

1

2

3

4

5

6

7

8

10

11

16

17

19

*Hal Varian is Chief Economist, Google Inc., Mountain View, California, and Emeritus Professor of Economics, University of California, Berkeley, California. Thanks to Jeffrey Oldham, Tom Zhang, Rob On, Pierre Grinspan, Jerry Friedman, Art Owen, Steve Scott, Bo Cowgill, Brock Noland, Daniel Stonehill, Robert Snedegar, Gary King, and the editors of this journal for helpful comments on earlier versions of this paper.

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

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

31 1 Tools to manipulate big data

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

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

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

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

54 A number of these tools are proprietary to Google, but have been de-
55 scribed in academic publications in sufficient detail that open-source imple-
56 mentations have been developed. Table 1 contains both the Google name
57 and the name of related open source tools. Further details can be found in
58 the Wikipedia entries associated with the tool names.

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

65 **2 Tools to analyze data**

66 The outcome of the big data processing described above is often a “small”
67 table of data that may be directly human readable or can be loaded into
68 an SQL database, a statistics package, or a spreadsheet. If the extracted
69 data is still inconveniently large, it is often possible to select a subsample
70 for statistical analysis. At Google, for example, I have found that random
71 samples on the order of 0.1 percent work fine for analysis of business data.

72 Once a dataset has been extracted it is often necessary to do some ex-
73 ploratory data analysis along with consistency and data-cleaning tasks. This
74 is something of an art which can be learned only by practice, but data clean-
75 ing tools such as OpenRefine and DataWrangler can be used to assist in data

Google name	Analog	Description
Google File System	Hadoop File System	This system supports files so large that they must be distributed across hundreds or even thousands of computers.
Bigtable	Cassandra	This is a table of data that lives in the Google File System. It too can stretch over many computers.
MapReduce	Hadoop	This is a system for accessing manipulating data in large data structures such as Bigtables. MapReduce allows you to access the data in parallel, using hundreds or thousands of machines to extract the data you are interested in. The query is “mapped” to the machines and is then applied in parallel to different shards of the data. The partial calculations are then combined (“reduced”) to create the summary table you are interested in.
Sawzall	Pig	This is a language for creating MapReduce jobs.
Go	None	Go is a flexible open-source general-purpose computer language that makes it easier to do parallel data processing.
Dremel, BigQuery	Hive, Drill, Impala	This is a tool that allows data queries to be written in a simplified form of SQL. With Dremel it is possible to run an SQL query on a petabyte of data (1000 terabytes) in a few seconds.

Table 1: Tools for manipulating big data.

76 cleansing.

77 Data analysis in statistics and econometrics can be broken down into four
78 categories: 1) prediction, 2) summarization, 3) estimation, and 4) hypothesis
79 testing. Machine learning is concerned primarily with prediction; the closely
80 related field of data mining is also concerned with summarization, and par-
81 ticularly in finding interesting patterns in the data. Econometricians, statis-
82 ticians, and data mining specialists are generally looking for insights that
83 can be extracted from the data. Machine learning specialists are often pri-
84 marily concerned with developing high-performance computer systems that
85 can provide useful predictions in the presence of challenging computational
86 constraints. Data science, a somewhat newer term, is concerned with both
87 prediction and summarization, but also with data manipulation, visualiza-
88 tion, and other similar tasks. Note that terminology is not standardized
89 in these areas, so these descriptions reflect general usage, not hard-and-fast
90 definitions. Other terms used to describe computer assisted data analysis
91 include knowledge extraction, information discovery, information harvesting,
92 data archaeology, data pattern processing, and exploratory data analysis.

93 Much of applied econometrics is concerned with detecting and summariz-
94 ing relationships in the data. The most common tool used to for summariza-
95 tion is (linear) regression analysis. As we shall see, machine learning offers
96 a set of tools that can usefully summarize more various sorts of nonlinear
97 relationships in the data. We will focus on these regression-like tools because
98 they are the most natural for economic applications.

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

103 In machine learning, the x -variables are usually called “predictors” or
104 “features.” The focus of machine learning is to find some function that
105 provides a good prediction of y as a function of x . Historically, most work

106 in machine learning has involved cross-section data where it is natural to
107 think of the data being independent and identically distributed (IID) or at
108 least independently distributed. The data may be “fat,” which means lots
109 of predictors relative to the number of observations, or “tall” which means
110 lots of observations relative to the number of predictors.

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

117 When confronted with a prediction problem of this sort an economist
118 would think immediately of a linear or logistic regression. However, there
119 may be better choices, particularly if a lot of data is available. These include
120 nonlinear methods such as 1) classification and regression trees (CART), 2)
121 random forests, and 3) penalized regression such as LASSO, LARS, and elas-
122 tic nets. (There are also other techniques such as neural nets, deep learning,
123 and support vector machines which I do not cover in this review.) Much
124 more detail about these methods can be found in machine learning texts; an
125 excellent treatment is available in Hastie et al. [2009], which can be freely
126 downloaded. Additional suggestions for further reading are given at the end
127 of this article.

128 **3 General considerations for prediction**

129 Our goal with prediction is typically to get good *out-of-sample predictions*.
130 Most of us know from experience that it is all too easy to construct a predictor
131 that works well in-sample, but fails miserably out-of-sample. To take a trivial
132 example, n linearly independent regressors will fit n observations perfectly
133 but will usually have poor out-of-sample performance. Machine learning

134 specialists refer to this phenomenon as the “overfitting problem.” They use
135 several different methods for dealing with this problem.

136 First, since simpler models tend to work better for out of sample forecasts,
137 machine learning experts have come up with various ways penalize models
138 for excessive complexity. In the machine learning world, this is known as
139 “regularization” and we will describe some examples below. Economists tend
140 to prefer simpler models for the same reason, but have not been as explicit
141 about quantifying complexity costs.

142 Second, it is conventional to divide the data into separate sets for the
143 purpose of training, testing and validation. You use the training data to
144 estimate a model, the validation data to choose your model, and the testing
145 data to evaluate how well your chosen model performs. (Often validation
146 and testing sets are combined.)

147 Third, it is often possible to capture the complexity of a model in terms
148 of a “tuning parameter.” One can then choose an appropriate value of this
149 tuning parameter by using *k-fold cross validation*.

- 150 1. Divide the data into k roughly equal subsets (folds) and label them by
151 $s = 1, \dots, k$. Start with subset $s = 1$.
- 152 2. Pick a value for the tuning parameter.
- 153 3. Fit your model using the $k - 1$ subsets other than subset s .
- 154 4. Predict for subset s and measure the associated loss.
- 155 5. Stop if $s = k$, otherwise increment s by 1 and go to step 2.

