

# QCB 508 – Week 10

*John D. Storey*

*Spring 2017*

## Contents

	<b>2</b>
<b>OLS Goodness of Fit</b>	<b>2</b>
Pythagorean Theorem . . . . .	2
OLS Normal Model . . . . .	3
Projection Matrices . . . . .	3
Decomposition . . . . .	3
Distribution of Projection . . . . .	4
Distribution of Residuals . . . . .	4
Degrees of Freedom . . . . .	5
Submodels . . . . .	5
Hypothesis Testing . . . . .	5
Generalized LRT . . . . .	5
Nested Projections . . . . .	6
$F$ Statistic . . . . .	6
$F$ Distribution . . . . .	7
$F$ Test . . . . .	7
Example: Davis Data . . . . .	7
Comparing Linear Models in R . . . . .	8
ANOVA (Version 2) . . . . .	8
Comparing Two Models with <code>anova()</code> . . . . .	8
When There's a Single Variable Difference . . . . .	9
Calculating the F-statistic . . . . .	9
Calculating the Generalized LRT . . . . .	10
ANOVA on More Distant Models . . . . .	10
Compare Multiple Models at Once . . . . .	11
<b>Generalized Linear Models</b>	<b>11</b>
Definition . . . . .	11
Exponential Family Distributions . . . . .	11
Natural Single Parameter EFD . . . . .	12
Dispersion EFDs . . . . .	12
Example: Normal . . . . .	12
EFD for GLMs . . . . .	13
Components of a GLM . . . . .	13
Link Functions . . . . .	13
Calculating MLEs . . . . .	14

Iteratively Reweighted Least Squares . . . . .	14
Estimating Dispersion . . . . .	15
CLT Applied to the MLE . . . . .	15
Approximately Pivotal Statistics . . . . .	15
Deviance . . . . .	16
Generalized LRT . . . . .	16
Example: Grad School Admissions . . . . .	17
glm() Function . . . . .	21
<b>Nonparametric Regression</b>	<b>22</b>
<b>Generalized Additive Models</b>	<b>22</b>
<b>Bootstrap for Statistical Models</b>	<b>22</b>
<b>Extras</b>	<b>22</b>
Source . . . . .	22
Session Information . . . . .	22

## OLS Goodness of Fit

### Pythagorean Theorem

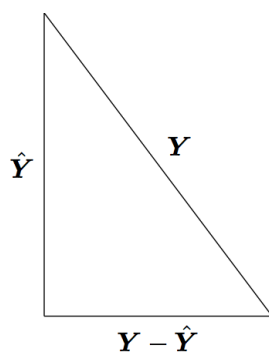


Figure 1: PythMod

Least squares model fitting can be understood through the Pythagorean theorem:  $a^2 + b^2 = c^2$ . However, here we have:

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

where the  $\hat{Y}_i$  are the result of a **linear projection** of the  $Y_i$ .

## OLS Normal Model

In this section, let's assume that  $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$  are distribution so that

$$\begin{aligned} Y_i &= \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + E_i \\ &= \mathbf{X}_i \boldsymbol{\beta} + E_i \end{aligned}$$

where  $\mathbf{E}|\mathbf{X} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I})$ . Note that we haven't specified the distribution of the  $\mathbf{X}_i$  rv's.

## Projection Matrices

In the OLS framework we have:

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

The matrix  $\mathbf{P}_{n \times n} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is a projection matrix. The vector  $\mathbf{Y}$  is projected into the space spanned by the column space of  $\mathbf{X}$ .

Project matrices have the following properties:

- $\mathbf{P}$  is symmetric
- $\mathbf{P}$  is idempotent so that  $\mathbf{P}\mathbf{P} = \mathbf{P}$
- If  $\mathbf{X}$  has column rank  $p$ , then  $\mathbf{P}$  has rank  $p$
- The eigenvalues of  $\mathbf{P}$  are  $p$  1's and  $n - p$  0's
- The trace (sum of diagonal entries) is  $\text{tr}(\mathbf{P}) = p$
- $\mathbf{I} - \mathbf{P}$  is also a projection matrix with rank  $n - p$

## Decomposition

Note that  $\mathbf{P}(\mathbf{I} - \mathbf{P}) = \mathbf{P} - \mathbf{P}\mathbf{P} = \mathbf{P} - \mathbf{P} = \mathbf{0}$ .

We have

$$\begin{aligned} \|\mathbf{Y}\|_2^2 &= \mathbf{Y}^T \mathbf{Y} = (\mathbf{P}\mathbf{Y} + (\mathbf{I} - \mathbf{P})\mathbf{Y})^T (\mathbf{P}\mathbf{Y} + (\mathbf{I} - \mathbf{P})\mathbf{Y}) \\ &= (\mathbf{P}\mathbf{Y})^T (\mathbf{P}\mathbf{Y}) + ((\mathbf{I} - \mathbf{P})\mathbf{Y})^T ((\mathbf{I} - \mathbf{P})\mathbf{Y}) \\ &= \|\mathbf{P}\mathbf{Y}\|_2^2 + \|(\mathbf{I} - \mathbf{P})\mathbf{Y}\|_2^2 \end{aligned}$$

where the cross terms disappear because  $\mathbf{P}(\mathbf{I} - \mathbf{P}) = \mathbf{0}$ .

Note: The  $\ell_p$  norm of an  $n$ -vector  $\mathbf{w}$  is defined as

$$\|\mathbf{w}\|_p = \left( \sum_{i=1}^n |w_i|^p \right)^{1/p}.$$

Above we calculated

$$\|\mathbf{w}\|_2^2 = \sum_{i=1}^n w_i^2.$$

## Distribution of Projection

Suppose that  $Y_1, Y_2, \dots, Y_n \stackrel{\text{iid}}{\sim} \text{Normal}(0, \sigma^2)$ . This can also be written as  $\mathbf{Y} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I})$ . It follows that

$$\mathbf{PY} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{PIP}^T).$$

where  $\mathbf{PIP}^T = \mathbf{PP}^T = \mathbf{PP} = \mathbf{P}$ .

