

HW6 q2

10/7/2019

a) predicting arrival delay

We are interested in predicting arrival delay (`ARR_DELAY`) using the variables that would be available before the flight takes off: departure time, departure airport (`ORD` or `DFW`), day of the month, and air carrier (airline). Also, make a “day of week” variable. First, for a baseline comparison, fit a linear model to the data using these variables. (Be naive: assume the effects of time of day or day of week are linear.) Show a table of the coefficients you find. Evaluate the model on the test set (using `predict`) and report its squared-error loss.

Table 1: Q2a

	Dependent variable:
	<code>ARR_DELAY</code>
<code>DEP_TIME</code>	0.015*** (0.001)
<code>ORIGINORD</code>	5.033*** (0.914)
<code>DAY_OF_MONTH</code>	−0.021 (0.043)
<code>CARRIERAS</code>	−5.150 (4.196)
<code>CARRIERB6</code>	1.547 (4.964)
<code>CARRIERDL</code>	−6.439*** (2.407)
<code>CARRIEREV</code>	3.838*** (1.165)
<code>CARRIERF9</code>	36.971*** (3.971)
<code>CARRIERNK</code>	0.969 (2.057)
<code>CARRIEROO</code>	7.317*** (1.431)
<code>CARRIERUA</code>	4.092*** (1.215)
<code>CARRIERVX</code>	17.159*** (5.718)
<code>DAY_OF_WEEK</code>	1.197*** (0.194)
Constant	−16.759*** (1.571)
Observations	16,639
R^2	0.036
Adjusted R^2	0.036
Residual Std. Error	49.628 (df = 16625)
F Statistic	48.222*** (df = 13; 16625)

Note: *p<0.1; **p<0.05; ***p<0.01

b) The relationships may be complicated and a linear model may not be appropriate, so

use `npreg` to fit a Nadaraya–Watson kernel regression model; allow `npregbw` to select all bandwidths with cross-validation. Report the kernel regression’s performance on the test set (again using `predict` and squared-error loss) and compare to the linear model. Does this method seem to do dramatically better?

c) The plot function for `npregression` objects (such as the fit returned by `npreg`) can plot

the marginal association of each variable with the response. If you set the `plot.errors.method = “bootstrap”` option, it will also plot bootstrap-based standard errors for these. Make the plots with standard errors and interpret the results. Which variables seem

strongly related with delay length? What do the plots suggest about the appropriate-ness of linear regression? If you saw major non-linearities in any variable, do these non-linearities appear to harm the predictions enough to make linear regression perform dramatically worse than kernel regression?

d) Repeat this analysis, but use a locally linear kernel regression (with `regtype = "ll"` provided to `npregbw`). Compare the test error from this model to that for the previous one. Discuss possible reasons for any difference you see.