# iCapital

Technical Manual for

# Asset Returns Projection

iCapital Portfolio Analytics Quantitative Research Team

## 1. Introduction and Scope

This project involves inputting historical asset class data to get return and risk projections for individual assets on Architect. We use Principal Components Regression to accomplish this task. There are two main **inputs** to this project:

1. The historical returns of asset classes;
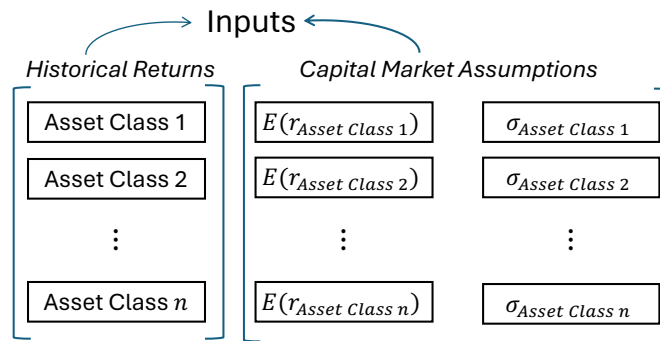2. The risk and return estimates provided by external CMA providers.



Fig. 1. Inputs to the process.

The **outputs** to the process will be estimates of *fund-level*:

- future expected return;
- and future expected risk.

We use **principal components regression (PCR)** to accomplish this task.
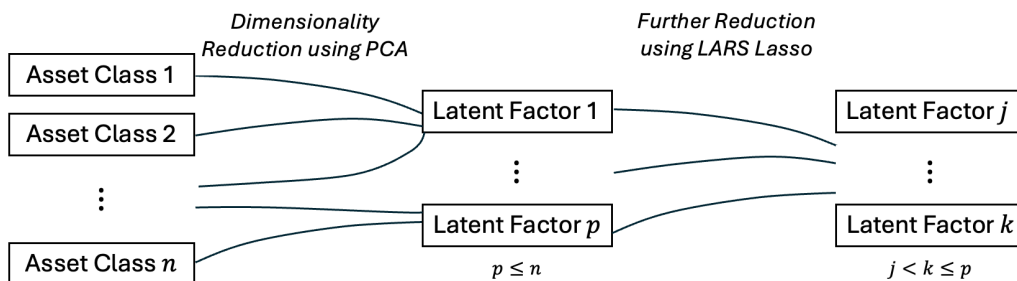
## 2. Methodology



Fig. 2. The first step in the methodology.

Principal components regression (PCR) is a technique used when the number of regressors is high, and are most likely correlated. Instead of regressing the dependent variable on all the explanatory variables, select principal components of the explanatory variables are used as regressors. One typically uses only a subset of all the principal components for regression; as a result, PCR could be interpreted as a kind of regularization procedure. It can also be viewed as a type of shrinkage estimator.

*Ordinary Least-Squares Regression*

$$\text{G3 Factor } 1 = \alpha_1 + \beta_{j,1}\boxed{\text{Latent Factor } j} + \cdots + \beta_{k,1}\boxed{\text{Latent Factor } k} + \epsilon_1$$

$$\vdots$$

$$\text{G3 Factor } m = \alpha_m + \beta_{j,m}\boxed{\text{Latent Factor } j} + \cdots + \beta_{k,m}\boxed{\text{Latent Factor } k} + \epsilon_m$$

$$m = 13 \qquad\qquad\qquad j < k \leq p$$

Fig. 3. The second step in the methodology.

*2.1. Summary*—The basic methodology can be summed up in four steps. We do this for each of the 13 factors:

1. Reduction of dimensionality of asset classes using principal components analysis, and selection of the optimal principal components using LARS Lasso.
2. Coefficient generation on reduced feature set via ordinary least-squares regression. Specifically, each G3 factor on the left-hand-side (LHS, dependent variable), and the optimal reduced feature set (the components from the LARS Lasso above) on the right-hand-side (RHS, independent variables).
3. Combine external CMAs with regression coefficients generated above to get future expected factor return and risk projections.
4. Use G3 factor loadings on future expected factor-projections to get asset-level projections.

*2.2. Potential Issues and Solutions*—Usually, the principal components which have a higher percent-variance-explained (the ones based on eigenvectors corresponding to the higher eigenvalues of the sample variance-covariance matrix of the explanatory variables) are selected as regressors. It should be noted, however, that for the purpose of predicting the outcome, the principal components with lower percent-variance-explained may also be important.[*] The logic behind this is simple – PCR does not consider the response variables (our factors) when deciding which components to drop; that decision is based only on the magnitude of the percent-variance-explained within the components themselves.

In this manner, we ensure that, both:
- the maximum variance of the independent variables; *and*
- the explanatory power of the independent variables *on* the dependent variables,

---

[*] When Kendall and Hotelling first proposed PCR in the 1950s, they proposed "complete" PCR, which means replacing the original variables by *all* the principal components. Which principal components are included in the final model is determined by looking at the explained variance of the parameter estimates. By the early 1980s, the term PCR had changed to mean "incomplete PCR."
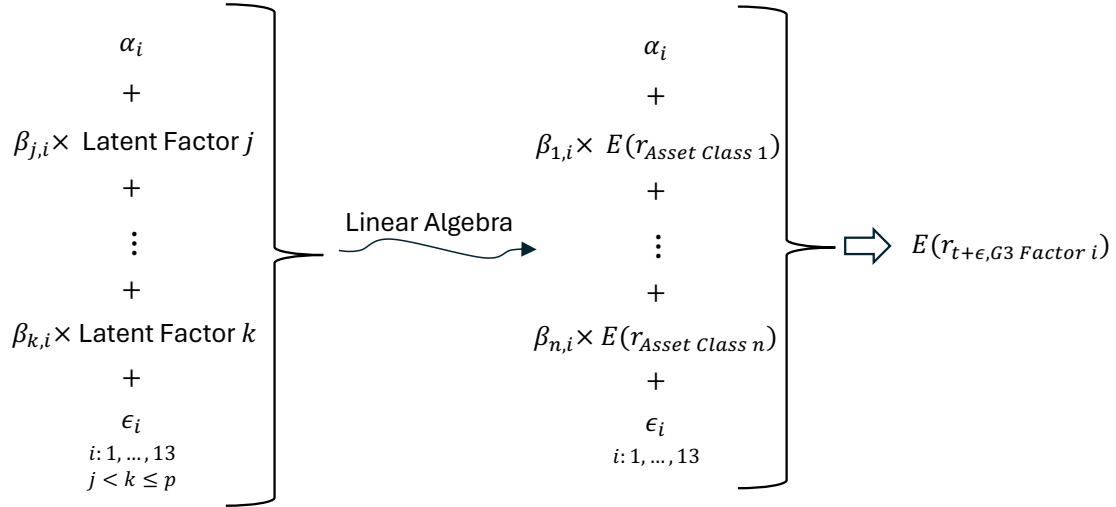
$$
\left.\begin{array}{c}
\alpha_i \\
+ \\
\beta_{j,i} \times \text{ Latent Factor } j \\
+ \\
\vdots \\
+ \\
\beta_{k,i} \times \text{ Latent Factor } k \\
+ \\
\epsilon_i \\
{\scriptstyle i:\, 1,\ldots,13} \\
{\scriptstyle j < k \le p}
\end{array}\right\}
\xrightarrow{\text{Linear Algebra}}
\left\{\begin{array}{c}
\alpha_i \\
+ \\
\beta_{1,i} \times E(r_{\text{Asset Class } 1}) \\
+ \\
\vdots \\
+ \\
\beta_{n,i} \times E(r_{\text{Asset Class } n}) \\
+ \\
\epsilon_i \\
{\scriptstyle i:\, 1,\ldots,13}
\end{array}\right\}
\Rightarrow E(r_{t+\epsilon,\, G3 \text{ Factor } i})
$$

