



Intro to R Programming for Biostatistics

Adam J Sullivan

by-nc-nd

Linear Regression in R

- We can use R to easily fit linear regressions for us.
- This section will explore the basic commands for linear regression as well as how to test assumptions.
- We will not teach linear regression, but only seek to display how R does it.

Linear Regression in R

- To fit Linear Regression models in R we use the `lm()` function.

```
lm(formula, data, subset, weights, na.action,  
method = "qr", model = TRUE, x = FALSE, y = TRUE,  
singular.ok = TRUE, contrasts = NULL, offset, ...)
```

- formula is the regression equation written as $y \sim x_1 + x_2 + \dots$
- data is the dataframe of interest.
- subset specific subset of data.
- weights for weighted data.

Regression Models in R

## # A tibble: 1,794 x 6	
##	country continent year lifeexp pop gdpPerCap
## 1	Afghanistan Asia 1952 28.861 8625333 779.4653
## 2	Afghanistan Asia 1957 38.332 5248934 828.8538
## 3	Afghanistan Asia 1962 31.997 10428843 851.1887
## 4	Afghanistan Asia 1967 34.828 13577966 836.1971
## 5	Afghanistan Asia 1972 36.888 13079468 739.0811
## 6	Afghanistan Asia 1977 38.438 14888172 786.8114
## 7	Afghanistan Asia 1982 39.854 138617957 852.3959
## 8	Afghanistan Asia 1987 40.822 14617921 649.3414
## 9	Afghanistan Asia 1992 41.674 16317921 635.3414
## 10	Afghanistan Asia 1997 41.763 22227415 635.3414
## # ... with 1,684 more rows	

Gapminder Data

8/64

library(Gapminder)

Gapminder Data

7/64

Variable	Meaning
country	Country Name
continent	Continent
year	Year Data Accounts For
lifeExp	Life Expectancy at Birth
pop	Total Population
gdpPerCap	per-Capita GDP

- Per-capita GDP (Gross domestic product) is given in units of international dollars, "a hypothetical unit of currency that has the same purchasing power parity that the U.S. dollar had in the United States at a given point in time" -- 2005, in this case.

- Worldwide data source.

- Contains 6 variables

Gapminder Data

[1] "coefficients" "residuals" "effects" "rank"
[5] "fitted.values" "assign" "qt"
[9] "xlevels" "call" "terms" "model"

names(kenya_model)

- We can find out what is contained in a list by using the names() function.
- We see that this is a list

Gapminder Regression

11/64

str(kenya_model)

- Basic Model Statement does not include much information
- If we look further we can see what type of object is returned.

Gapminder Regression

10/64

kenya <- gapminder %>% filter(country=="kenya")
kenya_model <- lm(lifeexp ~ year, data=kenya)
Call:
lm(formula = lifeexp ~ year, data = kenya)
##
Coefficients:
(Intercept) 0.2865
year -356.1818

Gapminder Regression

9/64

12/64

```
## Call:
## lm(formula = l1teexp ~ year, data = kenya)
## Residuals:
## Min 1Q Median 3Q Max
## -6.3554 -3.5739 -0.2819 3.0984 5.5687
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) -356.18098 145.08537 -2.454 0.0340 *
## year 0.28651 0.07329 2.818 0.0182 *
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## Residual standard error: 4.382 on 18 degrees of freedom
## Multiple R-squared: 0.4426, Adjusted R-squared: 0.3868
## F-statistic: 7.939 on 1 and 18 DF, p-value: 0.01823
```

Gapminder Regression

```
par(mfrow=c(2,2)) # optional 4 graphs/page
plot(kenya_model)
```

Diagnostic Model Plots

```
summary(kenya_model)
```

Gapminder Regression

- We can use other commands on a regression
- For example we can use the summary() function:

- Other useful functions are listed below:
- coefficients(kenya_model) # model coefficients
- confint(kenya_model, level=0.95) # CIs for model parameters
- fitted(kenya_model) # predicted values
- residuals(kenya_model) # residuals
- anova(kenya_model) # anova table
- vcov(kenya_model) # covariance matrix for model parameters
- influence(kenya_model) # regression diagnostics

Other Regression Functions

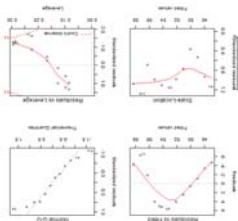
```
kenya_model$coefficients
## (Intercept) -356.1809769
## year 0.2865077
```

- Lets look at the coefficients
- We see different values that are listed here.

Gapminder Regression

```
kenya_summary <- summary(kenya_model)
names(kenya_summary)
```

Gapminder Regression



Diagnostic Model Plots

19/64

Comparing Models

- We can compare nested models using the `anova()` function.

```
kenya_model2 <- lm(lifeexp ~ year + pop, data=kenya)
anova(kenya_model1, kenya_model2)

## Analysis of Variance Table
##
## Model 1: lifeexp ~ year
## Res.Df    RSS of sum of sq    F    Pr(>F)
## 1    10 192.035
## 2     9  69.125  1 131.91 19.745 0.00015 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

20/64

broom Package: Easier to View Results

- There is a package called `broom` which makes results of regressions easier to view and compare.

- We will call this package and use the `tidy()` and `glance()` functions.

