

1

2

3

4

5

6

11

7

9

21 linear models. Machine learning techniques such as decision trees, support
22 vector machines, neural nets, deep learning and so on may allow for more
23 effective ways to model complex relationships.

24 In this essay I will describe a few of these tools for manipulating and an-
25 alyzing big data. I believe that these methods have a lot to offer and should
26 be more widely known and used by economists. In fact, my standard advice
27 to graduate students these days is “go to the computer science department
28 and take a class in machine learning.” There have been very fruitful collabo-
29 rations between computer scientists and statisticians in the last decade or so,
30 and I expect collaborations between computer scientists and econometricians
31 will also be productive in the future.

32 1 Tools for big data

33 Economists have historically dealt with data that fits in a spreadsheet, but
34 that is changing as new more detailed data becomes available; see Einav
35 and Levin [2013] for several examples and discussion. If you have more than
36 a million or so rows in a spreadsheet, you probably want to store it in a
37 relational database, such as MySQL. Relational databases offer a simple way
38 to store, manipulate and retrieve data using a Structured Query Language
39 (SQL) which is easy to learn and very useful for dealing with medium-sized
40 data sets.

41 However, if you have several gigabytes of data or several million observa-
42 tions, standard relational databases become unwieldy. Databases to manage
43 data of this size are generically known as “NoSQL” databases. The term is
44 used rather loosely, but is sometimes interpreted as meaning “not only SQL.”
45 NoSQL databases are more primitive than SQL databases in terms of data
46 manipulation capabilities but can handle larger amounts of data.

47 Due to the rise of computer mediated transactions, many companies have
48 found it necessary to develop systems to process billions of transactions per

49 day. For example, according to Sullivan [2012], Google has seen 30 trillion
50 URLs, crawls over 20 billion of those a day, and answers 100 billion search
51 queries a month. Analyzing even one day’s worth of data of this size is
52 virtually impossible with conventional databases. The challenge of dealing
53 with data sets of this size led to the development of several tools to manage
54 and analyze big data.

55 These tools are proprietary to Google, but have been described in aca-
56 demic publications in sufficient detail that open-source implementations have
57 been developed. The list below has both the Google name and the name of
58 related external tools. Further details can be found in the Wikipedia entries
59 associated with the tool names.

60 **Google File System** [Hadoop Distributed File System] This system sup-
61 ports files of to be distributed across hundreds or even thousands of
62 computers.

63 **Bigtable** [Cassandra] This is a table of data that lives in the Google File
64 System. It too can stretch over many computers.

65 **MapReduce** [Hadoop] This is a system for accessing manipulating data
66 in large data structures such as Bigtables. MapReduce allows you to
67 access the data in parallel, using hundreds or thousands of machines
68 to do the particular data extraction you are interested in. The query
69 is “mapped” to the machines and is then applied in parallel to dif-
70 ferent shards of the data. The partial calculations are then combined
71 (“reduced”) to create the summary table you are interested in.

72 **Go** [Pig] Go is an open-source general-purpose computer language that makes
73 it easier to do parallel data processing.

74 **Dremel** [Hive, Drill] This is a tool that allows data queries to be written
75 in a simplified form of SQL. With Dremel it is possible to run an SQL
76 query on a petabyte of data (1000 terabytes) in a few seconds.

77 Though these tools can be run on a single computer for learning purposes,
78 real applications use large clusters of computers such as those provided by
79 Amazon, Google, Microsoft and other cloud computing providers. The ability
80 to rent rather than buy data storage and processing has turned what was
81 previously a fixed cost into a variable cost and has lowered the barriers to
82 entry for working with big data.

83 **2 Analytic methods**

84 The outcome of the big data processing described above is often a “small”
85 table of data that may be directly human readable or can be loaded into an
86 SQL database, a statistics package, or a spreadsheet.

87 If the extracted data is still inconveniently large, it is often possible to
88 select a subsample for statistical analysis. At Google, for example, I have
89 found that random samples on the order of 0.1 percent work fine for analysis
90 of economic data.

91 Once a dataset has been extracted it is often necessary to do some ex-
92 ploratory data analysis along with consistency and data-cleaning tasks. This
93 is something of an art which can be learned only by practice, but there are
94 data cleaning software tools such as OpenRefine and DataWrangler that can
95 be used to assist in this task.

96 Data analysis in statistics and econometrics can be broken down into four
97 categories: 1) prediction, 2) summarization, 3) estimation, and 4) hypothesis
98 testing. Machine learning is concerned primarily with prediction; the closely
99 related field of data mining is also concerned with summarization. Econo-
100 metricians, statisticians, and data mining specialists are generally looking
101 for insights that can be extracted from the data. Machine learning special-
102 ists are often primarily concerned with developing computers systems that
103 can provide useful predictions and perform well in the presence of challeng-
104 ing computational constraints. Data science, a somewhat newer term, is

concerned with both prediction and summarization, but also with data manipulation, visualization, and other similar tasks. (Note: terminology is not standardized in these areas, so these statements reflect general usage, not hard-and-fast definitions.)

Much of applied econometrics is concerned with detecting and summarizing relationships in the data. The most common tool used to for summarization is (linear) regression analysis. As we shall see, machine learning offers a set of tools that can usefully summarize more complex relationships in the data. We will focus on these regression-like tools since those are the most natural for economic applications.

In the most general formulation of a statistical prediction problem, we are interested in understanding the conditional distribution of some variable y given some other variables $x = (x_1, \dots, x_P)$. If we want a point prediction we could use the mean or median of the conditional distribution.

In machine learning, the x -variables are usually called “predictors” or “features.” The focus of machine learning is to find some function that provides a good prediction of y as a function of x . Historically, most work in machine learning has involved cross-section data where it is natural to think of the data being IID or at least independently distributed. The data may be “fat,” which means lots of predictors relative to the number of observations, or “tall” which means lots of observations relative to the number of predictors.

We typically have some observed data on y and x and we want to compute a “good” prediction of y given new values of x . Usually “good” means it minimizes some loss function such as the sum of squared residuals, mean of absolute value of residuals, and so on. Of course, the relevant loss is that associated with *new* observations of x , not the observations used to fit the model.

When confronted with a prediction problem of this sort an economist would think immediately of a linear or logistic regression. However, there

135 may be better choices, particularly if a lot of data is available. These in-
136 clude nonlinear methods such as 1) neural nets, 2) support vector machines,
137 3) classification and regression trees, 4) random forests, and 5) penalized
138 regression such as lasso, lars, and elastic nets.

139 I will focus on the last three methods in the list above, since they seem
140 to work well on the type of data economists generally use. Neural nets and
141 support vector machines work well for many sorts of prediction problems, but
142 they are something of a black box. By contrast it is easy to understand the
143 relationships that trees and penalized regressions describe. Much more detail
144 about these methods can be found in machine learning texts; an excellent
145 treatment is available in Hastie et al. [2009], which can be freely downloaded.
146 Other suggestions for further reading are given at the end of this article.