Fig. 4. The third step in the methodology for a single factor. It is repeated for all factors.

are accounted for.

For details on $k$-fold Cross-Validation, and LARS Lasso, please refer to the Appendix.

## 3.   Results

*3.1.   Data*—The methodology in this paper was inspired by Haddad et al (2020).[†]

The asset classes used are: (i) Buyout and growth (ii) Commodity (iii) Emerging markets equities (iv) Emerging markets fixed income (v) Event driven (vi) Fund of funds (vii) Global macro (viii) Hedged equity (ix) Long/short credit (x) MultiStratHF (xi) MultiStratPE (xii) Non-US developed equities (xiii) Private debt (xiv) Private real asset (xv) Real estate (xvi) REITs (xvii) Secondaries (xviii) US equities (xix) US fixed income. Note that the asset class data is *desmoothed* (see below for details), and this ensures that all quarterly data is converted to monthly.

The factors used are: (i) Alt commodities (ii) Alt HF crowding (iii) Alt oil (iv) Alt trend (v) Emerging markets (vi) Equity market (vii) Equity momentum (viii) Equity quality (ix) Equity smallcap (x) Equity value (xi) Fixed credit (xii) Fixed duration (xiii) US dollar.

Correlations for all time series can be seen in Figure 5. Many correlations are positive, suggesting broad market factors tend to drive returns across asset classes. However, not all are positive, and even among the positive ones, the strength of correlations varies considerably. Equities, for example, tend to be highly correlated with each other: Alt HF Crowding, Emerging Markets, Equity Market, Equity Momentum, and Equity Quality. This suggests equity market returns are often driven by common underlying factors. Buyout and Growth is also highly correlated with many series, suggesting the multi-asset nature of buyout funds. Fixed Income shows very low correlation with most other asset classes, as do Commodities and Emerging

---

[†]  Haddad, V., S. Kozak, S. Santosh. 2020. Factor Timing. *The Review of Financial Studies*, 33(5):1980–2018.

Markets Fixed Income. As a final note, it should be noted that the asset classes are correlated, so Principal Components Analysis on them is justified.
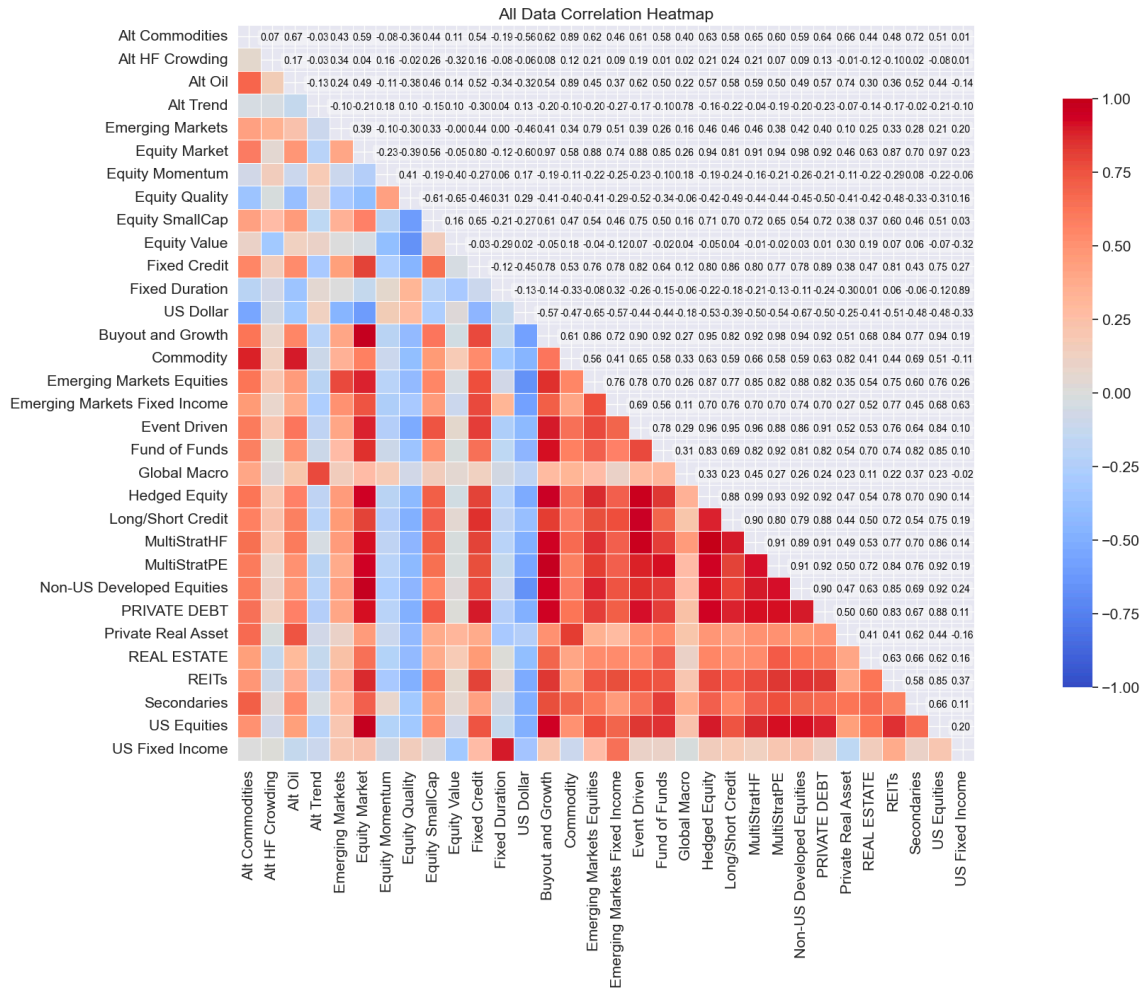
**All Data Correlation Heatmap**

