

Regression

Prof. Zhifeng Wang
Week 9, Spring 2023

Regression and Classification 一个是离散的 一个是连续的
的

Learn about...

- ▶ Linear and logistic regression.
- ▶ Regression trees.

Logistic regression 改造回归

Motivating Example

- ▶ **Regression** analysis is used to study the **relationship** between a **response** variable (y) and one or more **explanatory** variables (x_1, x_2, \dots).
 - ▶ Do people with higher **income** (x) spend more on **food** (y)? Or less?
 - ▶ Do people with higher **education** (x) receive higher **income** (y)? Or less?
- ▶ **Predict** value of a response variable (y) given the value of explanatory variables (x_1, x_2, \dots).
 - ▶ Given the house size (x), how do we **predict** the expected house price (y)?

(Simple) Linear Regression

1. ML 可以通过 LMS 去得到 β_1 和 β_2
1. ML 可以通过 SGD 去得到 β_1 和 β_2
Linear Regression 是一个 convex problem
可以用两种方法来解决

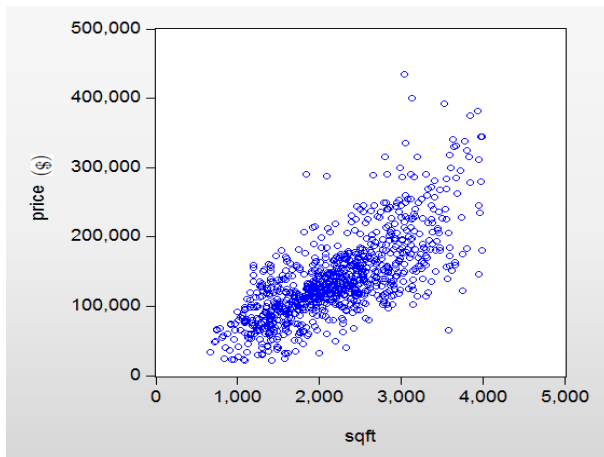
如果是非 convex 非凸 只能用 SGD 梯度下降来解决

- ▶ The simple linear (population) regression model:

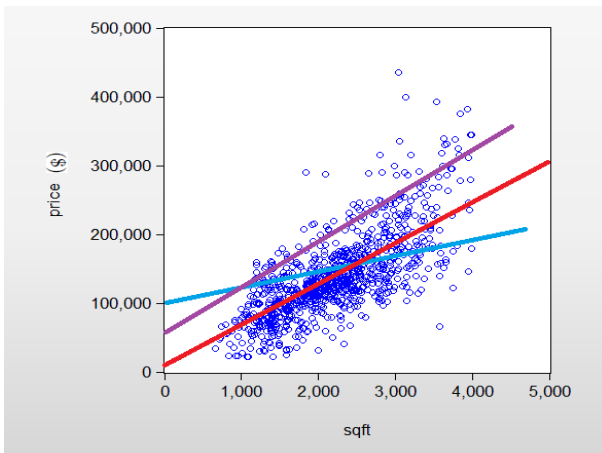
$$y = \beta_1 + \beta_2 x + e$$

- ▶ The random error e represents the unexplained part
 $e = y - (\beta_1 + \beta_2 x)$.
 - ▶ If the true relationship between y and x is linear, e can capture any unexplained variation in y .
 - ▶ e can be the effects of other variables not included in the model
 - ▶ e can be the effects of non-linearity in the relationship between y and x .
- ▶ There are two parameters β_1 and β_2 that need to be estimated given the sample data.

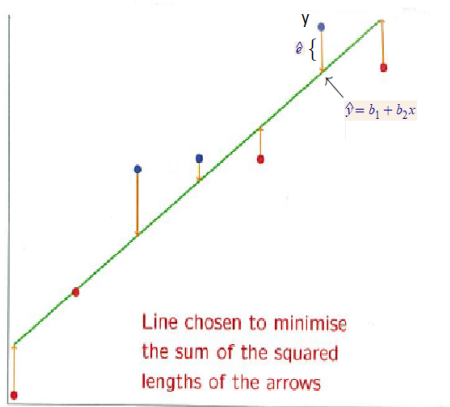
Motivating Example



Motivating Example

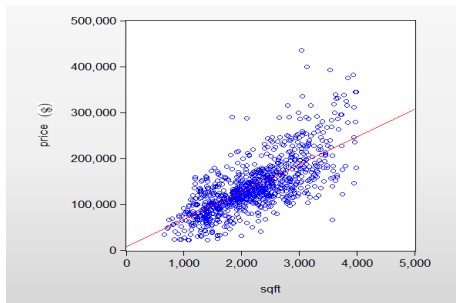


Least Square Method



- ▶ **Least squares** method choose a line to **minimise sum of square residuals**.
- ▶ The estimated model:
 $\hat{y} = \hat{\beta}_1 + \hat{\beta}_2x$
- ▶ \hat{y} is the fitted value of y .
- ▶ $\hat{\beta}_1$ and $\hat{\beta}_2$ are the **estimates** for β_1 and β_2 , respectively.
- ▶ **Residual**: $\hat{e} = y - \hat{y}$