156 Common choices for k are 10, 5, and the sample size minus 1 (“leave
157 one out”). After cross validation, you end up with k values of the tuning
158 parameter and the associated loss which you can then examine to choose
159 an appropriate value for the tuning parameter. Even if there is no tuning

160 parameter, it is useful to use cross validation to report goodness-of-fit mea-
161 sures since it measures out-of-sample performance which is generally more
162 meaningful than in-sample performance.

163 The test-train cycle and cross validation are very commonly used in ma-
164 chine learning and, in my view, should be used much more in economics,
165 particularly when working with large datasets. For many years, economists
166 have reported in-sample goodness-of-fit measures using the excuse that we
167 had small datasets. But now that larger datasets have become available,
168 there is no reason not to use separate training and testing sets. Cross-
169 validation also turns out to be a very useful technique, particularly when
170 working with reasonably large data. It is also a much more realistic measure
171 of prediction performance than measures commonly used in economics.

172 4 Classification and regression trees

173 Let us start by considering a discrete variable regression where our goal is to
174 predict a 0-1 outcome based on some set of features (what economists would
175 call explanatory variables or predictors.) In machine learning this is known as
176 a *classification problem*. A common example would be classifying email into
177 “spam” or “not spam” based on characteristics of the email. Economists
178 would typically use a generalized linear model like a logit or probit for a
179 classification problem.

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

186 Ironically, one of the earliest papers on the automatic construction of de-
187 cision trees was co-authored by an economist (Morgan and Sonquist [1963]).

features	predicted	actual/total
class 3	died	370/501
class 1-2, younger than 16	lived	34/36
class 2, older than 16	died	145/233
class 1, older than 16	lived	174/276

Table 2: Tree model in rule form.

188 However, the technique did not really gain much traction until 20 years later
189 in the work of Breiman et al. [1984] and his colleagues. Nowadays this predic-
190 tion technique is known as “classification and regression trees”, or “CART.”

191 To illustrate the use of tree models, I used the R package `rpart` to find
192 a tree that predicts Titanic survivors using just two variables, age and class
193 of travel.¹ The resulting tree is shown in Figure 1, and the rules depicted
194 in the tree are shown in Table 2. The rules fit the data reasonably well,
195 misclassifying about 30% of the observations in the testing set.

196 This classification can also be depicted in the “partition plot” shown in
197 Figure 2 which shows how the tree divides up the space of age and class
198 pairs into rectangular regions. Of course, the partition plot can only be used
199 for two variables while a tree representation can handle an arbitrarily large
200 number.

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

206 Trees tend to work well for problems where there are important nonlin-
207 earities and interactions. As an example, let us continue with the Titanic
208 data and create a tree that relates survival to age. In this case, the rule
209 generated by the tree is very simple: predict “survive” if age < 8.5 years.
210 We can examine the same data with a logistic regression to estimate the

¹All data and code used in this paper can be found in the online supplement.

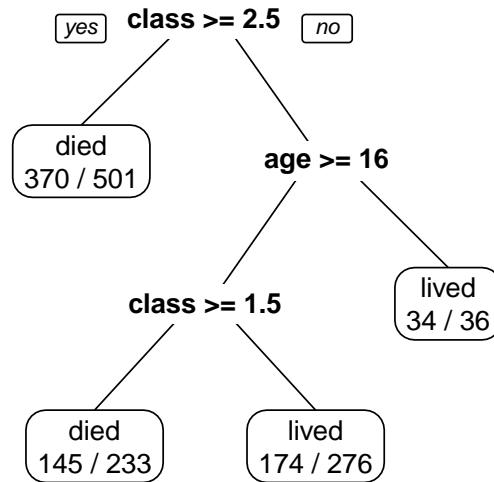


Figure 1: A classification tree for survivors of the Titanic. See text for interpretation.

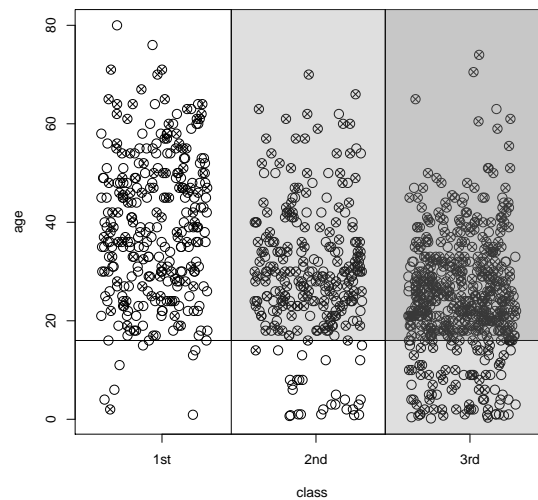


Figure 2: The simple tree model predicts death in shaded region. White circles indicate survival, black crosses indicate death.

Coefficient	Estimate	Std Error	t value	p value
Intercept	0.465	0.0350	13.291	0.000
age	-0.002	0.001	-1.796	0.072

Table 3: Logistic regression of survival vs age.

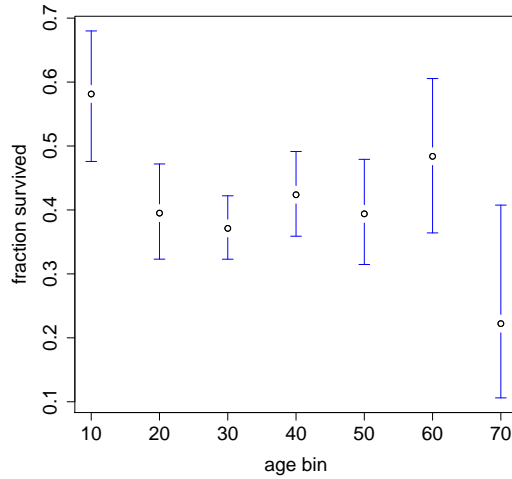


Figure 3: The figure shows the mean survival rates for different age groups along with confidence intervals. The lowest bin is “10 and younger”, the next is “older than 10, through 20” and so on.

211 probability of survival as a function of age, with results reported in Table 3.

212 The tree model suggests that age is an important predictor of survival
213 important, while the logistic model says it is barely important. This dis-
214 crepancy is explained in Figure 3 where we plot survival rates by age bins.
215 Here we see that survival rates for the youngest passengers were relatively
216 high and older passengers were relatively low. For passengers between these
217 two extremes, age didn’t matter much. It would be difficult to discover this
218 pattern from a logistic regression alone.²

²It is true that if you *knew* that there was a nonlinearity in age, you could use age dummies in the logit model to capture this effect. However the tree formulation made this nonlinearity immediately apparent.

