

1 Classification vs Regression

In class on October 16 (lecture 4), I went through a proof that “classification is easier than regression”. The math I sketched was

Let \hat{f} be any estimate of f_* . Let $\hat{g}(X) = \mathbf{1}(\hat{f}(X) > 1/2)$.

Then,

$$\begin{aligned} & \mathbb{P}(Y \neq \hat{g}(X)|X) - \mathbb{P}(Y \neq g_*(X)|X) \\ &= \dots = \\ &= (2f_*(X) - 1)(\mathbf{1}(g_*(X) = 1) - \mathbf{1}(\hat{g}(X) = 1)) \\ &= |2f_*(X) - 1|\mathbf{1}(g_*(X) \neq \hat{g}(X)) \\ &= 2 \left| f_*(X) - \frac{1}{2} \right| \mathbf{1}(g_*(X) \neq \hat{g}(X)) \end{aligned}$$

Now

$$g_*(X) \neq \hat{g}(X) \Rightarrow |\hat{f}(X) - f_*(X)| \geq |\hat{f}(X) - 1/2|$$

Therefore

$$\begin{aligned} & \mathbb{P}(Y \neq \hat{g}(X)) - \mathbb{P}(Y \neq g_*(X)) \\ &= \int (\mathbb{P}(Y \neq \hat{g}(X)|X) - \mathbb{P}(Y \neq g_*(X)|X)) d\mathbb{P}_X \\ &= \int 2 \left| \hat{f}(X) - \frac{1}{2} \right| \mathbf{1}(g_*(X) \neq \hat{g}(X)) d\mathbb{P}_X \\ &\leq 2 \int |\hat{f}(X) - f_*(X)| \mathbf{1}(g_*(X) \neq \hat{g}(X)) d\mathbb{P}_X \\ &\leq 2 \int |\hat{f}(X) - f_*(X)| d\mathbb{P}_X \end{aligned}$$

Fill in the part that is missing in the dots.

2 Another replication-like exercise

The repo includes an article from Journal of Money, Credit, and Banking as well as data nearly the same as theirs (it is from an earlier version). Load the data and produce out of sample forecasts for $h = 0, 1, 2, 3$. Consider the following methods:

1. Logistic regression
2. Lasso-logistic regression
3. Elnet-logistic regression ($\alpha = 0.5$)

For the regularized methods, try also adding all the squared predictors, and all the squared and cubic predictors. This means a total of 7 models. For the regularized models, use CV to choose λ in each case. Estimate the models using only data up to the year 2000. Produce a table with out-of-sample error rates for

each model and each forecast horizon. Use the QPS score on page 855. Do any of these predict the 2008-09 recession? Feel free to try any other classifiers you like.

The following code loads the data and does some processing. The `alldat` object contains the data for $h = 0$. It is easily adaptable for the other values of h .

```
states = read.csv('statelevel.csv')
national = read.csv('national.csv')

library(tidyverse)
library(lubridate)

gr <- function(x, horizon=3) x/lag(x, horizon)
national = national %>%
  mutate(termspread = GS10-TB3MS, DATE = ymd(DATE),
         SP500 = gr(SP500), PAYEMS = lag(gr(PAYEMS)),
         INDPRO = lag(gr(INDPRO)), GS10=NULL,
         TB3MS = NULL, RECPROUSM156N = NULL)
gri <- function(x, horizon=2) x*lag(x, horizon)
states = states %>% mutate(DATE = mdy(X), X=NULL) %>%
  mutate_at(vars(Alabama:Wyoming), funs(lag(gri(.))))
alldat = full_join(national,states,by='DATE') %>% na.omit()
```

For all seven models and the $h = 0$ horizon, produce a figure that plots (for the whole time period), the true recessions and the forecasts. Think about how you might make a nice looking graph. the true recession