Motivating Example



- ▶ The estimated model:

$$\widehat{price} = 7387.08 + 59.85SQFT$$

- ▶ The meaning of the **slope** coefficient is that for every extra square foot in size, house prices are **expected to increase** by around \$59.85.
- ▶ Bigger houses are expected to cost more.

Prediction

- ▶ The fitted regression equation

$$\widehat{price} = 7387.08 + 59.85SQFT$$

- ▶ We can use this to **predict** the expected price for a 3500 square feet house

$$\widehat{price} = 7387.08 + 59.85(3500) = \$216862.08.$$

Linear Regression

- ▶ Consider a data set of n observations:

$$\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$$

Usually \mathbf{x}_i consists of multiple attributes.

- ▶ Let \hat{y}_i denote the predicted (fitted) value for observation i .

Linear Regression

- ▶ *Linear regression* fits a simple equation of the form

$$\hat{y}_i = \mathbf{x}_i \hat{\boldsymbol{\beta}} \equiv \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \cdots + \hat{\beta}_p x_{i,p}$$

where \hat{y}_i denotes the predicted target variable and $\mathbf{x}_i = [1, x_{i,1}, \dots, x_{i,p}]$ denote the explanatory attributes, with $\hat{\boldsymbol{\beta}} = [\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p]^\top$.

- ▶ Note the “special” element for the intercept.
- ▶ For notational convenience, \mathbf{x}_i s and y_i s are often “stacked”:

$$\begin{aligned} \text{▶ } \mathbf{y} &= \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \text{ and } \hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} \\ \text{▶ } \mathbf{x} &= \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} \text{ so } \mathbf{x} \hat{\boldsymbol{\beta}} = \begin{bmatrix} \mathbf{x}_1 \hat{\boldsymbol{\beta}} \\ \vdots \\ \mathbf{x}_n \hat{\boldsymbol{\beta}} \end{bmatrix} \end{aligned}$$

$$\Rightarrow \hat{\mathbf{y}} = \mathbf{x} \hat{\boldsymbol{\beta}}$$

Regression

从几何上不是点到线的距离，SLM 是这个关注的是 y 的差值

- ▶ We estimate the parameters $\hat{\beta}$ via Least Squares (minimising sum of squared residuals):
- ▶ Very fast (compared to other methods); has a unique solution (if columns of \mathbf{x} are not “redundant”). If any columns of \mathbf{x} are “redundant”, in that they are a linear function of other columns. Then, regression can't be fit.
- ▶ Interpretable: β_k is the predicted effect of a unit change in x_k on the predicted value of Y , given other predictors fixed or constant.

- ▶ `lm()` is the workhorse function for fitting Linear Models.

Transforming x

- ▶ Suppose the effect of x on \hat{y} is believed to be nonlinear. Does it mean our model can't be linear?

No! \hat{y} only needs to be linear *in the parameters β* !

- ▶ $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \beta_2 x_i^2$ is still an Linear Model.
 - ▶ We don't have to stop at quadratic effects, but too high a power can be overfitting and unstable.
 - ▶ In `lm()`, powers need to be enclosed in `I()` (e.g., `I(x^2)`), or R will process them differently.
- ▶ Other nonlinear transformation of x also possible.
 - ▶ E.g., if x is right-skewed, taking $\sqrt{}$ or `log` can work better.

Categorical Predictors

- ▶ Categorical predictors are represented using *dummy* or *indicator* variables.

1. One category (level) is set as the baseline.

⇒ Categorical variable with I levels or categories, requires $I - 1$ dummy or indicator variables.

- ▶ R automatically does this for factor variables.
 - ▶ Beware categorical variables coded as numbers! “Process” them with `factor()` to let R know.
 - ▶ An ordinal factor can be identified using the function `ordered()`.
- ▶ You can create a dummy variable explicitly via `I(var == "val")` to add

$$x_{i,k} = \mathbb{I}\{\text{var}[i] \text{ is "val"}\} = \begin{cases} 1 & \text{if value of var for } i \text{ is "val"} \\ 0 & \text{otherwise} \end{cases}$$

Interaction

- ▶ Interaction occurs when the effect of one predictor variable depends on the level of another.