Also,  $(\mathbf{PY})^T(\mathbf{PY}) = \mathbf{Y}^T \mathbf{P}^T \mathbf{PY} = \mathbf{Y}^T \mathbf{PY}$ , a **quadratic form**. Given the eigenvalues of  $\mathbf{P}$ ,  $\mathbf{Y}^T \mathbf{PY}$  is equivalent in distribution to  $p$  squared iid Normal(0,1) rv's, so

$$\frac{\mathbf{Y}^T \mathbf{PY}}{\sigma^2} \sim \chi_p^2.$$

## Distribution of Residuals

If  $\mathbf{PY} = \hat{\mathbf{Y}}$  are the fitted OLS values, then  $(\mathbf{I} - \mathbf{P})\mathbf{Y} = \mathbf{Y} - \hat{\mathbf{Y}}$  are the residuals.

It follows by the same argument as above that

$$\frac{\mathbf{Y}^T (\mathbf{I} - \mathbf{P}) \mathbf{Y}}{\sigma^2} \sim \chi_{n-p}^2.$$

It's also straightforward to show that  $(\mathbf{I} - \mathbf{P})\mathbf{Y} \sim \text{MVN}_n(\mathbf{0}, \sigma^2(\mathbf{I} - \mathbf{P}))$  and  $\text{Cov}(\mathbf{PY}, (\mathbf{I} - \mathbf{P})\mathbf{Y}) = \mathbf{0}$ .

## Degrees of Freedom

The degrees of freedom,  $p$ , of a linear projection model fit is equal to

- The number of linearly dependent columns of  $\mathbf{X}$
- The number of nonzero eigenvalues of  $\mathbf{P}$  (where nonzero eigenvalues are equal to 1)
- The trace of the projection matrix,  $\text{tr}(\mathbf{P})$ .

The reason why we divide estimates of variance by  $n - p$  is because this is the number of effective independent sources of variation remaining after the model is fit by projecting the  $n$  observations into a  $p$  dimensional linear space.

## Submodels

Consider the OLS model  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}$  where there are  $p$  columns of  $\mathbf{X}$  and  $\boldsymbol{\beta}$  is a  $p$ -vector.

Let  $\mathbf{X}_0$  be a subset of  $p_0$  columns of  $\mathbf{X}$  and let  $\mathbf{X}_1$  be a subset of  $p_1$  columns, where  $1 \leq p_0 < p_1 \leq p$ . Also, assume that the columns of  $\mathbf{X}_0$  are a subset of  $\mathbf{X}_1$ .

We can form  $\hat{\mathbf{Y}}_0 = \mathbf{P}_0\mathbf{Y}$  where  $\mathbf{P}_0$  is the projection matrix built from  $\mathbf{X}_0$ . We can analogously form  $\hat{\mathbf{Y}}_1 = \mathbf{P}_1\mathbf{Y}$ .

## Hypothesis Testing

Without loss of generality, suppose that  $\boldsymbol{\beta}_0 = (\beta_1, \beta_2, \dots, \beta_{p_0})^T$  and  $\boldsymbol{\beta}_1 = (\beta_1, \beta_2, \dots, \beta_{p_1})^T$ .

How do we compare these models, specifically to test  $H_0 : (\beta_{p_0+1}, \beta_{p_0+2}, \dots, \beta_{p_1}) = \mathbf{0}$  vs  $H_1 : (\beta_{p_0+1}, \beta_{p_0+2}, \dots, \beta_{p_1}) \neq \mathbf{0}$ ?

The basic idea to perform this test is to compare the goodness of fits of each model via a pivotal statistic. We will discuss the generalized LRT and ANOVA approaches.

## Generalized LRT

Under the OLS Normal model, it follows that  $\hat{\boldsymbol{\beta}}_0 = (\mathbf{X}_0^T \mathbf{X}_0)^{-1} \mathbf{X}_0^T \mathbf{Y}$  is the MLE under the null hypothesis and  $\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{Y}$  is the unconstrained MLE. Also, the respective MLEs of  $\sigma^2$  are

$$\hat{\sigma}_0^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_{0,i})^2}{n}$$

$$\hat{\sigma}_1^2 = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_{1,i})^2}{n}$$

where  $\hat{\mathbf{Y}}_0 = \mathbf{X}_0 \hat{\boldsymbol{\beta}}_0$  and  $\hat{\mathbf{Y}}_1 = \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1$ .

The generalized LRT statistic is

$$\lambda(\mathbf{X}, \mathbf{Y}) = \frac{L(\hat{\boldsymbol{\beta}}_1, \hat{\sigma}_1^2; \mathbf{X}, \mathbf{Y})}{L(\hat{\boldsymbol{\beta}}_0, \hat{\sigma}_0^2; \mathbf{X}, \mathbf{Y})}$$

where  $2 \log \lambda(\mathbf{X}, \mathbf{Y})$  has a  $\chi_{p_1 - p_0}^2$  null distribution.

