



北京航空航天大学

— 经济管理学院 —

BEIHANG UNIVERSITY
SCHOOL OF ECONOMICS AND MANAGEMENT

Generalized Linear Models

Lecture 2: Linear models with R



- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model
- 4 Regression diagnostics
- 5 Variable selection
- 6 Conclusion
- 7 Assignment One

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1} + \varepsilon$$

- Response: y
- Predictors: x_1, \dots, x_{p-1}
- Error: ε

Let $\mathbf{y} = (y_1, \dots, y_n)'$, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$, $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})'$ and

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} & \cdots & x_{p-1,1} \\ 1 & x_{1,2} & x_{2,2} & \cdots & x_{p-1,2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{1,n} & x_{2,n} & \cdots & x_{p-1,n} \end{bmatrix} \quad \text{(the model matrix).}$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1} + \varepsilon$$

- Response: y
- Predictors: x_1, \dots, x_{p-1}
- Error: ε

Let $\mathbf{y} = (y_1, \dots, y_n)'$, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$, $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})'$ and

$$\mathbf{X} = \begin{bmatrix} 1 & x_{1,1} & x_{2,1} & \cdots & x_{p-1,1} \\ 1 & x_{1,2} & x_{2,2} & \cdots & x_{p-1,2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{1,n} & x_{2,n} & \cdots & x_{p-1,n} \end{bmatrix} \quad (\text{the model matrix}).$$

Then

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

Least squares estimation

Minimize: $(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$

Least squares estimation

Minimize: $(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$

Differentiate wrt β gives

The “normal” equations

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

If the errors are iid and normally distributed, then

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}).$$

If the errors are iid and normally distributed, then

$$\mathbf{y} \sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}).$$

So the likelihood is

$$L = \frac{1}{\sigma^T(2\pi)^{T/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

If the errors are iid and normally distributed, then

$$\mathbf{y} \sim \mathbf{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}).$$

So the likelihood is

$$L = \frac{1}{\sigma^T(2\pi)^{T/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

which is maximized when $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ is minimized.

If the errors are iid and normally distributed, then

$$\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}).$$

So the likelihood is

$$L = \frac{1}{\sigma^T(2\pi)^{T/2}} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

which is maximized when $(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$ is minimized.

So MLE = OLS.

```
fit <- lm(response ~ x1 + x2 + x3,  
          data=tibble)
```

$$\text{response} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$$

- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model
- 4 Regression diagnostics
- 5 Variable selection
- 6 Conclusion
- 7 Assignment One

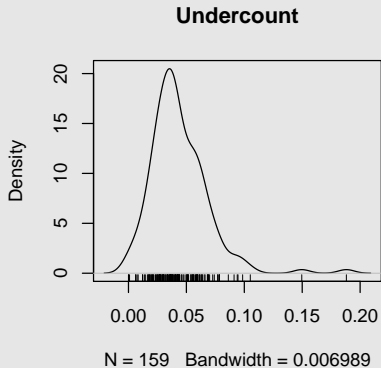
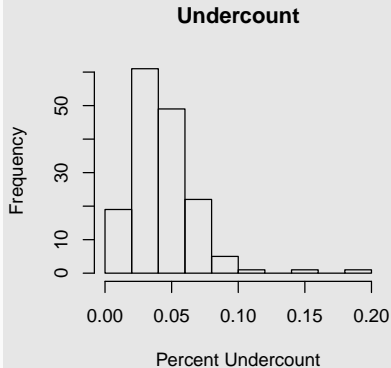
Have a look at the data

```
library(faraway)
data(gavote)
gavote
head(gavote)
?gavote
str(gavote)
```

```
summary(gavote)
gavote$undercount <-
  (gavote$ballots-gavote$votes)/gavote$ballots
summary(gavote$undercount)
with(gavote,
      sum(ballots-votes)/sum(ballots))
```

