# Generalised Linear Models (GLMs) II

# Deviance

# Deviance

$\log L$

$\leftarrow \log L(S)$

$\updownarrow$

$\leftarrow \log L(\hat{\beta}^A)$

From here on we will w.l.o.g. adopt the convention that $a(\phi) = \phi$.

**Definition:** the scaled deviance for model A is

$$\frac{D^A}{\phi} = -2\log\frac{\mathcal{L}(\hat{\beta}^A)}{\mathcal{L}(S)}$$

where $\hat{\beta}^A$ is the MLE of $\beta^A$, the parameters in the model A, and $\mathcal{L}(S)$ is the maximum likelihood for the saturated model

The deviance is just $D^A$.

$\phi \cdot$ scaled deviance.

## Example: normal

The saturated normal model uses $y_i$ to estimate $\mu_i$.

The deviance can be written as $D = \sum_i d_i$ where

$$d_i = (y_i - \hat{\mu}_i)^2$$

where $\hat{\mu}_i$ is the fitted mean using the MLE.

It is equivalent to SSE (sum of squared errors) in linear models.

Handwritten annotations (top):

model A $\quad \hat{\beta}^A \Rightarrow \hat{\theta}_i^A \quad \phi$ assume known

for all GLM $\quad \mu_i = b'(\theta_i) = g^{-1}(x_i^T \beta)$. model saturated s.
$\quad \hat{\beta}^s = \hat{\theta}_i^s \quad \phi$ assume known

$$\frac{D^A}{\phi} = 2\left[\log^{L(s)} - \log^{L(A)}\right]$$

$$= 2\left[\sum_{i=1}^{n}\left[\frac{y_i \hat{\theta}_i^s - b(\hat{\theta}_i^s)}{\phi} + c(y_i, \phi)\right]\right.$$

$$\left.- \sum_{i=1}^{n}\left[\frac{y_i \hat{\theta}_i^A - b(\hat{\theta}_i^A)}{\phi} + c(y_i, \phi)\right]\right]$$

$$= \frac{2}{\phi}\sum_{i=1}^{n}\left[y_i(\hat{\theta}_i^s - \hat{\theta}_i^A) - \left(b(\hat{\theta}_i^s) - b(\hat{\theta}_i^A)\right)\right]$$

Handwritten annotations (bottom left):

or normal $\quad \theta_i = \mu_i \quad b(\theta_i) = \frac{1}{2}\theta_i^2$

$\hat{\theta}_i^A = \hat{\mu}_i = x^T\beta$

or saturated $\quad \hat{\theta}_i^s = y_i$

$\frac{2}{\phi}\sum_{i=1}^{n}[y_i(y_i - \hat{\mu}_i) - \frac{1}{2}(y_i^2 - \hat{\mu}_i^2)] = \frac{1}{\phi}\sum_{i=1}^{n}(y_i - \hat{\mu}_i)^2 \rightarrow$ scaled deviance

$D^A = \phi \cdot \frac{1}{\phi}\sum_{i=1}^{n}(y_i - \hat{\mu}_i)^2 = \sum_{i=1}^{n}(y_i - \hat{\mu}_i)^2$

# Example: Poisson

*exercise*

The saturated Poisson model uses $y_i$ to estimate $\mu_i = \lambda_i$.

The deviance can be written as $D = \sum_i d_i$ where

$$d_i = -2 \left( y_i \log \frac{\hat{\mu}_i}{y_i} - (\hat{\mu}_i - y_i) \right)$$

where $\hat{\mu}_i$ is the fitted mean using the MLE.

# Example: binomial

*exercise*

The saturated binomial model uses $y_i$ to estimate $\mu_i$, that is $y_i/m_i$ to estimate $p_i$.

The deviance can be written as $D = \sum_i d_i$ where

$$d_i = -2 \left( y_i \log \frac{\hat{\mu}_i}{y_i} + (m_i - y_i) \log \frac{m_i - \hat{\mu}_i}{m_i - y_i} \right)$$

where $\hat{\mu}_i$ is the fitted mean using the MLE.