## Nested Projections

We can apply the Pythagorean theorem we saw earlier to linear subspaces to get:

$$\begin{aligned} \|\mathbf{Y}\|_2^2 &= \|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2 + \|\mathbf{P}_1\mathbf{Y}\|_2^2 \\ &= \|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2 + \|(\mathbf{P}_1 - \mathbf{P}_0)\mathbf{Y}\|_2^2 + \|\mathbf{P}_0\mathbf{Y}\|_2^2 \end{aligned}$$

We can also use the Pythagorean theorem to decompose the residuals from the smaller projection  $\mathbf{P}_0$ :

$$\|(\mathbf{I} - \mathbf{P}_0)\mathbf{Y}\|_2^2 = \|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2 + \|(\mathbf{P}_1 - \mathbf{P}_0)\mathbf{Y}\|_2^2$$

## F Statistic

The  $F$  statistic compares the improvement of goodness in fit of the larger model to that of the smaller model in terms of sums of squared residuals, and it scales this improvement by an estimate of  $\sigma^2$ :

$$\begin{aligned} F &= \frac{[\|(\mathbf{I} - \mathbf{P}_0)\mathbf{Y}\|_2^2 - \|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2] / (p_1 - p_0)}{\|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2 / (n - p_1)} \\ &= \frac{\left[ \sum_{i=1}^n (Y_i - \hat{Y}_{0,i})^2 - \sum_{i=1}^n (Y_i - \hat{Y}_{1,i})^2 \right] / (p_1 - p_0)}{\sum_{i=1}^n (Y_i - \hat{Y}_{1,i})^2 / (n - p_1)} \end{aligned}$$

Since  $\|(\mathbf{I} - \mathbf{P}_0)\mathbf{Y}\|_2^2 - \|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2 = \|(\mathbf{P}_1 - \mathbf{P}_0)\mathbf{Y}\|_2^2$ , we can equivalently write the  $F$  statistic as:

$$\begin{aligned}
F &= \frac{\|(\mathbf{P}_1 - \mathbf{P}_0)\mathbf{Y}\|_2^2 / (p_1 - p_0)}{\|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2 / (n - p_1)} \\
&= \frac{\sum_{i=1}^n (\hat{Y}_{1,i} - \hat{Y}_{0,i})^2 / (p_1 - p_0)}{\sum_{i=1}^n (Y_i - \hat{Y}_{1,i})^2 / (n - p_1)}
\end{aligned}$$

## **F Distribution**

Suppose we have independent random variables  $V \sim \chi_a^2$  and  $W \sim \chi_b^2$ . It follows that

$$\frac{V/a}{W/b} \sim F_{a,b}$$

where  $F_{a,b}$  is the  $F$  distribution with  $(a, b)$  degrees of freedom.

By arguments similar to those given above, we have

$$\frac{\|(\mathbf{P}_1 - \mathbf{P}_0)\mathbf{Y}\|_2^2}{\sigma^2} \sim \chi_{p_1 - p_0}^2$$

$$\frac{\|(\mathbf{I} - \mathbf{P}_1)\mathbf{Y}\|_2^2}{\sigma^2} \sim \chi_{n - p_1}^2$$

and these two rv's are independent.

## **F Test**

Suppose that the OLS model holds where  $\mathbf{E}|\mathbf{X} \sim \text{MVN}_n(\mathbf{0}, \sigma^2 \mathbf{I})$ .

In order to test  $H_0 : (\beta_{p_0+1}, \beta_{p_0+2}, \dots, \beta_{p_1}) = \mathbf{0}$  vs  $H_1 : (\beta_{p_0+1}, \beta_{p_0+2}, \dots, \beta_{p_1}) \neq \mathbf{0}$ , we can form the  $F$  statistic as given above, which has null distribution  $F_{p_1 - p_0, n - p_1}$ . The p-value is calculated as  $\Pr(F \geq F^*)$  where  $F$  is the observed  $F$  statistic and  $F^* \sim F_{p_1 - p_0, n - p_1}$ .

If the above assumption on the distribution of  $\mathbf{E}|\mathbf{X}$  only approximately holds, then the  $F$  test p-value is also an approximation.

## **Example: Davis Data**

```
> library("car")
> data("Davis", package="car")
```

```

> htwt <- tbl_df(Davis)
> htwt[12,c(2,3)] <- htwt[12,c(3,2)]
> head(htwt)
# A tibble: 6 × 5
   sex weight height repwt repht
  <fctr>   <int>   <int>   <int>   <int>
1     M     77    182     77    180
2     F     58    161     51    159
3     F     53    161     54    158
4     M     68    177     70    175
5     F     59    157     59    155
6     M     76    170     76    165

```

## Comparing Linear Models in R

Example: Davis Data

Suppose we are considering the three following models:

```

> f1 <- lm(weight ~ height, data=htwt)
> f2 <- lm(weight ~ height + sex, data=htwt)
> f3 <- lm(weight ~ height + sex + height:sex, data=htwt)

```

How do we determine if the additional terms in models **f2** and **f3** are needed?

## ANOVA (Version 2)

A generalization of ANOVA exists that allows us to compare two nested models, quantifying their differences in terms of goodness of fit and performing a hypothesis test of whether this difference is statistically significant.

A model is *nested* within another model if their difference is simply the absence of certain terms in the smaller model.

The null hypothesis is that the additional terms have coefficients equal to zero, and the alternative hypothesis is that at least one coefficient is nonzero.

Both versions of ANOVA can be described in a single, elegant mathematical framework.

