

Kernel Conditional Factor Models

Pierre Collin-Dufresne^{1,2} Damir Filipović^{1,2} **Urban Ulrych**^{1,2}

¹École Polytechnique Fédérale de Lausanne

²Swiss Finance Institute

EEA-ESEM 2024

Erasmus School of Economics, Rotterdam

Version: August 2024

1. Introduction

- *Factor models* are one of the fundamental tools in finance, used to uncover the relationship between *asset returns* and *underlying factors*:

$$r_{i,t+1} = \alpha_{i,t} + \beta_{i,t}f_{t+1} + \epsilon_{i,t+1}.$$

- Factor models traditionally perform *factor* (f_{t+1}) and *beta* ($\beta_{i,t}$) *learning* in a linear manner, which does *not* capture the *non-linear dynamics* of financial markets.
- Our paper, Collin-Dufresne, Filipović, and Ulrych (2024), proposes an innovative approach to enhance factor models by *incorporating non-linearity* through *(low-rank) kernel functions*.
- Utilize a reproducing kernel Hilbert space (RKHS) with the associated reproducing kernel as *hypothesis space* for modeling *factor loadings*, and employ *cross-sectional ridge regression* to directly learn *factors* (i.e., factor portfolios).
- The proposed framework *improves* the *accuracy* and *predictive power* of factor models in asset pricing.

2. Kernel Factors

- Let $x_{t,i} \in \mathcal{X}$ denote the *characteristics* of assets $i \in \mathcal{I}_t$, where \mathcal{X} denotes the set of observable asset characteristics.
- Let $m \geq 1$ represent the number of factors.
- *Excess returns* of assets $i \in \mathcal{I}_t$ are, in vector notation, given as

$$\mathbf{r}_{t+1} = g(\mathbf{x}_t)F_{t+1} + \epsilon_{t+1}.$$

- Objective: *Learn* the *factor loading* function $g = (g_1, \dots, g_m) : \mathcal{X} \rightarrow \mathbb{R}^{1 \times m} \cong \mathbb{R}^m$, and *factors* $F_{t+1} \in \mathbb{R}^m$ by a cross-sectional ridge regression, similar as in Kelly et al. (2019).
- As hypothesis space for g we select \mathbb{R}^m -valued reproducing kernel Hilbert space (RKHS) $\mathcal{G} = \mathcal{H} \otimes \mathbb{R}^m$ with operator-valued reproducing kernel $K(x, x') = k(x, x')I_m$ on \mathcal{X} .

2. Alternating Kernel Ridge Regression

- We then solve the *regularized optimization problem*

$$\min_{g \in \mathcal{G}, F \in \mathbb{R}^{m \times T}} \{ \mathcal{E}(g, F) + \lambda_1 \|g\|_{\mathcal{G}}^2 + \lambda_2 \|F\|_2^2 \}, \quad (1)$$

for an error function $\mathcal{E} : \mathcal{G} \times \mathbb{R}^{m \times T} \rightarrow \mathbb{R}_{>0}$ and some penalty parameters $\lambda_1, \lambda_2 > 0$.

- Using the (weighted) mean squared error function, (1) can be equivalently written as

$$\min_{g \in \mathcal{G}} \left\{ \sum_{t=0}^{T-1} \min_{F_{t+1} \in \mathbb{R}^m} \left\{ \sum_{i \in \mathcal{I}_t} \omega_{t,i} (r_{t+1,i} - g(x_{t,i}) F_{t+1})^2 + \lambda_2 \|F_{t+1}\|_2^2 \right\} + \lambda_1 \|g\|_{\mathcal{G}}^2 \right\}, \quad (2)$$

which reflects the *alternating kernel ridge regression algorithm* for solving it.

- This is a *generalization* of the linear alternating least-squares approach of Kelly et al. (2019) towards *non-linear factor loadings*.