Fig. 5. Correlations of all time series returns.

The data is monthly and the months for which all series have data is October 31 2007 to July 31 2023. Also, we break up the dataset into training and testing in an 80%-20% split, respectively. To preserve the temporal sequence (and any potential serial correlation), we do not shuffle the data before splitting it into the train and test set.

Finally, we use the *desmoothed* asset class returns, to "convert" the quarterly data into monthly data. For details on the desmoothing process, please refer to the Appendix.

Histograms for the standardized, desmoothed asset class returns can be seen in Figure 6. Some exhibit more positive skew (e.g. Global Macro, Real Estate) and others more negative skew (e.g. Event Driven). The dispersion also differs, with some sectors having tighter distributions (e.g. US Equities) compared to wider spreads in others (e.g. Long/Short Credit, Private Debt). REITs and Fund of Funds, are more concentrated around the mean. MultiStratHF and Secondaries, have a more even spread across the distribution. Hedged Equity has its peak slightly

to the right of center, indicating a bias towards positive returns, while Commodity's peak is left of center, suggesting a bias towards negative returns.
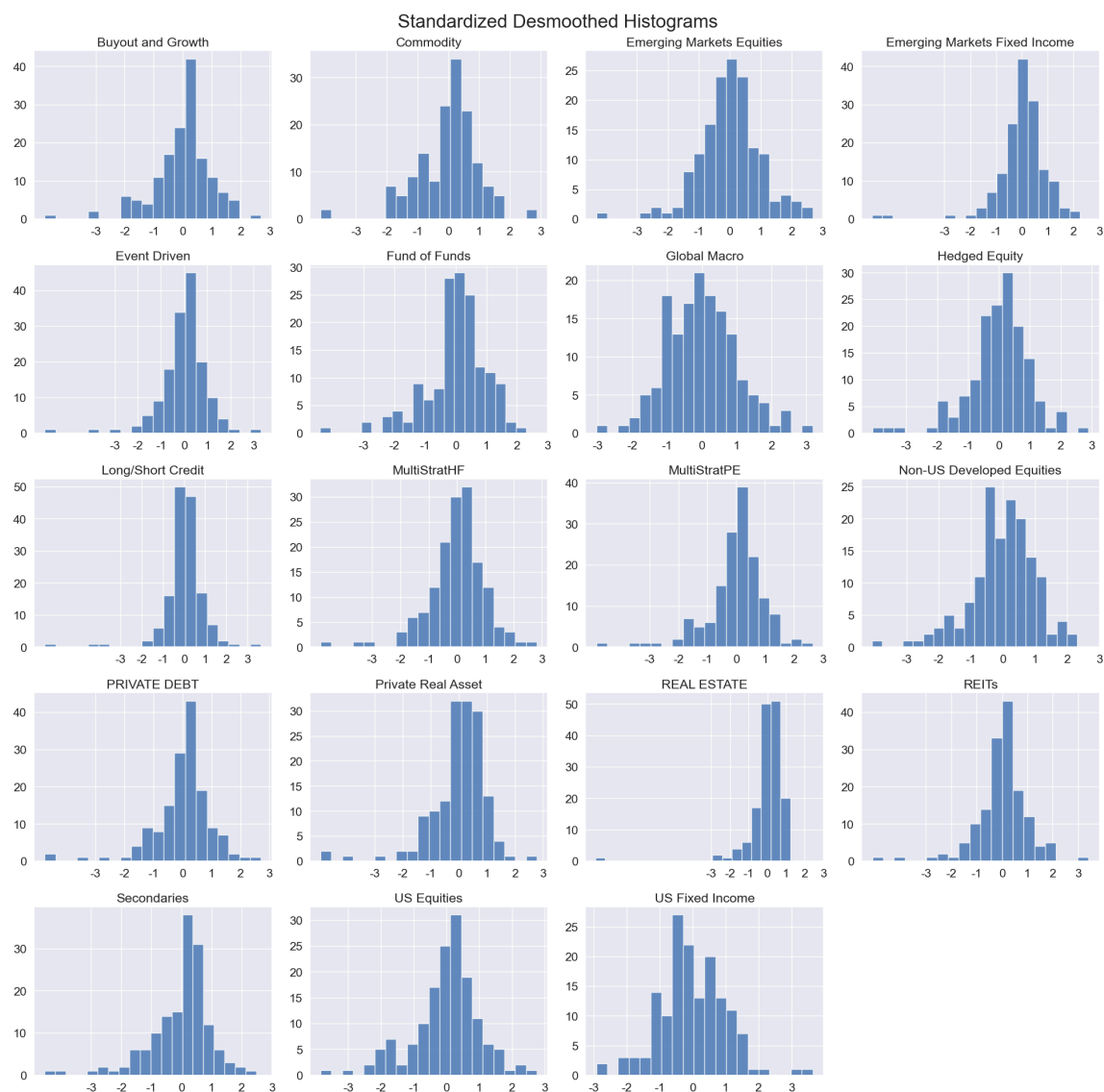


Fig. 6. Histograms of the standardized asset class returns.

*3.2.    Generating a Baseline RMSE*—We run a 10-fold cross-validation on the training set of the original (unmodified) features, to get a baseline RMSE. We compare this with the RMSE on the test set. The results are shown in Table 1 below.

