER-SATURATION

Linear models: Discrete strata



Linear models: Continuous predictor

- Now, suppose you have N observations, $(X_1, Y_1), \dots, (X_N, Y_N)$, with the values of X continuous.
- Consider fitting the following regression model:

$$E[Y|X = x] = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_{N-1} x^{N-1}.$$

- You may already know that this will fail beautifully:
 - ▶ Every observation will be predicted perfectly.
 - ▶ $\hat{\sigma}^2 = 0$.
- No inference on β possible.

Linear models: Polynomial interpolation

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Generate data
5 n <- 7
6 x <- rnorm(n, 5, 3)
7 y <- 6 + 0.5*x + rnorm(n, 0, 7)
```

SATURATION AND SUPER-SATURATION

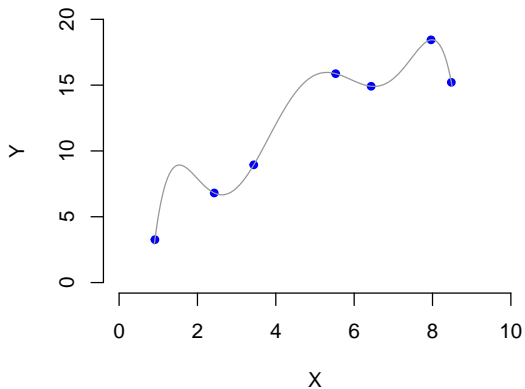
Linear models: Polynomial interpolation

```
1 ## Polynomial regression
2 zz <- lm(y ~ x + I(x^2) + I(x^3) + I(x^4) + I(x^5) + I(x^6))
3
4 > summary(zz)
5
6 Call:
7 lm(formula = y ~ x + I(x^2) + I(x^3) + I(x^4) + I(x^5) + I(x^6))
8
9 Residuals:
10 ALL 7 residuals are 0: no residual degrees of freedom!
11
12 Coefficients:
13             Estimate Std. Error t value Pr(>|t|)
14 (Intercept)  -60.5965         NaN     NaN     NaN
15 x             138.5121         NaN     NaN     NaN
16 I(x^2)       -103.8012         NaN     NaN     NaN
17 I(x^3)         37.0178         NaN     NaN     NaN
18 I(x^4)        -6.7027         NaN     NaN     NaN
19 I(x^5)         0.5969         NaN     NaN     NaN
20 I(x^6)        -0.0208         NaN     NaN     NaN
21
22 Residual standard error: NaN on 0 degrees of freedom
23 Multiple R-squared: 1, Adjusted R-squared: NaN
24 F-statistic: NaN on 6 and 0 DF, p-value: NA
25
26 ## Sad trombone noises
```

SATURATION AND SUPER-SATURATION

Linear models: Polynomial interpolation

Super-saturated (not so good)



SATURATION AND SUPER-SATURATION

Tabular data:

- Suppose I have the following data, represented by a 2×2 table:

	$X_1 = 0$	$X_1 = 1$
$X_2 = 0$	a	b
$X_2 = 1$	c	d

- We could convert the data to a binary-outcome representation and fit, say, a simple logistic regression model.
 - ▶ Sample size: $N = a + b + c + d$.
 - ▶ Such a model, having $K = 2$ degrees of freedom, will perfectly predict the relative proportions as seen in the table (no concerns).
- I'd like to challenge you to think about what might happen if we were to conceptualize the cell values as originating from Poisson counts.
- Sample size: $N = 4$?
- Such a model, having $K = 4$ degrees of freedom, will perfectly predict the exact *counts* seen in the table. Problem?

SATURATION AND SUPER-SATURATION

Tabular data:

```
1 ## Set seed for reproducibility
2 set.seed(7345)
3
4 ## Two-by-two tabular data
5 n <- 4
6
7 ## Columns
8 x1 <- c(0,0, 1,1)
9
10 ## Rows
11 x2 <- c(0,1, 0,1)
12
13 ## Counts
14 y <- rpois(n, lambda = exp(1 + x1 + x2 + x1*x2))
15
16 tab <- data.frame(matrix(y, nrow = 2))
17 names(tab) <- c("x1=0", "x1=1")
18 row.names(tab) <- c("x2=0", "x2=1")
19
20 > tab
21      x1=0 x1=1
22 x2=0    4   11
23 x2=1    2   48
24
25 > cbind(x1, x2, y)
26      x1 x2 y
27 [1,]  0  0  4
28 [2,]  0  1  2
29 [3,]  1  0  11
30 [4,]  1  1  48
```

SATURATION AND SUPER-SATURATION

Tabular data:

```
1 ## Model
2 zz <- glm(y ~ x1 + x2 + I(x1*x2), family = poisson(link = "log"))
3
4 > ## Should recover [1,1], [1,2], [2,1], and [2,2] entries perfectly
5 > as.numeric(exp(coef(zz)[1]))
6 [1] 4
7
8 > as.numeric(exp(coef(zz)[1] + coef(zz)[2]))
9 [1] 11
10
11 > as.numeric(exp(coef(zz)[1] + coef(zz)[3]))
12 [1] 2
13
14 > as.numeric(exp(coef(zz)[1] + coef(zz)[2] + coef(zz)[3] + coef(zz)[4])
15 [1] 48
16
17 ## Looks like a ``super-saturated`` model.
```

SATURATION AND SUPER-SATURATION

Tabular data:

```
1 > summary(zz)
2
3 Call:
4 glm(formula = y ~ x1 + x2 + I(x1 * x2), family = poisson(link = "log"))
5
6 Coefficients:
7             Estimate Std. Error z value Pr(>|z|)
8 (Intercept)   1.386      0.500   2.77  0.0056 **
9 x1            1.012      0.584   1.73  0.0832 .
10 x2           -0.693      0.866  -0.80  0.4235
11 I(x1 * x2)    2.166      0.928   2.33  0.0196 *
12 ---
13 (Dispersion parameter for poisson family taken to be 1)
14
15 Null deviance: 7.5800e+01 on 3 degrees of freedom
16 Residual deviance: 4.6629e-15 on 0 degrees of freedom
17 AIC: 23.84
18
19 Number of Fisher Scoring iterations: 3
20
21 ## No problem...surprisingly?
```

Tabular data:

- Super-saturation does not appear to pose a problem in this case.
 - ▶ There is no nuisance parameter.
 - ▶ Mean-variance relationship: $V(\mu) = \mu; \phi = 1$.
- Poisson standard errors rely on mean-variance relationship.
 - ▶ Analogue to OLS: if assuming $\sigma^2 = 1$ (bad), super-saturated linear model would have no problem.
- Recall that information about the nuisance comes from the *residuals*. If the residuals are all zero, so too will be the estimated nuisance.
- On the flip side: residuals are also where you assess assumptions. . .
 - ▶ Cannot use the data to assess the validity of the Poisson model in a “super-saturated” analysis of tabular data.

TABLE OF CONTENTS

- 1 A reason to delve into the theory
- 2 Generalized linear models (GLMs)
- 3 The canonical link function
- 4 The Gauss-Newton Algorithm
- 5 Variance
- 6 Examples: Hard-coding GLMs
- 7 Quick note: The Gauss-Newton algorithm and IRLS
- 8 Quick note: Binomial regression (compact representation)
- 9 Quick note: Handling offset terms
- 10 Quick note: Saturation and super-saturation
- 11 Quick note: Fixed vs. random design matrices**

Considerations:

- In this set of notes, we are treating GLMs as a likelihood-based approach, and we are assuming that the models put forth are correctly specified.
- It is mathematically convenient to think of \mathbf{X} as fixed and known in advance (that's what we have done this entire slide set).
- However, when the models are correctly specified, this makes no meaningful difference as key quantities described in these notes (such as $\mathbf{D}^T \mathbf{V}^{-1} \mathbf{D}$) are consistent for their their respective expectations over the marginal distribution of \mathbf{X} .
- We will cover various forms of GLM misspecification, at which point we will actually have to pay attention to whether \mathbf{X} is fixed or random and use the correct approach.

This unit:

- Generalized linear models.
- The canonical link.
- The Newton-Raphson and Gauss-Newton algorithms.
- Asymptotics.
- Offsets.
- Saturation.

SUMMARY: SO FAR

- Random vectors and matrices; multivariate normal theory.
- Ordinary least squares.
- Hypothesis testing and ANOVA.
- Weighted least squares.
- Misspecification.
- Confidence regions and prediction.
- Diagnostics.
- Regularization.
- Bayesian regression.
- Exponential families.
- Generalized linear models.

SUMMARY: COMING UP

- Sandwich and bootstrap.
- Quasi-likelihood.
- Hypothesis testing for GLMs.
- Diagnostics for GLMs.
- Further considerations for binary outcomes.
- Nonlinear least squares.