```
library(broom)
tidy(kenya_model1)
glance(kenya_model1)
```

21/64

broom Package: Easier to View Results

- We can also compare multiple models at the same time

- Using the commands we learned in data cleaning:

```
tidy1 <- tidy(kenya_model1)
tidy2 <- tidy(kenya_model2)
bind_rows(tidy1, tidy2)
```

23/64

broom Package: Easier to View Results

- We can also compare multiple models at the same time

- Using the commands we learned in data cleaning:

```
## term estimate std.error statistic p.value
## 1 (Intercept) -3.561810e+02 1.650854e+02 -2.45424 0.0399941
## 2 year 2.866977e-01 7.329116e-02 3.917614 0.00034640
## 3 (Intercept) -2.386902e+03 4.606624e+02 -5.137518 0.0005369503
## 4 year 1.250090e+00 2.387936e-01 5.24034 0.0005383790
## 5 pop -1.918461e-06 4.317401e-07 -4.44355 0.001515708
```

24/64

```
library(MASS)
fit <- lm(y~x1+x2+x3,data=mydata)
step <- stepAIC(fit, direction="both")
step$anova # display results
```

Variable Selection: Stepwise Regression

```
## r.squared adj.r.squared sigma statistic p.value df logLik
## 1 0.4425573 0.3868130 4.382174 7.939063 0.01824540 2 -33.66387
## 2 0.8254669 0.7868118 2.584812 21.283081 0.0003876588 3 -26.69040
## 3 1.73.32775 74.78247 192.03450 10
## 2 61.39280 63.33243 60.12522 9
```

brglm Package: Easier to View Results

```
glance2 <- glance(kenya_model2)
glance2 <- glance(kenya_model2)
blind_rows(glance2, glance2)
```

brglm Package: Easier to View Results

```
##
## No studentized residuals with Bonferroni p < 0.05
## Largest |rstudent|:
## student unadjusted p-value Bonferroni p
## 11 -2.582753 0.032477 0.38972
outlierest(kenya_model2) # Bonferroni p-value for most extreme obs
```

More Detailed Regression Diagnostics: Outliers

```
library(car)
outlierest(kenya_model2) # Bonferroni p-value for most extreme obs
qqplot(kenya_model2, main="QQ plot") qq plot for studentized resid
leveragePlots(kenya_model2) # leverage plots
```

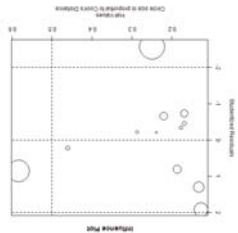
More Detailed Regression Diagnostics: Outliers

- We can see more regression diagnostics using the car package
- With this package we have the following functions

More Detailed Regression Diagnostics

```
# influential observations
# added variable plots
# Cook's D plot
# identify D values > 4/(n-k-1)
cutoff <- 4/((nrow(kenya_model2)$coefficients)-2)
plot(kenya_model2, which=1, cook.levels=cutoff)
# InfluencePlot(kenya_model2, id.method="identify", main="Influence Plot", sub="Circle size is proportional to Cook's D")
```

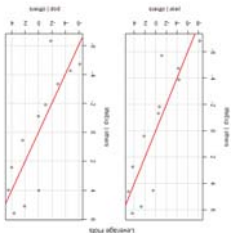
More Detailed Regression Diagnostics: Influential Observations



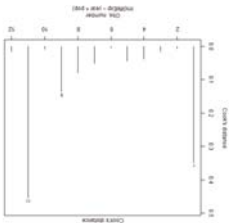
More Detailed Regression Diagnostics: Influential Observations

```
leveragePlots(kenya_model2) # leverage plots
```

More Detailed Regression Diagnostics: Outliers

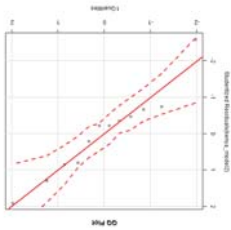


More Detailed Regression Diagnostics: Influential Observations

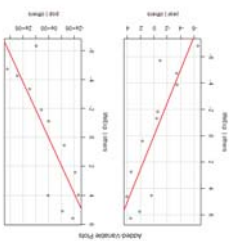


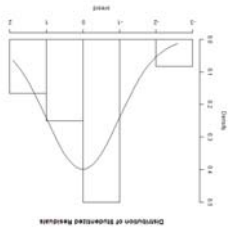
```
qqPlot(kenya_model2, main="QQ Plot") # qq plot for studentized resid
```

More Detailed Regression Diagnostics: Outliers

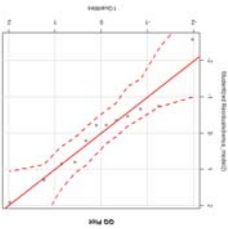


More Detailed Regression Diagnostics: Influential Observations





More Detailed Regression Diagnostics: Normality



More Detailed Regression Diagnostics: Normality

