

Linear Models in R, version 2

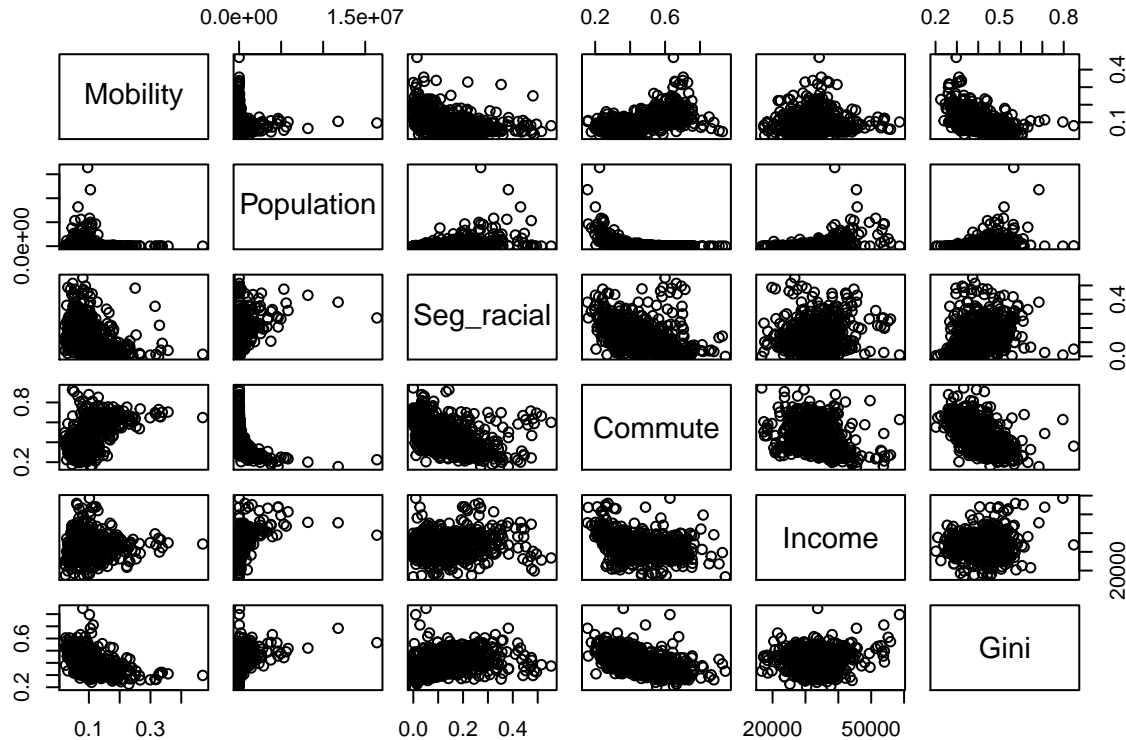
DJM, Revised: NAK

January 17, 2018

Linear models (Finally~!)

- R has lots of functions for working with different sorts of predictive models.
- We should review how they work with `lm`, and how they generalize to other sorts of models.
- We'll use the **Mobility** data from the book website:

```
mob <- read.csv("http://www.stat.cmu.edu/~cshalizi/uADA/15/hw/01/mobility.csv")
attach(mob)
mob <- data.frame(Mobility, Population, Seg_racial, Commute, Income, Gini)
pairs(mob)
```



Estimation Functions and Formulas

- To estimate a linear model in R: you use `lm`.

```
mob.lm1 <- lm(mob$Mobility ~ mob$Population + mob$Seg_racial + mob$Commute + mob$Income + mob$Gini)
```

- What `lm` returns is a complex object containing the estimated coefficients, the fitted values, a lot of diagnostic statistics, and a lot of information about exactly what work R did to do the estimation. We will come back to some of this later.

- The thing to focus on for now is the argument to `lm` in the line of code above, which tells the function exactly what model to estimate + it **specifies** the model. The R jargon term for that sort of specification is that it is the **formula** of the model.

The data argument

- While the line of code above works, it's not very elegant, because we have to keep typing `mob$` over and over.
- More abstractly, it runs specifying which variables we want to use (and how we want to use them) together with telling R where to look up the variables. This gets annoying if we want to, say, compare estimates of the same model on two different data sets (in this example, perhaps from different years).
- The solution is to separate the formula from the data source:

```
mob.lm2 <- lm(Mobility ~ Population + Seg_racial + Commute + Income + Gini, data=mob)
```

- The `data` argument tells `lm` to look up variable names appearing in the formula (the first argument) in a dataframe called `mob`.
- It therefore works even if there aren't variables in our workspace called `Mobility`, `Population`, etc., those just have to be column names in `mob`.
- In addition to being easier to write, read and re-use than our first effort, this format works better when we use the model for prediction, as explained below.