219 Trees also handle missing data well. Perlich et al. [2003] examined several
220 standard datasets and found that “logistic regression is better for smaller
221 datasets and tree induction for larger data sets.” Interestingly enough, trees
222 tend *not* to work very well if the underlying relationship really is linear,
223 but there are hybrid models such as RuleFit (Friedman and Popescu [2005])
224 which can incorporate both tree and linear relationships among variables.
225 However, even if trees may not improve on predictive accuracy compared to
226 linear models, the age example shows that they may reveal aspects of the
227 data that are not apparent from a traditional linear modeling approach.

228 4.1 Pruning trees

229 One problem with trees is that they tend to “overfit” the data. Just as a
230 regression with n observations and n variables will give you a good fit in
231 sample, a tree with many branches will also fit the training data well. In
232 either case, predictions using new data, such as the test set, could be very
233 poor.

234 The most common solution to this problem is to “prune” the tree by
235 imposing a cost for complexity. There are various measures of complexity,
236 but a common one is the number of terminal nodes (also known as “leafs.”
237 The cost of complexity is a tuning parameter that is chosen to provide the
238 best out-of-sample predictions, which is typically measured using the 10-fold
239 cross validation procedure mentioned earlier.

240 A typical tree estimation session might involve dividing your data into
241 ten folds, using nine of the folds to grow a tree with a particular complexity,
242 and then predict on the excluded fold. Repeat the estimation with different
243 values of the complexity parameter using other folds and choose the value
244 of the complexity parameter that minimizes the out-of-sample classification
245 error. (Some researchers recommend being a bit more aggressive and advo-
246 cate choosing the complexity parameter that is one standard deviation lower
247 than the loss-minimizing value.)

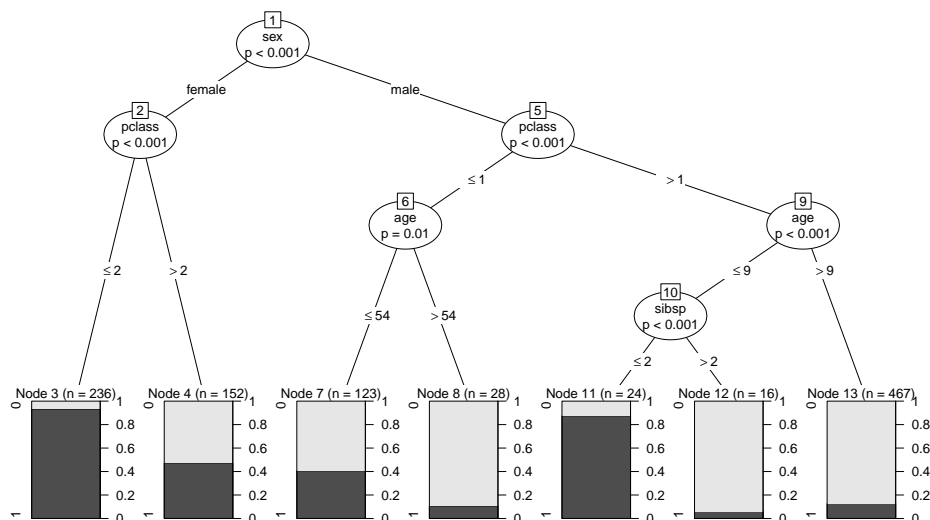


Figure 4: A ctree for survivors of the Titanic. The black bars indicate fraction of the group that survived.

Of course, in practice, the computer program handles most of these details for you. In the examples in this paper I mostly use default choices to keep things simple, but in practice these defaults will often be adjusted by the analyst. As with any other statistical procedure, skill, experience and intuition are helpful in coming up with a good answer. Diagnostics, exploration, and experimentation are just as useful with these methods as with regression techniques.

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

The first node divides by gender. The second node then divides by class. In the right-hand branches, the third node divides by age, and a fourth node divides by the number of siblings and spouses aboard. The bins at

the bottom of the figure show the total number of people in that leaf and a graphical depiction of their survival rate. One might summarize this tree by the following principle: “women and children first . . . particularly if they were traveling first class.” This simple example again illustrates that classification trees can be helpful in summarizing relationships in data, as well as predicting outcomes.³

4.2 Economic example: HMDA data

Munnell et al. [1996] examined mortgage lending in Boston to see if race played a significant role in determining who was approved for a mortgage. The primary econometric technique was a logistic regression where race was included as one of the predictors. The coefficient on race showed a statistically significant negative impact on probability of getting a mortgage for black applicants. This finding prompted considerable subsequent debate and discussion; see Ladd [1998] for an overview.

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

The tree fits pretty well, misclassifying 228 of the 2380 observations for an error rate of 9.6%. By comparison, a simple logistic regression does slightly better, misclassifying 225 of the 2380 observations, leading to an error rate of 9.5%. As you can see in Figure 5, the most important variable is `dmi` = “denied mortgage insurance”. This variable alone explains much of the variation in the data. The race variable (`black`) shows up far down the tree and seems to be relatively unimportant.

One way to gauge whether a variable is important is to exclude it from

³For two excellent tutorials on tree methods that use the Titanic data, see Stephens and Wehrley [2014].

