



Forecasting economic time series using targeted predictors

Jushan Bai^{a,b}, Serena Ng^{c,*}

^a Department of Economics, NYU, 19 W 4th Street, New York, NY 10012, United States

^b School of Economics and Management, Tsinghua University, Beijing, China

^c Department of Economics, Columbia University, 420 W. 118 St., New York 10027, United States

ARTICLE INFO

Article history:

Available online 28 August 2008

Keywords:

Diffusion index

Factor models

LASSO

LARS

Hard thresholding

ABSTRACT

This paper studies two refinements to the method of factor forecasting. First, we consider the method of quadratic principal components that allows the link function between the predictors and the factors to be non-linear. Second, the factors used in the forecasting equation are estimated in a way to take into account that the goal is to forecast a specific series. This is accomplished by applying the method of principal components to ‘targeted predictors’ selected using hard and soft thresholding rules. Our three main findings can be summarized as follows. First, we find improvements at all forecast horizons over the current diffusion index forecasts by estimating the factors using fewer but informative predictors. Allowing for non-linearity often leads to additional gains. Second, forecasting the volatile one month ahead inflation warrants a high degree of targeting to screen out the noisy predictors. A handful of variables, notably relating to housing starts and interest rates, are found to have systematic predictive power for inflation at all horizons. Third, the targeted predictors selected by both soft and hard thresholding changes with the forecast horizon and the sample period. Holding the set of predictors fixed as is the current practice of factor forecasting is unnecessarily restrictive.

© 2008 Elsevier B.V. All rights reserved.

1. Introduction

In recent years, the method of “diffusion index forecasts”, also known as factor augmented forecasts, has received the attention of both econometricians and practitioners. In the diffusion index forecasting methodology, the factors are first estimated from a large number of predictors, (X_{1t}, \dots, X_{Nt}) , by the method of principal components, and then augmented to a linear forecasting equation for y_{t+h} that consists of lags of y and other predictors. What makes the diffusion index (DI) methodology appealing is its capacity to incorporate information in a large number of predictors into the forecast in a simple and parsimonious way. However, this does not preclude refinements to the DI methodology. In particular, the methodology as it stands does not take into account the predictive ability of X_{it} for y_{t+h} when the factors are estimated. Furthermore, the framework is now confined to a linear relation between the predictors and the series to be forecasted.

Our goal is first to go beyond the linear principal components framework to permit a more flexible factor structure, and more importantly, to use only those predictors informative for y in estimating the factors. To this end, we consider two possible improvements to the DI framework. First, we allow the factors

to be non-linearly related to the predictors by expanding the set of predictors to include non-linear functions of the observed variables. Second, and independently of whether non-linearity is being considered, we take explicit account that the object of interest is ultimately the forecast of y . Accordingly, we form principal components using a subset of those predictors that are tested to have predictive power for y . As this set of predictors change with y , we refer to these as ‘targeted predictors’. Our approach therefore entertains more predictors than the current implementation of DI, but we will, in general, use fewer predictors to estimate the factors than the existing implementation of factor forecasting.

The primary focus of our analysis is how to reduce the influence of uninformative predictors for y within the confines of DI framework. We use ‘hard’ and ‘soft’ thresholding to determine which variables the factors are to be extracted from. The factors are the diffusion indices of the forecasting equation. Under hard thresholding, subset variable selection based on some pretest procedure is used to decide whether a predictor is ‘in’ or ‘out’. Under soft thresholding, the top ranked predictors are kept, where the ordering of the predictors depends on the particular soft-thresholding rule used.

We consider the LASSO and the Elastic Net soft-thresholding rules, which are special cases of the ‘Least Angle Regression’ (LARS) algorithm developed in Efron et al. (2004). These soft thresholding methods have been used in biostatistics to study whether groups

* Corresponding author.

E-mail addresses: Jushan.Bai@nyu.edu (J. Bai), serena.ng@columbia.edu (S. Ng).

of genes in a DNA microarray can be used to predict if a certain outcome (such as prostate cancer) occurs. Donoho and Johnstone (1994) provided many optimality results for soft-thresholding and showed that LASSO asymptotically comes close to being an ideal subset selector in terms of its function as an oracle. However, most of the theoretical and empirical analysis we are aware of assume iid data design. We are interested in the usefulness of soft-thresholding from the point of view of factor forecasting, which raises two specific issues. First, economic data are generally weakly dependent data and not iid, and it is not known how these methods perform. Furthermore, subset variable selection is only our intermediate object of interest, as ultimately, it is how the ordered predictors affect the factor estimates that determine forecast efficiency. To our knowledge, this is a new use of the soft-thresholding methodology. As we will see below, the results are encouraging.