147 3 General considerations for prediction

148 Our goal with prediction is typically to get good *out-of-sample predictions*.
149 Most of us know from experience that it is all too easy to construct a predictor
150 that works well in-sample, but fails miserably out-of-sample. To take a trivial
151 example, n linearly independent regressors will fit n observations perfectly
152 but will usually have poor out-of-sample performance. Machine learning
153 specialists refer to this phenomenon as the “overfitting problem.”

154 There are three major techniques for dealing with the overfitting problem
155 which are commonly used in machine learning.

156 First, since simpler models tend to work better for out of sample forecasts,
157 machine learning experts have come up with various ways penalize models for
158 excessive complexity. In the machine learning world, this is known as “reg-
159 ularization” and we will encounter a some examples later one. Economists
160 tend to prefer simpler models for the same reason, but have not been as
161 explicit about quantifying complexity costs.

162 Second, it is conventional to divide the data into separate sets for the

163 purpose of training, testing and validation. You use the training data to
164 estimate a model, the validation data to choose your model, and the testing
165 data to evaluate how well your chosen model performs. (Often validation
166 and testing sets are combined.)

167 Third, in the training stage, it may be necessary to estimate some “tuning
168 parameters” of the model. The conventional way to do this in machine
169 learning is to use *k-fold cross validation*.

- 170 1. Divide the data into k roughly equal subsets and label them by $s =$
171 $1, \dots, k$. Start with subset $s = 1$.
- 172 2. Pick a value for the tuning parameter.
- 173 3. Fit your model using the $k - 1$ subsets other than subset s .
- 174 4. Predict for subset s and measure the associated loss.
- 175 5. Stop if $s = k$, otherwise increment s by 1 and go to step 2.

176 Common choices for k are 10, 5, and the sample size minus 1 (“leave
177 one out”). After cross validation, you end up with k values of the tuning
178 parameter and the associated loss which you can then examine to choose
179 an appropriate value for the tuning parameter. Even if there is no tuning
180 parameter, it is useful to use cross validation to report goodness-of-fit mea-
181 sures since it measures out-of-sample performance which is what is typically
182 of interest.

183 Test-train and cross validation, are very commonly used in machine learn-
184 ing and, in my view, should be used much more in economics, particularly
185 when working with large datasets. For many years, economists have re-
186 ported in-sample goodness-of-fit measures using the excuse that we had small
187 datasets. But now that larger datasets have become available, there is no
188 reason not to use separate training and testing sets. Cross-validation also

189 turns out to be a very useful technique, particularly when working with rea-
190 sonably large data. It is also a much more realistic measure of prediction
191 performance than measures commonly used in economics.

192 4 Classification and Regression Trees

193 Let us start by considering a discrete variable regression where our goal is to
194 predict a 0-1 outcome based on some set of features (what economists would
195 call explanatory variables or predictors.) In machine learning this is known
196 as a *classification problem*. Economists would typically use a generalized
197 linear model like a logit or probit for a classification problem.

198 A quite different way to build a classifier is to use a decision tree. Most
199 economists are familiar with decision trees that describe a sequence of de-
200 cisions that results in some outcome. A tree classifier has the same general
201 form, but the decision at the end of the process is a choice about how to
202 classify the observation. The goal is to construct (or “grow”) a decision tree
203 that leads to good out-of-sample predictions.

204 Ironically, one of the earliest papers on the automatic construction of de-
205 cision trees was co-authored by an economist (Morgan and Sonquist [1963]).
206 However, the technique did not really gain much traction until 20 years later
207 in the work of Breiman et al. [1984] and his colleagues.

208 Consider the simple example shown in Figure 1, where we are trying to
209 predict survivors of the Titanic using just two variables, age and which class
210 of travel the passenger purchased.

211 Here is a set of rules that can be read off of this tree (more of a bush,
212 really):

- 213 • class 3: predict died (370 out of 501)
- 214 • class 1 or 2 and younger than 16: predict lived (34 out of 36)
- 215 • class 2 or 3 and older than 16: predict died (145 out of 233)

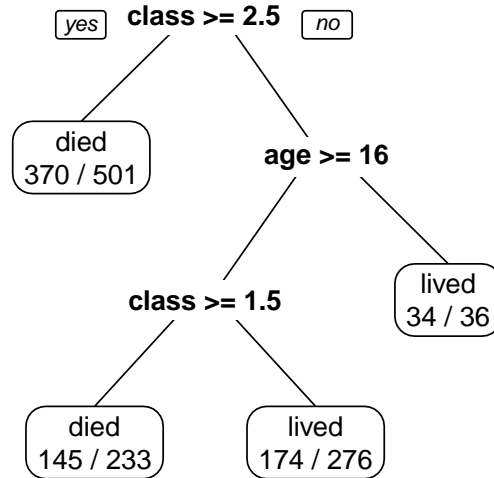


Figure 1: A classification tree for survivors of the Titanic. See text for interpretation.

- class 1, older than 16: predict lived: (174 out of 276)

The rules fit the data reasonably well, misclassifying about 30% of the observations in the testing set.

This classification can also be depicted in the “partition plot” shown in Figure 2 which shows how the tree divides up the space of (age, class) pairs. Of course, the partition plot can only be used for 2 variables while a tree representation can handle an arbitrarily large number.

It turns out that there are computationally efficient ways to construct classification trees of this sort. These methods generally are restricted to binary trees (two branches at each node). They can be used for classification with multiple outcomes (“classification trees”) , or with continuous dependent variables (“regression trees.”)

Trees tend to work well for problems where there are important nonlinearities and interactions. As an example, let us continue with the Titanic data and create a tree that relates survival to age. In this case, the rule

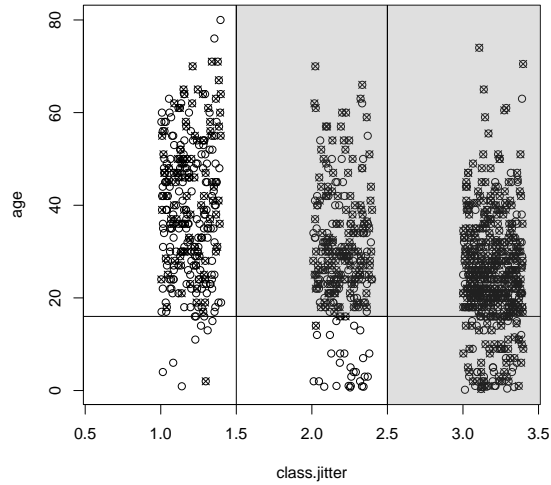


Figure 2: The simple tree model predicts death in shaded region. White circles indicate survival, black crosses indicate death.

231 generated by the tree is very simple: predict “survive” if age < 8.5 years.
 232 We can examine the same data with a logistic regression to estimate the
 233 probability of survival as a function of age:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.464813	0.034973	13.291	<2e-16 ***
age	-0.001894	0.001054	-1.796	0.0727 .