Transformations

```
mob.lm3 <- lm(Mobility ~ log(Population) + Seg_racial + Commute + Income + Gini, data=mob)
```

- Formulas are so important that R knows about them as a special data type.
- They *look* like ordinary strings, but they *act* differently, so there are special functions for converting strings (or potentially other things) to formulas, and for manipulating them.
- For instance, if we want to keep around the formula with log-transformed population, we can do as follows:

```
form.logpop <- "Mobility ~ log(Population) + Seg_racial + Commute + Income + Gini"
form.logpop <- as.formula(form.logpop)
mob.lm4 <- lm(form.logpop, data=mob)
```

Why formulas?

- Being able to turn strings into formulas is very convenient if we want to try out a bunch of different model specifications, because R has lots of tools for building strings according to regular patterns, and then we can turn all those into formulas.
- If we have already estimated a model and want the formula it used as the specification, we can extract that with the `formula` function:

```
formula(mob.lm3)
```

```
## Mobility ~ log(Population) + Seg_racial + Commute + Income +
##      Gini
```

```
formula(mob.lm3) == form.logpop
```

```
## [1] TRUE
```

Extracting Coefficients, Confidence Intervals, Fitted Values, Residuals, etc.

If we want the coefficients of a model we've estimated, we can get that with the `coefficients` function:

```
coefficients(mob.lm3)
```

```
##      (Intercept) log(Population)      Seg_racial      Commute  
##      8.338558e-02 -2.894236e-03 -5.656590e-02  1.450771e-01  
##           Income           Gini  
##      1.772105e-06 -1.621921e-01
```

```
mob.lm3$coefficients
```

```
##      (Intercept) log(Population)      Seg_racial      Commute  
##      8.338558e-02 -2.894236e-03 -5.656590e-02  1.450771e-01  
##           Income           Gini  
##      1.772105e-06 -1.621921e-01
```

Or even

```
summary(mob.lm3)$coef
```

```
##              Estimate  Std. Error  t value    Pr(>|t|)  
## (Intercept)  8.338558e-02 2.870373e-02  2.905044 3.784114e-03  
## log(Population) -2.894236e-03 1.874746e-03 -1.543802 1.230739e-01  
## Seg_racial    -5.656590e-02 1.713493e-02 -3.301203 1.009994e-03  
## Commute       1.450771e-01 1.934259e-02  7.500397 1.869467e-13  
## Income        1.772105e-06 2.878660e-07  6.156006 1.236337e-09  
## Gini          -1.621921e-01 2.225561e-02 -7.287695 8.277813e-13
```

Confidence Intervals

- If we want confidence intervals for the coefficients, we can use `confint`:

```
confint(mob.lm3,level=0.90) # default confidence level is 0.95
```

```
##              5 %          95 %  
## (Intercept)  0.036111577 1.306596e-01  
## log(Population) -0.005981875 1.934023e-04  
## Seg_racial    -0.084786513 -2.834528e-02  
## Commute       0.113220542 1.769336e-01  
## Income        0.000001298 2.246209e-06  
## Gini          -0.198846318 -1.255379e-01
```

- **WARNING:** This calculates confidence intervals assuming independent, constant-variance Gaussian noise everywhere, etc., etc., so it's not to be taken too seriously unless you've checked those assumptions somehow; see Chapter 2 of Shalizi's book, and Chapter 6 for alternatives.