Our primary evaluation of the different methods will be based on forecasts of inflation at different horizons and over different samples. The decision to focus on inflation is due in part to the fact that inflation forecasts are important to decision making for both private agents and government agencies. Inflation is chosen for this study also because inflation is well documented to be a challenging series to forecast, see, for example, Stock and Watson (1999). In particular, the reduction in mean-squared inflation forecast error from using the DI methodology has been found to be much smaller than forecasting real variables such as industrial production, see Boivin and Ng (2005). Testing the methodologies on inflation thus puts these methods considered to a non-trivial challenge. Additional results for other series support the main finding from analyzing inflation data that we can push the efficiency of DI forecasting one step further simply by forming the diffusion indices from targeted predictors.

2. Preliminaries

Suppose we are given data on a large number of predictors $X_t = (X_{1t}, \dots, X_{Nt})'$, $i = 1, \dots, N$ and $t = 1, \dots, T$. We are interested in forecasting y_{T+h}^h , the annualized value of the variable y_t in period $T+h$. The precise definition of y_{T+h}^h depends on whether y_t is stationary or not. We will consider both possibilities and make precise the definition of y_{T+h}^h in the empirical section. Note that y can even be one of the predictors. If $N < T$, a forecast that makes use of all available predictors is $\hat{y}_{T+h|T}^h = \hat{\alpha}'W_T + \hat{\Gamma}'X_T$, where W_T is a vector of predetermined variables such as a constant and lags of y_{t+h} , $\hat{\alpha}$ and $\hat{\Gamma}$ are obtained from least squares estimation of

$$y_{t+h}^h = \alpha'W_t + \Gamma'X_t + \epsilon_{t+h}. \tag{1}$$

Although $\hat{\alpha}$ and $\hat{\Gamma}$ are \sqrt{T} consistent, the mean-squared forecast error is increasing in N . Let $\hat{\sigma}_n^2 = \hat{\epsilon}^{n'}\hat{\epsilon}^n/T$ be the sum of squared residuals from estimating a model with n predictors, divided by T . In principle, we can use information criteria of the form

$$N^* = \min_n \left[\log(\hat{\sigma}_n^2) + n \frac{C_T}{T} \right],$$

to select the optimal number of predictors. The FPE with $C_T = 2$ is designed specifically for forecasting, but the BIC with $C_T = \log T$ is also widely used. However, when the predictors have no natural ordering, in theory, there are 2^N sets of predictors to consider, rendering the procedure impractical.

The factor approach to an h period-ahead forecast is to estimate the forecasting equation using data for $t = 1, \dots, T - h$:

$$y_{t+h}^h = \alpha'W_t + \beta(L)\hat{f}_t + \epsilon_{t+h} \tag{2}$$

where $\hat{f}_t \subset \hat{F}_t$, $\beta(L)$ are coefficients pertaining to f_t and p of its lags, \hat{F}_t are the principal component estimates of the $r \times 1$ vector F_t in the factor model

$$X_{it} = \lambda_i'F_t + e_{it} \tag{3}$$

or in matrix form

$$X_t = \Lambda F_t + e_t.$$

Eqs. (2) and (3) constitute the ‘diffusion index’ (DI) forecasting framework of Stock and Watson (2002). The DI forecast is $\hat{y}_{T+h|T}^h = \hat{\alpha}'W_T + \hat{\beta}'(L)\hat{f}_T$.

It is now understood that consistent estimation of the space spanned by F_t makes it possible to obtain \sqrt{T} consistent estimates of α and β and $\min[\sqrt{N}, \sqrt{T}]$ consistent forecasts of the conditional mean, $y_{T+h|T}^h$, if $\sqrt{T}/N \rightarrow 0$ as $N, T \rightarrow \infty$. As shown in Bai and Ng (2006), we can treat \hat{f}_t in the forecasting equation as though it is a vector of observed regressors. The forecasts generated by this methodology seem promising. Evaluations based on key macroeconomic variables find that the DI forecasts tend to do at least as well and often beat alternative methods such as forecast combination, empirical Bayes procedures, etc. See, for example, Stock and Watson (2006) and the references therein. An alternative method of factor forecasting, developed by Forni et al. (2005), also yield promising results. See Forni et al. (2001) and Boivin and Ng (2005).

One way of thinking about the DI methodology is that the factors provide a natural ranking for N mutually orthogonal linear combinations of X_t . If the bulk of the variation in X_t can be explained by a small number of these combinations, say, $N_{\max} \ll N$, the BIC or the FPE need only be evaluated $O(N_{\max})$ times (much smaller than 2^N) to arrive at f_t , the subset of F_t that best predicts y . As the principal component estimates of F_t are just linear combinations of X_t , the DI forecast can be written as

$$y_{t+h}^h = \alpha'W_t + \bar{\Gamma}'X_t + \epsilon_{t+h}$$

where $\bar{\Gamma}$ is a restricted version of Γ in (1). Viewed this way, the DI forecasts use all N of the predictors in forecasting to the extent that $\bar{\Gamma}$ has no element that equals zero exactly.