```
# Normality of residuals
# qq plot for studentized resid
qqPlot(kenya_model2, main="QQ plot")
# distribution of studentized residuals
library(MASS)
resid <- student(kenya_model2)
hist(resid, freq=F, ASE,
      main="Distribution of Studentized Residuals")
#fit <- lm(y ~ x)
#fit <- lm(y ~ x)
#fit <- lm(y ~ x)
```

More Detailed Regression Diagnostics: Normality

# Evaluate Collinearity vif(kenya_model2) # variance inflation factors	## year pop ## 30.51456 30.51456	## TRUE TRUE ## TRUE TRUE
sqrt(vif(kenya_model2)) > 2 # problems?		

More Detailed Regression Diagnostics: Multi-collinearity

Non-constant Variance Score Test ## Variance formula: ~ fitted.values ## ChiSquare = 1.111634 Df = 1 p = 0.2917271
--

More Detailed Regression Diagnostics: Error Variance

Evaluate homoscedasticity # non-constant error variance test # plot studentized residuals vs. fitted values spreadvslevelPlot(kenya_model2)
--

More Detailed Regression Diagnostics: Error Variance

Generalized Linear Models in R

- Generalized linear models are fit using the `glm()` function.
- The form of the glm function is

```
glm(formula, family=familyType(link=linkFunction), data=)
```

Family	Default Link Function
binomial	$\text{link} = \text{logit}$
poisson	$\text{link} = \text{log}$
inverse.gaussian	$\text{link} = 1/\text{mu}^2$
Gamma	$\text{link} = \text{inverse}$
quasibinomial	$\text{link} = \text{logit}$
quasi	$\text{link} = \text{identity}$, $\text{variance} = \text{constant}$
quasipoisson	$\text{link} = \text{log}$

Link Functions

Fitting Logistic Regression

```
fit <- glm(r~x1+x2+x3, data=mydata, family=binomial)
summary(fit) # display results
coef(fit) # 95% CI for the coefficients
exp(coef(fit)) # exponentiated coefficients
exp(confint(fit)) # 95% CI for exponentiated coefficients
predict(fit, type="response") # predicted values
residuals(fit, type="deviance") # residuals
```

Logistic Regression

- Logistic regression is useful when you are predicting a binary outcome from a set of continuous predictor variables.
- It is frequently preferred over discriminant function analysis because of its less restrictive assumptions.

Further GLM Help

- See `help(glm)` for other modeling options.
- See `help(family)` for other allowable link functions for each family.
- Two subtypes of generalized linear models will be covered here:
 - logistic regression
 - poisson regression

```
library(ResourceSelection)
hoslem.test(safesdead, fitted(mod.back.auc0), g=10)
```

Calibration: Hosmer-Lemeshew Test

- We usually determine the goodness of fit for logistic regression based on 1. **Calibration** - A model is well *calibrated* if the observed and predicted probabilities based on the model are reasonably close.
- 2. **Discrimination** - A model has good *discrimination* if the distribution of risk scores for cases and controls separate out. a. This means cases tend to have higher scores. b. This means controls tend to have lower scores. c. There is little overlap.

Testing Logistic Regression Models

- You can use `anova(fit1, fit2, test="Chisq")` to compare nested models.
- Additionally, `cplot(r~x, data=mydata)` will display the conditional density plot of the binary outcome `F` on the continuous `x` variable.

Comparing Nested Logistic Models

```
library(ggplot2)
library(ROCR)

prob <- predict(model)
pred <- predict(prob, data$F)
# I know, the following code is bizarre, just go with it.
perf <- performance(pred, "fpr", "fpr")
auc <- performance(pred, measure = "auc")
auc <- auc@y.values[[1]]
```

Discrimination: ROC Curve

- The **C Statistic** is given by $P(\hat{p}_1 > \hat{p}_2)$
- if the risk prediction is worthless we find that $C = 0.5$ or essentially the same as flipping a coin.
- if the risk is larger for all who are dead than all who are not dead than we have $C = 1$.
- We typically find this value with a Receiver Operating Characteristic (ROC) curve.

Discrimination: C-Statistic

- We then assess discrimination.
- To do this we use something called **Concordance** or **C Statistic**
- To understand what this is consider 2 different subjects
- 1. Subject 1 is dead
- 2. Subject 2 is not dead
- If we consider our model from above it predicts:
- 1. \hat{p}_1 the probability that subject 1 is dead.
- 2. \hat{p}_2 the probability that subject 2 is dead.

Discrimination: C-Statistic

Discrimination: ROC Curve

Graph .

[illegible]

Poisson Regression with Overdispersion

Poisson Regression

- Poisson regression is useful when predicting an outcome variable representing counts from a set of continuous predictor variables.

Poisson Regression in R

```
# where count is a count and
# x1-x3 are continuous predictors
fit <- glm(count ~ x1+x2+x3, data=mydata, family=poisson())
summary(fit) #display results
```

63/64

62/64

61/64

19/79