Fitted values and residuals

For every data point in the original data set, we have both a fitted value (\hat{y}) and a residual ($y - \hat{y}$). These are vectors, and can be extracted with the `fitted` and `residuals` functions:

```
head(fitted(mob.lm2)) # To let you know head() and tail() exist, jic
```

```
##           1           2           3           4           5           6
## 0.07048490 0.06299687 0.06926223 0.04927934 0.05791660 0.06455628
```

```
tail(residuals(mob.lm2))
```

```
##           736           737           738           739           740
## -0.045252255 -0.031707484  0.004026805  0.015472295 -0.025058476
##           741
##  0.007091485
```

Using bits of the lm output

- You may be more used to accessing all these things as parts of the estimated model — writing something like `mob.lm2$coefficients` to get the coefficients.
- This is fine as far as it goes, but we will work with many different sorts of statistical models in this course, and those internal names can change from model to model.
- If the people implementing the models did their job, however, functions like `fitted`, `residuals`, `coefficients` and `confint` will all, to the extent they apply, work, and work in the same way.

```
# Example of all the different parts of a lm() object
```

```
names(mob.lm2)
```

```
## [1] "coefficients" "residuals"      "effects"         "rank"
## [5] "fitted.values" "assign"         "qr"             "df.residual"
## [9] "na.action"    "xlevels"       "call"           "terms"
## [13] "model"
```

Methods and Classes (Extra details, but possibly important)

- In R things like `residuals` or `coefficients` are a special kind of function, called **methods**.
- Other methods, which you've used a lot without perhaps realizing it, are `plot`, `print` and `summary`.
- These are a sort of generic/meta function, which looks up the class of model being used, and then calls a specialized function which how to work with that class.
- The convention is that the specialized function is named `method.class`, e.g., `summary.lm`.
- If no specialized function is defined, R will try to use `method.default`.

Why methods?

- The advantage of methods is that you, as a user, don't have to learn a totally new syntax to get the coefficients or residuals of every new model class
- you just use `residuals(mdl)` whether `mdl` comes from a linear regression which could have been done two centuries ago, or is a Batrachian Emphasis Machine which won't be invented for another five years.

- (It also means that core parts of R don't have to be re-written every time someone comes up with a new model class.)
- The one draw-back is that the help pages for the generic methods tend to be pretty vague, and you may have to look at the help for the class-specific functions
- Compare `?summary` with `?summary.lm`.

(If you are not sure what the class of your model, `mdl`, is called, use `class(mdl)`.)

Making Predictions

- The point of a regression model is to do prediction, and the method for doing so is, naturally enough, called `predict`. It works like so:

```
predict(object, newdata)
```

- Here `object` is an already estimated model, and `newdata` is a data frame containing the new cases, real or imaginary, for which we want to make predictions.
- The output is (generally) a vector, with a predicted value for each row of `newdata`.
- If the rows of `newdata` have names, those will be carried along as names in the output vector.

```
predict(mob.lm2, newdata=mob[which(mob$State=="AL"),])
```

```
## numeric(0)
```

Subtleties of Predict!

- It is important to remember that making a prediction does *not* mean “changing the data and re-estimating the model”;
- It means taking the unchanged estimate of the model, and putting in new values for the covariates or independent variables.
 - In terms of the linear model, we change x , not $\hat{\beta}$.
- Notice that I used `mob.lm2` here, rather than the mathematically-equivalent `mob.lm1`.


```
-mob.lm1 <- lm(mob$Mobility ~ mob$Population + mob$Seg_racial + mob$Commute + mob$Income + mob$Gini)
-mob.lm2 <- lm(Mobility ~ Population + Seg_racial + Commute + Income + Gini, data=mob)
```
- Because I specified `mob.lm2` with a formula that just referred to column names, `predict` looks up columns with those names in `newdata`, puts them into the function estimated in `mob.lm2`, and calculates the predictions.
- Had I tried to use `mob.lm1`, it would have completely ignored `newdata`.
- This is one crucial reason why it is best to use clean formulas and a `data` argument when estimating the model.

Transformations

- If the formula specifies transformations, those will also be done on `newdata`;
- we don't have to do the transformations ourselves:

```
predict(mob.lm3, newdata=mob[which(mob$State=="AL"),])
```

```
## numeric(0)
```

- The `newdata` does not have to be a subset of the original data used for estimation, or related to it in any way at all

Fun with predict

- It just has to have columns whose names match those in the right-hand side of the formula.

```
predict(mob.lm3, newdata=data.frame(Population=1.5e6, Seg_racial=0,
                                     Commute=0.5, Income=3e4, Gini=median(mob$Gini)))
```

```
##          1
```

```
## 0.1033759
```

```
predict(mob.lm3, newdata=data.frame(Population=1.5e6, Seg_racial=0,
                                     Commute=0.5, Income=quantile(mob$Income,c(0.05,0.5,0.95)),
                                     Gini=quantile(mob$Gini,c(0.05,0.5,0.95))))
```

```
##          5%          50%          95%
```

```
## 0.1122663 0.1075794 0.1024651
```