In this paper, we consider refinements to the DI methodology using what we refer to as ‘targeted diffusion index forecasts’. The thrust of the refinement is to target the factor estimates to the objective of forecasting y_t . More precisely, we seek a model

$$y_{t+h}^h = \alpha'W_t + \gamma'x_t + \epsilon_{t+h}$$

where the $k^* \times 1$ vector x_t is a subvector of X_t . Written differently,

$$y_{t+h}^h = \alpha'W_t + \Gamma^*X_t + \epsilon_{t+h} \tag{4}$$

where the vector Γ^* effectively puts a zero weight on those predictors that are not useful in forecasting y . We will propose two ways of defining Γ^* . Before turning to such an analysis, we introduce a generalization of the method of principal components which can be used whether or not the predictors are targeted.

2.1. Quadratic principal components

By the method of principal component, the estimates of F_t are linear combinations of X_{it} that minimize the sum of squared residuals of the linear model, $X_{it} = \lambda_i'F_t + e_{it}$. This presupposes a linear link function between the data and the latent factors. A more flexible approach is to consider a non-linear link function, $g(\cdot)$, such that

$$g(X_{it}) = \phi_i'J_t + e_{it},$$

where J_t are the common factors, and ϕ_i is the vector of factor loadings. Define X_t^* to be X_t augmented by some or all of the unique

cross-products of the elements of X_{it} , and let $X^* = (X_1^*, \dots, X_T^*)$ be an $N^* \times T$ matrix. The second-order factor model is

$$X_t^* = \Phi J_t + e_t$$

where X_t^* is an $N^* \times 1$ vector. Estimation of J_t then proceeds by the usual method of principal components. If $X_t^* = \{X_{it}, X_{it}^2\}$, then $N^* = 2N$; we will henceforth refer to the procedure as SPC (squared principal components). In a previous version of this paper, the cross-product terms $X_{it}X_{jt}$, $i \neq j$ were also included, a method we referred to as QPC (quadratic principal components). The QPC is computationally demanding and was not noticeably better than the SPC. Results are therefore not included.

It is noteworthy that \hat{J}_t estimated from X_t^* is different from \hat{K}_t , where \hat{K}_t are estimates of the factors in the model $X_{it}^2 = \psi_i' K_t + \eta_{it}$. Whereas \hat{J}_t is a linear combination of the linear AND the quadratic terms of X_{it} , \hat{K}_t is a linear combination of X_{it}^2 . The latter is of interest when one is concerned with factors in the second moments of X_{it} . As our application concerns forecasting inflation, not its volatility, estimation of K_t will not be further considered.

Once the estimates of J_t are obtained using SPC, they are augmented to the forecasting equation in the same way as the standard DI. That is to say, \hat{J} is the matrix of eigenvectors corresponding to the r largest eigenvalues of the $T \times T$ matrix $X^{**}X^*$.¹ Note also that the forecasting equation remains linear whatever is g .

An alternative way of capturing non-linearity is to augment the forecasting equation to include functions of the factors. The simplest case, and one which we will refer to as PC^2 (squared factors), uses the following forecasting equation

$$y_{t+h}^h = \alpha' W_t + \beta_1' \hat{F}_t + \beta_2' \hat{F}_t^2 + \epsilon_{t+h}$$

The PC^2 is conceptually distinct from SPC. While the PC^2 forecasting model allows the volatility of factors estimated by linear principal components to have predictive power for y , the SPC model allows the factors to be possibly non-linear functions of the predictors while maintaining a linear relation between the factors and y . Ludvigson and Ng (2007) found that the square of the first factor estimated from a set of financial factors (i.e., volatility of the first factor) is significant in the regression model for the mean excess returns. In contrast, factors estimated from the second moment of data (i.e., volatility factors) are much weaker predictors of excess returns.

Clearly, the expanded matrix X^* can be much higher dimension than X because of the quadratic terms. But inclusion of noisy predictors can potentially lead to inferior factor estimates. Consideration of quadratic principal components leads naturally to the issue of what predictors should be used to form the diffusion indices.

3. Targeted predictors

A practical question in DI forecasting is how much data are really needed? In Boivin and Ng (2006), it was found that expanding the sample size simply by adding data that bear little information about the factor components does not necessarily improve forecasts. We need to take into account the properties of the idiosyncratic errors when constructing principal components. When the data are too noisy, we can be better off throwing away some data even though they are available. Results in Stock and Watson (2004a) suggest that the weighted principal components,

in the spirit of GLS, indeed provide better forecasts than the OLS-based principal component estimates.