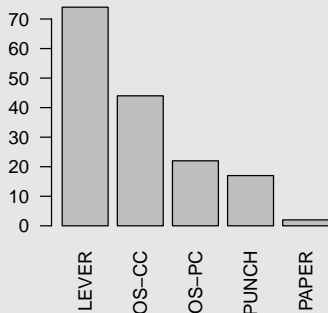
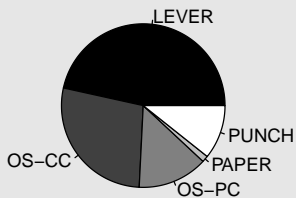
```
par(mfrow = c(1,2))
hist(gavote$undercount,main="Undercount",
      xlab="Percent Undercount")
plot(density(gavote$undercount),main="Undercount")
rug(gavote$undercount)
```

Initial data analysis - visualization



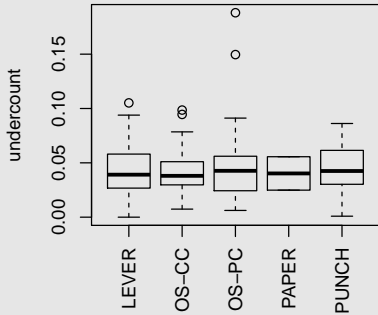
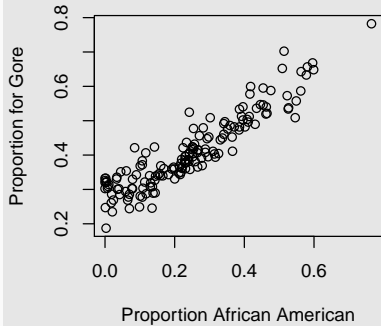
Initial data analysis - visualization

```
par(mfrow = c(1,2))  
pie(table(gavote$equip), col=gray(0:4/4))  
barplot(sort(table(gavote$equip), decreasing=TRUE), las=2)
```



```
par(mfrow = c(1,2))
gavote$pergore <- gavote$gore/gavote$votes
plot(pergore ~ perAA, gavote,
      xlab="Proportion African American",
      ylab="Proportion for Gore")
plot(undercount ~ equip, gavote, xlab="", las=3)
```

Initial data analysis - visualization



```
names(gavote)
names(gavote)[4] <- "usage"
nix <- c(3,10,11,12)
cor(gavote[,nix])
```

- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model**
- 4 Regression diagnostics
- 5 Variable selection
- 6 Conclusion
- 7 Assignment One

$$\text{undercount} = \beta_0 + \beta_1 \text{pergore} + \beta_2 \text{perAA} + \epsilon$$

```
lmod <- lm(undercount ~ pergore + perAA, gavote)
```

- `coef`: $\hat{\beta}$
- `fitted`: \hat{y}
- `predict`: predicted values based on linear model
- `residuals`: $y - \hat{y}$
- `deviance`: $RSS = \hat{\epsilon}^T \hat{\epsilon}$
- `df.residual`: degrees of freedom = $n - p$
- `summary`: a summary of the model

```
predict(lmod)
residuals(lmod)
deviance(lmod)
df.residual(lmod)
nrow(gavote)-length(coef(lmod))
lmodsum <- summary(lmod)
lmodsum$sigma
```


Fitting a linear model

- Estimation of σ

$$\hat{\sigma} = \sqrt{\text{RSS}/\text{df}}$$

- R^2 : coefficient of determination

$$R^2 = 1 - \frac{\sum(\hat{y}_i - y_i)^2}{\sum(y_i - \bar{y})^2} = 1 - \frac{\text{RSS}}{\text{TSS}} = r^2 = \text{cor}^2(\hat{y}, y)$$

- Adjusted R^2

$$R_a^2 = 1 - \frac{\text{RSS}/(n - p)}{\text{TSS}/(n - 1)}$$

```
sqrt(deviance(lmod)/df.residual(lmod))  
lmodsum$r.squared  
cor(predict(lmod),gavote$undercount)^2  
lmodsum$adj.r.squared
```