Problems w/ predict

- A very common programming error is to run `predict` and get out a vector whose length equals the number of rows in the original estimation data
- and which doesn't change no matter what you do to `newdata`.
- This is because if `newdata` is missing, or if R cannot find all the variables it needs in it, the default is the predictions of the model on the original data.
- An even more annoying form of this error consists of forgetting that the argument is called `newdata` and not `data`:

```
head(predict(mob.lm3)) # Equivalent to head(fitted(mob.lm3))
```

```
##          1          2          3          4          5          6
```

```
## 0.06707724 0.06499898 0.06773945 0.05266410 0.06632751 0.07133333
```

More problems

```
head(predict(mob.lm3,data=data.frame(Population=1.5e6, Seg_racial=0,
                                     Commute=0.5, Income=3e4, Gini=median(mob$Gini))))
```

```
##          1          2          3          4          5          6
```

```
## 0.06707724 0.06499898 0.06773945 0.05266410 0.06632751 0.07133333
```

```
# Don't do this!
```

- Returning the original fitted values when `newdata` is missing or messed up is not what I would have chosen, but nobody asked me.
- Because `predict` is a method, the generic help file is fairly vague, and many options are only discussed on the help pages for the class-specific functions

- compare `?predict` with `?predict.lm`.
- Common options include giving standard errors for predictions (as well point forecasts), and giving various sorts of intervals.

Using Different Model Classes

- All of this carries over to different model classes, at least if they've been well-designed.
- For instance, suppose we want to estimate a kernel regression (as in chapter 4) to the same data, using the same variables.

```
#
library(np) # Nonparametric methods library

## Nonparametric Kernel Methods for Mixed Datatypes (version 0.60-9)
## [vignette("np_faq",package="np") provides answers to frequently asked questions]
## [vignette("np",package="np") an overview]
## [vignette("entropy_np",package="np") an overview of entropy-based methods]
# This must be computed first. It is a computationally intese function.
mob.npbw <- npregbw(formula=formula(mob.lm2), data=mob, tol=1e-2, ftol=1e-2)

# Notice that no formula is needed here.
mob.np <- npreg(mob.npbw, data=mob)
```

(See chapter 4 on the `tol` and `ftol` settings.)

Why this is easy

- We can re-use the formula, because it's just saying what the input and target variables of the regression are, and we want that to stay the same.
- More importantly, both `lm` and `npreg` use the same mechanism, of separating the formula specifying the model from the data set containing the actual values of the variables.
- Of course, some models have variations in allowable formulas
 - interactions make sense for `lm` but not for `npreg`,
 - the latter has a special way of dealing with ordered categorical variables that `lm` doesn't
 - etc.
- After estimating the model, we can do most of the same things to it that we could do to a linear model.

Putting It All Together

- We can look at a summary:

```
summary(mob.np)

##
## Regression Data: 729 training points, in 5 variable(s)
##
## No. Complete Observations: 729
## No. Incomplete (NA) Observations: 12
## Observations omitted or excluded: 374 376 386 410 440 459 485 542 613 616 637 652
```

```
##           Population Seg_racial   Commute   Income     Gini
## Bandwidth(s):    1649603  0.1624437 0.03871639 2382.342 0.0318117
##
## Kernel Regression Estimator: Local-Constant
## Bandwidth Type: Fixed
## Residual standard error: 0.0302321
## R-squared: 0.6733646
##
## Continuous Kernel Type: Second-Order Gaussian
## No. Continuous Explanatory Vars.: 5
```

- We can look at fitted values and residuals:

```
head(fitted(mob.np))
```

```
## [1] 0.06430449 0.06742469 0.07513909 0.05630422 0.06187851 0.06751230
```

```
tail(residuals(mob.np))
```

```
##           736           737           738           739           740
## -4.472859e-02 -3.445805e-02 -6.568906e-08  2.774485e-02 -7.634712e-03
##           741
##  1.801038e-02
```

*We can make predictions:

```
predict(mob.np, newdata=data.frame(Population=1.5e6, Seg_racial=0,
                                   Commute=0.5, Income=3e4, Gini=median(mob$Gini)))
```

```
## [1] 0.09849096
```

- and we can plot things

```
par(mar=c(5,5,1,1),cex.lab=3,cex.axis=2,lwd=2,col=4,bty='n')
plot(mob.np,plot.errors.method='bootstrap')
```