As currently implemented, the factors are extracted from the same large data set, regardless of the series to be forecasted. It is conceivable that the series to be forecasted, y , is highly predictable by a subset of the N series, and this subset is different for different y . We now discuss two classes of procedures to isolate this subset of variables, which we call 'targeted predictors'.

3.1. Hard thresholding

The method of hard thresholding simply uses a statistical test to determine if the i -th predictor is significant without regard for the other predictors being considered. Boivin and Ng (2006) used the correlation coefficients of the errors to pick out the variables to be dropped. It thus exploits a particular hard-thresholding rule to decide which variables are to be used in factor analysis. However, the series to be forecasted was not taken into account.

Our implementation of hard thresholding is closest to Bair et al. (2006), who, like Boivin and Ng (2006), also suggested that the principal components estimated from a large group of variables (which in their analysis are genes) can be dominated by principal components estimated from a smaller set of predictors. They used the bivariate relation between y_{t+h} and X_{it} to screen the variables and referred to the resulting procedure as 'supervised principal components'.² 'Supervised learning' has been used to isolate out subsets of genes associated with certain disease when often, the number of genes (our N) is much larger than the cases (our T) under investigation. However, given the dependent nature of our data, our targeting (or supervising) cannot be based just on the bivariate relation between y_{t+h} and X_{it} . Instead, we need to consider this relation after controlling for other predictors W_t (such as lags of y_t) since a simple autoregressive forecast is always available as an alternative forecasting procedure. The details are:

- For each $i = 1, \dots, N$, perform a regression of y_t^h on W_{t-h} and X_{it-h} . In application, W_{t-h} includes a constant and four lags of y_t . Let t_i denote the t statistic associated with X_{it-h} .
- Obtain a ranking of the marginal predictive power of X_{it} by sorting $|t_1|, |t_2|, \dots, |t_N|$ in descending order.
- Let k_α^* be the number of series whose $|t_i|$ exceeds a threshold significance level, α ;
- Let $x_t(\alpha) = (x_{t[1]}, \dots, x_{t[k_\alpha^*]})$ be the corresponding set of k_α^* targeted predictors. Estimate F_t from $x_t(\alpha)$ by the method of principal components to yield \hat{F}_t .
- Estimate (2) using the BIC to select p and $\hat{f}_t \subset \hat{F}_t$.
- The h period ahead forecast is $\hat{y}_{T+h|T}^h = \hat{\alpha}' W_T + \hat{\beta}'(L)\hat{f}_T$.

Instead of including W_{t-h} as regressors in step 1, an equivalent method is to perform regressions on $M_w y$ and $M_w X_i$, where M_w is the matrix that projects onto the space orthogonal to W , making $M_w y$ and $M_w X_i$ the residuals associated with these projections. An alternative to step (c) is to let k_α^* be N minus the smallest j such that $p_{[j]} \geq \frac{\alpha}{N-j+1}$, where $p_{[j]}$ is the j -th ordered p -value of the test. This Bonferroni-type procedure, due to Holm (1979), is more powerful and generally selects fewer variables than step (c) discussed above, but the top variables selected are quite similar and results will not be reported.

The above algorithm essentially uses only those variables whose marginal predictive power for y is significant at some prescribed level, α , in the factor analysis. After the targeted principal components are estimated, steps (d)–(f) are standard in the DI framework.

¹ In practice, the data are always demeaned and standardized before forming principal components.

² Another use of hard thresholding is 'bagging'. Inoue and Kilian (2008) orthogonalized the data on about 30 variables and used hard thresholding at each bootstrap sample to reduce forecasting variance.

4. Soft thresholding

Hard thresholding can be sensitive to small changes in the data because of the discreteness of the decision rule. Another drawback of selecting predictors one at a time is that it does not take into account the information in other predictors. We may end up selecting variables that are too ‘similar’. It is well known that model averaging is effective only if we pool over variables that bear distinct information from each other.

We now consider ‘soft thresholding’ methods that perform subset selection and shrinkage simultaneously. In the context of (4), it estimates Γ^* and sets those elements corresponding to weak predictors to zero. It is in this sense that shrinkage and model selection are performed simultaneously. We now describe three procedures in this class.

4.1. LASSO

One way of dropping uninformative regressors is to use penalized regressions. Let $RSS(\alpha, \beta)$ be sum of squared residuals from a regression of y_{t+h}^h on all available regressors, W_t and X_{it} , $i = 1, \dots, N$. The solution to

$$\min_{\beta, \alpha} \text{RSS} + \lambda \sum_{j=1}^N \beta_j^2$$

for $0 \leq \lambda < \infty$ is the well-known ridge estimator that shrinks the least squares estimates of β_j towards zero. Note that $\sum_{j=1}^N \beta_j^2 = \|\beta\|_2^2$, the length of β given by the L_2 norm. By the nature of the L_2 penalty, the ridge estimates will almost never be zero exactly. In consequence, uninformative predictors can still inflate forecast error variance.