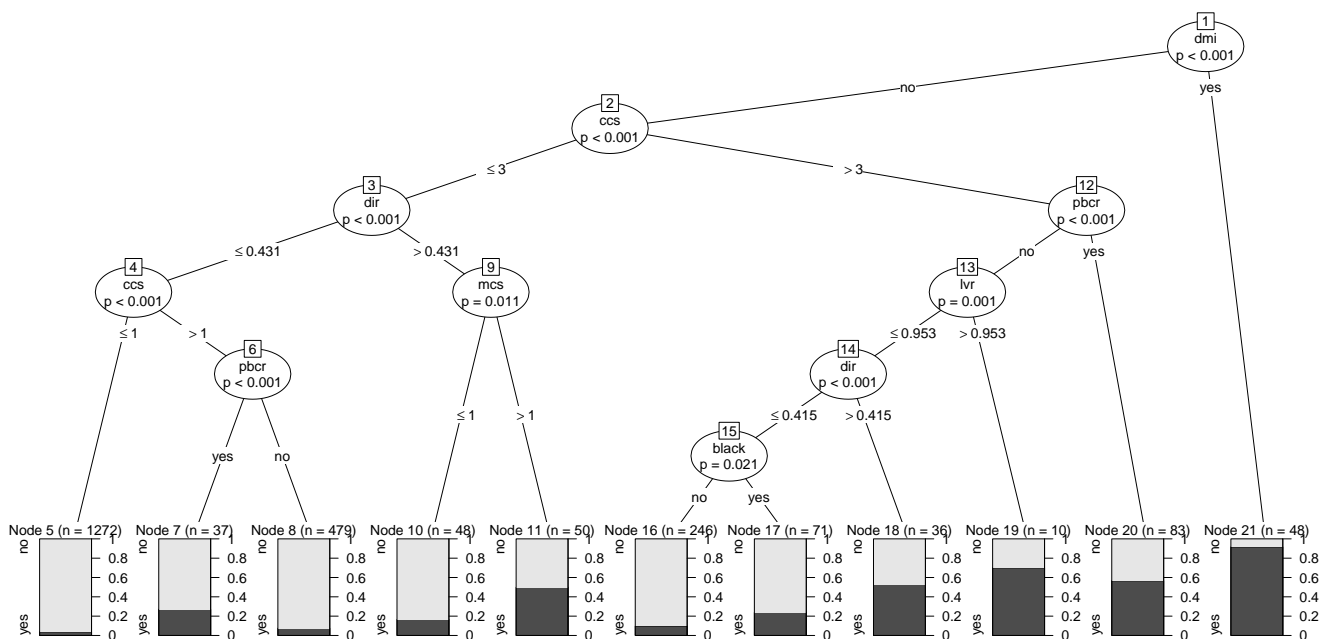


Figure 5: HMDA tree. The black bars indicate the fraction of each group that were denied mortgages. The most important determinant of this is the variable **dmi**, “denied mortgage insurance.”

290 the prediction and see what happens. When this is done, it turns out that
 291 the accuracy of the tree based model doesn’t change at all: exactly the same
 292 cases are misclassified. Of course, it is perfectly possible that there was
 293 racial discrimination elsewhere in the mortgage process, or that some of the
 294 variables included are highly correlated with race. But it is noteworthy that
 295 the tree model produced by standard procedures that omits race fits the
 296 observed data just as well as a model that includes race.

297 5 Boosting, bagging and bootstrap

298 There are several useful ways to improve classifier performance. Interestingly
 299 enough, some of these methods work by *adding* randomness to the data. This

300 seems paradoxical at first, but adding randomness turns out to be a helpful
301 way of dealing with the overfitting problem.

302 **Bootstrap** involves choosing (with replacement) a sample of size n from a
303 dataset of size n to estimate the sampling distribution of some statistic.
304 A variation is the “ m out of n bootstrap” which draws a sample of size
305 m from a dataset of size $n > m$.

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

308 **Boosting** involves repeated estimation where misclassified observations are
309 given increasing weight in each repetition. The final estimate is then a
310 vote or an average across the repeated estimates.⁴

311 Econometricians are well-acquainted with the bootstrap but rarely use the
312 other two methods. Bagging is primarily useful for nonlinear models such
313 as trees. (Friedman and Hall [2005].) Boosting tends to improve predictive
314 performance of an estimator significantly and can be used for pretty much
315 any kind of classifier or regression model, including logits, probits, trees, and
316 so on.

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

320 **Random forests** refers to a technique that uses multiple trees. A typical
321 procedure uses the following steps.

- 322 1. Choose a bootstrap sample of the observations and start to grow
323 a tree.

⁴Boosting is often used with decision trees, where it can dramatically improve their predictive performance.

- 324 2. At each node of the tree, choose a random sample of the predictors
325 to make the next decision. Do not prune the trees.
- 326 3. Repeat this process many times to grow a forest of trees
- 327 4. In order to determine the classification of a new observation, have
328 each tree make a classification and use a majority vote for the
329 final prediction

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

340 One defect of random forests is that they are a bit of a black box—they
341 don’t offer simple summaries of relationships in the data. As we have seen
342 earlier, a single tree can offer some insight about how predictors interact. But
343 a forest of a thousand trees cannot be easily interpreted. However, random
344 forests can determine which variables are “important” in predictions in the
345 sense of contributing the biggest improvements in prediction accuracy.

346 Note that random forests involves quite a bit of randomization; if you
347 want to try them out on some data, I strongly suggest choosing a particular
348 seed for the random number generator so that your results can be reproduced.
349 (See the online supplement for examples.)

350 I ran the random forest method on the HMDA data and found that it
351 misclassified 223 of the 2380 cases, a small improvement over the logit and
352 the ctree. I also used the importance option in random forests to see how

the predictors compared. It turned out that `dmi` was the most important predictor and race was second from the bottom which is consistent with the `ctree` analysis.

6 Variable selection

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

6.1 Lasso and friends

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

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

This estimation method is called *elastic net regression*; it contains three other methods as special cases. If there is no penalty term ($\lambda = 0$), this is *ordinary least squares*. If $\alpha = 1$ so that there is only the quadratic constraint, this is *ridge regression*. If $\alpha = 0$ this is called the *lasso*, an acronym for “least absolute shrinkage and selection operator.”

These penalized regressions are classic examples of regularization. In

374 this case, the complexity is the number and size of predictors in the model.
375 All of these methods tend to shrink the least squares regression coefficients
376 towards zero. The lasso and elastic net typically produces regressions where
377 some of the variables are set to be exactly zero. Hence this is a relatively
378 straightforward way to do variable selection.

379 It turns out that these estimators can be computed quite efficiently, so
380 doing variable selection on reasonably large problems is computationally fea-
381 sible. They also seem to provide good predictions in practice.

382 **6.2 Spike and slab regression**

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

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

396 Now we take a draw of γ from its prior distribution, which will just be a
397 list of variables in the regression. Conditional on this list of included vari-
398 ables, we take a draw from the prior distribution for the coefficients. We
399 combine these two draws with the likelihood in the usual way which gives
400 us a draw from posterior distribution on both probability of inclusion and
401 the coefficients. We repeat this process thousands of times using a Markov
402 Chain Monte Carlo (MCMC) technique which give us a table summarizing

the posterior distribution for γ (indicating variable inclusion), β (the coefficients), and the associated prediction of y . We can summarize this table in a variety of ways. For example, we can compute the average value of γ_p which shows the posterior probability that the variable p is included in the regressions.

6.3 Economic example: growth regressions

We illustrate these different methods of variable selection using data from Sala-i-Martin [1997]. This exercise involved examining a dataset of 72 countries and 42 variables in order to see which variables appeared to be important predictors of economic growth. Sala-i-Martin [1997] computed all possible subsets of regressors of manageable size and used the results to construct an importance measure he called CDF(0). Ley and Steel [2009] investigated the same question using Bayesian Model Averaging (BMA) a technique related to, but not identical with, spike-and-slab, while Hendry and Krolzig [2004] examined an iterative significance test selection method.