237 The tree model suggests that age is an important predictor of survival impor-
 238 tant, while the logistic model says it is barely important. This discrepancy is
 239 explained in Figure 3 where we plot survival rates by bins. Here we see that
 240 survival rates for those under 10 years old were elevated compared to older
 241 passengers, except for the very oldest group. So what mattered for survival
 242 is not so much age, but whether the passenger was a child or a senior. It
 243 would be difficult to discover this fact from a logistic regression alone.¹

¹It is true that if you knew that there was a nonlinearity in age, you use age dummies in

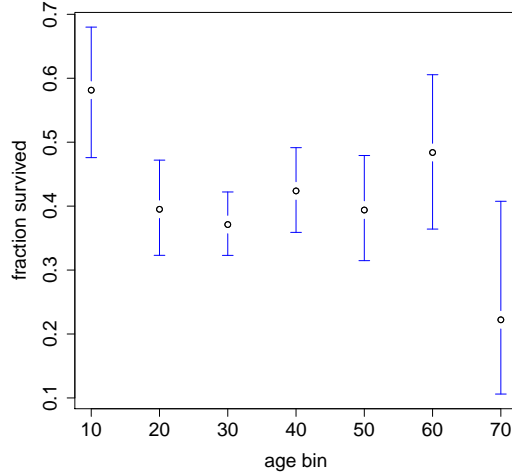


Figure 3: The figure shows the fraction of the population that survived for different age groups (0-10,10-20, and so on). The error bars are computed using the Wilson method.

244 Trees also handle missing data well. Perlich et al. [2003] examined several
 245 standard data sets and found that “logistic regression is better for smaller
 246 data sets and tree induction for larger data sets.” Interestingly enough, trees
 247 tend *not* to work very well if the underlying relationship really is linear,
 248 but there are hybrid models such as RuleFit (Friedman and Popescu [2005])
 249 which can incorporate both tree and linear relationships among variables.

250 However, even if trees may not improve on predictive accuracy compared
 251 to linear models, the age example shows that they may reveal aspects of the
 252 data that are not apparent from a traditional linear modeling approach.

the logit model to capture this effect. However the tree formulation made this nonlinearity quite apparent.

253 4.1 Pruning trees

254 One problem with trees is that they tend to overfit the data. The most
255 widely-used solution to this problem is to “prune” the tree by imposing some
256 complexity cost for having too many branches. This penalty for complexity
257 is a form of regularization, which was mentioned earlier.

258 So, a typical tree estimation session might involve dividing your data
259 into 10 folds, using 9 of the folds to grow a tree with a particular complexity,
260 and then predict on the excluded fold. Repeat the estimation with different
261 values of the complexity parameter using other folds and choose the value
262 of the complexity parameter that minimizes the out-of-sample classification
263 error. (Some researchers recommend being a bit more aggressive than that
264 and choosing the complexity parameter that is one standard deviation lower
265 than the loss-minimizing value.)

266 Of course, in practice, the computer program handles most of these details
267 for you. In the examples in this paper I mostly use default choices, but in
268 practices these default will often be tuned. As with any other statistical
269 procedure, skill, experience and intuition are helpful in coming up with a
270 good answer and diagnostics, exploration, and experimentation are just as
271 useful with these methods as with regression techniques.

272 There are many other approaches to creating trees, including some that
273 are explicitly statistical in nature. For example, a “conditional inference
274 tree,” or ctree for short, chooses the structure of the tree using a sequence
275 of hypothesis tests. The resulting trees tend to need very little pruning.
276 (Hothorn et al. [2006]) An example for the Titanic data is shown in Figure
277 4.

278 One might summarize this tree by the following principle: “women and
279 children first ... particularly if they were traveling first class.” This simple
280 example again illustrates that classification trees can be helpful in summa-
281 rizing relationships in data, as well as predicting outcomes.

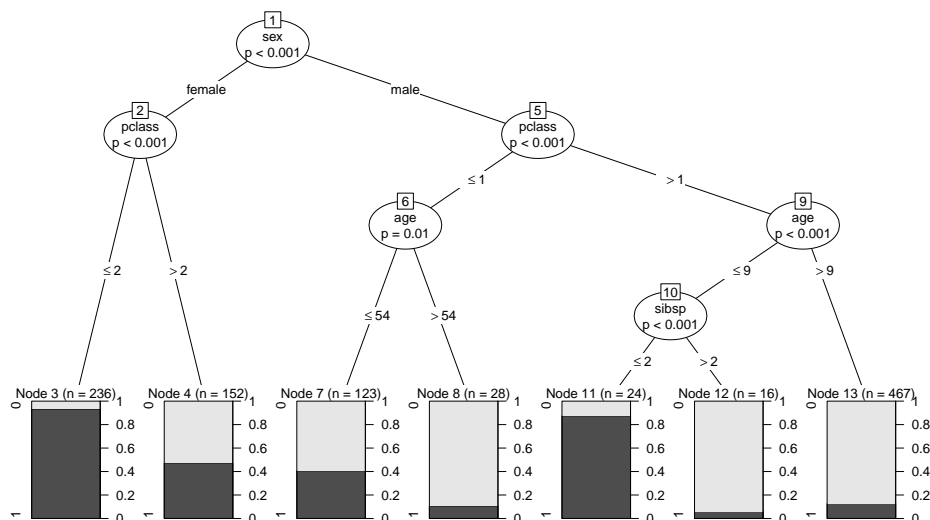


Figure 4: A ctree for survivors of the Titanic. The black bars indicate fraction of the group that survival.

282 4.2 Economic example: HMDA data

283 Munnell et al. [1996] examined mortgage lending in Boston to see if race
 284 played a significant role in determining who was approved for a mortgage.
 285 The primary econometric technique was a logistic regression where race was
 286 included as one of the predictors. The race effect indicated a statistically
 287 significant negative impact on probability of getting a mortgage for black
 288 applicants. This finding prompted lively subsequent debate and discussion,
 289 with 725 citations on Google Scholar as of July 2013.

290 Here I examine this question using the tree-based estimators described in
 291 the previous section. The data consists of 2380 observations of 12 predictors,
 292 one of which was race. Figure 5 shows a conditional tree estimated using the
 293 R package `party`. (For reasons of space, I have omitted variable descriptions
 294 which are readily available on the web site.)

295 The tree fits pretty well, misclassifying 228 of the 2380 observations for an
 296 error rate of 9.6%. By comparison, a simple logistic regression does slightly

297 better, misclassifying 225 of the 2380 observations, leading to an error rate
 298 of 9.5%. As you can see in Figure 5, the most important variable is `dmi`
 299 = “denied mortgage insurance”. This variable alone explains much of the
 300 variation in the data. The race variable (`black`) shows up far down the tree
 301 and seems to be relatively unimportant.