- R will create the required dummy variables from a categorical *factor*.
- The first level is used as the reference category.
- Use `relevel` to change the reference category

$$\text{undercount} = \beta_0 + \beta_1 \text{pergore} + \beta_2 \text{perAA} + \beta_3 d + \epsilon$$
$$d = \begin{cases} 0 & \text{rural} \\ 1 & \text{urban} \end{cases}$$

Interactions:

- Interactions are obtained by multiplying the relevant columns of the model matrix.
- Use $a:b$ for the interaction between a and b .
- Use $a*b$ to mean $a + b + a:b$

Limited order interactions:

- Interactions up to 2nd order can be specified using the \wedge operator.
- $(a+b+c)^2$ is identical to $(a+b+c)*(a+b+c)$

Nested factors:

- $a + b \%in\% a$ expands to $a + a:b$.

- Each coefficient gives the effect of a one unit increase of the predictor on the response variable, *holding all other variables constant*.
- Be careful with interactions: you cannot interpret the main effects when there is an interaction between them.

```
gavote$cpergore <- gavote$pergore - mean(gavote$pergore)
gavote$cperAA <- gavote$perAA - mean(gavote$perAA)
lmodi <- lm(undercount ~ cperAA+cpergore*usage+equip,
            gavote)
summary(lmodi)
```

- Use F-tests between models:
 - Model 1: p_1 parameters.
 - Model 2 (nested within Model 1): p_2 parameters.
 - H_0 : the smaller model is correct.

$$F = \frac{(RSS_2 - RSS_1)/(p_1 - p_2)}{RSS_1/(n - p_1)} \sim F_{p_1 - p_2, n - p_1}$$

```
anova(lmod, lmodi)
drop1(lmodi, test="F")
```

- Hierarchy principal: all lower-order terms corresponding to an interaction should be retained in the model.

- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model
- 4 Regression diagnostics
- 5 Variable selection
- 6 Conclusion
- 7 Assignment One

- Residuals vs Fitted: Check heteroskedasticity
- Scale-Location: standardized residuals vs fitted values.
- Normal QQ plot: Check for non-normality
- Residuals vs Leverage: Check for influential points

Produced using `ggfortify::autoplot` or `plot`

Fitted values

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is the “hat matrix”.

Fitted values

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is the “hat matrix”.

Theorem: Leave-one-out residuals

Let h_1, \dots, h_n be the diagonal values of \mathbf{H} . Then $e_{(i)} = e_i / (1 - h_i)$ is the prediction error that would be obtained if the i th observation was omitted.

Fitted values

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is the "hat matrix".

Theorem: Leave-one-out residuals

Let h_1, \dots, h_n be the diagonal values of \mathbf{H} . Then $e_{(i)} = e_i / (1 - h_i)$ is the prediction error that would be obtained if the i th observation was omitted.

Leverage values

h_i is called the "leverage" of observation i .

In R: `hatvalues(fit)`

Let $\mathbf{X}_{[i]}$ and $\mathbf{Y}_{[i]}$ be similar to \mathbf{X} and \mathbf{Y} but with the i th row deleted in each case. Let \mathbf{x}'_i be the i th row of \mathbf{X} and let

$$\hat{\beta}_{[i]} = (\mathbf{X}'_{[i]}\mathbf{X}_{[i]})^{-1}\mathbf{X}'_{[i]}\mathbf{Y}_{[i]}$$

be the estimate of β without the i th case. Then $e_{[i]} = y_i - \mathbf{x}'_i\hat{\beta}_{[i]}$.

Now $\mathbf{X}'_{[i]}\mathbf{X}_{[i]} = (\mathbf{X}'\mathbf{X} - \mathbf{x}_i\mathbf{x}'_i)$ and $\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i = h_i$.