# Scaled deviance for testing model adequacy

The binomial and Poisson models:

- $\phi = 1$ and the scaled deviance is just the deviance.
- When the individual responses are somewhat normal, if the model is adequate then the scaled deviance will be $\approx \chi^2_{n-p}$, where the saturated model has $n$ parameters (equal to the number of observations) and the fitted model has $p$ parameters.
- As a rule of thumb we need $m_i p_i$ and $m_i(1 - p_i) \geq 5$ for the binomial model and $\lambda_i \geq 5$ for the Poisson model.

The normal, gamma or inverse gaussian models:

- we have to estimate $\phi$ using $X^2/(n - p)$.
- we can't use the scaled deviance to test model adequacy.

*Handwritten annotations:*

if scaled deviance
> (1-a)th of $\chi^2_{n-p}$

fitted
model is not adequate

loss a lot of information

# Scaled deviance for model selection (nested models)

The difference between two scaled deviances can be used for model selection between two nested models (i.e., testing the significance of the extra parameters).

The binomial and Poisson models:

- For nested models, if the smaller model is correct then the difference between two scaled deviances is the log likelihood ratio and will be $\approx \chi^2_s$ where $s$ is the difference in the number of parameters.
- Example : See "LR test using deviance" in Bliss.pdf

The normal, gamma or inverse gaussian models:

- we have to estimate $\phi$ using $X^2/(n - p)$.
- see the next slide..

*if the difference is bigger than what we expect, then the additional parameter in the full model is useful, pick the full model*

For the normal model (linear model), if model A is nested within model B, under the null hypothesis that model A is correct we have

$$\frac{(D^A - D^B)/s}{X^2/(n-p)} \sim F_{s,n-p}$$

where $D^A$ and $D^B$ denote deviances for models A and B, respectively, we have $n$ observations, model A has $p - s$ parameters, and model B has $p$ parameters.

- $X^2$ (Pearson's chi-squared) is calculated using model B.
- It is equivalent to the F-test in linear model.

For other models, this distributional result only holds approximately, but it can still be used for comparing two nested models. You will use it for comparing gamma models in the lab problem for the week 4.

Handwritten annotations (left margin):

$(p-s \text{ parameters})$
model A is nested in model B.
$(p)$

$LR = \dfrac{D^A - D^B}{\phi}$

estimate $\phi$ with $\dfrac{X^2}{n-p}$

$\dfrac{D^A - D^B / s}{\dfrac{X^2}{n-p}} \sim F_{s,n-p}$

$\Rightarrow \dfrac{\dfrac{D^A - D^B}{\phi}/s}{\dfrac{X^2}{\phi}/n-p}$

$\phi$ is known

$\dfrac{D^A - D^B}{\phi} \sim \chi^2_s$    $\dfrac{X^2}{\phi} \sim \chi_{n-p}$

Handwritten annotation (below equation): estimate $\phi$.

# Scaled deviance for model selection (AIC)

The AIC is defined as

$$\text{AIC} = 2p - 2\log \mathcal{L}(\hat{\boldsymbol{\beta}})$$

where $\hat{\boldsymbol{\beta}}$ represent MLE of parameters in the model and $p$ is the number of the parameters. Given a choice, we prefer the model with the smaller AIC.

If model B has $s$ more parameters than model A (not necessarily nested within B), then

$$
\begin{aligned}
\text{AIC}^B - \text{AIC}^A &= 2s - 2\log \mathcal{L}(\hat{\boldsymbol{\beta}}^B) + 2\log \mathcal{L}(\hat{\boldsymbol{\beta}}^A) \\
&= 2s + \frac{D^B}{\phi} - \frac{D^A}{\phi}.
\end{aligned}
$$

Like the log likelihood ratio, the AIC needs an estimate of $\phi$.

