

# Machine Learning Methods for Demand Estimation

By PATRICK BAJARI, DENIS NEKIPELOV, STEPHEN P. RYAN, AND MIAOYU YANG\*

Over the past decade, there has been a high level of interest in modeling consumer behavior in the fields of computer science and statistics. These applications are motivated in part by the availability of large data sets, and are commonly used by firms in the retail, health care, and internet industries to improve business decisions. In this paper, we compare these methods to standard econometric models that are used by practitioners to study demand.<sup>1</sup> We are motivated by the problem of finding practical tools that would be of use to applied econometricians in estimating demand with large numbers of observations and covariates, such as in a scanner panel data set.

Many economists are unfamiliar with these methods, so briefly sketch several commonly-used techniques from the machine learning literature.<sup>2</sup> We consider eight different models that can be used for estimating demand: linear regression, the conditional logit, and six machine learning methods, all of which differ from standard approaches by combining an element of model selection into the estimation procedure. Several of these models can be seen as variants on regularization schemes, which reduce the number of covariates in a regression which receive non-zero coefficients, such as stepwise regression, forward stage-wise regression, LASSO, and support vector machines. We also consider two models

based on regression trees, which are flexible methods for approximating arbitrary functions: bagging and random forests. While these models may be unfamiliar to many economists, they are surprisingly simple and are based on underlying methods that will be quite familiar. Also, all of the methods that we use are supported in statistical packages. We perform our computations in the open source software package R. Therefore, application of these methods will not require writing complex code from scratch. However, applied econometricians may have to familiarize themselves with alternative software.

We apply our method to a canonical demand estimation problem. We use data from IRI Marketing Research (Bronnenberg, Kruger and Mela, 2008) via an academic license at the University of Chicago. It contains scanner panel data from grocery stores within one grocery store chain for six years. We used sales data on salty snacks, which is one of the categories provided in the IRI data. The number of observations are 837,460, which includes 3,149 unique products.

If we allow for product and store level fixed effects, our model effectively has many thousands of explanatory variables. Therefore, variable selection will be an important problem. If we included all of these variables in a standard regression model, the parameters would be poorly estimated. Also, many of the regressors will be multicollinear which will make the models predict poorly out of sample.

In our results, we find that the six models we use from the statistics and computer science literature predict demand out of sample in standard metrics much more accurately than a panel data or logistic model. We do not claim that these models dominate all methods proposed in the voluminous demand estimation literature. Rather,

\* Bajari: University of Washington and NBER, Department of Economics, 331 Savery Hall, Seattle, WA 98195; bajari@uw.edu. Nekipelov: University of Virginia, Department of Economics, 254 Monroe Hall, Charlottesville, VA 22904; dn4w@virginia.edu. Ryan: University of Texas at Austin and NBER, Department of Economics, 2225 Speedway Stop C3100, BRB 3.134D, Austin, TX 78712; sryan@utexas.edu. Yang: University of Washington, Department of Economics, 331 Savery Hall, Seattle, WA 98195; yangmiao@uw.edu.

<sup>1</sup>Hastie, Tibshirani and Friedman (2009) is a comprehensive reference.

<sup>2</sup>See, for example, Belloni, Chernozhukov and Hansen (2014) and Varian (2014).

we claim that as compared to common methods an applied econometrician might use in off the shelf statistical software, these methods are considerably more accurate. Also, the methods that we propose are all available in the well documented, open software package R as well as commercially-available software.

Finally, we propose using an idea dating back at least to Bates and Granger (1969). We treat each of these eight independent predictions as regressors and form a combined model by regressing the dependent variable on to the prediction of each component model. We use a three-way cross validation to avoid overfitting the models in practice. We split the sample into three disjoint sets; we use the first set to fit all eight models, we use the second set to fit our regression on the eight independent model predictions, and we use the third set of the data to test the fit out of sample. We find that this combination procedure can lead to substantial improvements in fit with little additional work. And, as we detail in Bajari et al. (2014), the combined model exhibits standard asymptotic behavior, even though the component models may not, which simplifies the construction of standard errors.