2. Cross-Sectional Ridge Regression

- For a given function $g \in \mathcal{G}$ and for each $t = 0, \dots, T - 1$, we solve the *cross-sectional ridge regression*:

$$\min_{F_{t+1} \in \mathbb{R}^m} \left\{ (\mathbf{r}_{t+1} - g(\mathbf{x}_t)F_{t+1})^\top \boldsymbol{\Omega}_t (\mathbf{r}_{t+1} - g(\mathbf{x}_t)F_{t+1}) + \lambda_2 \|F_{t+1}\|_2^2 \right\}. \quad (3)$$

- The solution to (3) is *unique* and, for each t , *explicitly* given by

$$\hat{F}_{t+1} = (g(\mathbf{x}_t)^\top \boldsymbol{\Omega}_t g(\mathbf{x}_t) + \lambda_2 \mathbf{I}_m)^{-1} g(\mathbf{x}_t)^\top \boldsymbol{\Omega}_t \mathbf{r}_{t+1},$$

reflecting that \hat{F}_{t+1} is a *factor portfolio*.

2. Time-Series Kernel Ridge Regression

- Conversely, for given factors $F_{t+1} \in \mathbb{R}^m$, solving the outer optimization of (2) for g amounts to *time-series kernel ridge regression*:

$$\min_{g \in \mathcal{G}} \left\{ \sum_{t=0}^{T-1} (\mathbf{r}_{t+1} - g(\mathbf{x}_t)F_{t+1})^\top \boldsymbol{\Omega}_t (\mathbf{r}_{t+1} - g(\mathbf{x}_t)F_{t+1}) + \lambda_1 \|g\|_{\mathcal{G}}^2 \right\}. \quad (4)$$

- By the *representer theorem*, this would lead to a ridge regression of dimension $M = \sum_{t=0}^{T-1} M_t$, which is too *computationally costly*.
- Following Filipović et al. (2023), we compute a *low-rank approximation* of the *kernel function*

$$k(x, x') \approx \sum_{j=1}^d \phi_j(x) \phi_j(x'),$$

for the orthonormal functions $\phi(\cdot) = (\phi_1(\cdot), \dots, \phi_d(\cdot))$ in \mathcal{H} .

2. Low-Rank Approximation

- Accordingly, we replace the full hypothesis space \mathcal{G} by the subspace \mathcal{G}_ϕ spanned by $\phi_j(\cdot)v_j$, for $v_j \in \mathbb{R}^m$ and $j = 1, \dots, d$. Any g in this subspace is of the form

$$g(\cdot) = \sum_{j=1}^d \phi_j(\cdot)v_j, \quad v_j \in \mathbb{R}^m, \quad (5)$$

and problem (4) becomes *quadratic* in $\mathbf{v} = [v_1; \dots; v_d] \in \mathbb{R}^{dm}$.

- Differentiating the objective function (4) in \mathbf{v} yields the FOC with the *unique solution*:

$$\hat{\mathbf{v}} = \left(\left(\sum_{t=0}^{T-1} (\phi(\mathbf{x}_t)^\top \Omega_t \phi(\mathbf{x}_t)) \otimes (F_{t+1} F_{t+1}^\top) \right) + \lambda_1 I_{dm} \right)^{-1} \left(\sum_{t=0}^{T-1} ((\phi(\mathbf{x}_t)^\top \Omega_t) \otimes F_{t+1}) \mathbf{r}_{t+1} \right).$$

