

## Question 1: Linear models and causal interpretation

Jake Bowers and Paul Testa

January 27, 2014

Imagine that the US State Department contacted you with the following table:

	ols	lad	lasso
(Intercept)	-1630935.30	-34789.04	-1472511.35
treaty	38265.01	468.64	21553.88
polity	-3584.88	-139.81	-2409.05
racecorrel	-145137.88	-177.65	-131756.45
apc	1236.86	-8.20	1060.24
waraims2	207998.41	9503.18	196505.30
duration	60.29	32.25	61.60
finalprop	278245.15	2129.97	229147.10
lnnewpop	85944.76	1727.78	79421.59

Table 1: Fits: least squares (smoothed conditional means), least absolute deviations (smoothed conditional medians), and lasso (tuning parameter chosen to minimize sum of squared error after 10-fold cross-validation).

The senior diplomat is very smart, but knows nothing about statistics. She would like to act on the idea that humanitarian treaties cause increases in civilian casualties, but she needs to understand what is happening here. The analyst who made the table just submitted a bunch of stuff to wikileaks and is in hiding. So, the diplomat doesn't know what to do.

Your job is to explain where these coefficients came from, what one can say based on them, and what kinds of statements and beliefs ought not to be strongly supported by them.

Here are a few details:

The data and model specification come from [Valentino et al. \(2006\)](#).<sup>1</sup>

Here is a link to the data which was used to produce the above table. You can load the data directly into R this way:

```
load(url("http://jakebowers.org/PS531Data/warttreatydeath.rda"))
```

The data contains all interstate wars from 1900 to 2003. The outcome (`noncomdead`) is number of civilians intentionally killed by one side or another of a war (the rows in the dataset are the sides in the war, usually two sides). The key explanatory variable (`treaty`) records whether the side ratified an international treaty about the protection of civilians (the Hague Convention of 1899 and 1907, or the Geneva Convention of 1949). The authors also added other variables to their data model: polity score, whether the conflict was racial or religious, whether the wars involved an attrition or counterinsurgency strategy, the war aims, the war duration, relative capabilities, adversary population size.

Using what you learned from reading Gerber and Green, Holland, and Berk, your classwork last term, and the web, help this diplomat understand this table. You'll need to specify formally (i.e. with notation and arithmetic) what a "causal relationship" might *mean* here, how such a relationship might relate to the fit of a linear model to some data, and interrogate this linear model given these data. For example, I suspect that diplomat will want to know why we get three different numbers for the `treaty` coefficient, what each of those numbers mean, and which, if any, she should use (for some purpose). In some sense, this person is asking, "Can I trust this model enough to act on it?" It is not enough to say, "The first model chooses coefficients based on the least squares criterion" since the diplomat does not know what the "least squares criterion" is, or why changes in US policy ought to be made based on such a criteria for choosing a best fitting plane. Nor is it enough to say, "We see the effects of treaty controlling for polity, etc...." because the diplomat does not understand what "controlling for" means. If you say, "holding constant", you need to show what you mean. Is there evidence, for example, that you have treaty and non-treaty wars for every relevant combination of the "control" variables (if you don't, in what sense are these scholars "controlling for" something)? Finally, but importantly, do not

<sup>1</sup>But they say "The results . . . indicate that the laws of war do not provide strong protections for civilians in times of war." (368) and they show a coefficient of  $-0.876$  using a logged outcome measure. And so the diplomat is even more confused. For this class, let's focus on the unlogged version.

talk about  $p$ -values or confidence intervals here. Statistical inference does not answer any questions that the diplomat cares to ask right now.

Finally, the diplomat shows you this torn piece of paper in case it might help you:

```

• No statistical inference...
read(url("http://jakobovars.org/P8531Data/"))

#ols <- lm(lmnoncom ~ treaty + polity + racecorrel + apc + waraims2 + duration + finalprop + lnewpop, data = wartreedydeath)
#l1 <- lm(noncomdead ~ treaty + polity + racecorrel + apc + waraims2 + duration + finalprop + lnewpop, data = wartreedydeath)
#l1coef <- coef(l1)

library(quantreg)
rq1 <- rq(noncomdead ~ treaty + polity + racecorrel + apc + waraims2 + duration + finalprop + lnewpop, data = wartreedydeath)
rqcoef <- coef(rq1)

library(glmnet)
set.seed(20140126)
X <- model.matrix(noncomdead ~ treaty + polity + racecorrel + apc + waraims2 + duration + finalprop + lnewpop - 1, data = wartreedydeath)
y <- wartreedydeath$noncomdead
cvglmnet1 <- cv.glmnet(X, y, alpha = 1, standardize = TRUE)

lassocoef <- as.numeric(coef(cvglmnet1, s = "lambda.min"))

thetab <- cbind(ols = l1coef, l1d = rqcoef, lasso = lassocoef)
round(thetab, 4)

      ols      l1d      lasso
(Intercept) -1630936.3 -34789.0 -1472511
treaty       38265.0   468.6   21553
polity       -3584.9  -139.8   -2407
racecorrel   -145137.9 -177.6  -13177
apc          1236.9   -8.2    1/
waraims2     207998.4  8509.2   19
duration      60.3    32.3
finalprop    278245.2  2130.0
lnewpop      85944.8  1727.8

library(xtable)

xtab <- xtable(thetab, caption = "Fits: least squares, l1, and lasso")

```

## References

Valentino, B., Huth, P., and Croco, S. (2006). Covenants without the sword: International law and the protection of civilians in times of war. *World Politics*, 58(3):339.