Table 4 shows 10 predictors that were chosen by Ley and Steel [2009], Sala-i-Martin [1997], `lasso`, and `spike-and-slab`. The table is based on that in Ley and Steel [2009] but metrics used are not strictly comparable across the various 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 would likely be preferred to exhaustive search. All 4 of these variable selection methods give similar results for the first 4 or 5 variables, after which they diverge. In this particular case, the dataset appears to be too small to resolve the question of what is “important” for economic growth.

predictor	BMA	CDF(0)	lasso	spike-slab
GDP level 1960	1.000	1.000	-	0.9992
Fraction Confucian	0.995	1.000	2	0.9730
Life expectancy	0.946	0.942	-	0.9610
Equipment investment	0.757	0.997	1	0.9532
Sub-Saharan dummy	0.656	1.000	7	0.5834
Fraction Muslim	0.656	1.000	8	0.6590
Rule of law	0.516	1.000	-	0.4532
Open economy	0.502	1.000	6	0.5736
Degree of Capitalism	0.471	0.987	9	0.4230
Fraction Protestant	0.461	0.966	5	0.3798

Table 4: Comparing variable selection algorithms. See text for discussion.

7 Time series

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

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

The challenge is that there are billions of queries so it is hard to determine

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

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

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

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

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

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

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

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

463 It is easy to add an additional state variable for seasonality if that is ap-
 464 propriate. The parameters to estimate are the regression coefficients β and
 465 the variances of (e_{it}) for $i = 1, \dots, 3$. We can then use these estimates to
 466 construct the optimal forecast based on techniques drawn from the literature
 467 on Kalman filters.

468 For the regression we use the spike-and-slab variable choice mechanism
 469 described above. A draw from the posterior distribution now involves a draw
 470 of variances of (e_{1t}, e_{2t}, e_{3t}) , a draw of the vector γ that indicates which vari-
 471 ables are in the regression, and a draw of the regression coefficients β for the
 472 included variables. The draws of μ_t , b_t , and β can be used to construct esti-
 473 mates of y_t and forecasts for y_{t+1} . We end up with an (estimated) posterior
 474 distribution for each parameter of interest. If we seek a point prediction, we

475 can average over these draws, which is essentially a form of Bayesian model
 476 averaging.

477 As an example, consider the non-seasonally adjusted data for new homes
 478 sold in the U.S. (HSN1FNSA) from the St. Louis Federal Reserve Economic
 479 Data. This time series can be submitted to Google Correlate, which then
 480 returns the 100 queries that are the most highly correlated with the series.
 481 We feed that data into the BSTS system which identifies the predictors with
 482 the largest posterior probabilities of appearing in the housing regression are
 483 shown in Figure 6. In these figures, black bars indicate a negative relation-
 484 ship and white bars indicate a positive relationship. Two predictors, [oldies
 485 lyrics] and [www.mail2web] appear to be spurious so we remove them and
 486 re-estimate, yielding the results in Figure 7.

487 The fit is shown in Figure 8 which shows the incremental contribution
 488 of the trend, seasonal, and two two of the regressors. Even with only two
 489 predictors, queries on [appreciate rate] and queries on [irs 1031], we get a
 490 pretty good fit.⁵

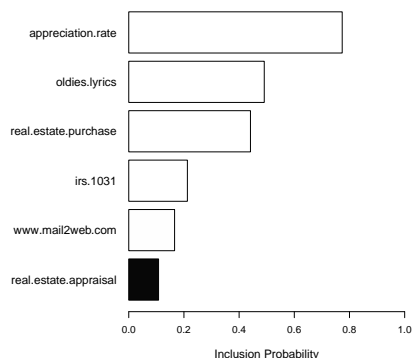


Figure 6: Initial predictors.

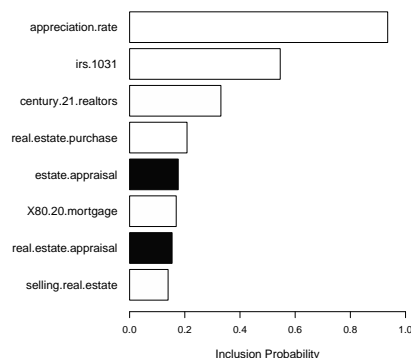


Figure 7: Final predictors.

⁵IRS section 1031 has to do with deferring capital gains on certain sorts of property exchange.

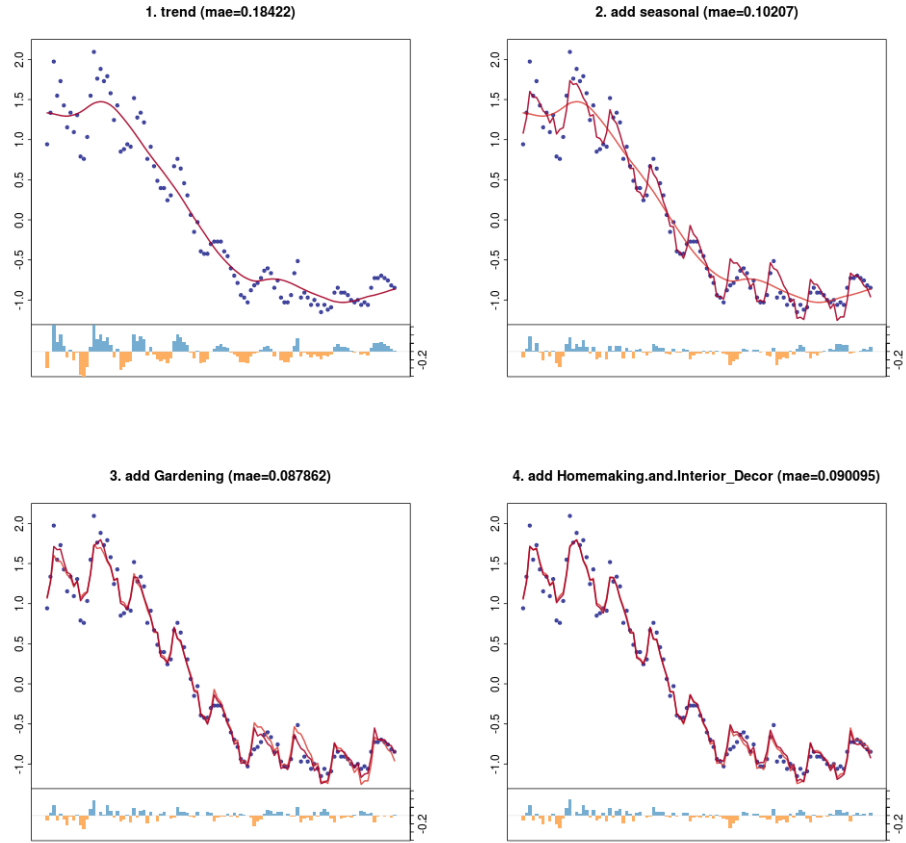


Figure 8: Incremental plots. The plots show the impact of the trend, seasonal, and a few individual regressors. The residuals are shown on the bottom.

491 8 Econometrics and machine learning

492 There are a number of areas where there would be opportunities for fruitful
493 collaboration between econometrics and machine learning. I mentioned above
494 that most machine learning uses IID data. However, the BSTS model shows
495 that some of these techniques can be adopted for time series models. It is
496 also possible to use machine learning techniques to look at panel data and
497 there has been some work in this direction.