It should be stated, the absolute values of the RMS errors are small, and the results should be taken with a grain of salt. That said, there is considerable disparity between results. Equity Momentum, Equity Market, and Equity Value show significant improvement from the training set to the test set. These equity market factors have very similar drivers, and can be predicted easily using our unmodified asset classes. Fixed Duration shows considerable improvement as well. Alt Commodities, Alt Oil, and US Dollar exhibit negative percentage differences. This should not be surprising, as these factors are volatile.

Table 1. RMSE Results for OLS Regression for the Baseline.

| Dependent Var | Train | Test | %Δ |
|---|---|---|---|
| Alt Commodities | 0.0248 | 0.0213 | -14.21% |
| Alt HF Crowding | 0.0372 | 0.0468 | 25.83% |
| Alt Oil | 0.0539 | 0.0518 | -3.85% |
| Alt Trend | 0.0172 | 0.0154 | -10.03% |
| Emerging Markets | 0.0074 | 0.0093 | 26.38% |
| Equity Market | 0.0016 | 0.0039 | 142.51% |
| Equity Momentum | 0.0326 | 0.0826 | 153.02% |
| Equity Quality | 0.0220 | 0.0219 | -0.61% |
| Equity SmallCap | 0.0131 | 0.0193 | 47.20% |
| Equity Value | 0.0124 | 0.0242 | 95.38% |
| Fixed Credit | 0.0114 | 0.0141 | 22.87% |
| Fixed Duration | 0.0053 | 0.0095 | 80.60% |
| US Dollar | 0.0177 | 0.0146 | -17.96% |

*3.3.*   *Principal Components Analysis*—We run a PCA on the training part of the RHS dataset. The results of the variance explained can be seen in Figure 7.
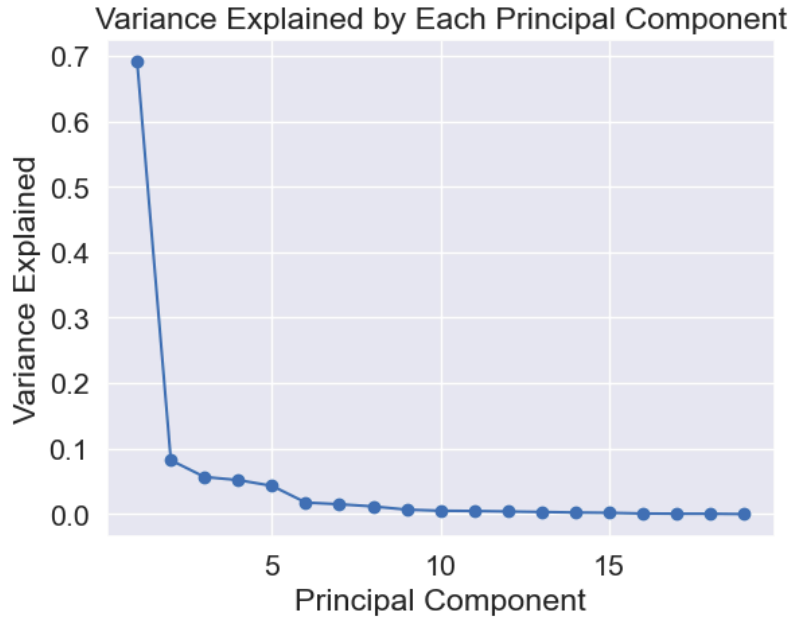


Fig. 7. The variance explained by each principal component. The first component explains ∼69% of the variance, with the rest of the components explaining 8% or less, each.

*3.4.*   *LARS Lasso*—The LARS Lasso helps us decide which principal components to keep, for each factor, since the percent-variance-explained metric only utilizes the asset class data, and not the factor data. In other words, there could be principal components with low percent-

**6**

variance-explained, but they could be significant in explaining one of the factors. Table 2 shows the principal components that are kept for each dependent variable.

Table 2. Principal Components Kept for Each Dependent Variable.

| Factor | Principal Components | | | | | | | | | | | | | | | | | | | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | |
| Alt Commodities | | x | x | x | x | | x | x | x | | | x | x | x | | x | x | | x | 13 |
| Alt HF Crowding | x | x | x | x | x | x | x | x | x | x | x | x | x | x | x | x | x | x | x | 19 |
| Alt Oil | x | x | x | x | x | x | x | | | x | x | x | x | x | x | x | x | x | | 16 |
| Alt Trend | | x | x | | | | x | | | | | | | | | | | x | | 4 |
| Emerging Markets | | x | x | x | | | | | x | | | x | | | x | x | | x | | 8 |
| Equity Market | x | x | x | | | x | x | x | x | | | x | x | | x | x | | x | | 12 |
| Equity Momentum | | | | | | | | | | | | | | | | | x | | | 1 |
| Equity Quality | | x | | x | | x | | | | | | x | | | x | x | | x | x | 8 |
| Equity SmallCap | | x | x | x | x | x | x | x | | | x | x | x | | x | x | | x | | 13 |
| Equity Value | x | | x | x | x | | | x | | | | x | x | x | x | x | | x | x | 12 |
| Fixed Credit | | x | x | | | | | | x | | | | | | | x | | | | 4 |
| Fixed Duration | | x | x | | x | | | | | | | x | | | | | | x | x | 6 |
| US Dollar | | x | x | x | | | | | x | | | x | | x | | | | x | x | 8 |
| Total | 4 | 10 | 11 | 7 | 7 | 4 | 7 | 5 | 7 | 2 | 3 | 10 | 6 | 5 | 7 | 10 | 3 | 10 | 6 | |

There is a significant variation in the number of principal components retained across different dependent variables. For instance, Alt HF Crowding retains all 19 components, suggesting that variance across these components is evenly spread, or the model relies heavily on capturing as much information as possible from the data. In contrast, Equity Momentum retains only one component, indicating that most of the useful variance can be captured by a single principal component.

The number of components retained can also suggest the complexity of the underlying patterns within each variable. A higher number of components might imply more complex interactions or more nuanced underlying structures. For example, Emerging Markets and US Dollar, retaining 7 and 8 components respectively, may indicate more complex financial behaviors or influences requiring multiple factors for adequate representation.

The total counts of each PC across all variables (shown in the last row) provide a useful overview of how often each PC is utilized. PCs with higher counts are more universally useful across multiple dependent variables, suggesting they capture more generalizable features of the data. PCs utilized less frequently may capture more specific, unique, or less significant aspects of the data.