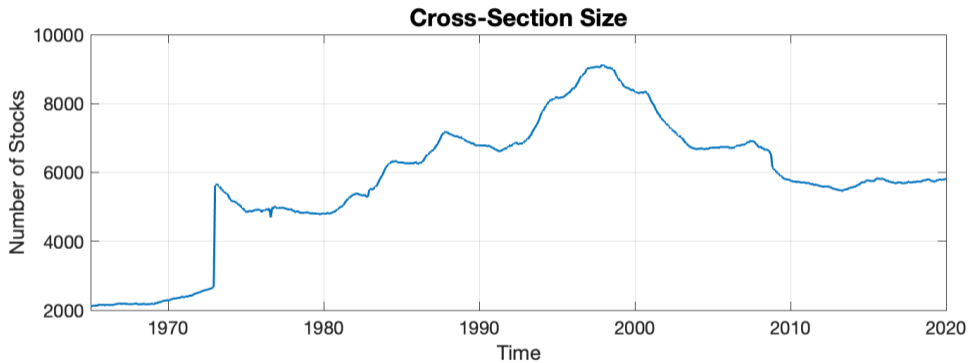
3. Kernel Selection

- Analyzed kernels:
 - ▶ Linear: $k(\mathbf{x}, \mathbf{y}) = 1 + \frac{\mathbf{x}^T \mathbf{y}}{\rho^2}$
 - ▶ Quadratic: $k(\mathbf{x}, \mathbf{y}) = (1 + \frac{\mathbf{x}^T \mathbf{y}}{\rho^2})^2$
 - ▶ Gaussian: $k(\mathbf{x}, \mathbf{y}) = e^{-\frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{2\rho^2}}$
- Extend the approach by *including industry classification* as an additional characteristic, encoded by $a \in \{1, \dots, A\}$.
- At any t and for any asset i , the original characteristics $x_{t,i}$ are thus extended to $(x_{t,i}, a_{t,i})$.
- A simple separable approach for *incorporating industries*:
 - ▶ $\bar{k}((\mathbf{x}, \mathbf{a}), (\mathbf{y}, \mathbf{b})) = k(\mathbf{x}, \mathbf{y})k_a(\mathbf{a}, \mathbf{b})$, with $k_a(\mathbf{a}, \mathbf{b}) = \begin{cases} 1, & \mathbf{a} = \mathbf{b} \\ \rho, & \text{else} \end{cases}$
for $\rho \in [0, 1]$.

3. Empirical Analysis

- Conducted an *Out-of-Sample* (OOS) *analysis* of the proposed smart kernel factor model.
- Analyzed *US stocks* from January 1965 to December 2019 with *monthly cross-sections* and 94 characteristics per stock.
- Rank-normalized all characteristics into the interval $(-1, 1)$ for each month t .
- Incorporated *industry classifications* ($A = 11$) based on Standard Industrial Classification (SIC) two-digit codes.
- Utilized a *training period* of 10 years, followed by *hyperparameter validation* on the subsequent 5 years.
- Performed a *rolling window out-of-sample backtest* with 10 years of training and 1 year of testing.

3. Empirical Analysis

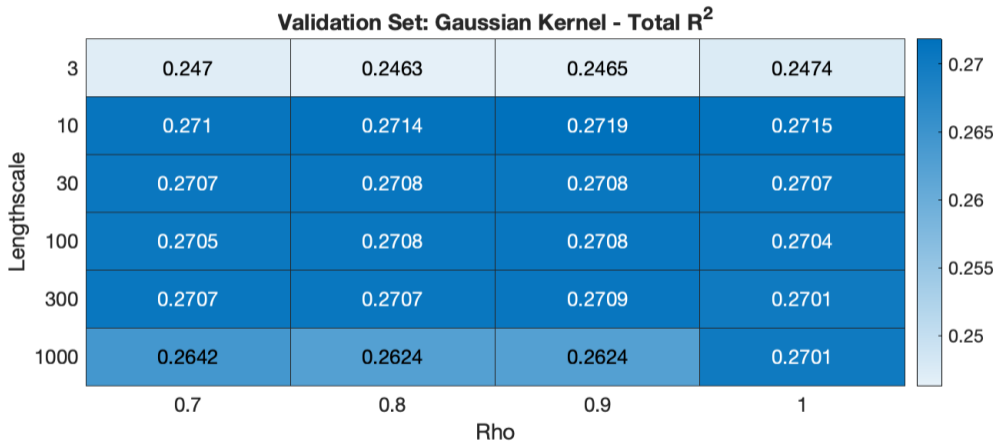


3. Empirical Performance

- $R_{Total}^2 = 1 - \frac{\sum_{(i,t) \in OOS} (r_{i,t} - \hat{g}(x_{i,t-1}) \hat{F}_t)^2}{\sum_{(i,t) \in OOS} r_{i,t}^2}$, $R_{Pred}^2 = 1 - \frac{\sum_{(i,t) \in OOS} (r_{i,t} - \hat{g}(x_{i,t-1}) \hat{\lambda}_{t-1})^2}{\sum_{(i,t) \in OOS} r_{i,t}^2}$, where $\hat{\lambda}_{t-1}$ is the prevailing sample average of \hat{F} up to month $t - 1$.
- Comparison of OOS R_{Total}^2 and R_{Pred}^2 for different methods and $m = 5$ factors:

Method ($m = 5$)	R_{Total}^2 (%)	R_{Pred}^2 (%)
1/ N Portfolio	8.80	0.02
Kelly et al.	14.11	< 0
Linear Kernel ($d = 94$)	14.13	0.65
Quadratic Kernel ($d = 200$)	14.42	0.22
Gaussian Kernel ($d = 200$)	14.38	0.36
Linear with Industry ($d = 95$)	14.20	0.74
Quadratic with Industry ($d = 200$)	14.43	0.11
Gaussian with Industry ($d = 200$)	14.40	0.06

3. Gaussian Kernel: Validation & Effect of Industry



4. Conclusion

- This paper presents a novel approach for learning *factors* and *factor loadings* in a *non-linear* manner using *kernel-based methodology*.
- We *extend* existing linear-learning-based approaches by introducing:
 - ① *Non-linear dependence* on characteristics, allowing for greater flexibility in modeling factor relationships. Notably, our linear kernel specification *nests* Kelly et al. (2019).
 - ② *Regularization*, allowing for more factors and *improving OOS performance* in terms of explained variation (R^2).
 - ③ Additional characteristics, such as *industries*, potentially enhancing model accuracy.
- A preliminary *empirical analysis* demonstrates that our kernel-based extension *outperforms* current linear-learning-based models in terms of OOS R^2 , demonstrating the *effectiveness* of *non-linear learning* and *regularization*.
- An *interpretable* approach with *fast* computation speed \implies *practical* applicability!

Thank you for your attention!

References I

- P. Collin-Dufresne, D. Filipović, and U. Ulrych. Smart kernel factors. *Work in Progress*, 2024.
- D. Filipović, M. Multerer, and P. Schneider. Kernel conditional distribution machines. *Working Paper*, 2023.
- B. T. Kelly, S. Pruitt, and Y. Su. Characteristics are covariances: A unified model of risk and return. *Journal of Financial Economics*, 134(3):501–524, 2019.

Appendix: Low-Rank Approximation

- We assume that the kernel matrix $\mathbf{K} := k(\mathbf{X}, \mathbf{X}^\top)$ admits a *low-rank approximation* $\mathbf{K} \approx \mathbf{L}\mathbf{L}^\top$ for some $M \times d$ -matrix \mathbf{L} and such that there exists a bi-orthogonal $M \times d$ matrix \mathbf{B} with
 - ▶ $\mathbf{K}\mathbf{B} = \mathbf{L}$,
 - ▶ $\mathbf{B}^\top \mathbf{L} = \mathbf{I}_d$, and
 - ▶ $\text{Im } \mathbf{B} = \text{span}\{\mathbf{e}_{\pi_1}, \dots, \mathbf{e}_{\pi_d}\}$,for some pivot indices $\{\mathbf{e}_{\pi_1}, \dots, \mathbf{e}_{\pi_d}\} \subseteq \{1, \dots, M\}$. That is, only d rows of \mathbf{B} are different from zero.
- Matrices \mathbf{B} and \mathbf{L} can be computed recursively, see Filipović et al. (2023).
- This yields the *low-rank approximation* of the *kernel function*

$$k(x, x') \approx \sum_{j=1}^d \phi_j(x) \phi_j(x'),$$

for the orthonormal functions $\phi(\cdot) = (\phi_1(\cdot), \dots, \phi_d(\cdot))$ in \mathcal{H} given by $\phi(\cdot) := \mathbf{B}^\top k(\cdot, \mathbf{X})$.