498 However, the most important area for collaboration involves causal infer-
499 ence. Econometricians have developed several tools for causal inference such
500 as instrumental variables, regression discontinuity, difference-in-differences
501 and various forms of natural and designed experiments. (Angrist and Krueger
502 [2001].) Machine learning work has, for the most part, dealt with pure pre-
503 diction. In a way this is ironic, since theoretical computer scientists, such
504 as Pearl [2009a,b] have made significant contributions to causal modeling.
505 However, it appears that these theoretical advances have not as yet been
506 incorporated into machine learning practice to a significant degree.

507 8.1 Causality and prediction

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

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

518 To enlarge on this point, let us consider an experiment (natural or de-

signed) that attempts to estimate the impact of some policy, such as adding police to precincts. There are two critical questions.

- How will police be assigned to precincts in both the experiment and the policy implementation? 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, 6) based on asking for volunteers, and so on. Ideally the assignment procedure in the experiment will be similar to that used in the policy. Developing accurate predictions about which precincts will receive additional police under the proposed policy based on the experimental data can clearly be helpful in predicting the expected 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 want to estimate the *causal* impact of some treatment we need to compare the outcome with the intervention to what *would have happened* without the intervention. But this counterfactual cannot be observed, so it must be predicted by some model. The better predictive model you have for the counterfactual, the better you will be able to estimate the causal effect, an observation that is true for both pure experiments and natural experiments.

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

To state this in a slightly more formal way, consider the identity from Angrist and Pischke [2008], page 11:

$$\begin{aligned} \text{observed difference in outcome} &= \text{average treatment effect on the treated} \\ &+ \text{selection bias} \end{aligned}$$

545 If you want to model the average treatment effect as a function of other vari-
546 ables, you will usually need to model both the observed difference in outcome
547 and the selection bias. The better your predictive model for those compo-
548 nents, the better predictions you can make about the average treatment ef-
549 fect. Of course, if you have a true randomized treatment-control experiment,
550 selection bias goes away and those treated are an unbiased random sample
551 of the population.

552 To illustrate these points, let us consider the thorny problem of estimat-
553 ing the causal effect of advertising on sales. (Lewis and Rao [2013].) The
554 difficulty is that there are many confounding variables, such as seasonality or
555 weather, that cause both increased ad exposures and increased purchases by
556 consumers. For example, consider the (probably apocryphal) story about an
557 advertising manager who was asked why he thought his ads were effective.
558 “Look at this chart,” he said. “Every December I increase my ad spend and,
559 sure enough, purchases go up.” Of course, in this case seasonality can be
560 included in the model. However, generally there will be other confounding
561 variables that affect both exposure to ads and the propensity of purchase,
562 which makes causal interpretations of observed relationships problematic.

563 The ideal way to estimate advertising effectiveness is, of course, to run a
564 controlled experiment. In this case the control group provides an estimate
565 of the counterfactual: what would have happened without ad exposures.
566 But this ideal approach can be quite expensive, so it is worth looking for
567 alternative ways to predict the counterfactual. One way to do this is to use
568 the Bayesian Structural Time Series method described earlier.

569 Suppose a given company wants to determine the impact of an advertising
570 campaign on its sales. It first uses BSTS (or some other technique) to build
571 a model predicting the time series of sales as a function its past history,
572 seasonal effects and other possible predictors such as Google queries on its
573 company name, its competitors’ names, or products that it produces. Since
574 there are many possible choices for predictors, it is important to use some

575 variable selection mechanism such as those described earlier.

576 It next runs an ad campaign for a few weeks and records sales during
577 this period. Finally, it makes a forecast of what sales *would have been* in
578 the absence of the ad campaign using the model developed in the first stage.
579 Comparing the actual outcome to the counterfactual outcome gives us an
580 estimate of causal effect of advertising.

581 Figure 9 shows the outcome of such a procedure. It is based on the
582 approach proposed in Brodersen et al. [2013], but where the covariates are
583 chosen automatically from Google Trends categories using BSTS. Panel *a*
584 shows the actual sales and the prediction of what the sales would have been
585 without the campaign based on the BSTS forecasting model. Panel *b* shows
586 the difference between actual and predicted sales, and Panel *c* shows the
587 cumulative difference. It is clear from this figure that there was a significant
588 causal impact of advertising which can then be compared to the cost of the
589 advertising to evaluate the campaign.

590 This procedure does not use a control group in the conventional sense.
591 Rather it uses a general time series model based on trend extrapolation,
592 seasonal effects, and relevant covariates to forecast the what would have
593 happened without the ad campaign.

594 A good predictive model can be *better* than a randomly-chosen control
595 group, which is usually thought to be the gold standard. For example, sup-
596 pose that you run an ad campaign in 100 cities and retain 100 cities as a
597 control. After the experiment is over, you discover the weather was dramat-
598 ically different across the cities in the study. Should you add weather as a
599 predictor of the counterfactual? Of course! If weather affects sales (which it
600 does) then you will get a more accurate prediction of the counterfactual and
601 thus a better estimate of the causal effect of advertising.

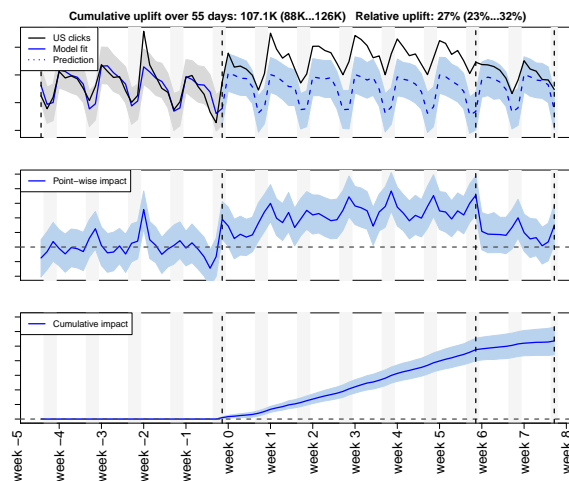


Figure 9: Actual and predicted sales.

9 Model uncertainty

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

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

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

However, I think that model uncertainty has crept in to applied econo-

619 metrics through the back door. Many papers in applied econometrics present
620 regression results in a table with several different specifications: which vari-
621 ables are included in the controls, which variables are used as instruments,
622 and so on. The goal is usually to show that the estimate of some interesting
623 parameter is not very sensitive to the exact specification used.

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

629 In this period of “big data” it seems strange to focus on *sampling uncer-*
630 *tainty*, which tends to be small with large datasets, while completely ignoring
631 *model uncertainty* which may be quite large. One way to address this is to
632 be explicit about examining how parameter estimates vary with respect to
633 choices of control variables and instruments.