## Comparing Two Models with `anova()`

This provides a comparison of the improvement in fit from model **f2** compared to model **f1**:



```
> anova(f1, f2)
Analysis of Variance Table

Model 1: weight ~ height
Model 2: weight ~ height + sex
  Res.Df  RSS Df Sum of Sq    F    Pr(>F)
1     198 14321
2     197 12816  1    1504.9 23.133 2.999e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## When There's a Single Variable Difference

Compare above `anova(f1, f2)` p-value to that for the `sex` term from the `f2` model:

```
> library(broom)
> tidy(f2)
  term      estimate std.error statistic    p.value
1 (Intercept) -76.6167326 15.71504644 -4.875374 2.231334e-06
2 height      0.8105526  0.09529565  8.505662 4.499241e-15
3 sexM        8.2268893  1.71050385  4.809629 2.998988e-06
```

## Calculating the F-statistic

```
> anova(f1, f2)
Analysis of Variance Table

Model 1: weight ~ height
Model 2: weight ~ height + sex
  Res.Df  RSS Df Sum of Sq    F    Pr(>F)
1     198 14321
2     197 12816  1    1504.9 23.133 2.999e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

How the F-statistic is calculated:

```
> n <- nrow(htwt)
> ss1 <- (n-1)*var(f1$residuals)
> ss1
[1] 14321.11
> ss2 <- (n-1)*var(f2$residuals)
> ss2
```

```
[1] 12816.18
> ((ss1 - ss2)/anova(f1, f2)$Df[2])/(ss2/f2$df.residual)
[1] 23.13253
```

## Calculating the Generalized LRT

```
> anova(f1, f2, test="LRT")
Analysis of Variance Table

Model 1: weight ~ height
Model 2: weight ~ height + sex
  Res.Df  RSS Df Sum of Sq  Pr(>Chi)
1     198 14321
2     197 12816   1    1504.9 1.512e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

> library(lmtest)
> lrtest(f1, f2)
Likelihood ratio test

Model 1: weight ~ height
Model 2: weight ~ height + sex
  #Df LogLik Df  Chisq Pr(>Chisq)
1    3 -710.9
2    4 -699.8   1 22.205   2.45e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

These tests produce slightly different answers because `anova()` adjusts for degrees of freedom when estimating the variance, whereas `lrtest()` is the strict generalized LRT. See [here](#).

## ANOVA on More Distant Models

We can compare models with multiple differences in terms:

```
> anova(f1, f3)
Analysis of Variance Table

Model 1: weight ~ height
Model 2: weight ~ height + sex + height:sex
  Res.Df  RSS Df Sum of Sq    F    Pr(>F)
1     198 14321
```

```

2      196 12567  2      1754 13.678 2.751e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