Consider replacing the L_2 penalty by an L_1 penalty $\|\beta\|_1 = \sum_{j=1}^N |\beta_j|$. The solution to

$$\min_{\beta, \alpha} \text{RSS} + \lambda \sum_{j=1}^N |\beta_j|$$

is the LASSO estimator (‘least absolute shrinkage selection operator’) of Tibshirani (1996). The dual to this problem is

$$\min_{\beta, \alpha} \text{RSS subject to } \sum_{j=1}^N |\beta_j| \leq c$$

where the parameter $c \geq 0$ controls the amount of shrinkage. An important feature of the L_1 penalty is that some coefficient estimates can be exactly zero. The shrinkage under LASSO depends only on λ (or c) and the value of the unrestricted estimates, but not on the correlation of the predictors as is the case under ridge estimation. As shown in Fan and Li (2001), LASSO enjoys a ‘sparsity property’; it estimates zero components of the true parameter vector exactly as zero with probability approaching one as the sample size increases. The asymptotic distribution of the estimator is the same whether or not the zero restrictions are imposed. LASSO thus possesses the ‘oracle property’ in the sense of Fan and Li (2001). That is, the asymptotic distribution of the estimator based on the overall model and the one based on the more parsimonious model coincide.

The LASSO estimator and the ridge estimator are special cases of bridge estimators which are solutions to

$$\min_{\beta} \text{RSS} + \lambda \sum_j |\beta_j|^\gamma.$$

As discussed in Fu (1998), ‘bridge’ estimators have a Bayesian interpretation. The bridge penalty function $\sum_j |\beta_j|^\gamma$ can be thought of as the log prior distribution of the parameter vector, β . The prior

distribution with $\gamma = 2$ is Gaussian, and the prior distribution with $\gamma = 1$ is a Laplace (or double exponential). A small γ favors models either with many parameters set to zero, or parameters with large absolute values from long tailed density. Large values of γ favor models with regression parameters of small but non-zero values from a normal like, or short tailed density. Mol et al. (2006) considered penalized regression models as an alternative to DI forecast and analyzed the problem from a Bayesian perspective. We stay within the DI framework and are interested in which regressors to use in the estimation of the factors.

If the regressors are orthogonal, the LASSO estimates, denoted $\tilde{\beta}$, are

$$\tilde{\beta}_i = \text{sign}\{\hat{\beta}_i\}(|\hat{\beta}_i| - \lambda/2)_+$$

where $\hat{\beta}_i$ is the unrestricted OLS estimate of β_i , $z_+ = z$ if $z > 0$ and 0 otherwise. Therefore, when the least squares coefficients are too small in absolute value, LASSO sets them to zero. Clearly, the LASSO estimate is a non-linear and non-differentiable function of the data. Fu (1998) proposed a shooting algorithm that iteratively solves for the LASSO estimates without using quadratic programming, but the method is unstable when $N > T$. Using convex theory, Osborne et al. (2000) showed that the solution path for $\tilde{\beta}$ is piecewise linear in c . More efficient algorithms are available by exploiting this feature. Our implementation of LASSO will be discussed below.

4.2. The elastic net

The LASSO estimator is an improvement over the ridge estimator when there are many zero coefficients in the true model, since the ridge estimator will never set the coefficients to zero exactly. However, LASSO is not without its drawback. Empirically, it seems that when there is high correlation in the predictors, LASSO is dominated by the ridge. Conceptually there are two problems as highlighted by Zou and Hastie (2005). First, if $N > T$, LASSO can select at most T variables. Second, if there is a group of variables with high pairwise coefficients, LASSO tends to select only one variable from the group and does not care which one. These concerns suggest that a convex combination of ridge and LASSO estimation might be desirable. The result is the ‘elastic net’ (EN) estimator of Zou and Hastie (2005).

The idea of the elastic net is to stretch the fishing net that retains all the ‘big fish’. Like LASSO, the EN simultaneously shrinks the estimates and performs model selection. The LASSO penalty is convex, but not strictly convex. Strict convexity enforces the grouping effect so that predictors with similar properties will have similar coefficients. The EN objective function is

$$\min_{\beta} \text{RSS} + \lambda_1 \sum_{j=1}^N |\beta_j| + \lambda_2 \sum_{j=1}^N \beta_j^2.$$

The EN penalty is thus a convex combination of the LASSO and the ridge penalty and is strictly convex when $\frac{\lambda_2}{\lambda_1 + \lambda_2} > 0$. A computationally appealing property of the EN is that it can be reformulated as a LASSO problem and hence solved using algorithms for LASSO. To see this, define new variables (when the W variables are absent) as follows:

$$X^+ = (1 + \lambda_2)^{-1/2} \begin{pmatrix} X \\ \sqrt{\lambda_2} I_N \end{pmatrix} \quad y^+ = \begin{pmatrix} y \\ 0_N \end{pmatrix}.$$