- ▶ Trivial example: gender vs. height for children.
- ▶ The increase in height depends on the gender
- ▶ Slope for boys will be higher than slope for girls.

⇒ Interaction between height and gender.

- ▶ Represented in LMs as a product between predictor variables (and indicators).

- ▶ In `lm()`, use `x1:x2` to add $\beta_{12}x_{i,1}x_{i,2}$.

- ▶ *Principle of marginality* says that if you include $x_{i,1}x_{i,2}$ in the model, you should also include $x_{i,1}$ and $x_{i,2}$.

⇒ Use `x1*x2` to add $\beta_1x_{i,1} + \beta_2x_{i,2} + \beta_{12}x_{i,1}x_{i,2}$.

Transforming y

- ▶ LMs work best when residuals are symmetric and don't have extreme outliers.
 - ▶ $\sqrt{}$ transformations are often used for y a count
 - ▶ log transformations often used for other strictly positive measurements
- ▶ Transforming y changes interpretation:
 - ▶ $\log(y) = \beta_0 + \beta_1 x$
 - ▶ a 1 unit increase in x leads to a $100\beta_1\%$ change in y .
 - ▶ $\log(y) = \beta_0 + \beta_1 \log x$
 - ▶ a 1% change in x results in a $\beta_1\%$ in y .

Performance Evaluation

Performance of a regression task can be evaluated by looking at the prediction error.

Mean Squared Error (In-Sample):

$$\text{MSE} = \frac{\text{SSE}}{n-p-1}, \text{ where } \text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- ▶ Most common; easiest to work with.
- ▶ Very sensitive to outliers.
- ▶ Often reported as $\sqrt{\text{MSE}}$, the *Root Mean Squared Error (RMSE)* which is on the same scale as the data.
- ▶ Loosely, p is the number of explanatory variables in the model for \hat{y}_i .

Mean Absolute Error(In Sample): $\text{MAE} = \frac{1}{n-p-1} \sum_{i=1}^n |y_i - \hat{y}_i|$

- ▶ More resistant to outliers.

R^2 and Adjusted R^2

- ▶ High value of R^2 means the model explains a high proportion of total variation, $0 \leq R^2 \leq 1$.
- ▶ More complexity in a model almost always increases $R^2 \implies$ overfitting.

\implies *Adjusted R^2*

$$R_{\text{adj}}^2 = R^2 - (1 - R^2) \frac{p}{n - p - 1}$$

is more directly comparable across linear models of differing complexity: where p is number of explanatory variables.

Example: Iris data

- ▶ What if we wanted to predict petal length from species?

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \hat{\beta}_2 x_{i,2}$$

where

$$x_{i,1} = \begin{cases} 1 & \text{if species is versicolor} \\ 0 & \text{otherwise} \end{cases}, \quad x_{i,2} = \begin{cases} 1 & \text{if species is virginica} \\ 0 & \text{otherwise} \end{cases}$$

- ▶ Then,

$\hat{\beta}_0$ is the predicted mean for setosa

$\hat{\beta}_1$ is how much higher the predicted mean for versicolor is than that for setosa

$\hat{\beta}_2$ is how much higher the predicted mean for virginica is than that for setosa

Example: Iris data

```
data(iris)
summary(lm(Petal.Length ~ Species, data = iris))
```

```
## Call:
## lm(formula = Petal.Length ~ Species, data = iris)
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.260 -0.258  0.038  0.240  1.348
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      1.4620     0.0609   24.0   <2e-16 ***
## Speciesversicolor  2.7980     0.0861   32.5   <2e-16 ***
## Speciesvirginica   4.0900     0.0861   47.5   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Residual standard error: 0.43 on 147 degrees of freedom
## Multiple R-squared:  0.941, Adjusted R-squared:  0.941
## F-statistic: 1.18e+03 on 2 and 147 DF, p-value: <2e-16
```

Regression Output

- ▶ The standard errors ('Std. Error') show how accurately the regression coefficients have been estimated given the sample size; larger values indicate less accuracy.
- ▶ The asterisks (*) indicate which attributes have a statistically significant effect upon the response, after controlling for other predictors in the model.
- ▶ The 'Multiple R-squared' (R^2) value shows the proportion of variation in the response which is collectively explained by the explanatory attributes.
- ▶ As R^2 continues to increase as more variables are included in a regression equation, 'Adjusted R-squared' involves a penalty for the number of predictors.