*3.5. Principal Components Regression*—The results of the Principal Components Regression using the principal components selected from the LARS Lasso above are shown in Table 3. Note that the "Diff" column shows the difference in the percent change of RMSE from the training to the test data, from OLS with the original features, to PCR with the reduced,

transformed feature set.

Table 3. Tain and Test Scores for Principal Components Regression, with Percent Differences for all Factors. OLS Results Included.

| Dependent Variable | Train Score | Test Score | %Δ PCR | %Δ OLS | Diff |
|---|---|---|---|---|---|
| Alt Commodities | 0.053297 | 0.054787 | 2.80% | -14.21% | 17.01% |
| Alt HF Crowding | 0.037208 | 0.073891 | 98.59% | 25.83% | 72.76% |
| Alt Oil | 0.053003 | 0.143000 | 169.80% | -3.85% | 173.65% |
| Alt Trend | 0.020795 | 0.039424 | 89.58% | -10.03% | 99.61% |
| Emerging Markets | 0.065726 | 0.069887 | 6.33% | 26.38% | -20.05% |
| Equity Market | 0.009693 | 0.095859 | 888.97% | 142.51% | 746.46% |
| Equity Momentum | 0.034990 | 0.049986 | 42.86% | 153.02% | -110.16% |
| Equity Quality | 0.026716 | 0.029138 | 9.06% | -0.61% | 9.67% |
| Equity SmallCap | 0.030698 | 0.036976 | 20.45% | 47.20% | -26.75% |
| Equity Value | 0.013929 | 0.028793 | 106.72% | 95.38% | 11.34% |
| Fixed Credit | 0.026211 | 0.022726 | -13.30% | 22.87% | -36.17% |
| Fixed Duration | 0.012558 | 0.021283 | 69.48% | 80.60% | -11.12% |
| US Dollar | 0.026205 | 0.019955 | -23.85% | -17.96% | -5.89% |

Equity Market shows an extraordinarily high "Diff" of 746.46%. Such a massive jump could indicate that the model exceptionally overfits the training data but somehow aligns better with the testing data characteristics or that specific features influential in the test data are heavily weighted in the model.and conversely, Equity Momentum has a negative "Diff" of -110.16%, suggesting a degradation in model performance from training to testing, highlighting potential overfitting or instability in the model when exposed to new data.

The wide range in "Diff" across variables suggests that no one-size-fits-all model tweak or adjustment will universally enhance model performance. Instead, tailored adjustments, perhaps incorporating more sophisticated regularization techniques or feature engineering specific to each variable's characteristics, might be necessary to stabilize and improve model outputs across different conditions and datasets.

*3.6.    Calculating Factor Expected Returns*—The final step in the project is to calculate factor expected returns. We will use the asset class expected returns provided by external CMA providers for this. For now, we use the anomaly dataset historical returns.

*3.6.1.    Issues and Solution*—It's important to note that we are using the coefficients from the regression run after two modifications on the original feature set, i.e., the asset classes: (i) standardization or normalization; and (ii) principal components analysis.We will be using the returns provided by an external CMA provider (in our case, Envestnet), along with the coefficients and intercept from the PCR. These externally-provided returns are the 10-year projections of the original feature set. We must somehow find a way to adjust the coefficients so that they match up with the original feature set, and not the standardized feature set in "principal components

space".‡ To do this, we rely on the equivalence between PCA and singular value decomposition. Please see the Appendix for details on the methodology used and the corresponding linear algebra.

Another issue faced is the mapping of iCapital's Asset Classes to the Envestnet Asset Classes. This is shown in Table 4.

Table 4. Asset Class Mappings between iCapital and Envestnet, with Expected Returns from Envestnet.

| iCapital | $E(r)$ | Envestnet |
|---:|---|---|
| Commodity | 4.9% | Commodity |
| Emerging Markets Equities | 8.3% | Int'l Emerging Mkts |
| Emerging Markets Fixed Income | 5.6% | Emerging-Markets Bond |
| REAL ESTATE | 6.2% | Private Real Estate |
| REITs | 6.5% | REITs |
| PRIVATE DEBT | 8.8% | Private Credit |
| Global Macro | 5.9% | Global Macro |
| Hedged Equity | 7.0% | Hedged Equity |
| Long/Short Credit | 6.9% | Long/Short Credit |
| Event Driven | 5.7% | Event Driven |
| Non-US Developed Equities | 7.3% | Int'l Developed Mkts |
| MultiStratPE | 9.8% | Private Equity |
| US Equities | 6.5% | Large-Cap Core |
| US Fixed Income | 4.4% | Intermediate Bond |
| MultiStratHF | 5.7% | Multi-Strategy |
| Fund of Funds | 7.0% | Hedged Equity |
| Private Real Asset | 5.6% | Other |
| Buyout and Growth | 6.0% | Alternative |
| Secondaries | 9.8% | Private Equity |

*3.6.2.   Final Results*—Table 5 shows the projections and the historical means of the factor returns across the entire time period (October 2007 to December 2023). Data has been left as monthly here for factual clarity, but should be annualized in the final version of the implementation.

11 out of the 13 factors have historical mean values that fall within the 95% confidence window as defined by their regression standard errors. This suggests that the projections for these variables are aligned with their historical performances, indicating consistency in their behavior over time or accurate modeling.

However, Alt HF Crowding and Equity Value do not fit within their respective projected ranges, which might suggest that these projections were outliers in historical terms or that unusual market conditions affected these variables during the period under consideration. Fortunately, the window for Equity Value is small, at 0.05%. Furthermore, Alt HF Crowding can be volatile.

---

‡ Recall that principal components analysis transforms the original features into a new set of orthogonal features, capturing maximum variance in a lower-dimensional space.