# Diagnostics

# Residuals

**Response residuals:**

$$y_i - \hat{\mu}_i$$

Unless $v(\mu)$ is constant, as in the normal case, the response residuals are not homoskedastic and hence not very useful.

*same variance across samples*

```
> residuals(glmfit, type="response")
```

**Pearson residuals:**

$$r_P(i) = \frac{y_i - \hat{\mu}_i}{\sqrt{v(\hat{\mu}_i)}}$$

Pearson residuals are (approximately) homoskedastic, and $\sum_i r_P(i)^2 = X^2$.

```
> residuals(glmfit, type="pearson")
```

*Handwritten annotations:*

$r_i = y_i - \hat{\mu}_i$

$var(r_i) = var(y_i) = v(\mu_i) a(\phi)$.

for linear model $var(r_i)$ is constant across samples

$v(\mu_i) = 1$

$var\left(\frac{y_i - \mu_i}{\sqrt{v(\mu_i)}}\right) = a(\phi) \rightarrow$ should be homoskedastic

since $\mu_i$ unknown, use estimate $\hat{\mu}_i$

defined to estimate $\phi$

# Residuals

**Deviance residuals:**

$$r_D(i) = \text{sign}(y_i - \hat{\mu}_i)\sqrt{d_i}$$

where the deviance is $D = \sum_i d_i = \sum_i r_D(i)^2$.

```
> residuals(glmfit)
```

- As for linear models, patterns in the residuals indicate structure in the data that has not been captured by the model.

- We can plot the residuals against predictor variables, the responses, or the fitted means. Often a plot against $\hat{\eta}_i = x_i^T\hat{\beta}$ works well.

- With count data residual plots exhibit banding due to the discrete nature of the responses, and this can make it hard to see other patterns. In this case we can use a smoothed fit of the residuals to help spot trends/patterns.
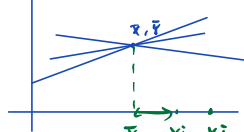
*banding pattern dominate* ?

## Leverage

The leverage measures the potential influence of a point on the fitted model. We borrow the definition of leverage from the theory of linear models, and use the hat matrix from the IWLS algorithm.

In the IWLS algorithm

$$\hat{Z} = X\hat{\beta} = X(X^T\hat{\Sigma}^{-1}X)^{-1}X^T\hat{\Sigma}^{-1}Z,$$

where $\hat{\Sigma}_{ii} = (g'(\hat{\mu}_i))^2 v(\hat{\mu}_i)\phi$ (assuming $a(\phi) = \phi$). Then, the hat matrix for $\hat{Z}$ is

$$H = X(X^T\hat{\Sigma}^{-1}X)^{-1}X^T\hat{\Sigma}^{-1}.$$

The leverage for the $i$-th observation is $H_{ii}$, the $i$-th diagonal element of $H$.

Note that $H$ does not depend on $\phi$ (it cancels out).

A large leverage does not necessarily mean a point *has* influenced the fit.

> influence(glmfit)$hat

---

*Handwritten annotations:*

linear model

$\hat{y} = X\hat{\beta} = X(X^TX^{-1})X^Ty$

$= Hy$ ← $n \times n$

$h_{ii}$: distance between $\bar{x}$ and $x_i$

all fitted lines should go through $(\bar{x}, \bar{r})$

if $h_{ii} > h_{jj}$

the point $(x_i)$ far away from $\bar{x}$ has a potentially larger influence than the point close to $\bar{x}$ $(x_j)$

measures the potential influence → doesn't measure the direct influence

The variance of $\hat{Z}$ is

$$\text{var}(\hat{z}) = H \text{ var}(z) H^T$$

$$H\text{Var }ZH^T$$
$$= [X(X^T\Sigma^{-1}X)^{-1}X^T\Sigma^{-1}]\Sigma[\Sigma^{-1}X(X^T\Sigma^{-1}X)^{-1}X^T]$$
$$= X(X^T\Sigma^{-1}X)^{-1}X^T = H\Sigma.$$

Whence $\text{Var}(Z - \hat{Z}) = (I - H)\Sigma$.