Interpretation

$$\hat{y}_i = 1.4620 + 2.7980\mathbb{I}\{i \text{ is versicolor}\} + 4.0900\mathbb{I}\{i \text{ is virginica}\}$$

*** all three β_k s are highly *statistically significant*:

β_0 there is enough evidence to believe that population mean petal length for setosa is different (higher) from 0 (a trivial statement); it's about 1.4620

β_1 there is enough evidence to believe that population mean petal length for versicolor is different (higher) from that of setosa (the baseline); it's about 4.2600

β_2 there is enough evidence to believe that population mean petal length for virginica is different (higher) from that of setosa (the baseline); it's about 5.5520

R^2 : The model explains about 94.1% of the variation in the data, and the value of adjusted R^2 is also about 94.1%.

Model selection

- ▶ When there are many potential predictor variables and interaction terms, prediction performance for future data will often deteriorate if a very complex model is fitted.

Stepwise regression aims to select the most important terms for inclusion in the final model:

Forward selection: Start with the minimal model, and add one at a time. Stop when nothing can be added to improve the criterion.

Backwards elimination: Start with the maximal model, and remove one at a time. Stop when nothing can be removed to improve the criterion.

Bidirectional elimination: Start with some initial model, and try to add or remove one at a time. Stop when nothing can be changed to improve the criterion.

All-subsets regression: Try every single possible combination of terms. Takes a very long time!

Common criteria

R^2_{adj} : $1 - (1 - R^2) \times (n - 1) / (n - p - 1)$ (bigger is better)

Mallows C_p : $\text{SSE} / \hat{\sigma}^2_{\text{max}} - (n - 2p - 2)$, where

$\hat{\sigma}^2_{\text{max}} = \text{SSE}_{\text{max}} / (n - p_{\text{max}} - 1)$ (smaller = better)

Akaike Information Criterion (AIC): (smaller is better)

Bayesian Information Criterion (BIC): (smaller is better)

Correlation

- ▶ Another commonly used performance measure of a numeric prediction model is the *correlation* R between observed and predicted response.

$$R = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

- ▶ Ideally, plot of predicted versus actual response should almost straight with positive slope.
- ▶ So R close to 1 means good prediction.

Example: Swiss Fertility data

```
library(datasets)
data(swiss)
```

- Data about 47 French-speaking provinces of Switzerland around 1888.

Fertility standardised fertility measure

Agriculture % of males involved in agriculture as occupation

Examination % of draftees receiving highest mark on army exam

Education % with education beyond primary school for draftees

Catholic % Catholic (as opposed to Protestant)

Infant.Mortality % live births who lives less than a year

Regression Output

```
summary(swiss.fit <- lm(Fertility ~ ., data = swiss))
```

```
## Call:
## lm(formula = Fertility ~ ., data = swiss)
## Residuals:
##      Min       1Q   Median       3Q      Max
## -15.274  -5.262   0.503   4.120  15.321
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    66.9152    10.7060     6.25 1.9e-07 ***
## Agriculture    -0.1721     0.0703    -2.45 0.0187 *
## Examination    -0.2580     0.2539    -1.02 0.3155
## Education      -0.8709     0.1830    -4.76 2.4e-05 ***
## Catholic        0.1041     0.0353     2.95 0.0052 **
## Infant.Mortality 1.0770     0.3817     2.82 0.0073 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Residual standard error: 7.17 on 41 degrees of freedom
## Multiple R-squared:  0.707, Adjusted R-squared:  0.671
## F-statistic: 19.8 on 5 and 41 DF, p-value: 5.59e-10
```

Regression Equation

The fitted regression equation is

$$\widehat{Fertility} = 66.915 - 0.172 \text{ Agriculture} \\ - 0.258 \text{ Examination} - 0.871 \text{ Education} \\ + 0.104 \text{ Catholic} + 1.077 \text{ Infant.Mortality}$$

- ▶ E.g., for every additional percentage point of draftees with education beyond primary school, the predicted fertility measure decreases by 0.8709 units, keeping other predictors fixed.

Stepwise regression

```
step(swiss.fit, data = swiss)
```

```
## Start: AIC=190.7
## Fertility ~ Agriculture + Examination + Education + Catholic +
## Infant.Mortality
##           Df Sum of Sq  RSS AIC
## - Examination      1      53 2158 190
## <none>                2105 191
## - Agriculture      1     308 2413 195
## - Infant.Mortality  1     409 2514 197
## - Catholic         1     448 2553 198
## - Education        1    1163 3268 209
## Step: AIC=189.9
## Fertility ~ Agriculture + Education + Catholic + Infant.Mortality
##           Df Sum of Sq  RSS AIC
## <none>                2158 190
## - Agriculture      1     264 2422 193
## - Infant.Mortality  1     410 2568 196
## - Catholic         1     957 3115 205
## - Education        1    2250 4408 221
## Call:
## lm(formula = Fertility ~ Agriculture + Education + Catholic +
## Infant.Mortality, data = swiss)
## Coefficients:
## (Intercept)      Agriculture      Education      Catholic
##      62.101         -0.155         -0.980          0.125
## Infant.Mortality
##      1.078
```