634 10 Summary and further reading

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

641 I have already mentioned Hastie et al. [2009] which has detailed descrip-
642 tions of all the methods discussed here but at a relatively advanced level.
643 James et al. [2013] describes many of the same topics at an undergraduate-
644 level, along with R code and many examples.⁶ Murphy [2012] examines ma-

⁶There are several economic examples in the book where the tension between predictive modeling and causal inference is apparent.

chine learning from a Bayesian point of view.

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

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>.
- Joshua D. Angrist and Jörn-Steffen Pischke. *Mostly Harmless Econometrics*. Princeton University Press, 2008.
- 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. Inferring causal impact using Bayesian structural

671 time series models. Technical report, Google, Inc., 2013. URL [http:](http://research.google.com/pubs/pub41854.html)
672 [//research.google.com/pubs/pub41854.html](http://research.google.com/pubs/pub41854.html).

673 Yan Carrière-Swallow and Felipe Labbé. Nowcasting with Google Trends in
674 an emerging market. *Journal of Forecasting*, 2011. doi: 10.1002/for.1252.
675 URL <http://ideas.repec.org/p/chb/bcchwp/588.html>. Working Pa-
676 pers Central Bank of Chile 588.

677 Rich Caruana and Alexandru Niculescu-Mizil. An empirical comparison of
678 supervised learning algorithms. In *Proceedings of the 23rd International*
679 *Conference on Machine Learning*, Pittsburgh, PA, 2006.

680 Jennifer L. Castle, Xiaochuan Qin, and W. Robert Reed. How to pick the
681 best regression equation: A review and comparison of model selection algo-
682 rithms. Technical Report 13/2009, Department of Economics, University
683 of Canterbury, 2009. URL [http://www.econ.canterbury.ac.nz/RePEc/](http://www.econ.canterbury.ac.nz/RePEc/cbt/econwp/0913.pdf)
684 [cbt/econwp/0913.pdf](http://www.econ.canterbury.ac.nz/RePEc/cbt/econwp/0913.pdf).

685 Hyunyoung Choi and Hal Varian. Predicting the present with Google Trends.
686 Technical report, Google, 2009. URL [http://google.com/googleblogs/](http://google.com/googleblogs/pdfs/google_predicting_the_present.pdf)
687 [pdfs/google_predicting_the_present.pdf](http://google.com/googleblogs/pdfs/google_predicting_the_present.pdf).

688 Hyunyoung Choi and Hal Varian. Predicting the present with Google Trends.
689 *Economic Record*, 2012. URL [http://people.ischool.berkeley.edu/](http://people.ischool.berkeley.edu/~hal/Papers/2011/ptp.pdf)
690 [~hal/Papers/2011/ptp.pdf](http://people.ischool.berkeley.edu/~hal/Papers/2011/ptp.pdf).

691 Pedro Domingos. A few useful things to know about machine learning. *Com-*
692 *munications of the ACM*, 55(10), October 2012. URL [http://homes.cs.](http://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf)
693 [washington.edu/~pedrod/papers/cacm12.pdf](http://homes.cs.washington.edu/~pedrod/papers/cacm12.pdf).

694 Liran Einav and Jonathan Levin. The data revolution and economic analysis.
695 Technical report, NBER Innovation Policy and the Economy Conference,
696 2013.

- 697 Andrey Feuerverger, Yu He, and Shashi Khatri. Statistical significance of
698 the Netflix challenge. *Statistical Science*, 27(2):202–231, 2012. URL <http://arxiv.org/abs/1207.5649>.
699
- 700 Jerome Friedman. Stochastic gradient boosting. Technical report, Stan-
701 ford University, 1999. URL [http://www-stat.stanford.edu/~jhf/ftp/](http://www-stat.stanford.edu/~jhf/ftp/stobst.pdf)
702 [stobst.pdf](http://www-stat.stanford.edu/~jhf/ftp/stobst.pdf).
- 703 Jerome Friedman and Peter Hall. On bagging and nonlinear estimation.
704 Technical report, Stanford University, 2005. URL [http://www-stat.](http://www-stat.stanford.edu/~jhf/ftp/bag.pdf)
705 [stanford.edu/~jhf/ftp/bag.pdf](http://www-stat.stanford.edu/~jhf/ftp/bag.pdf).
- 706 Jerome Friedman and Bogdan E. Popescu. Predictive learning via rule
707 ensembles. Technical report, Stanford University, 2005. URL [http://www-stat.](http://www-stat.stanford.edu/~jhf/R-RuleFit.html)
708 [stanford.edu/~jhf/R-RuleFit.html](http://www-stat.stanford.edu/~jhf/R-RuleFit.html).
- 709 Sharad Goel, Jake M. Hofman, Sbastien Lahaie, David M. Pennock, and
710 Duncan J. Watts. Predicting consumer behavior with web search. *Pro-*
711 *ceedings of the National Academy of Sciences*, 2010. URL [http://www.](http://www.pnas.org/content/107/41/17486.full)
712 [pnas.org/content/107/41/17486.full](http://www.pnas.org/content/107/41/17486.full).
- 713 Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of*
714 *Statistical Learning: Data Mining, Inference, and Prediction*. Springer-
715 Verlag, 2 edition, 2009. URL [http://www-stat.stanford.edu/~tibs/](http://www-stat.stanford.edu/~tibs/ElemStatLearn/download.html)
716 [ElemStatLearn/download.html](http://www-stat.stanford.edu/~tibs/ElemStatLearn/download.html).
- 717 Rebecca Hellerstein and Menno Middelorp. Forecasting with
718 internet search data. *Liberty Street Economics Blog of the*
719 *Federal Reserve Bank of New York*, January 2012. URL
720 [http://libertystreeteconomics.newyorkfed.org/2012/01/](http://libertystreeteconomics.newyorkfed.org/2012/01/forecasting-with-internet-search-data.html)
721 [forecasting-with-internet-search-data.html](http://libertystreeteconomics.newyorkfed.org/2012/01/forecasting-with-internet-search-data.html).
- 722 David F. Hendry and Hans-Martin Krolzig. We ran one regression. *Oxford*
723 *Bulletin of Economics and Statistics*, 66(5):799–810, 2004.