## Compare Multiple Models at Once

We can compare multiple models at once:

```

> anova(f1, f2, f3)
Analysis of Variance Table

Model 1: weight ~ height
Model 2: weight ~ height + sex
Model 3: weight ~ height + sex + height:sex
  Res.Df  RSS Df Sum of Sq    F    Pr(>F)
1     198 14321
2     197 12816  1    1504.93 23.4712 2.571e-06 ***
3     196 12567  1     249.04  3.8841  0.05015 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

## Generalized Linear Models

### Definition

The generalized linear model (GLM) builds from OLS and GLS to allow for the case where  $Y|\mathbf{X}$  is distributed according to an exponential family distribution. The estimated model is

$$g(E[Y|\mathbf{X}]) = \mathbf{X}\boldsymbol{\beta}$$

where  $g(\cdot)$  is called the **link function**. This model is typically fit by numerical methods to calculate the maximum likelihood estimate of  $\boldsymbol{\beta}$ .

### Exponential Family Distributions

Recall that if  $Y$  follows an EFD then it has pdf of the form

$$f(y; \boldsymbol{\theta}) = h(y) \exp \left\{ \sum_{k=1}^d \eta_k(\boldsymbol{\theta}) T_k(y) - A(\boldsymbol{\eta}) \right\}$$

where  $\boldsymbol{\theta}$  is a vector of parameters,  $\{T_k(y)\}$  are sufficient statistics,  $A(\boldsymbol{\eta})$  is the cumulant generating function.

The functions  $\eta_k(\boldsymbol{\theta})$  for  $k = 1, \dots, d$  map the usual parameters  $\boldsymbol{\theta}$  (often moments of the rv  $Y$ ) to the *natural parameters* or *canonical parameters*.

$\{T_k(y)\}$  are sufficient statistics for  $\{\eta_k\}$  due to the factorization theorem.

$A(\boldsymbol{\eta})$  is sometimes called the *log normalizer* because

$$A(\boldsymbol{\eta}) = \log \int h(y) \exp \left\{ \sum_{k=1}^d \eta_k(\boldsymbol{\theta}) T_k(y) \right\}.$$

## Natural Single Parameter EFD

A natural single parameter EFD simplifies to the scenario where  $d = 1$  and  $T(y) = y$

$$f(y; \eta) = h(y) \exp \{ \eta(\theta)y - A(\eta(\theta)) \}$$

where without loss of generality we can write  $E[Y] = \theta$ .

## Dispersion EFDs

The family of distributions for which GLMs are most typically developed are dispersion EFDs. An example of a dispersion EFD that extends the natural single parameter EFD is

$$f(y; \eta) = h(y, \phi) \exp \left\{ \frac{\eta(\theta)y - A(\eta(\theta))}{\phi} \right\}$$

where  $\phi$  is the dispersion parameter.

## Example: Normal

Let  $Y \sim \text{Normal}(\mu, \sigma^2)$ . Then:

$$\theta = \mu, \eta(\mu) = \mu$$

$$\phi = \sigma^2$$

$$A(\mu) = \frac{\mu^2}{2}$$

$$h(y, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{y^2}{\sigma^2}}$$

## EFD for GLMs

There has been a very broad development of GLMs and extensions. A common setting for introducing GLMs is the dispersion EFD with a general link function  $g(\cdot)$ .

See the classic text *Generalized Linear Models*, by McCullagh and Nelder, for such a development.

## Components of a GLM

1. *Random*: The particular exponential family distribution.

$$Y \sim f(y; \eta, \phi)$$

2. *Systematic*: The determination of each  $\eta_i$  from  $\mathbf{X}_i$  and  $\beta$ .

$$\eta_i = \mathbf{X}_i \beta$$

3. *Parametric Link*: The connection between  $E[Y_i | \mathbf{X}_i]$  and  $\mathbf{X}_i \beta$ .

$$g(E[Y_i | \mathbf{X}_i]) = \mathbf{X}_i \beta$$

## Link Functions

Even though the link function  $g(\cdot)$  can be considered in a fairly general framework, the **canonical link function**  $\eta(\cdot)$  is often utilized.

The canonical link function is the function that maps the expected value into the natural parameter.

In this case,  $Y | \mathbf{X}$  is distributed according to an exponential family distribution with

$$\eta(E[Y | \mathbf{X}]) = \mathbf{X} \beta.$$

## Calculating MLEs

Given the model  $g(E[Y|\mathbf{X}]) = \mathbf{X}\boldsymbol{\beta}$ , the EFD should be fully parameterized. The Newton-Raphson method or Fisher's scoring method can be utilized to find the MLE of  $\boldsymbol{\beta}$ .

### Newton-Raphson

1. Initialize  $\boldsymbol{\beta}^{(0)}$ . For  $t = 1, 2, \dots$
2. Calculate the score  $s(\boldsymbol{\beta}^{(t)}) = \nabla \ell(\boldsymbol{\beta}; \mathbf{X}, \mathbf{Y}) |_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(t)}}$  and observed Fisher information
$$H(\boldsymbol{\beta}^{(t)}) = -\nabla \nabla^T \ell(\boldsymbol{\beta}; \mathbf{X}, \mathbf{Y}) |_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(t)}}$$

. Note that the observed Fisher information is also the negative Hessian matrix.
3. Update  $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + H(\boldsymbol{\beta}^{(t)})^{-1} s(\boldsymbol{\beta}^{(t)})$ .
4. Iterate until convergence, and set  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^{(\infty)}$ .

### Fisher's scoring

1. Initialize  $\boldsymbol{\beta}^{(0)}$ . For  $t = 1, 2, \dots$
2. Calculate the score  $s(\boldsymbol{\beta}^{(t)}) = \nabla \ell(\boldsymbol{\beta}; \mathbf{X}, \mathbf{Y}) |_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(t)}}$  and expected Fisher information
$$I(\boldsymbol{\beta}^{(t)}) = -E \left[ \nabla \nabla^T \ell(\boldsymbol{\beta}; \mathbf{X}, \mathbf{Y}) |_{\boldsymbol{\beta}=\boldsymbol{\beta}^{(t)}} \right]$$
3. Update  $\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} + I(\boldsymbol{\beta}^{(t)})^{-1} s(\boldsymbol{\beta}^{(t)})$ .
4. Iterate until convergence, and set  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^{(\infty)}$ .

When the canonical link function is used, the Newton-Raphson algorithm and Fisher's scoring algorithm are equivalent.

Exercise: Prove this.

## Iteratively Reweighted Least Squares

For the canonical link, Fisher's scoring method can be written as an iteratively reweighted least squares algorithm, as shown earlier for logistic regression. Note that the Fisher information is

$$I(\boldsymbol{\beta}^{(t)}) = \mathbf{X}^T \mathbf{W} \mathbf{X}$$

where  $\mathbf{W}$  is an  $n \times n$  diagonal matrix with  $(i, i)$  entry equal to  $\text{Var}(Y_i|\mathbf{X}; \boldsymbol{\beta}^{(t)})$ . The score function is

$$s(\boldsymbol{\beta}^{(t)}) = \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^{(t)})$$

and the current coefficient value  $\boldsymbol{\beta}^{(t)}$  can be written as

$$\boldsymbol{\beta}^{(t)} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{X} \boldsymbol{\beta}^{(t)}.$$

Putting this together we get

$$\boldsymbol{\beta}^{(t)} + I(\boldsymbol{\beta}^{(t)})^{-1} s(\boldsymbol{\beta}^{(t)}) = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}^{(t)}$$

where

$$\mathbf{z}^{(t)} = \mathbf{X}\boldsymbol{\beta}^{(t)} + \mathbf{W}^{-1} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}^{(t)}).$$

This is a generalization of the iteratively reweighted least squares algorithm we showed earlier for logistic regression.

## Estimating Dispersion

For the simple dispersion model above, it is typically straightforward to calculate the MLE  $\hat{\phi}$  once  $\hat{\boldsymbol{\beta}}$  has been calculated.

## CLT Applied to the MLE

Given that  $\hat{\boldsymbol{\beta}}$  is a maximum likelihood estimate, we have the following CLT result on its distribution as  $n \rightarrow \infty$ :

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{D} \text{MVN}_p(\mathbf{0}, \phi(\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1})$$

## Approximately Pivotal Statistics

The previous CLT gives us the following two approximations for pivotal statistics. The first statistic facilitates getting overall measures of uncertainty on the estimate  $\hat{\boldsymbol{\beta}}$ .

$$\hat{\phi}^{-1}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})^T (\mathbf{X}^T \hat{\mathbf{W}} \mathbf{X}) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \sim \chi_1^2$$

This second pivotal statistic allows for performing a Wald test or forming a confidence interval on each coefficient,  $\beta_j$ , for  $j = 1, \dots, p$ .

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\hat{\phi}[(\mathbf{X}^T \hat{\mathbf{W}} \mathbf{X})^{-1}]_{jj}}} \sim \text{Normal}(0, 1)$$

## Deviance

Let  $\hat{\boldsymbol{\eta}}$  be the estimated natural parameters from a GLM. For example,  $\hat{\boldsymbol{\eta}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  when the canonical link function is used.

Let  $\hat{\boldsymbol{\eta}}_n$  be the **saturated model** where  $Y_i$  is directly used to estimate  $\eta_i$  without model constraints. For example, in the Bernoulli logistic regression model  $\hat{\boldsymbol{\eta}}_n = \mathbf{Y}$ , the observed outcomes.

The **deviance** for the model is defined to be

$$D(\hat{\boldsymbol{\eta}}) = 2\ell(\hat{\boldsymbol{\eta}}_n; \mathbf{X}, \mathbf{Y}) - 2\ell(\hat{\boldsymbol{\eta}}; \mathbf{X}, \mathbf{Y})$$

## Generalized LRT

Let  $\mathbf{X}_0$  be a subset of  $p_0$  columns of  $\mathbf{X}$  and let  $\mathbf{X}_1$  be a subset of  $p_1$  columns, where  $1 \leq p_0 < p_1 \leq p$ . Also, assume that the columns of  $\mathbf{X}_0$  are a subset of  $\mathbf{X}_1$ .

Without loss of generality, suppose that  $\boldsymbol{\beta}_0 = (\beta_1, \beta_2, \dots, \beta_{p_0})^T$  and  $\boldsymbol{\beta}_1 = (\beta_1, \beta_2, \dots, \beta_{p_1})^T$ .

Suppose we wish to test  $H_0 : (\beta_{p_0+1}, \beta_{p_0+2}, \dots, \beta_{p_1}) = \mathbf{0}$  vs  $H_1 : (\beta_{p_0+1}, \beta_{p_0+2}, \dots, \beta_{p_1}) \neq \mathbf{0}$

We can form  $\hat{\boldsymbol{\eta}}_0 = \mathbf{X}\hat{\boldsymbol{\beta}}_0$  from the GLM model  $g(\mathbb{E}[Y|\mathbf{X}_0]) = \mathbf{X}_0\boldsymbol{\beta}_0$ . We can analogously form  $\hat{\boldsymbol{\eta}}_1 = \mathbf{X}\hat{\boldsymbol{\beta}}_1$  from the GLM model  $g(\mathbb{E}[Y|\mathbf{X}_1]) = \mathbf{X}_1\boldsymbol{\beta}_1$ .

The  $2\log$  generalized LRT can then be formed from the two deviance statistics

$$2\log \lambda(\mathbf{X}, \mathbf{Y}) = 2\log \frac{L(\hat{\boldsymbol{\eta}}_1; \mathbf{X}, \mathbf{Y})}{L(\hat{\boldsymbol{\eta}}_0; \mathbf{X}, \mathbf{Y})} = D(\hat{\boldsymbol{\eta}}_0) - D(\hat{\boldsymbol{\eta}}_1)$$

where the null distribution is  $\chi^2_{p_1-p_0}$ .



## Example: Grad School Admissions

Let's revisit a logistic regression example now that we know how the various statistics are calculated.