All-Subsets regression: regsubsets() in leaps

- Gives best model for each number of predictor.

```
library(leaps)
regsub <- regsubsets(Fertility ~ ., data = swiss)
summary(regsub)
```

Subset selection object
Call: regsubsets.formula(Fertility ~ ., data = swiss)
5 Variables (and intercept)
##

		Forced in	Forced out
## Agriculture		FALSE	FALSE
## Examination		FALSE	FALSE
## Education		FALSE	FALSE
## Catholic		FALSE	FALSE
## Infant.Mortality		FALSE	FALSE

1 subsets of each size up to 5
Selection Algorithm: exhaustive

##		Agriculture	Examination	Education	Catholic	Infant.Mortality
## 1	(1)	" "	" "	"*"	" "	" "
## 2	(1)	" "	" "	"*"	"*"	" "
## 3	(1)	" "	" "	"*"	"*"	"*"
## 4	(1)	"*"	" "	"*"	"*"	"*"
## 5	(1)	"*"	"*"	"*"	"*"	"*"

Selecting best overall model

```
summary(regsub)$cp # Mallow's cp
## [1] 35.205 18.486 8.178 5.033 6.000

with(summary(regsub), which[which.min(cp), ])

##      (Intercept)      Agriculture      Examination      Education
##              TRUE              TRUE              FALSE              TRUE
##      Catholic Infant.Mortality
##              TRUE              TRUE
```

- I.e., best Mallows C_P measure is for 4 predictors, which are Agriculture, Education, Catholicism, and Infant Mortality.

Interactions

```
summary(swiss.fit <- lm(Fertility ~ (Agriculture +  
  Examination + Education + Catholic + Infant.Mortality)^2,  
  data = swiss))
```

```
## Call:  
## lm(formula = Fertility ~ (Agriculture + Examination + Education +  
##   Catholic + Infant.Mortality)^2, data = swiss)  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -8.76  -3.89  -0.68   3.14  14.10   
## Coefficients:  
##                                Estimate Std. Error t value Pr(>|t|)      
## (Intercept)                   253.97615    67.99721   3.74  0.00076 ***  
## Agriculture                   -2.10867     0.70163  -3.01  0.00522 **  
## Examination                   -5.58074     2.75010  -2.03  0.05109 .  
## Education                     -3.47089     2.68377  -1.29  0.20547   
## Catholic                      -0.17693     0.40653  -0.44  0.66642   
## Infant.Mortality              -5.95748     3.08963  -1.93  0.06303 .  
## Agriculture:Examination        0.02137     0.01377   1.55  0.13091   
## Agriculture:Education          0.01906     0.01523   1.25  0.22009   
## Agriculture:Catholic           0.00263     0.00285   0.92  0.36387   
## Agriculture:Infant.Mortality   0.06370     0.02981   2.14  0.04060 *  
## Examination:Education          0.07517     0.03634   2.07  0.04703 *  
## Examination:Catholic          -0.00153     0.01079  -0.14  0.88791   
## Examination:Infant.Mortality   0.17101     0.12907   1.33  0.19485   
## Education:Catholic            -0.00713     0.01018  -0.70  0.48865   
## Education:Infant.Mortality     0.03359     0.12420   0.27  0.78863   
## Catholic:Infant.Mortality       0.00992     0.01617   0.61  0.54409   
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
## Residual standard error: 6.47 on 31 degrees of freedom  
## Multiple R-squared:  0.819, Adjusted R-squared:  0.731  
## F-statistic: 9.35 on 15 and 31 DF, p-value: 1.08e-07
```

Stepwise selection of interactions

```
swiss.fit2 <- lm(Fertility ~ (Agriculture + Examination +  
  Education + Catholic + Infant.Mortality)^2, data = swiss)  
step(swiss.fit2, data = swiss)
```

```
## Start: AIC=188  
## Fertility ~ (Agriculture + Examination + Education + Catholic +  
##   Infant.Mortality)^2  
##
```

	Df	Sum of Sq	RSS	AIC
## - Examination:Catholic	1	0.8	1300	186
## - Education:Infant.Mortality	1	3.1	1302	186
## - Catholic:Infant.Mortality	1	15.8	1315	187
## - Education:Catholic	1	20.6	1320	187
## - Agriculture:Catholic	1	35.6	1335	187
## <none>			1299	188
## - Agriculture:Education	1	65.6	1365	188
## - Examination:Infant.Mortality	1	73.6	1373	189
## - Agriculture:Examination	1	100.9	1400	190
## - Examination:Education	1	179.3	1478	192
## - Agriculture:Infant.Mortality	1	191.4	1491	192

Note that removal of Examination:Catholic has the lowest AIC (removed from model).