Let $\gamma = \frac{\lambda_1}{\sqrt{1 + \lambda_2}}$. Then the EN estimator can be reformulated as

$$\beta^{++} = \underset{\beta}{\text{argmin}} \text{RSS}^+ + \gamma \sum_{j=1}^N |\beta_j|$$

with RSS^+ is the sum of squared residuals from a regression of y^+ on X^+ . To remove a double shrinkage effect (which is in both LASSO and ridge), the EN estimator that proposed by Zou and Hastie (2005) is $\beta^+ = (1 + \lambda_2)\beta^{++}$. As will be discussed below, our main interest is not so much in the point estimates, but the ordering of variables provided by the EN. With this in mind, we now turn to the implementation of LASSO and EN.

4.3. Least angle regressions

A widely-used variable selection method is the forward selection regression whereby the $(k + 1)$ -th predictor is added to the ‘in’ set if it has the maximum correlation with the residual vector from the k -step. The residual vector is then projected on the remaining predictors and a new predictor is found. Forward selection regressions tend to be too aggressive in the sense of eliminating too many predictors correlated with the ones included. Another popular method is forward stagewise regression, which is more cautious than forward selection regressions as it takes smaller steps towards the final model. Briefly, if $\hat{\mu}_k$ is the current estimate of y with k predictors and $\hat{c} = X'(y - \hat{\mu}_k)$ is the ‘current correlation’ (assuming each column of X is standardized), there exists a j such that $|\hat{c}_j|$ is maximized. Consider the updating rule $\hat{\mu}_{k+1} = \hat{\mu}_k + \hat{\gamma} \text{sign}(\hat{c}_j)X_j$. Forward selection sets $\hat{\gamma} = |\hat{c}_j|$ whereas forward stagewise regression sets $\hat{\gamma}$ to a small constant. As we will see below LASSO uses yet another $\hat{\gamma}$ and replaces X_j by some other quantity.

Efron et al. (2004) showed that LASSO and forward stagewise regressions are in fact special cases of what is known as LARS, or least angle regressions. At each step, the $\hat{\gamma}$ in LARS is endogenously chosen so that the algorithm proceeds equiangularly between the variables in the most correlated set (hence the ‘least angle direction’) until the next variable is found. After k steps, there are k variables in the active set. If we end after k steps, we will have an active set of k predictors, or in other words, the coefficients corresponding to the remaining $N - k$ predictors will be set to zero. If we continue until $k = N$, we will have a set K of indices of predictors ordered according to when they join the active set. How many coefficients to set to zero is thus re-cast into a stopping rule for k .

Formally, the LARS algorithm begins at $\hat{\mu}_0 = 0$. Suppose $\hat{\mu}$ is the current estimate and let $\hat{c} = X'(y - \hat{\mu})$. Define K as the set of indices corresponding to variables with the largest absolute correlations,

$$\hat{C} = \max_j |\hat{c}_j| \quad K = \{j : |\hat{c}_j| = |\hat{C}|\}.$$

Let $s_j = \text{sign}(\hat{c}_j)$ and define the active matrix corresponding to K as

$$X_K = (s_j X_j)_{j \in K}.$$

Let $G_K = X'_K X_K$ and $A_K = (1'_K G_K^{-1} 1_K)^{-1/2}$, where 1_K is a vector of ones equaling the size of K . A unit equiangular vector with columns of the active set matrix X_K can be defined as

$$u_K = X_K w_K, \quad w_K = A_K G_K^{-1} 1_K, \quad a_K = X' u_K,$$

so that $X'_K u_K = A_K 1_K$. LARS then updates $\hat{\mu}$ as

$$\hat{\mu}^{\text{new}} = \hat{\mu} + \hat{\gamma} u_K$$

where

$$\hat{\gamma} = \min_{j \in K^c}^+ \left(\frac{\hat{C} - \hat{c}_j}{A_K - a_j}, \frac{\hat{C} + \hat{c}_j}{A_K + a_j} \right)$$

where the minimum is taken over only the positive components.

LARS has several advantages. First, it gives us a ranking of the predictors when the presence of other predictors are taken into account, which is unlike the case of hard thresholding. Second,

the algorithm implicitly avoids strongly correlated predictors, since if one of the correlated predictors is already included, the new residual will have a low correlation with variables strongly correlated with the variable just included. Third, LARS is not as ‘greedy’ as forward regressions which, when a good direction is found, it exploits the direction to a maximum. Fourth, LARS is fast; the computation cost is of the same order as the usual OLS. Indeed, starting from zero, the LARS solution paths grow piecewise linearly in a predictable way.