### I. Summary of Machine Learning Methods

We explore using machine learning techniques to predict demand. We briefly discuss each method in turn before applying them to estimation of demand using a scanner panel data set.

A typical specification for demand of product  $j$  in group  $h$  in market  $m$  at time  $t$  would be:

$$(1) \quad Y_{jht} = f(\mathbf{X}, \mathbf{D}, \mathbf{p})' \beta + \zeta_{hm} + \eta_{mt} + \epsilon_{jmt},$$

where  $f$  generates arbitrary interactions between the observables ( $\mathbf{X}$ ), demographics ( $\mathbf{D}$ ), and prices ( $\mathbf{p}$ ). Such a model may have thousands of right-hand side variables; an extreme example from Rajaraman and Ullman (2011) notes Google estimates the demand for a given webpage by using a model of the network structure of literally billions

of other webpages on the right-hand side. Dummy variables on nests are captured by  $\zeta_{hm}$ . Seasonality is captured by the term  $\eta_{mt}$ , which varies by time (say, quarters) across markets.

In ordinary least squares (OLS), the parameters of Equation 1 are typically estimated using the closed-form formula  $\beta = (\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{Y})$ . This formula requires an inversion of  $(\mathbf{X}'\mathbf{X})$ , imposing a rank and order condition on the matrix  $\mathbf{X}$ . We highlight this because in many settings, the number of right-hand side variables can easily exceed the number of observations. Even in the simplest univariate model, one can saturate the right-hand side by using a series of basis functions of  $\mathbf{X}$ . This restriction requires the econometrician to make choices about which variables to include in the regression. We will return to this below, as some of machine learning methods we discuss below allow the econometrician to skirt the order condition by combining model selection and estimation simultaneously.

A large literature on differentiated products has focused on logit-type models, where the idiosyncratic error term is assumed to be distributed as a Type I Extreme Value. Under that restriction, market shares are given by:

$$(2) \quad s_{jht} = \frac{\exp(\theta' \mathbf{X}_{jht})}{\sum_{k \in J} \exp(\theta' \mathbf{X}_{kht})}.$$

Quantities are then computed by multiplying through by market size.<sup>3</sup>

Stepwise regression starts with the intercept as the base model on a set of demeaned covariates. The algorithm then searches over the set of covariates, selects the one with the highest correlation with the residual, and adds that variable to the next model. The method then estimates OLS using that subset of covariates, then repeats the search for the covariate with the next highest correlation. The procedure produces a series of nested models and runs

<sup>3</sup>Berry, Levinsohn and Pakes (1995) extends this model to deal with unobserved heterogeneity and vertical characteristics that are observed to both the firm and consumers.

until no covariates have a sufficiently high correlation with the error term.

Forward stagewise regression is a variant on the stepwise regression. Whereas all of the coefficients can change at each step in the stepwise regression, forward stagewise regression only updates one coefficient at each step. The method finds the variable with the highest correlation with the error term and adds that covariance to the coefficient. This continues until none of the covariates have any correlation with the error term.

These methods build up the model over time in addition to estimating a fit. One advantage of this approach is that the methods can recover the true data-generating process when the number of covariates is larger than the number of observations and the true model is sparse, e.g. the number of coefficients with true non-zero values is less than the number of observations.

Support vector machines (SVM) are a penalized method of regression, using the following:

$$(3) \quad \min_{\beta} \sum_{i=1}^n V(y_i - \mathbf{X}_i' \beta) + \frac{\lambda}{2} \|\beta\|,$$

where the penalty function is:

$$(4) \quad V_{\epsilon}(r) = \begin{cases} 0 & \text{if } |r| < \epsilon, \\ |r| - \epsilon & \text{otherwise.} \end{cases}$$

The tuning parameter,  $\epsilon$ , controls which errors are included in the regression. Errors of sufficiently small size are treated as zeros. Typically only a partial set of the covariates are assigned a non-zero value in SVM regression.