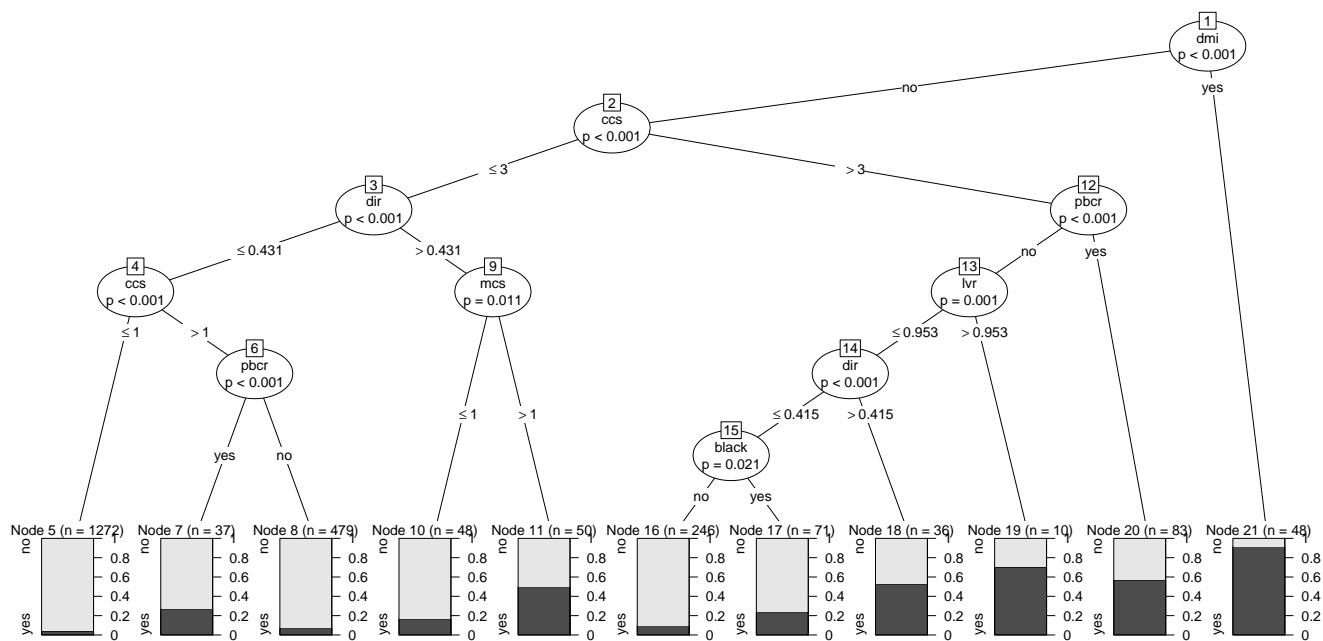


Figure 5: HMDA tree. The black bars indicate the fraction of each group that were denied mortgages. The most important determinant of this is the variable `dmi`, “denied mortgage insurance.”

302 One way to gauge whether a variable is important is to exclude it from
 303 the prediction and see what happens. When this is done, it turns out that
 304 the accuracy of the tree based model doesn’t change at all: exactly the same
 305 cases are misclassified. So there is a plausible decision tree model that ignores
 306 race that fits the observed data just as well as a model that includes race.

307 5 Boosting, bagging and bootstrap

308 There are several useful ways to improve classifier performance. Interestingly
309 enough, the some of these methods work by *adding* randomness to the data.
310 This seems paradoxical at first, but adding randomness turns out to be a
311 helpful way of dealing with the overfitting problem.

312 **Bootstrap** involves choosing (with replacement) a sample of size n from a
313 data set of size n to estimate the sampling distribution of some statistic.
314 A variation is the “ m out of n bootstrap” which draws a sample of size
315 m from a dataset of size $n > m$.

316 **Bagging** involves averaging across models estimated with several different
317 bootstrap samples in order to improve the performance of an estimator.

318 **Boosting** involves repeated estimation where misclassified observations are
319 given increasing weight in each repetition. The final estimate is then a
320 vote or an average across the repeated estimates.

321 Econometricians are well-acquainted with the bootstrap rarely use the
322 other two methods. Bagging is primarily useful for nonlinear models such
323 as trees. (Friedman and Hall [2005].) Boosting tend to improve predictive
324 performance of an estimator significantly and can be used for pretty much
325 any kind of classifier or regression model, including logits, probits, trees, and
326 so on.

327 It is also possible to combine these techniques and create a “forest” of
328 trees that can often significantly improve on single-tree methods. Here is a
329 rough description of how such “random forests” work.

330 **Random forests** refers to a technique that uses multiple trees. A typical
331 procedure uses the follow steps.

- 332 1. Choose a bootstrap sample of the observations and start to grow
333 a tree

- 334 2. At each node of the tree, choose a random sample of the predictors
335 to make the next decision. Do not prune the trees.
- 336 3. Repeat this process many times to grow a forest of trees
- 337 4. The final classification is then determined by majority vote among
338 all the trees in the forest

339 This method produces surprisingly good out-of-sample fits, particularly
340 with highly nonlinear data. In fact, Howard [2013] claims “ensembles of
341 decision trees (often known as Random Forests) have been the most successful
342 general-purpose algorithm in modern times.” He goes on to indicate that
343 “the algorithm is very simple to understand, and is fast and easy to apply.”
344 See also Caruana and Niculescu-Mizil [2006] who compare several different
345 machine learning algorithms and find that ensembles of trees perform quite
346 well. There are a number variations and extensions of the basic “ensemble of
347 trees” model such as Friedman’s “Stochastic Gradient Boosting” (Friedman
348 [1999]).

349 One defect of random forests is that they are a bit of a black box—
350 they don’t offer simple summaries of the data. However, they can be used
351 to determine which variables are “important” in predictions in the sense of
352 contributing the biggest improvements in prediction accuracy.

353 Note that random forests involves quite a bit of randomization; if you
354 want to try them out on some data, I strongly suggest choosing a particular
355 seed for the random number generator so that your results can be reproduced.

356 I ran the random forest method on the HMDA data and found that it
357 misclassified 223 of the 2380 cases, a small improvement over the logit and
358 the ctree. I also used the importance option in random forests to see how
359 the predictors compared. It turned out that `dmi` was the most important
360 predictor and race was second from the bottom which is consistent with the
361 ctree analysis.

362 6 Variable selection

363 Let us return to the familiar world of linear regression and consider the prob-
364 lem of variable selection. There are many such methods available, including
365 stepwise regression, principal component regression, partial least squares,
366 AIC and BIC complexity measures and so on. Castle et al. [2009] describes
367 and compares 21 different methods.

368 6.1 Lasso and friends

369 Here we consider a class of estimators that involves penalized regression.
370 Consider a standard multivariate regression model where we predict y_t as a
371 linear function of a constant, b_0 , and P predictor variables. We suppose that
372 we have standardized all the (non-constant) predictors so they have mean
373 zero and variance one.

Consider choosing the coefficients (b_1, \dots, b_P) for these predictor variables
by minimizing the sum of squared residuals plus a penalty term of the form

$$\lambda \sum_{p=1}^P [(1 - \alpha)|b_p| + \alpha|b_p|^2]$$

374 This estimation method is called *elastic net regression*; it contains three other
375 methods as special cases. If there is no penalty term (i.e., $\lambda = 0$), this is
376 *ordinary least squares*. If $\alpha = 1$ so that there is only the quadratic constraint,
377 this is *ridge regression*. If $\alpha = 0$ this is called the *lasso*.

378 These penalized regressions are classic examples of regularization. In
379 this case, the complexity is the number and size of predictors in the model.
380 All of these methods tend to shrink the least squares regression coefficients
381 towards zero. The lasso and elastic net typically produces regressions where
382 some of the variables are set to be exactly zero. Hence this is a relatively
383 straightforward way to do variable selection.

