

Lecture 15

Motivating example: air pollution data

- Data on Chicago air quality and death between 1987 and 2000
- Variables include:
 - deaths
 - ozone concentration
 - sulphur dioxide concentration
 - temperature

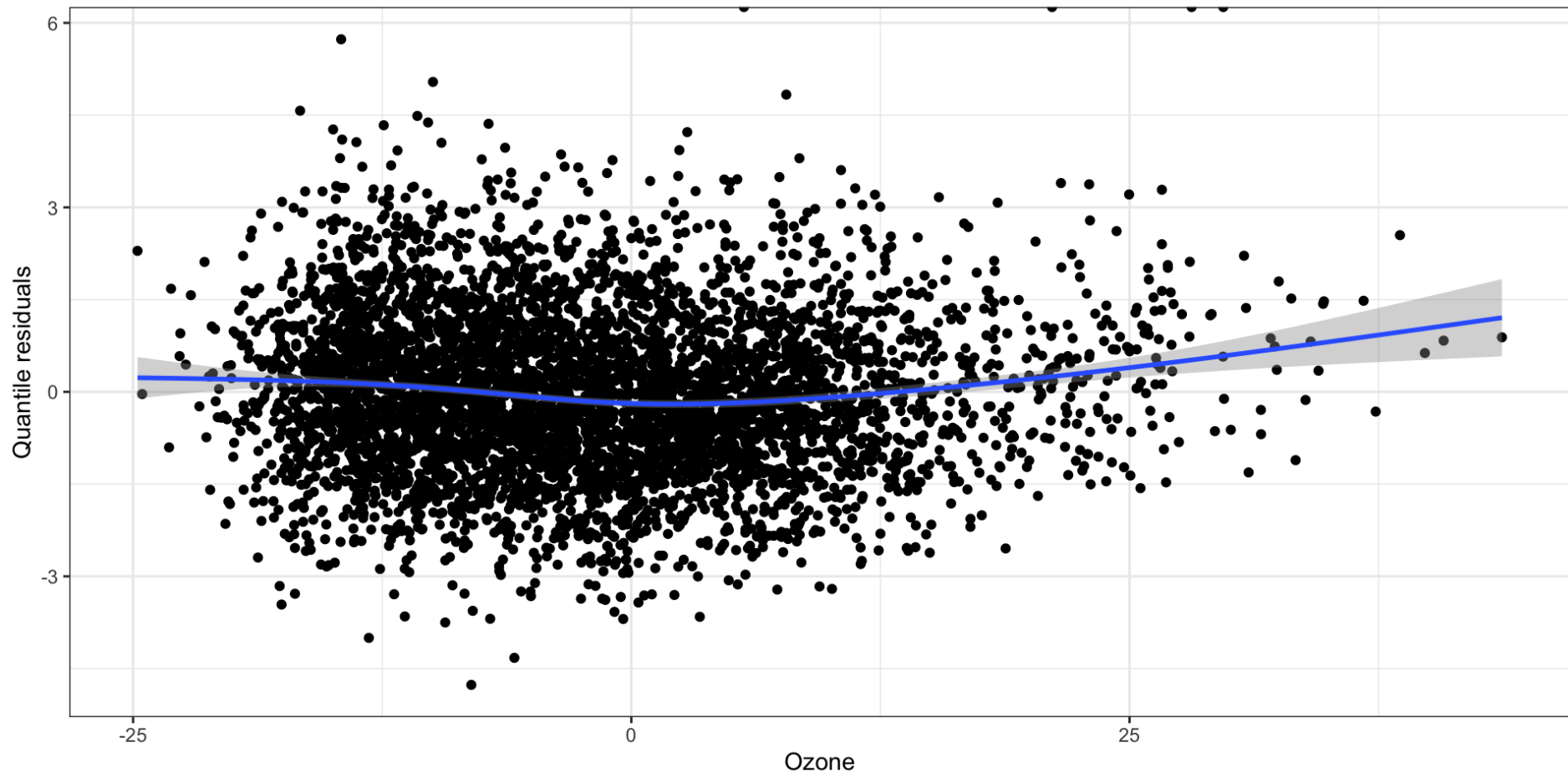
Motivating example: air pollution data

$$\text{Deaths}_i \sim \text{Poisson}(\lambda_i)$$

$$\log(\lambda_i) = \beta_0 + \beta_1 \text{Ozone}_i$$

Quantile residual plot

- variance of quantile residuals is
 \approx constant, but too high



- nice random scatter around 0 (maybe slight pattern for high values Ozone)
- maybe some more variability than we want

GOF test

```
1 m1$deviance
```

```
[1] 9551.836
```

```
1 m1$df.residual
```

```
[1] 5112
```

```
1 pchisq(m1$deviance, m1$df.residual, lower.tail=F)
```

```
[1] 6.362106e-273
```

=> maybe Poisson distribution is not appropriate
for Y_i

Overdispersion

Overdispersion occurs when the response variable Y_i has greater variability than the model accounts for