Sherman-Morrison-Woodbury formula

Suppose \mathbf{A} is a square matrix, and \mathbf{u} and \mathbf{v} are column vectors of the same dimension. Then

$$(\mathbf{A} + \mathbf{u}\mathbf{v}')^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{u}\mathbf{v}'\mathbf{A}^{-1}}{1 + \mathbf{v}'\mathbf{A}^{-1}\mathbf{u}}.$$

So by SMW,

$$(\mathbf{X}'_{[i]}\mathbf{X}_{[i]})^{-1} = (\mathbf{X}'\mathbf{X})^{-1} + \frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}'_i(\mathbf{X}'\mathbf{X})^{-1}}{1 - h_i}.$$

Also note that $\mathbf{X}'_{[i]}\mathbf{Y}_{[i]} = \mathbf{X}'\mathbf{Y} - \mathbf{x}_i y_i$.

Therefore

$$\begin{aligned}\hat{\beta}_{[i]} &= \left[(\mathbf{X}'\mathbf{X})^{-1} + \frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i\mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}}{1 - h_i} \right] (\mathbf{X}'\mathbf{Y} - \mathbf{x}_iy_i) \\ &= \hat{\beta} - \left[\frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i}{1 - h_i} \right] \left[y_i(1 - h_i) - \mathbf{x}_i'\hat{\beta} + h_iy_i \right] \\ &= \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_ie_i/(1 - h_i)\end{aligned}$$

Thus

$$\begin{aligned}e_{[i]} &= y_i - \mathbf{x}'_i \hat{\boldsymbol{\beta}}_{[i]} \\&= y_i - \mathbf{x}'_i \left[\hat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i e_i / (1 - h_i) \right] \\&= e_i + h_i e_i / (1 - h_i) \\&= e_i / (1 - h_i),\end{aligned}$$

Cross-validation statistic

$$CV = \frac{1}{T} \sum_{i=1}^n [e_i / (1 - h_i)]^2,$$

- Measures MSE of out-of-sample prediction
- Asymptotically equivalent to AIC (up to monotonic transformation)

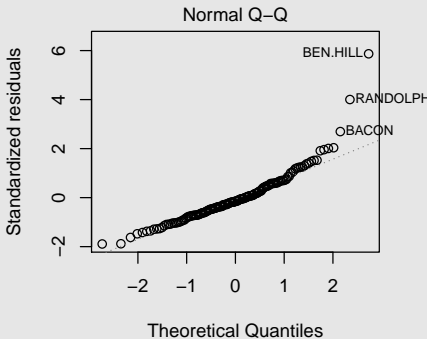
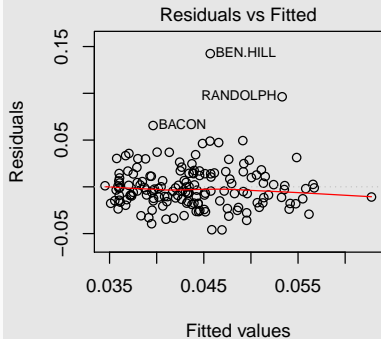
Cook distances

$$D_i = \frac{e_i^2 h_i}{\hat{\sigma}^2 p (1 - h_i)}$$

- Measures change in fit if observation i dropped.

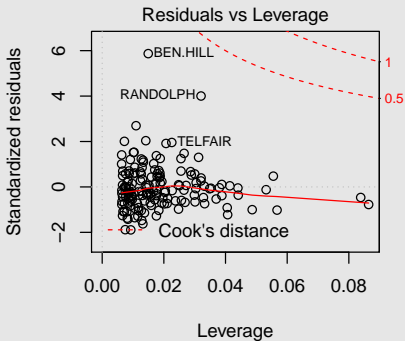
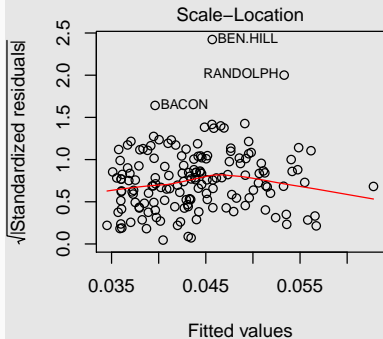
Regression diagnostics

```
par(mfrow = c(1,2))  
plot(lmod, which = c(1:2))
```