```
> mydata <-  
+   read.csv("http://www.ats.ucla.edu/stat/data/binary.csv")  
> dim(mydata)  
[1] 400  4  
> head(mydata)  
  admit gre  gpa rank  
1     0 380 3.61   3  
2     1 660 3.67   3  
3     1 800 4.00   1  
4     1 640 3.19   4  
5     0 520 2.93   4  
6     1 760 3.00   2
```

Fit the model with basic output. Note the argument `family = "binomial"`.

```
> mydata$rank <- factor(mydata$rank, levels=c(1, 2, 3, 4))  
> myfit <- glm(admit ~ gre + gpa + rank,  
+             data = mydata, family = "binomial")  
> myfit
```

```
Call:  glm(formula = admit ~ gre + gpa + rank, family = "binomial",  
          data = mydata)
```

Coefficients:

(Intercept)	gre	gpa	rank2
-3.989979	0.002264	0.804038	-0.675443
rank3	rank4		
-1.340204	-1.551464		

Degrees of Freedom: 399 Total (i.e. Null); 394 Residual

Null Deviance: 500

Residual Deviance: 458.5 AIC: 470.5

This shows the fitted coefficient values, which is on the link function scale – logit aka log odds here. Also, a Wald test is performed for each coefficient.

```
> summary(myfit)
```

Call:

```
glm(formula = admit ~ gre + gpa + rank, family = "binomial",  
    data = mydata)
```

Deviance Residuals:

```

      Min       1Q   Median       3Q      Max
-1.6268 -0.8662 -0.6388  1.1490  2.0790

Coefficients:
              Estimate Std. Error z value Pr(>|z|)
(Intercept) -3.989979   1.139951  -3.500 0.000465 ***
gre          0.002264   0.001094   2.070 0.038465 *
gpa          0.804038   0.331819   2.423 0.015388 *
rank2       -0.675443   0.316490  -2.134 0.032829 *
rank3       -1.340204   0.345306  -3.881 0.000104 ***
rank4       -1.551464   0.417832  -3.713 0.000205 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 499.98  on 399  degrees of freedom
Residual deviance: 458.52  on 394  degrees of freedom
AIC: 470.52

Number of Fisher Scoring iterations: 4

```

Here we perform a generalized LRT on each variable. Note the rank variable is now tested as a single factor variable as opposed to each dummy variable.