Table 5. Factor Projections with Historical Means.

| Dependent Variable | Projection $\pm$ 95% Conf. Int. | Historical Mean | In Interval |
|---|---|---|---|
| Alt Commodities | $-0.45 \pm 1.15$ | -0.15 | ✓ |
| Alt HF Crowding | $0.47 \pm 0.47$ | -0.14 | |
| Alt Oil | $1.33 \pm 0.49$ | 0.78 | ✓ |
| Alt Trend | $-0.07 \pm 0.70$ | 0.33 | ✓ |
| Emerging Markets | $0.26 \pm 1.98$ | -0.58 | ✓ |
| Equity Market | $0.17 \pm 0.06$ | 0.47 | ✓ |
| Equity Momentum | $1.23 \pm 3.61$ | 0.11 | ✓ |
| Equity Quality | $-0.18 \pm 1.17$ | 0.47 | ✓ |
| Equity SmallCap | $0.12 \pm 0.29$ | -0.02 | ✓ |
| Equity Value | $0.11 \pm 0.05$ | -0.30 | |
| Fixed Credit | $-0.07 \pm 0.65$ | 0.43 | ✓ |
| Fixed Duration | $-0.33 \pm 1.26$ | 0.20 | ✓ |
| US Dollar | $-0.41 \pm 0.94$ | 0.16 | ✓ |

Finally, the interval sizes vary considerably among the dependent variables, with Equity Momentum showing the widest interval, indicating the highest uncertainty, while Equity Value has the narrowest, suggesting a high level of confidence in these projections.

## 4.   Conclusion and Further Research

We applied Principal Components Regression (PCR) coupled with cross-validation and LARS Lasso to generate projections for our factor set. By utilizing PCR *and* LARS Lasso, we optimally reduced the dimensionality of our dataset, enabling us to focus on the most significant components, minimize multicollinearity, and enhance predictive power.

Our approach involved standardizing the feature set, applying PCA for dimensionality reduction, using cross-validation and LARS Lasso for optimal principal component choice, and then fitting linear regression models to the transformed data.

The findings underscore the potential of PCR in enhancing decision-making in multiple areas of alternative investment management, providing a scalable and interpretable framework for predictive modeling.

Next steps involve:
- Using G3 factor loadings to project out asset-level returns; and
- Devising a technique for covariance matrix projection.

# Appendix

## Linear Algebra for "Converting" PCR Coefficients

Let $R_{T \times N}$ be the historical returns of $N$ indices corresponding to CMA asset classes for $T$ periods. Using Singular Value Decomposition (SVD):

$$R_{T \times N} = U_{T \times T} \Sigma_{T \times N} V_{N \times N}^T \tag{1}$$

where $U$, $V$ are the left and right orthonormal eigenvectors and $\Sigma$ has the singular values along its diagonal. We may reduce the dimensionality by limiting to the first $K$ eigenvalues (principal components):

$$R_{T \times N} \approx U_{T \times K} \Sigma_{K \times K} V_{K \times N}^T = X_{T \times K} V_{K \times N}^T \tag{2}$$

Let $F_{T \times P}$ be the historical returns of $P$ risk factors, we can propose a linear relationship between the factors and the CMAs as:

$$F_{T \times P} = R_{T \times N} \beta_{N \times P}^T \approx X_{T \times K} V_{K \times N}^T \beta_{N \times P}^T = X_{T \times K} \beta_{K \times P}^T \tag{3}$$

So we can estimate $\beta'_{K \times P} = f(F_{T \times P}, X_{T \times K})$. From the above equations we have –

$$r_N \approx V_{N \times K} x_K \tag{4}$$

$$f_P \approx \beta_{P \times K} x_K \approx B_{P \times K} V_{K \times N}^T r_N \tag{5}$$

Thus, given a vector of projected returns for the CMAs $r_N$, we can estimate a vector of projected returns for risk factors $f_P$.

## What is $k$-fold Cross-Validation?

$k$-fold cross-validation is a widely used technique for assessing model performance. One splits the dataset into $k$ equal-sized subsets or folds. The model is then trained $k$ times, each time using $k-1$ folds for training and the remaining fold for testing. This process is repeated $k$ times, with each fold being used exactly once for testing. For time series, we must preserve the sequence of datapoints, so while we use all $k$ folds, it is done on an expanding-window basis. In other words, train on $k^{\text{th}}$ fold, and test on the $k+1^{\text{th}}$ fold. This process ensures that every data point is used for testing exactly once. After each training iteration, the model's performance is evaluated using the testing fold, typically by computing a performance metric which, in our case, is RMSE. The RMSEs from each iteration are then averaged to obtain a final performance metric for the model. Figure 8 below illustrates the process.

$k$-fold cross-validation provides a more robust estimate of the model's performance compared to a single train-test split, as it evaluates the model on multiple subsets of the data. It also helps in detecting overfitting by assessing the model's generalization ability across different subsets of data.[§]

---

[§] For documentation, the code used in the `KFold` function from the `scikit-learn` package, please refer to https://scikit-learn.org/stable/modules/cross_validation.html.
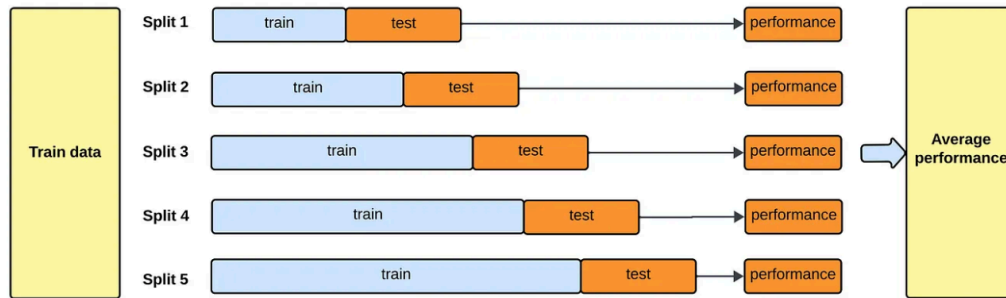
Fig. 8. An illustration of the cross-validation process for time series analysis.

We also consider leave-one-out cross-validation (LOOCV), which is a special case of $k$-fold cross-validation. Here, $k$ is equal to the number of samples in the dataset. In LOOCV, the model is trained $k$ times, each time leaving out one sample for testing and training on the remaining samples. This process is repeated for each sample in the dataset, and the performance metrics are averaged across all iterations to obtain the final evaluation metric.

$k$-fold cross-validation strikes a balance between computational efficiency and robustness by dividing the dataset into smaller subsets, whereas LOOCV exhaustively trains the model on nearly all possible combinations of training and testing samples, making it computationally expensive for larger datasets. $k$-fold cross-validation provides a more efficient estimate of model performance and is more commonly used in practice.