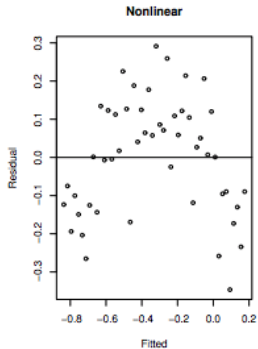
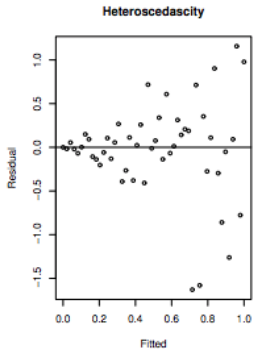
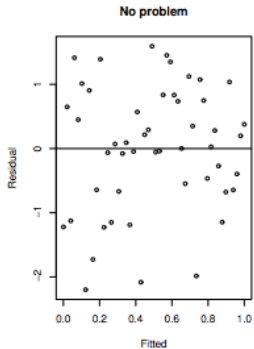


Regression diagnostics

```
par(mfrow = c(1,2))  
plot(lmod, which = c(3,5))
```



Which one is not expected?



- OLS works well for normal errors
- How about when we have outliers?
 - 1 outliers are incorrect observations: drop them
 - 2 outliers are real observations: robust regression
- Robust regression can downweight the effects of large errors
- `r1m` in R

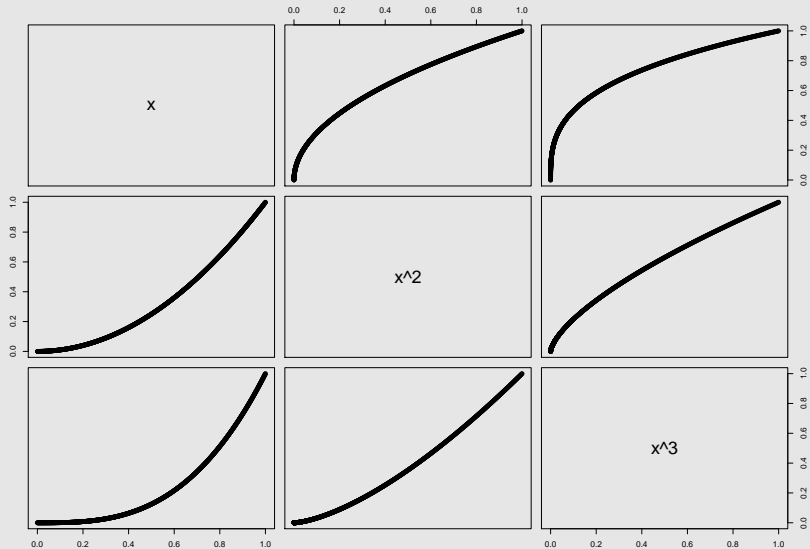
Simple approach: Add x , x^2 , x^3 , ... to model.

```
fit <- lm(y ~ x + I(x^2) + I(x^3))
```

Problems:

- predictors are highly correlated, so difficult to separate effects of terms.
- numerical instability of coefficients