$H\Sigma$ is symmetric

$(H\Sigma)^T = H\Sigma$

also $\Sigma$ is diagonal

$(H\Sigma)^T = \Sigma^T H^T = \Sigma H^T$

$$\text{var}(z - \hat{z}) = \text{var}(z - Hz) = \text{var}\left((I-H)z\right) = (I-H)\text{var}(z)(I-H)^T$$

$$= \Sigma - H\Sigma - \Sigma H^T + H\Sigma H^T$$

$$= \Sigma - H\Sigma$$

$$\xrightarrow{\quad} H\Sigma$$

$$= (I-H)\Sigma$$

$$= (I-H)\Sigma(I-H)^T$$

$$= (\Sigma - H\Sigma)(I-H^T)$$

# Studentised residuals

From the above we have $\mathrm{Var} \dfrac{Z_i - \hat{Z}_i}{\sqrt{(1-H_{ii})g'(\mu_i)^2 v(\mu_i)\phi}} = 1$. Also

$$
\begin{aligned}
Z_i - \hat{Z}_i &= g(\mu_i) + (Y_i - \mu_i)g'(\mu_i) - [g(\mu_i) + (\hat{Y}_i - \mu_i)g'(\mu_i)] \\
&= (Y_i - \hat{Y}_i)g'(\mu_i)
\end{aligned}
$$

*will be cancelled out.*

so, noting that $\hat{Y}_i = \hat{\mu}_i$,

$$
\mathrm{Var} \frac{Y_i - \hat{\mu}_i}{\sqrt{(1 - H_{ii})v(\hat{\mu}_i)\hat{\phi}}} \approx 1
$$

*Pearson residual*

The LHS is just $r_P(i)/\sqrt{(1 - H_{ii})\hat{\phi}} =: r_{SP}(i)$, which we call the *i*-th **studentised Pearson residual**.

# Studentised residuals

**Studentised Pearson residual:**

$$r_{SP}(i) = \frac{r_P(i)}{\sqrt{(1 - H_{ii})\hat{\phi}}}$$

By analogy we define **Studentised deviance residual:**

$$r_{SD}(i) = \frac{r_D(i)}{\sqrt{(1 - H_{ii})\hat{\phi}}}$$

# Jack-knife residuals

"direct influence"

A direct measure of the influence of a point is the **jack-knife residual**, which is the change in $\hat{\mu}_i$ when you remove $y_i$ from the set of observations, then scaled to standardise the variance.
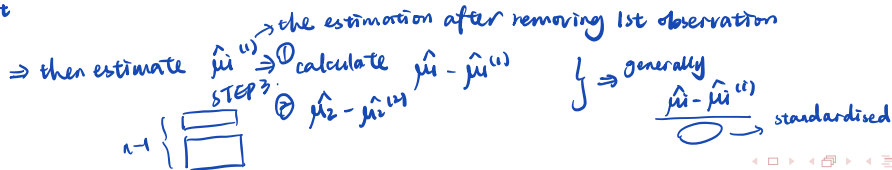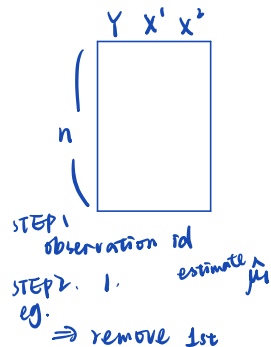
The jack-knife residual can be approximated by

$$\text{sign}(y_i - \hat{\mu}_i)\sqrt{(1 - H_{ii})r_{SD}^2(i) + H_{ii}r_{SP}^2(i)}.$$

```
> rstudent(glmfit)
```

weighted sum

Y X' X²

n

STEP 1
observation id

STEP 2. 1.
eg.
⇒ remove 1st

Y X' X²

n-1

⇒ then estimate $\hat{\mu}_i^{(i)}$ ⇒ ① calculate $\hat{\mu}_i - \hat{\mu}_i^{(i)}$