384 It turns out that these estimators can be computed quite efficiently, so

385 doing variable selection on reasonably large problems is computationally fea-
386 sible. They also seem to provide good predictions in practice.

387 **6.2 Spike and slab regression**

388 Another approach to variable selection that is novel to most economists is
389 spike-and-slab regression, a Bayesian technique. Suppose that you have P
390 possible predictors in some linear model. Let γ be a vector of length P
391 composed of zeros and ones that indicate whether or not a particular variable
392 is included in the regression.

393 We start with a Bernoulli prior distribution on γ ; for example, initially
394 we might think that all variables have an equally likely chance of being in
395 the regression. Conditional on a variable being in the regression, we specify a
396 prior distribution for the regression coefficient associated with that variable.
397 For example, we might use a Normal prior with mean 0 and a large variance.
398 These two priors are the source of the method's name: the "spike" is the
399 probability of a coefficient being non-zero; the "slab" is the (diffuse) prior
400 describing the values that the coefficient can take on.

401 Now we take a draw of γ from its prior distribution, which will just
402 be a list of variables in the regression. Conditional on this list of included
403 variables, we take a draw from the prior distribution for the coefficients. We
404 combine these two draws with the likelihood in the usual way which gives
405 us a draw from posterior distribution on both γ and the coefficients. We
406 repeat this process thousands of times which give us a table summarizing the
407 posterior distribution for γ and the coefficients and the associated prediction
408 of y .

409 We end up with a table of thousands of draws from the posterior distri-
410 butions of γ , β , and y which we can summarize in a variety of ways. For
411 example, we can compute the average value of γ_p which shows the posterior
412 probability variable p is included in the regressions.

predictor	BMA	CDF(0)	lasso	spike-slab
GDP level 1960	1.000	1.000	-	0.9992
Fraction Confucian	0.995	1.000	6	0.9730
Life expectancy	0.946	0.942	5	0.9610
Equipment investment	0.757	0.997	1	0.9532
Sub-Saharan dummy	0.656	1.000	-	0.5834
Fraction Muslim	0.656	1.000	-	0.6590
Rule of law	0.516	1.000	-	0.4532
Open economy	0.502	1.000	3	0.5736
Degree of Capitalism	0.471	0.987	-	0.4230
Fraction Protestant	0.461	0.966	-	0.3798

Table 1: Comparing variable selection algorithms. See text for discussion.

6.3 Economic example: growth regressions

We illustrate the lasso and spike and slab regression with an example from Sala-i-Martin [1997]. This involves examining a multi-country set of predictors of economic growth in order to see which variables appeared to be the most important. Sala-i-Martin [1997] looked at all possible subsets of regressors of manageable size. Ley and Steel [2009] investigated the same question using Bayesian techniques related to, but not identical with, spike-and-slab, while Hendry and Krolzig [2004] examined an iterative significance test selection method.

Table 1 shows 10 predictors that were chosen by Sala-i-Martin [1997], Ley and Steel [2009], `lasso`, and `spike-and-slab`. The table is based on that in Ley and Steel [2009] but metrics used are not strictly comparable across models. The “BMA” and “spike-slab” columns are posterior probabilities of inclusion; the “lasso” column is just the ordinal importance of the variable with a dash indicating that it was not included in the chosen model; and the CDF(0) measure is defined in Sala-i-Martin [1997].

The `lasso` and the Bayesian techniques are very computationally efficient and on this ground would likely be preferred to exhaustive search. All 4 of these variable selection methods give similar results for the first 4 or 5

432 variables, after which they diverge. In this particular case, the data set
433 appears to be too small to resolve the question of what is “important” for
434 economic growth.

435 7 Time series

436 The machine learning techniques described up until now are generally applied
437 to cross-sectional data where independently distributed data is a plausible
438 assumption. However, there are also techniques that work with time series.
439 Here we describe an estimation method which we call Bayesian Structural
440 Time Series (BSTS) that seems to work well for variable selection problems
441 in time series applications.

442 Our research in this area was motivated by Google Trends data which
443 provides an index of the volume of Google queries on specific terms. One
444 might expect that queries on [file for unemployment] might be predictive
445 of the actual rate of filings for initial claims, or that queries on [Orlando
446 vacation] might be predictive of actual visits to Orlando. Indeed, Choi and
447 Varian [2009, 2012], Goel et al. [2010], Carrière-Swallow and Labbé [2011],
448 McLaren and Shanbhoge [2011], Arola and Galan [2012], Hellerstein and
449 Middeldorp [2012] and many others have shown that Google queries do have
450 significant short-term predictive power for various economic metrics.

451 The challenge is that there are billions of queries so it is hard to determine
452 exactly which queries are the most predictive for a particular purpose. Google
453 Trends classifies the queries into categories, which helps a little, but even then
454 we have hundreds of categories as possible predictors so that overfitting and
455 spurious correlation are a serious concern. BSTS is designed to address these
456 issues. We offer a very brief description here; more details are available in
457 Scott and Varian [2012a,b].

458 Consider a classic time series model with *constant* level, linear time trend,
459 and regressor components:

460 • $y_t = \mu + bt + \beta x_t + e_t$.

461 The “local linear trend” is a stochastic generalization of this model where
 462 the level and time trend can vary through time.

463 • Observation: $y_t = \mu_t + z_t + e_{1t} = \text{level} + \text{regression}$

464 • State 1: $\mu_t = \mu_{t-1} + b_{t-1} + e_{2t} = \text{random walk} + \text{trend}$

465 • State 2: $z_t = \beta x_t = \text{regression}$

466 • State 3: $b_t = b_{t-1} + e_{3t} = \text{random walk for trend}$

467 It is easy to add an additional state variable for seasonality if that is ap-
 468 propriate. The parameters to estimate are the regression coefficients β and
 469 the variances of (e_{it}) for $i = 1, \dots, 3$. We can then use these estimates to
 470 construct the optimal Kalman forecast.

471 For the regression we use the spike-and-slab variable choice mechanism
 472 described above. A draw from the posterior distribution now involves a draw
 473 of variances of (e_{1t}, e_{2t}) , a draw of the vector γ that indicates which vari-
 474 ables are in the regression, and a draw of the regression coefficients β for
 475 the included variables. The draws of μ_t , b_t , and β can be used to construct
 476 estimates of y_t and forecasts for y_{t+1} . We end up with an (estimated) pos-
 477 terior distribution for the metric of interest. If we seek a point prediction,
 478 we could average over these draws, which is essentially a form of Bayesian
 479 model averaging.

480 As an example, consider the non-seasonally adjusted housing data (HSN1FNSA)
 481 from the Federal Reserve Economic Data provided by the St. Louis Fed. This
 482 time series can be submitted to Google Correlate, which then returns the 100
 483 queries that are the most highly correlated with the series. We feed that data
 484 into the BSTS system which identifies the predictors with the largest poste-
 485 rior probabilities of appearing in the housing regression are shown in Figure
 486 6. Two predictors, [oldies lyrics] and [www.mail2web] appear to be spurious

so we remove them and re-estimate, yielding the results in Figure 7. The fit is shown in Figure 8 which shows the incremental contribution of the trend, seasonal, and individual regressors components.