## What is LARS Lasso Regression?

LARS Lasso regression is a regularized linear regression technique that combines the Least Angle Regression (LARS) algorithm with the Lasso (Least Absolute Shrinkage and Selection Operator) regularization. It is useful when dealing with high-dimensional data or when feature selection is desired.

The main idea behind LARS Lasso is to efficiently solve the Lasso regression problem, which aims to minimize the sum of squared residuals while imposing an L1 penalty on the regression coefficients. The L1 penalty encourages sparsity in the coefficients, effectively shrinking some of them to exactly zero. This property makes Lasso regression well-suited for feature selection, as it automatically identifies and selects the most relevant features for the model.

The LARS algorithm is used to efficiently compute the Lasso solution path, which is the sequence of regression coefficients as the regularization parameter varies. It starts with all coefficients equal to zero and iteratively adds or removes variables based on their correlation with the residual. The algorithm proceeds in a stepwise manner, gradually increasing the magnitude of the coefficients until all variables are included in the model or a certain stopping criterion is met.

One of the key advantages of LARS Lasso is its computational efficiency. It can handle

high-dimensional data with a large number of features relatively quickly compared to other regularization techniques. Additionally, the Lasso regularization helps to prevent overfitting by constraining the magnitude of the coefficients and promoting a simpler and more interpretable model.

When using LARS Lasso, the regularization parameter (often denoted as $\lambda$) controls the strength of the L1 penalty. A higher value of $\lambda$ leads to more regularization and sparser coefficients, while a lower value allows for more complex models. The optimal value of $\lambda$ is determined through techniques such as cross-validation, as described above.
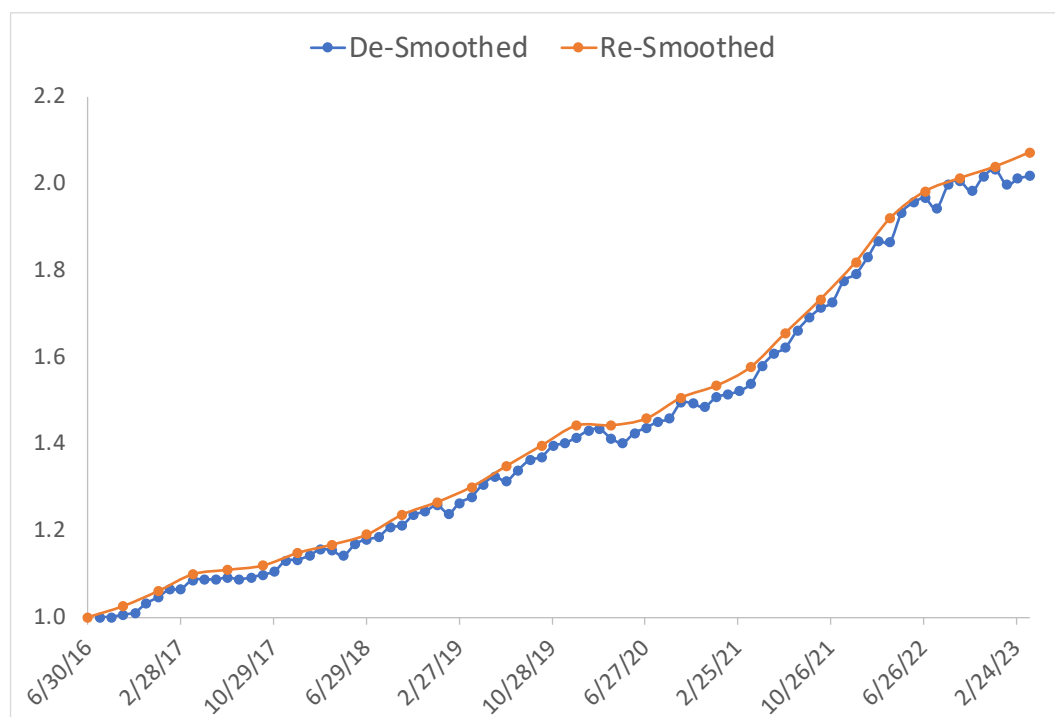
## What is Desmoothing?



Fig. 9. An example of a fund's de-smoothed ($x$), and re-smoothed ($y$) returns as generated by GMAM 3. The reported re-smoothed returns are quarterly, and the de-smoothed returns are monthly. Shown here is the cumulative return of $1 invested in each of the set of returns.

It is well-known in the finance practitioner literature that returns from alternative investment such as PE, hedge funds, or real estate funds are highly serially correlated (see, for example, Getmansky, et al. (2004).[¶] In other words, past values correlate with present values. The serial correlation occurs because of the lack of liquidity in the fund itself or some of the assets held within it. For instance, these illiquid assets may not trade frequently, leading to subjective and otherwise noisy valuations. The effect is such that when funds contain illiquid assets, their reported returns may seem steadier than their actual economic returns (returns that consider

---

[¶] Getmansky, M., A. Lo, and I. Makarov. 2004. An Econometric Model of Serial Correlation and Illiquidity in Hedge Fund Returns. *Journal of Financial Economics*, 74(3):529–609.

all available market information about those securities). The positive serial return correlation commonly leads to a downward bias in estimated return variance, or a smoothing effect.

The effect extends to the reported returns of real estate funds(see, for example, Geltner (1993).[††] Investors typically demand monthly or quarterly reporting. But valuations of many properties included in the funds are effectively updated only annually. Each quarter some properties have their valuations updated, and others do not. For some properties, the lack of a new valuation within a quarter might result in a carry-over of their last known value into the current quarter.

Finally, Financial Accounting Standard 157, released by the FASB in 2006 during the run-up to the financial crisis, and now called Accounting Standards Code Topic 820, requires companies to mark their assets to market. The rule was a radical change from historic cost accounting and required general partners to periodically mark the assets to market. This may also result in a managerial bias towards smoothing asset values.