```

> anova(myfit, test="LRT")
Analysis of Deviance Table

Model: binomial, link: logit

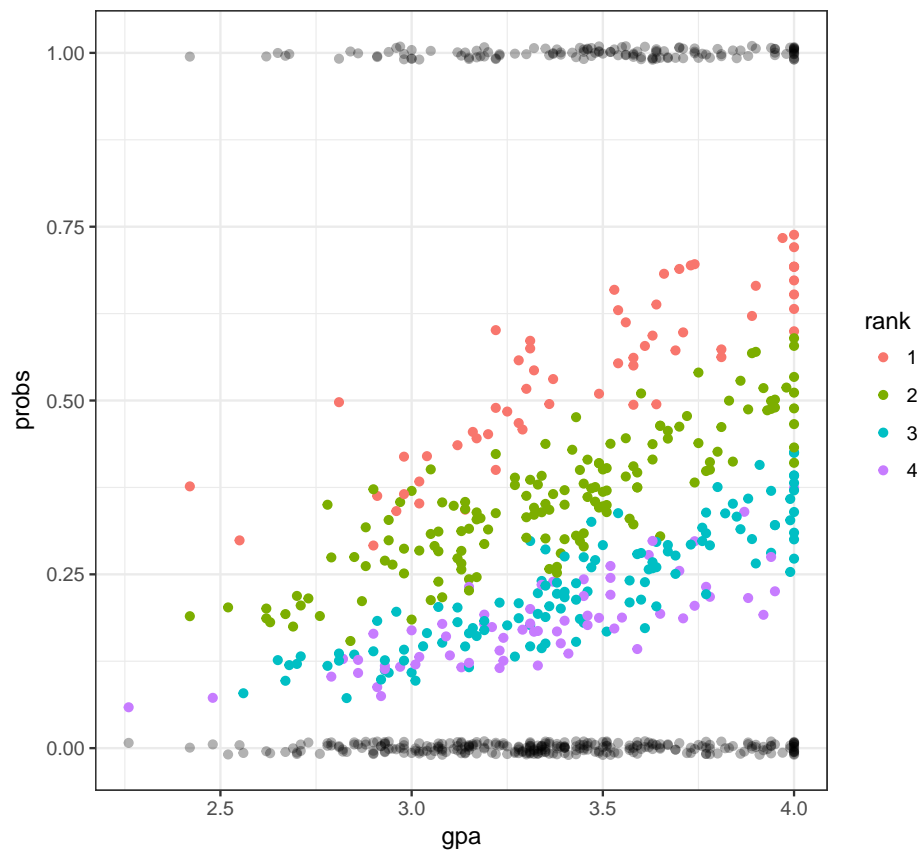
Response: admit

Terms added sequentially (first to last)

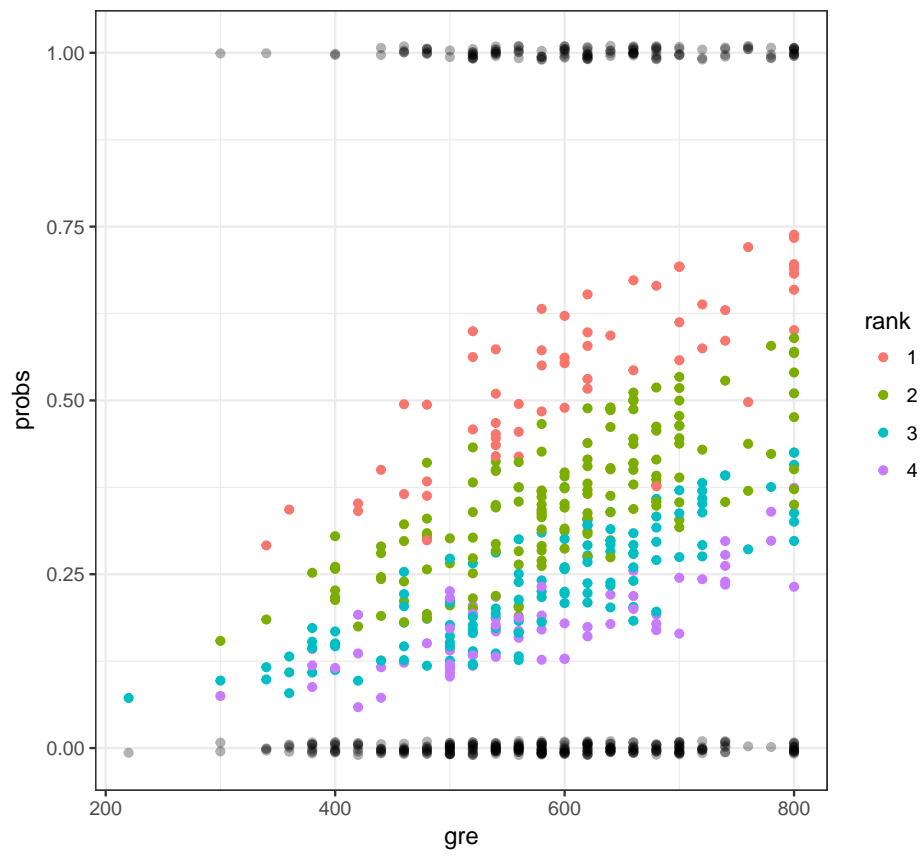
      Df Deviance Resid. Df Resid. Dev  Pr(>Chi)
NULL                                399    499.98
gre   1  13.9204    398    486.06 0.0001907 ***
gpa   1   5.7122    397    480.34 0.0168478 *
rank  3  21.8265    394    458.52 7.088e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

> mydata <- data.frame(mydata, probs = myfit$fitted.values)
> ggplot(mydata) + geom_point(aes(x=gpa, y=probs, color=rank)) +
+   geom_jitter(aes(x=gpa, y=admit), width=0, height=0.01, alpha=0.3)