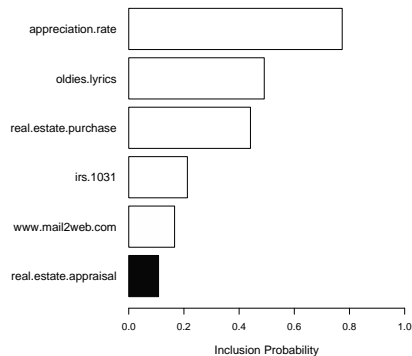


Figure 6: Initial predictors.

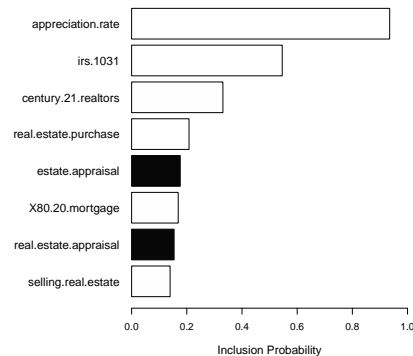


Figure 7: Final predictors.

8 Econometrics and machine learning

There are a number of areas where there would be opportunities for fruitful collaboration between econometrics and machine learning. I mentioned above that most machine learning uses IID data. However, the BSTS model shows that some of these techniques can be adopted for time series models. It should also be possible to incorporate panel data into this framework and there has been some work in this direction.

Econometricians have developed several tools for causal modeling such as instrumental variables, regression discontinuity, and various forms of experiments. (Angrist and Krueger [2001].) Machine learning work has, for the most part, dealt with pure prediction. In a way this is ironic, since theoretical computer scientists, such as Pearl [2009a,b] have made significant contributions to causal modeling. However, it appears that these theoretical

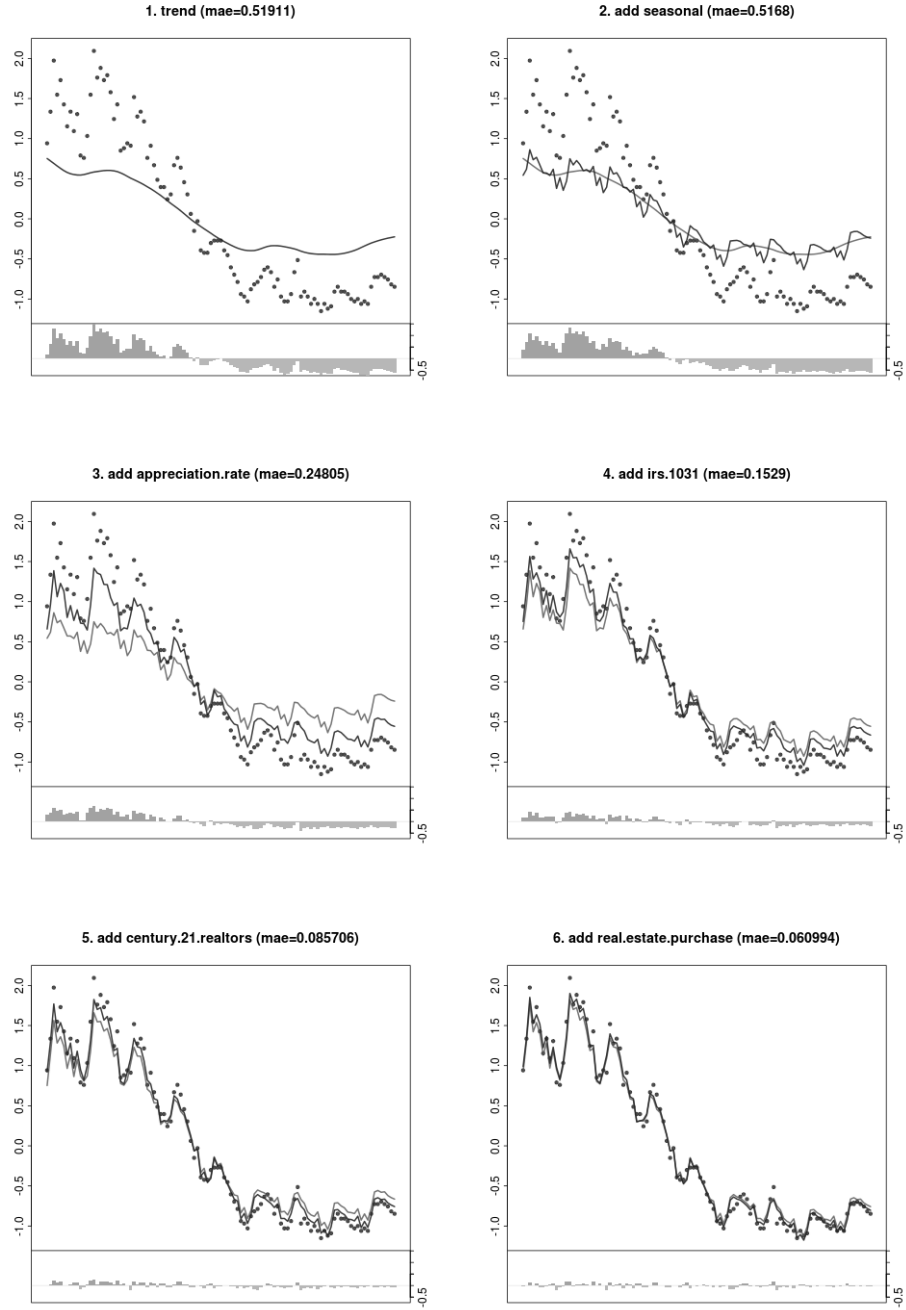


Figure 8: Incremental plots. The plots show the impact of the trend, seasonal, and a few individual regressors. The residuals are shown on the bottom.

advances have not as yet been incorporated into machine learning practice to a significant degree.

8.1 Causality and prediction

As economists know well there is a big difference between correlation and causation. A classic example: there are often more police in precincts with high crime, but that does not imply that increasing the number of police in a precinct would increase crime.

The machine learning models we have described so far have been entirely about prediction. If our data was generated by policymakers who assigned police to areas with high crime, then the observed relationship between police and crime rates could be highly predictive *within* sample, but not useful in predicting the causal impact of explicitly *assigning* additional police to a precinct.

To enlarge on this point, let us consider an experiment (natural or designed) that attempts to estimate the impact of some policy, such as adding police to neighborhoods. There are two critical questions.

- Which neighborhoods will receive additional police in the experiment and policy implementation and how will this be determined? Possible assignment rules could be 1) random, 2) based on perceived need, 3) based on cost of providing service, 4) based on resident requests, 5) based on a formula or set of rules, and so on. Ideally the assignment procedure in the experiment will be similar to that used in the policy. A good model for predicting which neighborhoods will receive additional police can clearly be helpful in estimating the impact of the policy.
- What will be the impact of these additional police in both the experiment and the policy? As Rubin [1974] and many subsequent authors have emphasized, when we consider the causal impact of some treatment we need to compare the outcome with the intervention to what

531 *would have happened* without the intervention. But this counterfactual
532 cannot be observed, so it must be predicted by some model. The better
533 predictive model you have for the counterfactual, the better you will be
534 able to estimate the causal effect, an observation that is true for both
535 pure experiments and natural experiments.