The latent returns $x_{i,t}$, henceforth just $x$ for simplicity of notation, generated above are not smoothed. They are the true economic returns. Since GMAM 3 attempts to mimic true reported returns, we must re-smooth the latent returns to reflect the reported returns so that investors can appropriately compare their investments in Alts with publicly available investments. In GMAM 3, this is done using a moving average (MA) process with smoothing parameters estimated using a Bayesian linear regression. Note that the MA is defined in an econometric sense, and not in a literal sense, although it may be interpreted as such. More specifically, the final returns presented to the user are re-smoothed using the following process:

$$y = \Phi x \tag{6}$$

where $y$ are the reported returns; $\Phi$ is a matrix consisting of smoothing coefficients, $\phi$ and $\tilde{\phi}$ defined in more detail below); and $\mathbf{x}$ is a vector of $x$ values. When both $x$ and $y$ have the same reporting frequency:

$$\Phi = \begin{bmatrix} \phi_1 & \cdots & \phi_P & \tilde{\phi}_{P+1} & 0 & 0 & 0 & 0 & 0 \\ 0 & \phi_1 & \cdots & \phi_P & \tilde{\phi}_{P+1} & \cdots & \cdots & \cdots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & \cdots & \phi_1 & \cdots & \phi_P & \tilde{\phi}_{P+1} & 0 \\ 0 & \cdots & \cdots & \cdots & \cdots & \phi_1 & \cdots & \phi_P & \tilde{\phi}_{P+1} \end{bmatrix} \tag{7}$$

$\phi$ is a vector of $P$ unrestricted coefficients, and $\tilde{\phi}$ is a vector of $P + \Delta t$ restricted and unrestricted coefficients. $P$ and $\Delta t$ represent terms that map the observation period of the reported, smoothed returns $y$ to the latent, de-smoothed returns, $x$.

The reported returns are $y_s$ s.t. $s \in 1 : S$, and the latent returns are $x_t$ s.t. $t \in 1 : T$. In equation (7) above, $S = T - P$, and the reported returns and de-smoothed returns have the same frequency.

[††] Geltner, D. 1993. Estimating Market Values from Appraised Values Without Assuming an Efficient Market. *The Journal of Real Estate Research*, 8(3):324–345.

However, $S$ may not always equal $T$; for example, when the reported returns are quarterly, and the de-smoothed returns are monthly. To resolve this, the following formula maps $t$ to $s$:

$$t(s) = P + (s \times \Delta t) \tag{8}$$

For example, if the reported returns $y$ are quarterly, and the latent returns $x$ are monthly, $P = 3$, and $\Delta t = 3$ (since each quarter consists of three months). And so, for the first quarter, $s = 1$, and $t = 6$. When $\Delta t > 1$, the values in $\tilde{\phi}_j$ must add up to 1 for each month of the quarter (first, second, third). If not, months that fall earlier in the quarter will have a different long-run impact on NAV than months later in the quarter. For instance if we have a single quarter lag, the first month of the quarter plus the first month of the previous quarter must add to 1. This implies that the sum total of values in $\tilde{\phi}_j$ is 3 for quarterly data and 1 for monthly data..[‡‡]

In particular,

$$y_s = \left(\tilde{\phi}_{1:(P+\Delta t)}\right)' x_{(t[s]-P-\Delta t+1):t[s]} + \varepsilon_t^y \tag{9}$$
$$= \left(\tilde{\phi}_{(P+1):(P+\Delta t)}\right)' x_{(t[s]-\Delta t+1):t[s]} + \phi' x_{(t[s]-P):(t[s]-\Delta t)} + \varepsilon_t^y \tag{10}$$

The complete posterior distribution of $y$ is given by:

$$p(y \mid rest) \sim MN\left(\Phi x, \frac{1}{\tau_y} I\right)) \tag{11}$$

where $\tau_y$ is a global precision parameter estimated for independent measurement variance. The precision parameter has Gamma-distributed prior: $p(\tau_y) \sim Gamma\left(\alpha_{y_0}, \zeta_{y_0}\right)$ with given shape hyperparameter, $\alpha_{y_0}$, and given inverse scale hyperparameter, $\zeta_{y_0}$; $\tau_\phi$ is a global precision multiplier parameter, distributed as $p(\tau_\phi) \sim Gamma\left(\alpha_{\phi_0}, \zeta_{\phi_0}\right)$ with given shape and inverse scale hyperparameters; $M_0$ is a $P \times P$ matrix consisting of precision hyperparameters for $\phi_0$; and $\phi_0$ is a hyperparameter (i.e., fed into the model).

The economic assumption behind the smoothing process is simple: that the observed fund returns ($y$) are a weighted average of the fund's economic returns $x$ over the most recent $(P+\Delta t)$ periods, inclusive of the current period. The econometric implications of this are that under the given assumption, the observed fund returns follow a MA process of order $P+\Delta t$.

This restriction is similar to Getmansky, et al. (2004) where the observed return ($R_t^0$) for some period $t$, is a weighted average of the "true" returns ($R_t^C$) over the most recent $k+1$ periods: $R_t^0 = \theta_0 R_t^C + \cdots + \theta_k R_{(t-k)}^C$, with $\sum_{i=0}^k \theta_i = 1$ to ensure that all information is eventually incorporated into observed returns, and $\theta_i \in [0,1]$ for $i = 1, \ldots, k$. In our case the observed returns are the reported returns, $y$, the "true" returns are the latent returns generated using the factors, $x$, and the $\theta_i$ terms from Getmansky, et al. (2004) are denoted as $\phi_i$ in our estimation process. They are generated using a multivariate normal distribution in the hierarchical Bayesian model, as follows:

---

[‡‡] Note that the vector $\phi \in \tilde{\phi}$. For quarterly reported data, say $\phi = [a,b,c]'$ where $a,b,c \in \mathbb{R}$. Then, $\tilde{\phi} = [a,b,c,1-a,1-b,1-c]'$ and the sum of all terms in $\tilde{\phi}$ is 3. For monthly reported data, say $\phi = [a,b,c,d,e]'$ where $a,b,c,d,e \in \mathbb{R}$. Then $\tilde{\phi} = [a,b,c,d,e,1-a-b-c-d-e]'$ and the sum of all terms in $\tilde{\phi}$ here, is 1. This process is also described in more detail in the Appendix.

$$p\left(\phi|rest\right) \sim MN\left(\phi_0, \frac{1}{\tau_y \tau_\phi} M_0^{-1}\right) \tag{12}$$

where $\tau_y$ is as defined above; $\tau_\phi$ is a global precision multiplier parameter, distributed as $p(\tau_\phi) \sim Gamma\left(\alpha_{\phi_0}, \zeta_{\phi_0}\right)$ with given shape and inverse scale hyperparameters; $M_0$ is a $P \times P$ matrix consisting of precision hyperparameters for $\phi_0$; and $\phi_0$ is a hyperparameter (i.e., fed into the model).