Superficially, LASSO and LARS seem quite different. However, as also shown in Efron et al. (2004), LASSO is in fact a special case of LARS that imposes the sign restriction that, if $\hat{\beta}^k$ is the vector of estimates at the k -th step, the sign of $\hat{\beta}_j^k$ must agree with the sign \hat{c}_j for those j in the active set. Variables in the active set that fail the sign restriction can be ‘kicked out’ of the active set under LASSO. Therefore, unlike LARS, the size of the active set under LASSO need not be monotonically increasing. Under LASSO, the tuning parameter, λ , determines the severity of the penalty and thus how many parameters are set to zero. The LARS implementation of LASSO turns the choice of this tuning parameter into the choice of k , or in other words, the size of the active set. To determine k^* , one can use an information criterion such as the BIC. That is, have

$$k^* = \underset{k}{\text{argmin}} \text{BIC}(k) = \log(\hat{\sigma}_k^2) + k \frac{\log T}{T}.$$

By choosing k^* , the BIC also sets coefficients on the $k^* + 1$ to N predictors as ordered by LARS/LASSO/EN to zero. The BIC can be replaced by the AIC, which is very similar to generalized cross-validation (GCV), and which Efron et al. (2004) considered.³

Our interests in soft thresholding was initiated by the sparsity property arising from the L_1 penalty of LASSO. If the number of non-zero coefficients is indeed very small, then one can simply use this sparse set of predictors for forecasting. However, when the model structure is not sufficiently sparse, i.e., when k^* is not sufficiently small, one is again faced with the problem that including too many predictors will inject excess sampling variability to the forecasts. In this situation, it seems reasonable to resort to dimension reduction by constructing diffusion indices from a selected subset of significant predictors. This being the case, the particular aspect of soft thresholding that turns out to be more useful for us is the ordering of the predictors provided by the LARS, and/or LARS implementation of LASSO. Specifically, the ordered set of variables will be used to construct diffusion indices. Of course, if k^* is too small, the principal component estimates will be imprecise. As we cannot know a priori how big is k^* , we always estimate the factors using a fixed number of series but these series will be ordered according to LARS/LASSO. On the other hand, if k^* is in fact small, we also entertain a forecasting model that uses k^* non-standardized predictors directly without forming factors. But it should be kept in mind that this is not a DI forecast.

In the empirical analysis to follow, we consider three methods of using LARS-ordered variables: (i) estimate principal components \hat{F}_t from the first 30 series that LARS/LASSO/EN select; (ii) enter the first five ($k = 5$) predictors to the forecasting equation directly; (iii) enter $k = k^*$ predictors to the forecasting equation directly. Put differently, (ii) and (iii) use a small number of selected variables as predictors (i.e., as our \hat{F}_t without principal components analysis), while the \hat{F}_t in (i) are principal components estimated

³ Comments on the LARS article reflect concerns by many that the GCV will overfit. The BIC evaluates models using in-sample errors, but an out-of-sample variant can also be considered. As discussed below, k^* is not of primary importance given that we have to estimate the factors. For this reason, alternative methods for determining k^* was not pursued.

from the first 30 series selected by LARS. Then the BIC is used to determine p and the corresponding \hat{f}_t (a subset of \hat{F}_t) that enters the forecasting equation (2). For (i), we use 30 series because our experience has been that for well-behaved data, the principal component estimates have reasonably good properties when the number of cross-section units exceed 30. As well, even the tightest hard-thresholding criterion tends to pick out over 30 series. We therefore want to see if the DI methodology works well with as few as 30 predictors. It is important to emphasize that (i) is a DI forecast and as such, neither the LARS/LASSO parameter estimates, nor k^* per se, will be used directly because the variables selected by LARS/LASSO enter the forecast only via the factors.

As we will see, some of our subsamples considered has $N > T$ and some has N in the same order as T , even though for the whole sample $T > N$. In practice, we always use the X^+ and Y^+ data matrices so that we effectively implement what amounts to LARS-EN and LASSO-EN, with the difference between LARS-EN and LASSO-EN being whether or not to impose a sign restriction. As a matter of notation, we simply refer to the three methods as LA(5), LA(PC), and LA(k^*) and the use of EN is implicit. The LASSO results can be similarly defined. We also use the LARS algorithm to produce an ordering of the variables when x_{it} and x_{it}^2 are included. A complete list of the methods considered is given below:

PC	f_t estimated from all 132 of the X_{it} available;
SPC	f_t estimated from all 132 of X_{it} and X_{it}^2 available;
TPC	f_t are the targeted principal components estimated from a subset of available X_{it} where the subset is selected by hard thresholding;
TSPC	f_t are the targeted principal components estimated from X_{it} and X_{it}^2 ;
TSTPC	f_t are targeted principal components estimated from a subset of X_{it} and X_{it}^2 ;
PC ²	f_t estimated from all 132 of the X_{it} available, $[f_t, f_t^2]$ used in the forecasting equation;
TPC ²	f_t estimated from a subset of X_{it} and X_{it}^2 , $[f_t, f_t^2]$ used in the forecasting equation;
LA(5)	f_t is the vector of 5 best predictors selected by LARS;
LA(PC)	f_t are the factors estimated from the 30 best predictors in $\{x_{it}\}$ as selected by LARS;
LA(k^*)	f_t is the vector of k^* best predictors selected by LARS;
LA(SPC)	f_t are the factors estimated from the 30 best predictors in $\{x_{it}, x_{it}^2\}$ as selected by LARS.

5. Results

The variable we are interested in forecasting, y_{t+h} , is the logarithm of PUNEW, or CPI all items, using factors estimated from some or all of the 132 predictors and or their cross-products.⁴ These are monthly time series available from 1960:1 to 2003:12 for a total of $T = 528$ observations. As argued by Nelson and Plosser (1982) and Beveridge and Nelson (1981), many of those series are I(1) non-stationary or contain an I(1) components, the data are therefore transformed by taking logs, first or second differences when necessary, as in Stock and Watson (2006). In particular, the logarithm of CPI is assumed to be integrated of order two. Following Stock and Watson (2002), define

$$y_{t+h}^h = \frac{1200}{h} \cdot (y_{t+h} - y_t) - 1200 \cdot (y_t - y_{t-1})$$

⁴The data are taken from Mark Watson's web site <http://www.princeton.edu/~mwatson>.

and let

$$z_t = 1200 \cdot (y_t - y_{t-1}) - 1200 \cdot (y_{t-1} - y_{t-2}).$$

For $h = 1, 6, 12,$ and 24 , the factor augmented forecast given information in time t is

$$\hat{y}_{t+h|t}^h = \hat{\alpha}_0 + \hat{\alpha}'_1(L)z_t + \hat{\beta}'_1(L)\hat{f}_t$$

where the number of lags of z_t and \hat{f}_t are determined by the BIC with the maximum number of lags set to six when the sample size permits, and is reduced to four otherwise. In the notation of the preceding discussion, W_t consists of a constant and lags of z_t . To simplify notation, \hat{f}_t generically denotes estimated factors used for forecasting, where $\hat{f}_t \subset \hat{F}_t$. In all cases, \hat{F}_t is a 10×1 vector. That is, we select the factors used for forecasting from the first ten estimated factors. It is understood that \hat{F}_t are estimated using different number of series. Although we are forecasting the change in inflation, we will continue to refer to the forecasts as inflation forecasts.

Our main interest is in figuring out n , the number of variables used to estimate F_t . Given that W_t are lags of inflation, the question more precisely phrased is what variables have predictive power for inflation after controlling for lags of inflation themselves. We do not restrict the optimal predictors to be the same for every time period. Instead, the predictors are selected at each t , and the forecasting equation is re-estimated after new factors are estimated. Our first estimation consists of ten years of data (120 data points) starting in 1960:3; the sample is extended one month at a time. There is one forecasting equation for each h . The last observation used in estimation when $h = 12$ is 2002:12. For each h , we have about 400 out-of-sample forecasts. We use the average of the forecast errors to evaluate the different procedures. We will refer to the ratio of the MSE for a given method to the MSE of an AR(4) as RMSE (relative mean-squared error). That is,

$$RMSE(\text{method}) = \frac{MSE(\text{method})}{MSE(AR(4))}.$$

An entry less than one indicates that the specified method is superior to the simple AR(4) forecast.

Because parameter instability is prevalent in economic time series, a method that forecasts well in one sample is not guaranteed to forecast well in another sample period. Therefore, in addition to the full sample analysis, we also consider seven forecast subsamples: 70:3–80:12, 80:3–90:12, 90:3–00:12, 70:3–90:12, 70:3–00:12, 80:3–00:12, and 70:3–03:12. For example, in the case of 70:3–00:12, the first forecast of 70:3 is based on estimation up to 60:3–70:3-h. The last forecast is for 00:12, and it uses parameters estimated for the sample 60:3–00:12-h. Table A.1 provides summary statistics for both y^h and the level of inflation over these samples. Notably, the mean of inflation over the estimation sample can be higher or lower than the forecast sample. Table A.2 then shows that five of the seven forecast samples considered had decelerating inflation. Note that inflation is the most volatile when $h = 1$. This feature will help understand the results to follow.

5.1. Number of variables chosen

We use the t statistic and three cutoff points for hard thresholding:

1. hard thresholding, $|t| > 1.28$;
2. hard thresholding, $|t| > 1.65$;
3. hard thresholding, $|t| > 2.58$;