LASSO is another penalized regression method. The regression is given by:

$$(5) \quad \min_{\beta} \frac{1}{2} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \left( t - \sum_{j=1}^p |\beta_j| \right),$$

where  $t$  is the tuning parameter governing

how strictly additional regressors are penalized. LASSO typically results in a number of covariates being given zero weights.

Regression trees approximate functions by partitioning the characteristic space into a series of hyper-cubes and reporting the average value of the function in each of those partition. Regression trees generalize fixed effects to allow them to depend on values of  $\mathbf{X}$ . In the limit as the hypercubes grow infinitesimally small, the tree reports the average value  $Y = f(\mathbf{X} = \mathbf{x})$ , which is a perfect reconstruction of the underlying function  $f$ . In practice, the tree is expanded until the reduction in squared prediction error falls under some threshold. Often, the tree is grown until a specific number of splits are achieved.

The literature has proposed several variations on the regression tree estimator. One is bagging (Breiman, 1996), which uses resampling and model combination to obtain a predictor. The idea is to sample the data with replacement  $B$  times, train a regression tree on each resampled set of data, and then predict the outcome at each  $x$  through a simple average of the prediction under each of the  $B$  trees.

A second approach, which we have found to work exceptionally well in practice, are random forests, as in Breiman (2001). Random forests expand on the idea of using collections of predictors to make predictions by introducing randomness into the set of variables which are considered at node level for splitting. Before each split, only  $m \leq p$  of the explanatory variables are included in the split search. Repeating this across  $B$  trees results in a forest of random trees. The regression predictor for the true function is then:

$$(6) \quad \hat{f}_{rf}^B(\mathbf{x}) = \frac{1}{B} \sum_{b=1}^B T_b(\mathbf{x}).$$

Trees of sufficient size can be unbiased but exhibit high variance, and therefore may benefit from averaging.

## II. Empirical Application

This section compares econometric models with machine learning ones using a typical demand estimation scenario—grocery store sales. We find that the machine learning models in general produce better out-of-sample fits than linear models without loss of in-sample goodness of fit. If we combine all the models linearly with non-negative weights, the resulting combination of models produces better out-of-sample fit than any model in the combination.

The data we use is provided by IRI Marketing Research via an academic license at the University of Chicago. It contains scanner panel data from grocery stores within one grocery store chain for six years. We used sales data on salty snacks, which is one of the categories in the IRI data. A unit of observation is product  $j$ , uniquely defined by a UPC (Universal Product Code), in store  $m$  at week  $t$ . The number of observations are 1,510,563, which includes 3,149 unique products. Let  $q_{jmt}$  be the number of bags of salty snack  $j$  sold in store  $m$  at week  $t$ . If  $q_{jmt} = 0$ , we do not know if it is due to no sale or out-of-stock and the observation is not filled in. The price  $p_{jmt}$  is defined as the quantity weighted average of prices for product  $j$  in store  $m$  at week  $t$ . Therefore if  $q_{jmt} = 0$ , the weight is also 0. In addition to price and quantity, the data contains attributes of the products (such as brand, volume, flavor, cut type, cooking method, package size, fat and salt levels) and promotional variables (promotion, display and feature).

The response variable is log of quantity sold per week. The covariates are log of price, product attributes variables, promotional variables, store fixed effects, and week fixed effects. We provide the same covariate matrix to all of the models except for the logit model, where all the fixed effects are excluded.<sup>4</sup>

In order to estimate and compare models, we split our data into three sets: train-

ing, validation, and holdout. We estimate the model using the training set, and then use the validation set to assign weights to each model when building the combined model. This approach mitigates overfitting in the training model; for example, the linear model tends to get a good in-sample fit but a bad out-of-sample fit, and granting these model a large in-sample weight would produce poor predictions in the holdout sample. Finally, we each model to predict fit in a holdout sample. 25 percent of the data is used as the holdout sample, 15 percent is used as the validate set, and the remaining 60 percent is used as the training set.