Stepwise selection of interactions (continued)

```
## Step: AIC=186
## Fertility ~ Agriculture + Examination + Education + Catholic +
## Infant.Mortality + Agriculture:Examination + Agriculture:Education +
## Agriculture:Catholic + Agriculture:Infant.Mortality + Examination:Educa
## tion +
## Examination:Infant.Mortality + Education:Catholic + Education:Infant.Mo
## rtality +
## Catholic:Infant.Mortality
##           Df Sum of Sq  RSS AIC
## - Education:Infant.Mortality  1      3.9 1304 184
## - Catholic:Infant.Mortality  1     17.3 1317 185
## - Agriculture:Catholic       1     37.1 1337 185
## <none>                        1300 186
## - Education:Catholic        1     56.8 1357 186
## - Agriculture:Education      1     69.5 1369 186
## - Examination:Infant.Mortality 1     86.0 1386 187
## - Agriculture:Examination    1    114.3 1414 188
## - Examination:Education      1    178.4 1478 190
## - Agriculture:Infant.Mortality 1    205.3 1505 191
```

Note that removal of `Education:Infant.Mortality` has the lowest AIC (removed from model).

Stepwise selection of interactions (continued 2)

```
## Step: AIC=184.2
## Fertility ~ Agriculture + Examination + Education + Catholic +
## Infant.Mortality + Agriculture:Examination + Agriculture:Education +
## Agriculture:Catholic + Agriculture:Infant.Mortality + Examination:Educa
## tion +
## Examination:Infant.Mortality + Education:Catholic + Catholic:Infant.Mor
## tality
##
## Df Sum of Sq RSS AIC
## - Catholic:Infant.Mortality 1 25.8 1330 183
## - Agriculture:Catholic 1 36.4 1340 184
## <none> 1304 184
## - Agriculture:Education 1 79.2 1383 185
## - Education:Catholic 1 79.3 1383 185
## - Agriculture:Examination 1 116.3 1420 186
## - Examination:Education 1 185.9 1490 188
## - Agriculture:Infant.Mortality 1 219.8 1524 190
## - Examination:Infant.Mortality 1 230.5 1534 190
```

Note that removal of Catholic:Infant.Mortality has the lowest AIC (removed from model).

Stepwise selection of interactions (continued 3)

```
## Step: AIC=183.1
## Fertility ~ Agriculture + Examination + Education + Catholic +
## Infant.Mortality + Agriculture:Examination + Agriculture:Education +
## Agriculture:Catholic + Agriculture:Infant.Mortality + Examination:Educa
## tion +
## Examination:Infant.Mortality + Education:Catholic
## Df Sum of Sq RSS AIC
## - Agriculture:Catholic 1 26.7 1356 182
## <none> 1330 183
## - Education:Catholic 1 91.7 1421 184
## - Agriculture:Education 1 92.2 1422 184
## - Agriculture:Examination 1 121.2 1451 185
## - Examination:Education 1 197.2 1527 188
## - Examination:Infant.Mortality 1 210.7 1540 188
## - Agriculture:Infant.Mortality 1 220.4 1550 188
## Step: AIC=182
## Fertility ~ Agriculture + Examination + Education + Catholic +
## Infant.Mortality + Agriculture:Examination + Agriculture:Education +
## Agriculture:Infant.Mortality + Examination:Education + Examination:Infan
## nt.Mortality +
## Education:Catholic
## Df Sum of Sq RSS AIC
## <none> 1356 182
## - Agriculture:Education 1 75.0 1431 183
## - Agriculture:Examination 1 99.7 1456 183
## - Examination:Education 1 174.6 1531 186
## - Education:Catholic 1 216.6 1573 187
## - Agriculture:Infant.Mortality 1 271.1 1627 189
## - Examination:Infant.Mortality 1 272.9 1629 189
```

Stepwise selection of interactions (continued 4)

```
## Call:
## lm(formula = Fertility ~ Agriculture + Examination + Education +
##      Catholic + Infant.Mortality + Agriculture:Examination + Agriculture:Edu
##      cation +
##      Agriculture:Infant.Mortality + Examination:Education + Examination:Infa
##      nt.Mortality +
##      Education:Catholic, data = swiss)
## Coefficients:
##              (Intercept)              Agriculture
##              225.9101                -1.9067
##              Examination              Education
##              -5.1202                -2.4735
##              Catholic              Infant.Mortality
##              0.2112                -5.2693
##      Agriculture:Examination      Agriculture:Education
##              0.0149                0.0191
##      Agriculture:Infant.Mortality      Examination:Education
##              0.0635                0.0639
##      Examination:Infant.Mortality      Education:Catholic
##              0.1722                -0.0124
```