536 So even though a predictive model will not necessarily allow one to con-
537 clude anything about causality by itself, such a model may help in estimating
538 the causal impact of an intervention when it occurs.

539 For example, let us consider the problem of estimating the causal effect of
540 advertising. The problem here is that there are many confounding variables,
541 such as seasonality or weather, that cause both increased ad exposures and
542 increased purchases by consumers.

543 This is illustrated by the (probably apocryphal) story about an advertis-
544 ing manager who was asked why he thought his ads were effective. “Look
545 at this chart,” he said. “Every December I increase my ad spend and, sure
546 enough, purchases go up.” Of course, seasonality can be observed and in-
547 cluded in the model. However, generally there will be other confounding
548 variables that affect both exposure to ads and the propensity of purchase,
549 which makes causal interpretations of relationships problematic.

550 The ideal way to estimate advertising effectiveness is, of course, to run a
551 controlled experiment. In this case the control group provides an estimate of
552 what would have happened without ad exposures. But this ideal approach
553 can be quite expensive, so it is worth looking for alternative ways to predict
554 the counterfactual. One way to do this is to use the Bayesian Structural Time
555 Series method described earlier. In this case, a model based on historical time
556 series data can, in some cases, be used to estimate what *would have happened*
557 in the absence of the advertising intervention. See Brodersen et al. [2013] for
558 details of this approach.

559 9 Model uncertainty

560 An important insight from machine learning is that averaging over many
561 small models tends to give better out-of-sample prediction than choosing a
562 single model.

563 In 2006, Netflix offered a million dollar prize to researchers who could
564 provide the largest improvement to their existing movie recommendation
565 system. The winning submission involved a “complex blending of no fewer
566 than 800 models” though they also point out that “predictions of good quality
567 can usually be obtained by combining a small number of judiciously chosen
568 methods.” (Feuerverger et al. [2012].) It also turned out that a blend of the
569 best and second best model outperformed both of them.

570 Ironically, it was recognized many years ago that averages of macroeco-
571 nomic model forecasts outperformed individual models, but somehow this
572 idea was rarely exploited in traditional econometrics. The exception is the
573 literature on Bayesian model averaging which has seen a steady flow of work;
574 see Steel [2011] for a survey.

575 However, I think that model uncertainty has crept in to applied econo-
576 metrics through the back door. Many papers in applied econometrics present
577 regression results in a table with several different specifications: which vari-
578 ables are included in the controls, which variables are used as instruments,
579 and so on. The goal is usually to show that the estimate of some interesting
580 parameter is not very sensitive to the exact specification used.

581 One way to think about it is that these tables illustrate a simple form of
582 model uncertainty: how an estimated parameter varies as different models are
583 used. In these papers the authors tend to examine only a few representative
584 specifications, but there is no reason why they couldn’t examine many more
585 if the data were available.

586 In this period of “big data” it seems strange to focus on *sampling un-*
587 *certainty*, which tends to be small with large data sets, while completely
588 ignoring *model uncertainty* which may be quite large. One way to address

589 this is to be explicit about examining how parameter estimates vary with
590 respect to choices of control variables and instruments.

591 10 Summary and further reading

592 Since computers are now involved in many economic transactions, big data
593 will only get bigger. Data manipulation tools and techniques developed for
594 small datasets will become increasingly inadequate to deal with new prob-
595 lems. Researchers in machine learning have developed ways to deal with
596 large data sets and economists interested in dealing with such data would be
597 well advised to invest in learning these techniques.

598 I have already mentioned Hastie et al. [2009] which has detailed descrip-
599 tions of all the methods discussed here but at a relatively advanced level.
600 James et al. [2013] describes many of the same topics at an undergraduate-
601 level, along with R code and many examples.²

602 Venables and Ripley [2002] contains good discussions of these topics with
603 emphasis on applied examples. Leek [2013] presents a number of YouTube
604 videos with gentle and accessible introductions to several tools of data anal-
605 ysis. Howe [2013] provides a somewhat more advanced introduction to data
606 science that also includes discussions of SQL and NoSQL databases. Wu
607 and Kumar [2009] gives detailed descriptions and examples of the major al-
608 gorithms in data mining, while Williams [2011] provides a unified toolkit.
609 Domingos [2012] summarizes some important lessons which include “pitfalls
610 to avoid, important issues to focus on and answers to common questions.”

²There are several economic examples in the book where the tension between predictive modeling and causal modeling is apparent.

References

- Joshua D. Angrist and Alan B. Krueger. Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic Perspectives*, 15(4):69–85, 2001. URL <http://www.aeaweb.org/articles.php?doi=10.1257/jep.15.4.69>.
- Concha Arola and Enrique Galan. Tracking the future on the web: Construction of leading indicators using internet searches. Technical report, Bank of Spain, 2012. URL <http://www.bde.es/webbde/SES/Secciones/Publicaciones/PublicacionesSeriadas/DocumentosOcasionales/12/Fich/do1203e.pdf>.
- L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth and Brooks/Cole, Monterey, 1984.
- Kay H. Brodersen, Nicolas Remy, Fabian Gallusser, Steven L. Scott, Jim Koehler, and Penny Chu. Measuring the causal impact of marketing campaigns. Technical report, Google, Inc., 2013. URL TBA.
- Yan Carrière-Swallow and Felipe Labbé. Nowcasting with Google Trends in an emerging market. *Journal of Forecasting*, 2011. doi: 10.1002/for.1252. URL <http://ideas.repec.org/p/chb/bcchwp/588.html>. Working Papers Central Bank of Chile 588.
- Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of supervised learning algorithms. In *Proceedings of the 23rd International Conference on Machine Learning*, Pittsburgh, PA, 2006.
- Jennifer L. Castle, Xiaochuan Qin, and W. Robert Reed. How to pick the best regression equation: A review and comparison of model selection algorithms. Technical Report 13/2009, Department of Economics, University of Canterbury, 2009. URL <http://www.econ.canterbury.ac.nz/RePEc/cbt/econwp/0913.pdf>.