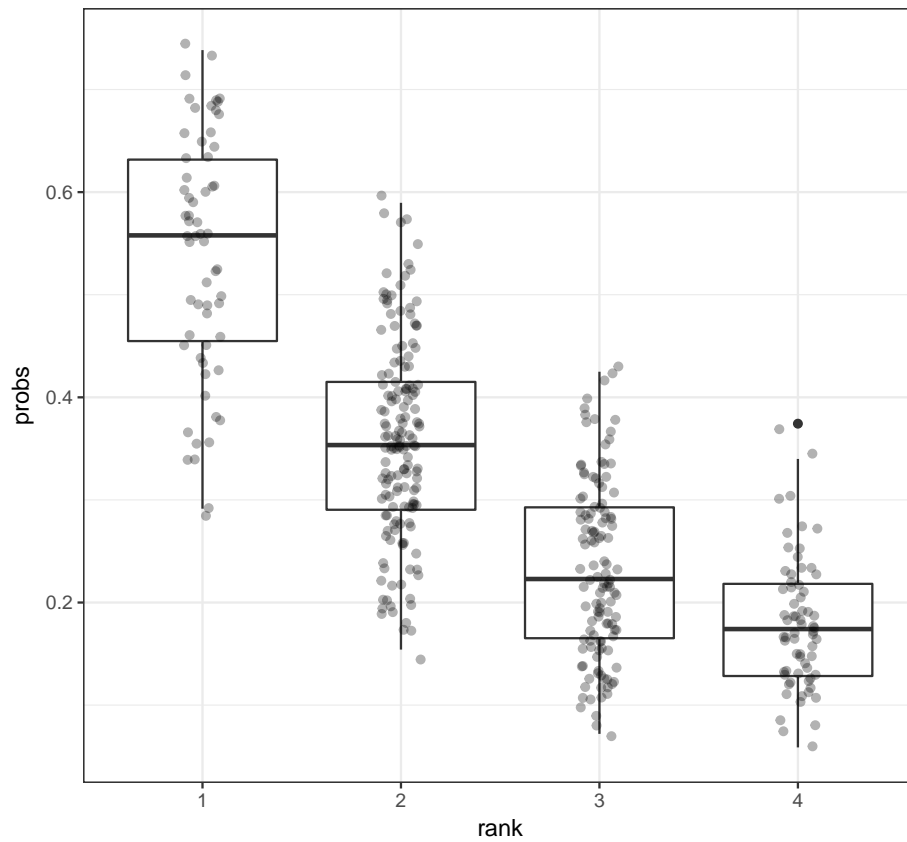
```



```
> ggplot(mydata) + geom_point(aes(x=gpa, y=probs, color=rank)) +  
+   geom_jitter(aes(x=gpa, y=admit), width=0, height=0.01, alpha=0.3)
```



```
> ggplot(mydata) + geom_boxplot(aes(x=rank, y=probs)) +
+   geom_jitter(aes(x=rank, y=probs), width=0.1, height=0.01, alpha=0.3)
```



## glm() Function

The `glm()` function has many different options available to the user.

```
glm(formula, family = gaussian, data, weights, subset,
    na.action, start = NULL, etastart, mustart, offset,
    control = list(...), model = TRUE, method = "glm.fit",
    x = FALSE, y = TRUE, contrasts = NULL, ...)
```

To see the different link functions available, type:

```
help(family)
```

## Nonparametric Regression

## Generalized Additive Models

## Bootstrap for Statistical Models

## Extras

### Source

License

Source Code

### Session Information

```
> sessionInfo()
R version 3.3.2 (2016-10-31)
Platform: x86_64-apple-darwin13.4.0 (64-bit)
Running under: macOS Sierra 10.12.4

locale:
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods
[7] base

other attached packages:
[1] lmtest_0.9-35  zoo_1.8-0      car_2.1-4
[4] broom_0.4.2    dplyr_0.5.0    purrr_0.2.2
[7] readr_1.0.0    tidyr_0.6.1    tibble_1.2
[10] ggplot2_2.2.1  tidyverse_1.1.1 knitr_1.15.1
[13] magrittr_1.5   devtools_1.12.0

loaded via a namespace (and not attached):
[1] reshape2_1.4.2  splines_3.3.2  haven_1.0.0
[4] lattice_0.20-34 colorspace_1.3-2 htmltools_0.3.5
[7] yaml_2.1.14     mgcv_1.8-17    nloptr_1.0.4
[10] foreign_0.8-67  withr_1.0.2    DBI_0.5-1
[13] modelr_0.1.0    readxl_0.1.1   plyr_1.8.4
[16] stringr_1.1.0   MatrixModels_0.4-1 munsell_0.4.3
```

[19]	<a href="#">gtable_0.2.0</a>	<a href="#">rvest_0.3.2</a>	<a href="#">psych_1.6.12</a>
[22]	<a href="#">memoise_1.0.0</a>	<a href="#">evaluate_0.10</a>	<a href="#">labeling_0.3</a>
[25]	<a href="#">forcats_0.2.0</a>	<a href="#">SparseM_1.74</a>	<a href="#">quantreg_5.29</a>
[28]	<a href="#">pbkrtest_0.4-6</a>	<a href="#">parallel_3.3.2</a>	<a href="#">Rcpp_0.12.9</a>
[31]	<a href="#">scales_0.4.1</a>	<a href="#">backports_1.0.5</a>	<a href="#">jsonlite_1.2</a>
[34]	<a href="#">lme4_1.1-12</a>	<a href="#">mnormt_1.5-5</a>	<a href="#">hms_0.3</a>
[37]	<a href="#">digest_0.6.12</a>	<a href="#">stringi_1.1.2</a>	<a href="#">grid_3.3.2</a>
[40]	<a href="#">rprojroot_1.2</a>	<a href="#">tools_3.3.2</a>	<a href="#">lazyeval_0.2.0</a>
[43]	<a href="#">MASS_7.3-45</a>	<a href="#">Matrix_1.2-8</a>	<a href="#">xml2_1.1.1</a>
[46]	<a href="#">lubridate_1.6.0</a>	<a href="#">assertthat_0.1</a>	<a href="#">minqa_1.2.4</a>
[49]	<a href="#">rmarkdown_1.3</a>	<a href="#">httr_1.2.1</a>	<a href="#">R6_2.2.0</a>
[52]	<a href="#">nnet_7.3-12</a>	<a href="#">nlme_3.1-131</a>	