Final fit with interactions

```
summary(swiss.fit2.steps)

##
## Call:
## lm(formula = Fertility ~ Agriculture + Examination + Education +
##      Catholic + Infant.Mortality + Agriculture:Examination + Agriculture:Education +
##      Agriculture:Infant.Mortality + Examination:Education + Examination:Infant.Mortality +
##      Education:Catholic, data = swiss)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.608 -3.665 -0.564  2.922 13.736
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    225.91005     52.45757    4.31  0.00013 ***
## Agriculture     -1.90665     0.56282   -3.39  0.00176 **
## Examination     -5.12020     1.57831   -3.24  0.00259 **
## Education       -2.47350     1.20277   -2.06  0.04725 *
## Catholic         0.21116     0.05418    3.90  0.00042 ***
## Infant.Mortality -5.26935     2.28727   -2.30  0.02729 *
## Agriculture:Examination  0.01488     0.00928    1.60  0.11771
## Agriculture:Education   0.01908     0.01372    1.39  0.17301
## Agriculture:Infant.Mortality 0.06353     0.02402    2.64  0.01216 *
## Examination:Education   0.06389     0.03010    2.12  0.04092 *
## Examination:Infant.Mortality 0.17219     0.06489    2.65  0.01189 *
## Education:Catholic     -0.01238     0.00524   -2.36  0.02374 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 6.23 on 35 degrees of freedom
## Multiple R-squared:  0.811, Adjusted R-squared:  0.752
## F-statistic: 13.7 on 11 and 35 DF, p-value: 1.28e-09
```

Linear Regression for Big Data

- ▶ Computational issues arise for “big” data sets, either in the sense of many rows (observations) or many columns (attributes).
 - ▶ For many rows, the data may need to be read in chunks, and stored economically.
 - ▶ For many columns, solution may be slow and numerically unstable despite being non-iterative.
 - ▶ For many columns, attribute selection is particularly important but backwards stepwise regression may be infeasible.

Linear Regression for Big Data

- ▶ The `biglm` function extends the capabilities of `lm` for linear regression.
- ▶ Data can be read in chunks, and fitted models can be updated with additional data by the `update` function.
- ▶ The `bigglm` function extends the capabilities of `glm` for generalised linear models.

Logistic Regression

- ▶ Designed when you have binary response.
- ▶ Logistic regression involves fitting an equation of the form

$$\Pr(Y_i = 1) = \text{squash}(\beta_0 + \beta_1 x_{i,1} + \cdots + \beta_i x_{i,p})$$

where $\text{squash}(x) = 1/(1 + e^{-x})$

- ▶ Statisticians call it the *logistic* function:

$$\text{logit}\{\Pr(Y_i = 1)\} = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_i x_{i,p},$$

where $\text{logit}(q) = \log \frac{q}{1-q}$, the log of the odds associated with probability q .

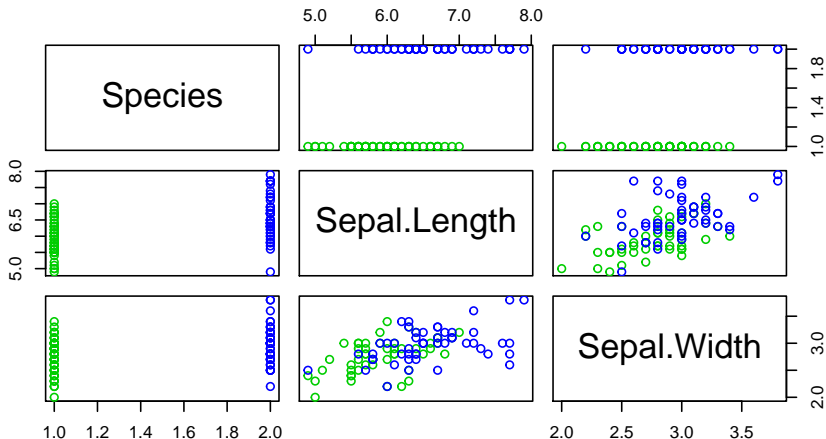
Logistic

- ▶ Implemented in R by `glm`.
 - ▶ To fit logistic regression, specify
`family = binomial("logit")`