- 638 Hyunyoung Choi and Hal Varian. Predicting the present with Google Trends.
 639 Technical report, Google, 2009. URL [http://google.com/googleblogs/
 640 pdfs/google_predicting_the_present.pdf](http://google.com/googleblogs/pdfs/google_predicting_the_present.pdf).
- 641 Hyunyoung Choi and Hal Varian. Predicting the present with Google Trends.
 642 *Economic Record*, 2012. URL [http://people.ischool.berkeley.edu/
 643 ~hal/Papers/2011/ptp.pdf](http://people.ischool.berkeley.edu/~hal/Papers/2011/ptp.pdf).
- 644 Pedro Domingos. A few useful things to know about machine learning. *Com-
 645 munications of the ACM*, 55(10), October 2012. URL [http://homes.cs.
 646 washington.edu/~pedrod/papers/cacm12.pdf](http://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf).
- 647 Liran Einav and Jonathan Levin. The data revolution and economic analysis.
 648 Technical report, NBER Innovation Policy and the Economy Conference,
 649 2013.
- 650 Andrey Feuerverger, Yu He, and Shashi Khatri. Statistical significance of
 651 the Netflix challenge. *Statistical Science*, 27(2):202–231, 2012. URL [http:
 652 //arxiv.org/abs/1207.5649](http://arxiv.org/abs/1207.5649).
- 653 Jerome Friedman. Stochastic gradient boosting. Technical report, Stan-
 654 ford University, 1999. URL [http://www-stat.stanford.edu/~jhf/ftp/
 655 stobst.pdf](http://www-stat.stanford.edu/~jhf/ftp/stobst.pdf).
- 656 Jerome Friedman and Peter Hall. On bagging and nonlinear estimation.
 657 Technical report, Stanford University, 2005. URL [http://www-stat.
 658 stanford.edu/~jhf/ftp/bag.pdf](http://www-stat.stanford.edu/~jhf/ftp/bag.pdf).
- 659 Jerome Friedman and Bogdan E. Popescu. Predictive learning via rule
 660 ensembles. Technical report, Stanford University, 2005. URL [http:
 661 //www-stat.stanford.edu/~jhf/R-RuleFit.html](http://www-stat.stanford.edu/~jhf/R-RuleFit.html).
- 662 Sharad Goel, Jake M. Hofman, Sbastien Lahaie, David M. Pennock, and
 663 Duncan J. Watts. Predicting consumer behavior with web search. *Pro-*

664 *ceedings of the National Academy of Sciences*, 2010. URL <http://www.pnas.org/content/107/41/17486.full>.
665

666 Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of*
667 *Statistical Learning: Data Mining, Inference, and Prediction*. Springer-
668 Verlag, 2 edition, 2009. URL <http://www-stat.stanford.edu/~tibs/ElemStatLearn/download.html>.
669

670 Rebecca Hellerstein and Menno Middelorp. Forecasting with
671 internet search data. *Liberty Street Economics Blog of the*
672 *Federal Reserve Bank of New York*, January 2012. URL
673 <http://libertystreeteconomics.newyorkfed.org/2012/01/forecasting-with-internet-search-data.html>.
674

675 David F. Hendry and Hans-Martin Krolzig. We ran one regression. *Oxford*
676 *Bulletin of Economics and Statistics*, 66(5):799–810, 2004.

677 Torsten Hothorn, Kurt Hornik, and Achim Zeileis. Unbiased recursive par-
678 titioning: A conditional inference framework. *Journal of Computational*
679 *and Graphical Statistics*, 15(3):651–674, 2006.

680 Jeremy Howard. The two most important algorithms in predictive mod-
681 eling today. Conference presentation, February 2013. URL <http://strataconf.com/strata2012/public/schedule/detail/22658>.
682

683 Bill Howe. Introduction to data science. Technical report, University of
684 Washington, 2013. URL <https://class.coursera.org/datasci-001/lecture/index>.
685

686 Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An*
687 *Introduction to Statistical Learning with Applications in R*. Springer, New
688 York, 2013.

689 Jeff Leek. Data analysis, 2013. URL <http://blog.revolutionanalytics.com/2013/04/coursera-data-analysis-course-videos.html>.
690

- 691 Eduardo Ley and Mark F. J. Steel. On the effect of prior assumptions in
692 Bayesian model averaging with applications to growth regression. *Journal of Applied Econometrics*, 24(4):651–674, 2009. URL [http://ideas.
693 repec.org/a/jae/japmet/v24y2009i4p651-674.html](http://ideas.repec.org/a/jae/japmet/v24y2009i4p651-674.html).
694
- 695 Nick McLaren and Rachana Shanbhoge. Using internet search data
696 as economic indicators. *Bank of England Quarterly Bulletin*,
697 June 2011. URL [http://www.bankofengland.co.uk/publications/
698 quarterlybulletin/qb110206.pdf](http://www.bankofengland.co.uk/publications/quarterlybulletin/qb110206.pdf).
- 699 James N. Morgan and John A. Sonquist. Problems in the analysis of survey
700 data, and a proposal. *Journal of the American Statistical Association*, 58
701 (302):415–434, 1963. URL <http://www.jstor.org/stable/2283276>.
- 702 Alicia H. Munnell, Geoffrey M. B. Tootell, Lynne E. Browne, and James
703 McEneaney. Mortgage lending in boston: Interpreting HDMA data. *Amer-
704 ican Economic Review*, pages 25–53, 1996.
- 705 Judea Pearl. *Causality*. Cambridge University Press, 2009a.
- 706 Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*,
707 4:96–146, 2009b.
- 708 Claudia Perlich, Foster Provost, and Jeffrey S. Simonoff. Tree induction vs.
709 logistic regression: A learning-curve analysis. *Journal of Machine Learning
710 Research*, 4:211–255, 2003. URL [http://machinelearning.wustl.edu/
711 mlpapers/paper_files/PerlichPS03.pdf](http://machinelearning.wustl.edu/mlpapers/paper_files/PerlichPS03.pdf).
- 712 Donald Rubin. Estimating causal effects of treatment in randomized and non-
713 randomized studies. *Journal of Educational Psychology*, 66(5):689, 1974.
- 714 Xavier Sala-i-Martin. I just ran two million regressions. *American Economic
715 Review*, 87(2):178–83, 1997.

- 716 Steve Scott and Hal Varian. Bayesian variable selection for nowcasting
 717 economic time series. Technical report, Google, 2012a. URL <http://www.ischool.berkeley.edu/~hal/Papers/2012/fat.pdf>. Presented
 718 at JSM, San Diego.
 719
- 720 Steve Scott and Hal Varian. Predicting the present with Bayesian structural
 721 time series. Technical report, Google, 2012b. URL <http://www.ischool.berkeley.edu/~hal/Papers/2013/pred-present-with-bsts.pdf>.
 722
- 723 Mark F. J. Steel. Bayesian model averaging and forecasting. *Bulletin*
 724 *of E.U. and U.S. Inflation and Macroeconomic Analysis*, 200:30–41,
 725 2011. URL http://www2.warwick.ac.uk/fac/sci/statistics/staff/academic-research/steel/steel_homepage/publ/bma_forecast.pdf.
 726
- 727 Danny Sullivan. Google: 100 billion searches per month, search to integrate
 728 gmail, launching enhanced search app for iOS. *Search Engine Land*, 2012.
 729 URL <http://searchengineland.com/google-search-press-129925>.
- 730 W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer-
 731 Verlag, New York, 4 edition, 2002.
- 732 Graham Williams. *Data Mining with Rattle and R*. Springer, New York,
 733 2011.
- 734 Xindong Wu and Vipin Kumar, editors. *The Top Ten Algorithms in*
 735 *Data Mining*. CRC Press, 2009. URL <http://www.cs.uvm.edu/~icdm/algorithms/index.shtml>.
 736