- 724 Torsten Hothorn, Kurt Hornik, and Achim Zeileis. Unbiased recursive par-
725 titioning: A conditional inference framework. *Journal of Computational*
726 *and Graphical Statistics*, 15(3):651–674, 2006.
- 727 Jeremy Howard. The two most important algorithms in predictive mod-
728 eling today. Conference presentation, February 2013. URL [http://](http://strataconf.com/strata2012/public/schedule/detail/22658)
729 strataconf.com/strata2012/public/schedule/detail/22658.
- 730 Bill Howe. Introduction to data science. Technical report, University of
731 Washington, 2013. URL [https://class.coursera.org/datasci-001/](https://class.coursera.org/datasci-001/lecture/index)
732 [lecture/index](https://class.coursera.org/datasci-001/lecture/index).
- 733 Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An*
734 *Introduction to Statistical Learning with Applications in R*. Springer, New
735 York, 2013.
- 736 Helen F. Ladd. Evidence on discrimination in mortgage lending. *Journal of*
737 *Economic Perspectives*, 12(2):41–62, 1998.
- 738 Jeff Leek. Data analysis, 2013. URL [http://blog.revolutionanalytics.](http://blog.revolutionanalytics.com/2013/04/coursera-data-analysis-course-videos.html)
739 [com/2013/04/coursera-data-analysis-course-videos.html](http://blog.revolutionanalytics.com/2013/04/coursera-data-analysis-course-videos.html).
- 740 Randall A. Lewis and Justin M. Rao. On the near impossibility of mea-
741 suring the returns to advertising. Technical report, Google, Inc. and
742 Microsoft Research, 2013. URL [http://justinmrao.com/lewis_rao_](http://justinmrao.com/lewis_rao_nearimpossibility.pdf)
743 [nearimpossibility.pdf](http://justinmrao.com/lewis_rao_nearimpossibility.pdf).
- 744 Eduardo Ley and Mark F. J. Steel. On the effect of prior assumptions in
745 Bayesian model averaging with applications to growth regression. *Jour-*
746 *nal of Applied Econometrics*, 24(4):651–674, 2009. URL [http://ideas.](http://ideas.repec.org/a/jae/japmet/v24y2009i4p651-674.html)
747 [repec.org/a/jae/japmet/v24y2009i4p651-674.html](http://ideas.repec.org/a/jae/japmet/v24y2009i4p651-674.html).
- 748 Nick McLaren and Rachana Shanbhoge. Using internet search data
749 as economic indicators. *Bank of England Quarterly Bulletin*,

750 June 2011. URL [http://www.bankofengland.co.uk/publications/](http://www.bankofengland.co.uk/publications/quarterlybulletin/qb110206.pdf)
751 [quarterlybulletin/qb110206.pdf](http://www.bankofengland.co.uk/publications/quarterlybulletin/qb110206.pdf).

752 James N. Morgan and John A. Sonquist. Problems in the analysis of survey
753 data, and a proposal. *Journal of the American Statistical Association*, 58
754 (302):415–434, 1963. URL <http://www.jstor.org/stable/2283276>.

755 Alicia H. Munnell, Geoffrey M. B. Tootell, Lynne E. Browne, and James
756 McEneaney. Mortgage lending in Boston: Interpreting HDMA data. *Amer-*
757 *ican Economic Review*, pages 25–53, 1996.

758 Kevin P. Murphy. *Machine Learning A Probabalistic Perspective*. MIT Press,
759 2012. URL <http://www.cs.ubc.ca/~murphyk/MLbook/>.

760 Judea Pearl. *Causality*. Cambridge University Press, 2009a.

761 Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*,
762 4:96–146, 2009b.

763 Claudia Perlich, Foster Provost, and Jeffrey S. Simonoff. Tree induction vs.
764 logistic regression: A learning-curve analysis. *Jounral of Machine Learning*
765 *Research*, 4:211–255, 2003. URL [http://machinelearning.wustl.edu/](http://machinelearning.wustl.edu/mlpapers/paper_files/PerlichPS03.pdf)
766 [mlpapers/paper_files/PerlichPS03.pdf](http://machinelearning.wustl.edu/mlpapers/paper_files/PerlichPS03.pdf).

767 Donald Rubin. Estimating causal effects of treatment in randomized and non-
768 randomized studies. *Journal of Educational Psychology*, 66(5):689, 1974.

769 Xavier Sala-i-Martín. I just ran two million regressions. *American Economic*
770 *Review*, 87(2):178–83, 1997.

771 Steve Scott and Hal Varian. Bayesian variable selection for nowcasting
772 economic time series. Technical report, Google, 2012a. URL [http:](http://www.ischool.berkeley.edu/~hal/Papers/2012/fat.pdf)
773 [//www.ischool.berkeley.edu/~hal/Papers/2012/fat.pdf](http://www.ischool.berkeley.edu/~hal/Papers/2012/fat.pdf). Presented
774 at JSM, San Diego.

775 Steve Scott and Hal Varian. Predicting the present with Bayesian structural
776 time series. Technical report, Google, 2012b. URL [http://www.ischool.](http://www.ischool.berkeley.edu/~hal/Papers/2013/pred-present-with-bsts.pdf)
777 [berkeley.edu/~hal/Papers/2013/pred-present-with-bsts.pdf](http://www.ischool.berkeley.edu/~hal/Papers/2013/pred-present-with-bsts.pdf).

778 Mark F. J. Steel. Bayesian model averaging and forecasting. *Bulletin*
779 *of E.U. and U.S. Inflation and Macroeconomic Analysis*, 200:30–41,
780 2011. URL [http://www2.warwick.ac.uk/fac/sci/statistics/staff/](http://www2.warwick.ac.uk/fac/sci/statistics/staff/academic-research/steel/steel_homepage/publ/bma_forecast.pdf)
781 [academic-research/steel/steel_homepage/publ/bma_forecast.pdf](http://www2.warwick.ac.uk/fac/sci/statistics/staff/academic-research/steel/steel_homepage/publ/bma_forecast.pdf).

782 Revor Stephens and Curt Wehrley. Getting started with R. *Kaggle*, 2014.
783 URL [https://www.kaggle.com/c/titanic-gettingStarted/details/](https://www.kaggle.com/c/titanic-gettingStarted/details/new-getting-started-with-r)
784 [new-getting-started-with-r](https://www.kaggle.com/c/titanic-gettingStarted/details/new-getting-started-with-r).

785 Danny Sullivan. Google: 100 billion searches per month, search to integrate
786 gmail, launching enhanced search app for iOS. *Search Engine Land*, 2012.
787 URL <http://searchengineland.com/google-search-press-129925>.

788 W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer-
789 Verlag, New York, 4 edition, 2002.

790 Graham Williams. *Data Mining with Rattle and R*. Springer, New York,
791 2011.

792 Xindong Wu and Vipin Kumar, editors. *The Top Ten Algorithms in*
793 *Data Mining*. CRC Press, 2009. URL [http://www.cs.uvm.edu/~icdm/](http://www.cs.uvm.edu/~icdm/algorithms/index.shtml)
794 [algorithms/index.shtml](http://www.cs.uvm.edu/~icdm/algorithms/index.shtml).