If $Y_i \sim \text{Poisson}(\lambda_i)$ then $\text{Var}(Y_i) = \lambda_i$

But if $Y_i \not\sim \text{Poisson}$, and we incorrectly assume it is Poisson, then $\text{Var}(Y_i)$ may be larger than the model accounts for

Recap: sandwich estimator for GLMs

$$\hat{\beta} \text{ solves } u(\beta) = \frac{x^T (y - \mu)}{n} = \frac{1}{n} \sum_{i=1}^n (y_i - \mu_i) x_i = 0$$

$$\beta^* \text{ solves } E\left[\frac{1}{n} (y_i - \mu_i) x_i\right] = 0$$

$$\sqrt{n} (\hat{\beta} - \beta^*) \xrightarrow{d} N(0, \underbrace{J_n(\beta^*)^{-1}}_{\text{bread}} \underbrace{V_n(\beta^*)}_{\text{meat}} \underbrace{J_n(\beta^*)^{-1}}_{\text{bread}})$$

$$\Rightarrow \hat{\beta} \approx N(\beta^*, J_n(\beta^*)^{-1} V_n(\beta^*) J_n(\beta^*)^{-1})$$

$$J_n(\beta^*) = -E[U'(\beta^*)] \quad V_n(\beta^*) = \text{Var}(U(\beta^*))$$

• If model is correct (both shape and distribution) then

$$J_n(\beta^*) = V_n(\beta^*)$$

• If model is incorrect, $\hat{\beta}$ still has this asymptotic variance

$$\hat{\beta} \xrightarrow{P} \beta^*$$

β^* = the coefficients for relationship if
mean $\mu_i = g(\beta^T x_i)$ is correct

Assumptions about both mean and variance

$$\frac{\partial \mu}{\partial \theta} = v(\mu)$$

Suppose we assume

$$\mathbb{E}[\gamma_i] = \mu_i = g^{-1}(\beta^T X_i)$$

$$\text{var}(\gamma_i) = \phi v(\mu_i)$$

} true if $\gamma_i \sim \text{EDM}(\mu_i, \phi)$

But we don't have to assume a distribution for γ_i ,
just the first 2 moments

$$\hat{\beta} \text{ solves } u(\beta) = \frac{X^T(\gamma - \mu)}{\phi} = \frac{1}{\phi} \sum_{i=1}^n (\gamma_i - \mu_i) X_i = 0$$

$$\text{var}(u(\beta)) = \sum_{i=1}^n \frac{1}{\phi^2} \cdot \text{var}(\gamma_i) X_i X_i^T$$

$$\parallel = \sum_{i=1}^n \frac{1}{\phi} \cdot v(\mu_i) X_i X_i^T$$

$$- \mathbb{E}[u'(\beta)] = - \frac{1}{\phi} \sum_{i=1}^n \mathbb{E}\left[\frac{\partial}{\partial \beta} (\gamma_i - \mu_i) X_i\right] = \frac{1}{\phi} \sum_{i=1}^n v(\mu_i) X_i X_i^T$$

If we correctly assume

$$E(Y_i) = \mu_i = g^{-1}(\beta^T X_i)$$

$$\text{var}(Y_i) = \phi v(\mu_i)$$

then

$$\hat{\beta} \approx N(\beta^*, \phi (X^T W X)^{-1})$$

\uparrow
 $w = \text{diag}(v(\mu_i))$

E.g. for

$$\log(\mu_i) = \beta^T X_i \quad (\text{like in Poisson})$$

$$\text{var}(Y_i) = \phi \mu_i \quad (\text{like in Poisson})$$

then

$$\hat{\beta} \approx N(\beta^*, \phi (X^T W X)^{-1})$$

$\swarrow \quad \uparrow$
 $\text{diag}(\mu_i)$

$$\phi \stackrel{||}{=} \text{var}(\hat{\beta}) \text{ Poisson model}$$

Quasi-Poisson models

$$\log(\mu_i) = \beta^T X_i \quad \leftarrow \text{same as Poisson regression}$$

$$\text{Var}(Y_i) = \phi \mu_i \quad \phi \text{ is allowed to be } > 1$$

ϕ Var Poisson

$$\Rightarrow \hat{\beta}_{QP} \equiv \hat{\beta}_{\text{Poisson}}$$

$$\hat{\text{Var}}(\hat{\beta}_{QP}) = \hat{\phi} \hat{\text{Var}}(\hat{\beta}_{\text{Poisson}})$$

↑
need to estimate ϕ ...

Example: Chicago air quality

Poisson model:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	4.743277988	0.0013382057	3544.50583	0.000000e+00
o3median	-0.002301345	0.0001285909	-17.89664	1.252641e-71

Quasi-Poisson model:

Call:

glm(formula = death ~ o3median, family = quasipoisson, data = chicago)

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	4.7432780	0.0018822	2520.02	<2e-16 ***
o3median	-0.0023013	0.0001809	-12.72	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for quasipoisson family taken to be 1.978347)

Null deviance: 9873.8 on 5113 degrees of freedom

$$0.0001809 = 0.00012859 \sqrt{1.978}$$

Class activity

https://sta712-f23.github.io/class_activities/ca_lecture_15.html