Table 1 reports the root mean squared prediction error (RMSE) across the validation and out-of-sample data sets, along with the estimated weights of each model in the combined model. In the scenario of out-of-sample prediction error, the best two models are random forest and support vector machine. The combined model, where we regress the actual value of the response variable on a constrained linear model of the predictions from eight models, outperforms all the eight models, which follows the optimal combination of forecasts in Bates and Granger (1969). Random forest receives the largest weight in the combined model (65.6 percent), and the stepwise and SVM models receive the majority of the rest. It is interesting to observe the combined model does not simply choose the submodel with the best RMSE; there are important covariances among the models which generate better fit in combination than any one given submodel.

## III. Conclusion

In this paper, we review and apply several popular methods from the machine learning literature to the problem of demand estimation. Machine learning models bridge the gap between parametric models with user-selected covariates and completely non-parametric approaches. We demonstrate that these methods can produce superior predictive accuracy as compared to a standard linear regression or

<sup>4</sup>For brevity, we have minimized the description of the data set and details of implementation of the machine learning methods; Bajari et al. (2014) provides more details about the construction of the data set.

TABLE 1—MODEL COMPARISON: PREDICTION ERROR

	Validation		Out-of-Sample		Percent Weight
	RMSE	Std. Err.	RMSE	Std. Err.	
Linear	1.169	0.022	1.193	0.020	6.62
Stepwise	0.983	0.012	1.004	0.011	12.13
Forward Stagewise	0.988	0.013	1.003	0.012	0.00
LASSO	1.178	0.017	1.222	0.012	0.00
Random Forest	0.943	0.017	0.965	0.015	65.56
SVM	1.046	0.024	1.068	0.018	15.69
Bagging	1.355	0.030	1.321	0.025	0.00
Logit	1.190	0.020	1.234	0.018	0.00
Combined	0.924		0.946		100.00

logit model. We also show that a linear combination of the underlying models can improve fit even further with very little additional work. While these methods are not yet commonly used in economics, we think that practitioners will find value in the flexibility, ease-of-use, and scalability of these methods to a wide variety of applied settings.

One concern has been the relative paucity of econometric theory for machine learning models. In related work (Bajari et al., 2014), we provide asymptotic theory results for rates of convergence of the underlying machine learning models. We show that while several of the machine learning models have non-standard asymptotics, with slower-than-parametric rates of convergence, the model formed by combining estimates retains standard asymptotic properties. This simplifies the construction of standard errors for both parameters and predictions, making the methods surveyed here even more accessible for the applied practitioner.

## REFERENCES

- Bajari, Patrick, Denis Nekipelov, Stephen P. Ryan, and Miaoyu Yang.** 2014. “Demand Estimation with Machine Learning and Model Combination.” Working Paper, University of Texas at Austin.
- Bates, John M, and Clive WJ Granger.** 1969. “The combination of forecasts.” *Or*, 451–468.
- Belloni, Alexandre, Victor Chernozhukov, and Christian Hansen.** 2014. “High-Dimensional Methods and Inference on Structural and Treatment Effects.” *The Journal of Economic Perspectives*, 28(2): 29–50.
- Berry, Steven, James Levinsohn, and Ariel Pakes.** 1995. “Automobile Prices in Equilibrium.” *Econometrica*, 63(4): 841–90.
- Breiman, Leo.** 1996. “Bagging predictors.” *Machine learning*, 24(2): 123–140.
- Breiman, Leo.** 2001. “Random forests.” *Machine learning*, 45(1): 5–32.
- Bronnenberg, Bart J, Michael W Kruger, and Carl F Mela.** 2008. “Database paper-The IRI marketing data set.” *Marketing Science*, 27(4): 745–748.
- Hastie, Trevor, Robert Tibshirani, and Jerome Friedman.** 2009. *The elements of statistical learning*. Vol. 2, Springer.
- Rajaraman, Anand, and Jeffrey D Ullman.** 2011. “Mining of Massive Datasets.” *Lecture Notes for Stanford CS345A Web Mining*, 67(3): 328.
- Varian, Hal R.** 2014. “Big data: New tricks for econometrics.” *The Journal of Economic Perspectives*, 28(2): 3–27.