the estimation after removing 1st observation

STEP 3. ② $\hat{\mu}_2 - \hat{\mu}_2^{(i)}$

n-1

estimate $\hat{\mu}_i$

⇒ generally

$\dfrac{\hat{\mu}_i - \hat{\mu}_i^{(i)}}{\bigcirc}$ → standardised

influence of ith observation on "ith" fitted mean

Y x1 x2  observation
1   $\hat{\mu}_1$
2   $\hat{\mu}_2$
... ...
n   $\hat{\mu}_n$

⇓

Y x1 x2
1   $\hat{\mu}_1^{(i)}$
2   $\hat{\mu}_2^{(i)}$
... ...
n   $\hat{\mu}_n^{(i)}$

⇓

look at

$$\frac{\sum_{j=1}^n (\hat{\mu}_j - \hat{\mu}_j^{(i)})^2}{\quad}$$ → standardize.

the difference between $\hat{\mu}_j$
and $\hat{\mu}_j^{(i)}$ (after removing ith instance).

# Cook's distance

" influence of ith observation on ALL fitted mean "

Another measure of the influence of the $i$-th observation is **Cook's distance**:

$$\frac{(\hat{\beta}^{(i)} - \hat{\beta})^T X^T \hat{\Sigma}^{-1} X (\hat{\beta}^{(i)} - \hat{\beta})}{p}$$

where $\hat{\beta}^{(i)}$ is the estimate of $\beta$ obtained when $y_i$ is omitted.

```
> cooks.distance(glmfit)
```

# Plotting leverage/jack-knife residuals/cook's distance

- When plotting it is helpful to consider them ordered by absolute size.

- If we plot the ordered absolute values against the percentage points of a half-normal distribution, then it is easier to see if the largest values are in keeping with the others. That is, plot the $i$-th ordered absolute residual against $\Phi^{-1}((n+i)/(2n+1))$, for $i = 1, \ldots, n$.

- If all is well we expect to see a smooth plot, while a jump or kink in the tail indicates a potential problem.

- Note that leverage/jack-knife residuals/cook's distance will not be normal, so don't expect a straight line.

# Checking linearity

- A non-linear link $g$ (anything except the identity) makes it harder to check the assumption that $\underbrace{g(\mu_i) = x_i^T \beta.}$   $g(\hat{\mu}) \to g(y_i).$

- The easiest thing to do is to plot $g(y_i)$ against $\underbrace{\{x_{ij}\}_{i=1}^n}$ for each $j$ and look for linear relationships.

*iterated weighted least square estimator.*

- A more sophisticated approach is to plot $\boxed{z_i = g(\mu_i) + (y_i - \mu_i)g'(\mu_i)}$ against $\{x_{ij}\}_{i=1}^n$, where we use $\hat{\mu}_i$ for $\mu_i$. From the IWLS algorithm we know that these plots should be linear.

- If there are non-linearities present then we can consider transforming $\{x_{ij}\}_{i=1}^n$. Transforming the responses $y_i$ is often not a good idea for a glm, as this can break assumptions made about the distribution of $Y_i$.

# Example

See Gala.pdf.

# Learning goals

- Be able to define Generalised Linear Model (GLM).
- Be able to obtain the cannonical link function for GLM.
- (Challenging) Be able to explain ideas for the iterated weighted least squares (IWLS) algorithm.
- Be able to compute the variance of parameter estimates in GLM using IWLS algorithm.
- Understand (scaled) deviance
  - Be able to define (scaled) deviance.
  - Be able to compute (scaled) deviance.
  - Be able to use it to test model adequacy, perform model selection for nested models and non-nested models.
- Be able to perform diagnostics for GLM using R.
  - Understand what different 'diagnostics measurements' aim to measure.
  - Be able to obtain those diagnostic measurements from the *glm* output.
  - Be able to perform diagnostics for GLM using R.