Example: Iris Data

- Recall the scenario from the SVM lecture.

```
iris2 <- transform(subset(iris, Species != "setosa",  
  c("Species", "Sepal.Length", "Sepal.Width")),  
  Species = factor(Species))
```



Logistic Regression

```
summary(glm(I(Species == "virginica") ~ ., data = iris2,  
            family = binomial("logit")))
```

```
## Call:  
## glm(formula = I(Species == "virginica") ~ ., family = binomial("logit"),  
##      data = iris2)  
## Deviance Residuals:  
##      Min       1Q   Median       3Q      Max   
## -1.874  -0.895  -0.055   0.961   2.357   
## Coefficients:  
##              Estimate Std. Error z value Pr(>|z|)      
## (Intercept)   -13.046     3.097   -4.21  2.5e-05 ***  
## Sepal.Length    1.902     0.517    3.68  0.00023 ***  
## Sepal.Width     0.405     0.863    0.47  0.63908   
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
## (Dispersion parameter for binomial family taken to be 1)  
##      Null deviance: 138.63  on 99  degrees of freedom  
## Residual deviance: 110.33  on 97  degrees of freedom  
## AIC: 116.3  
## Number of Fisher Scoring iterations: 4
```

Interpretation

- ▶ Only sepal length is significant.
- ▶ *In the presence of sepal length*, sepal width is not.
- ▶ For every unit increase in sepal length, the predicted odds of it being a virginica are multiplied by $e^{1.9024} = 6.7018$.
- ▶ Standard functions (like `predict()`) are available. However, you must specify the type:
 - ▶ By default, predicts $\text{logit}\{\hat{\text{Pr}}(Y_{\text{new}} = 1)\}$.
 - ▶ Specify `type="response"` to predict $\hat{\text{Pr}}(Y_{\text{new}} = 1)$.

Regression Trees

- ▶ The same ideas can be applied to regression trees as in the classification trees
- ▶ A regression tree is similar to a decision tree, except that the predicted value at a terminal leaf is given by the mean or median response variable of observations allocated to that leaf.
- ▶ Splits can be chosen to minimise sum of squared errors rather than Entropy measures.

Example: Swiss fertility data

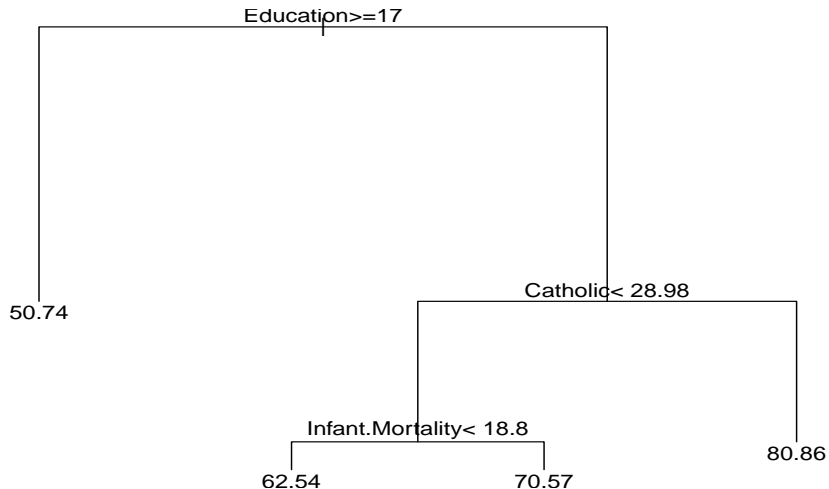
- ▶ When given a quantitative response variable, **rpart** and others automatically change to regression tree mode:

```
library(rpart)
(swiss.tree <- rpart(Fertility ~ ., data = swiss))

## n= 47
##
## node), split, n, deviance, yval
##      * denotes terminal node
##
##  1) root 47 7178.0 70.14
##    2) Education>=17 7  628.3 50.74 *
##    3) Education< 17 40 3454.0 73.54
##      6) Catholic< 28.98 23  827.5 68.13
##        12) Infant.Mortality< 18.8 7  167.3 62.54 *
##        13) Infant.Mortality>=18.8 16  346.6 70.57 *
##      7) Catholic>=28.98 17 1042.0 80.86 *
```


Visualisation of regression trees

```
plot(swiss.tree)  
text(swiss.tree, digits = 4)
```



Performance Evaluation

Mean Squared Error (Out-Sample):

$$\text{MSE} = \frac{\text{SSE}}{n_{\text{test}}}, \text{ where } \text{SSE} = \sum_{i=1}^{n_{\text{test}}} (y_i - \hat{y}_i)^2$$

Mean Absolute Error(Out Sample): $\text{MAE} = \frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} |y_i - \hat{y}_i|$