Polynomial regression



$$z_1 = a_1 + b_1x$$

$$z_2 = a_2 + b_2x + c_2x^2$$

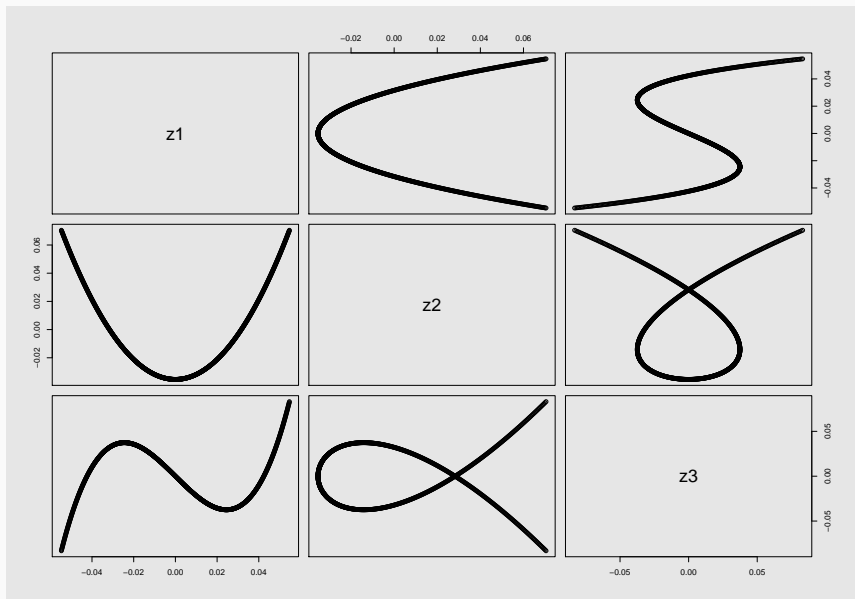
$$z_3 = a_3 + b_3x + c_3x^2 + d_3x^3$$

where coefficients such that $z_i'z_j = 0$ when $i \neq j$.

```
fit <- lm(y ~ poly(x,3))
```

- Because z_i involve constants a_1, a_2, \dots , the intercept will be affected.
- Coefficients of lower order terms are unchanged when higher order terms are added.

Orthogonal polynomials



```
plmodi <- lm(undercount ~ poly(cperAA,4)+  
             cpergore*usage+equip, gavote)  
summary(plmodi)
```

- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model
- 4 Regression diagnostics
- 5 Variable selection**
- 6 Conclusion
- 7 Assignment One

$$\text{AIC} = -2 \text{ maximum log likelihood} + 2p$$

- step will minimize AIC using backwards selection
- R330: : allpossregs will fit all possible regressions
- Do not use coefficient t-tests for variable selection
- Hierarchy principle: don't eliminate lower terms while retaining higher order terms.
- If explanation is the goal, you may not rely on completely automated variable selection methods.

```
biglm <- lm(undercount ~
            (equip+econ+usage+atlanta)^2+
            (equip+econ+usage+atlanta)*(perAA+pergore),
            gavote)
smallm <- step(biglm,trace=FALSE)
drop1(smallm,test="F")
finalm <- lm(undercount~equip + econ + perAA +
            equip:econ + equip:perAA, gavote)
summary(finalm)
```

- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model
- 4 Regression diagnostics
- 5 Variable selection
- 6 Conclusion
- 7 Assignment One

What is your data analysis like?

- An initial data analysis that explores the numerical and graphical characteristics of the data.
- Variable selection to choose the best model.
- An exploration of transformations to improve the fit of the model.
- Diagnostics to check the assumptions of your model.
- Some predictions of future observations for interesting values of the predictors.
- An interpretation of the meaning of the model with respect to the particular area of application.

- 1 Linear models
- 2 Looking at data
- 3 Fiting a linear model
- 4 Regression diagnostics
- 5 Variable selection
- 6 Conclusion
- 7 Assignment One

The exercise on Page 24 of our textbook.

- Choose any one of the data sets and perform linear data analysis.
- Your analysis should at least consist of:
 - an initial data analysis
 - fitting of a linear model
 - interpretation
 - testing
 - diagnostics
 - variable selection
 - conclusion