# Structural Deep Learning in Conditional Asset Pricing

Authors: Jianqing Fan, Zheng Tracy Ke, Yuan Liao, Andreas Neuhierl

Presented by Guilherme Piantino

# Content

# Conteúdo

# Introduction

**Motivation**

- Deep learning methods have been very successful in asset pricing, but they are often criticized as black boxes.
- The objective is to open the black box by joining rigorous asymptotic theory with finance theory. Thereby providing economic understanding for why deep learning models produce successful predictions.

**Contribution**

- They develop new nonparametric methods to obtain an economic interpretation of asset return predictions obtained from deep neural networks (DNN).
- The main objective is to understand deep learning prediction for the cross section of expected returns. To achieve this, they embed the neural network predictions in a factor pricing framework
- Finally they rigorously derive a theoretical framework of deep learning methods in asset pricing

# Related literature

- Most applications are based on "pooled machine learning", i.e., a single neural network function trained using data pooled cross-sectionally and over time.
- Differently, this paper study "period-by-period machine learning.
- Regarding the literature of conditional (or time-varying) asset pricing, the paper´s approach is different because it allows not only the mapping to characteristics but also to betas to vary over time.
- Considering the literature related to PCA, they applied local-PCA, which, differently from ordinary PCA, can be applied to conditional factor models, that is, with time-varying betas.

# Conteúdo

## The conditional asset pricing model

The time-varying factor model with intercepts:

$$y_{i,t} = \alpha_{i,t-1} + \beta'_{i,t-1}\lambda_{t-1} + \beta'_{i,t-1}(f_t\mathbb{E}(f_t|\mathcal{F}_{t-1})) + u_{i,t}$$

where $y_{i,t}$ is the excess return; $f_t$ is a K×1 vector of latent factors; $\alpha_{i,t-1}$ and $\beta_{i,t-1}$ respectively denote the (possibly) time-varying alpha and beta; $\lambda_{t-1}$ is the vector of factor risk premia; and $u_{i,t}$ is the idiosyncratic return.

Let $x_{i,t-1}$ be a d-dimensional vector of observed characteristics of stock i.

$$\alpha_{i,t-1} = g_{\alpha,t}(x_{i,t-1}) + \gamma_{\alpha,i,t-1} \qquad \mathbb{E}(\gamma_{\alpha,i,t-1}|x_{i,t-1}, f_t) = 0$$
$$\beta_{i,t-1} = g_{\beta,t}(x_{i,t-1}) + \gamma_{\beta,i,t-1} \qquad \mathbb{E}(\gamma_{\beta,i,t-1}|x_{i,t-1}, f_t) = 0$$

where $g_{\alpha,t}(\cdot)$ and $g_{\beta,t}(\cdot)$ are time-varying nonparametric functions of characteristics; $\gamma_{\alpha,i,t-1}$ and $\gamma_{\beta,i,t-1}$ respectively represent the source of alphas and betas not explained by characteristics

## Machine learning method

To obtain a prediction function at period t on the cross-sectional data $\{(y_{i,t}, x_{i,t1}) : i = 1, ..., N\}$, they apply DNN to solve for $\hat{m}_t(\cdot)$:

$$\hat{m}_t(\cdot) = \underset{m_t \in \mathcal{M}_{J,L}}{\arg\min} \sum_{i=1}^{N} (y_i - m_t(x_{i,t-1}))^2$$

Using $\hat{m}_t(\cdot)$, one can then compute both in-sample and out-of-sample expected returns:

$$\hat{y}_{i,t} := \hat{m}_t(x_{i,t-1}), \qquad \text{in-sample}$$
$$\hat{y}_{i,t+1|t} := \hat{m}_t(x_{i,t}), \qquad \text{out-of-sample}$$

The out-of-sample predictor is often used to predict $y_{i,t+1}$. But little interpretation has been given regarding the source of predictability in these models.

# Structural machine learning

The decompositions are obtained for a generic machine learning method, although in this paper only DNN estimation is used.

**In-sample decomposition**

$$y_{i,t} = \underbrace{g_{\alpha,t}(x_{i,t-1}) + \overbrace{\beta_{i,t-1}'\lambda_{t-1}}^{g_{riskP,t}(x_{i,t-1})} + g_{\beta,t}(x_{i,t-1})'(f_t - \mathbb{E}(f_t|\mathcal{F}_{t-1}))}_{\approx \hat{y}_t} + e_{i,t}$$

**Out-of-sample decomposition**

$$y_{i,t+1} = g_{\alpha,t}(x_{i,t}) + g_{riskP,t}(x_{i,t}) + \underbrace{g_{\beta,t}(x_{i,t})'(f_t - \mathbb{E}(f_t|\mathcal{F}_{t-1})) + e_{i,t}}_{\xi_{i,t+1}} + Op(1)$$

where $Op(1)$ converges to zero when $N \to \infty$. The mispricing component and the risk premium are indeed the only predictable parts. $g_{\alpha,t}(\cdot)$ and $g_{riskP,t}(\cdot)$ are assumed to change slowly over time.

## Estimation of components

The method for estimating the individual components of the return decomposition, that is the mispricing component, the factors, risk exposures and risk premia, is divided in three steps as below:

- **First:** They run period-by-period cross-sectional DNNs to estimate the nonparametric spot returns $\mathbb{E}(Y_s|X_s - 1, f_s)$.
- **Second:** They apply local PCA, running time-domain smoothing with a sequence of kernel-based weights, due to the possible nonlinearity of expected returns and time-varying of alphas/betas. Then they estimate betas $G_{\beta,t}(X_{t-1})$
- **Third:** They apply DNN to separately predict out-of-sample alphas and risk premium. The out-of-sample predictions can be constructed by plugging in $X_T$ to these estimated functions.

# Conteúdo

## Data and model parameters

- Data set of Chen and Zimmermann (2021), which contains monthly data of a set of firm specific characteristics.
- They delete all cases for which book-to-market is not observed and the first 24 months of observations to avoid forward looking biases.
- The remaining data set has 2,343,844 firm months observation starting in January 1955.
- Throughout, they use a 60 months window for estimation, sliding forward by one month, after each estimation.
- In the implementation, they use a one, two and three layer network with four nodes on each hidden layer. They use a learning rate of 0.001, 2000 epochs and a constant bandwidth of h = 0.75.

# In-sample return decomposition

- They decompose realized returns into a mispricing component $g_{\alpha,t}$, a risk-based component driven by the risk premium component $g_{\beta,t}(X)'\lambda_{t-1}$, and the exposure to the factor shock $g_{\beta,t}(X)'(f_t - \mathbb{E}(f_t|\mathcal{F}_{t-1}))$.

- They establish a benchmark of how much of realized return can be explained, relying on the the following measure of $R^2$:
$$R^2 = 1 - \frac{\sum_{i,t}(y_{i,t} - prediction_{i,t})^2}{\sum_{i,t} y_{i,t}^2}$$

- where $prediction_{it}$ is equal to the different parts of the return decomposition, to assess their explanatory strength, that is prediction is the same as fitting in this case.

## Table I: In-Sample Decomposition - Realized Returns (Full Sample)

This table shows empirical estimates for the in-sample decomposition of realized returns (equation (2.6)). $R^2_{\hat{y}}$ measures the quality of the in-sample fit from the period-by-period DNN regressions of excess returns onto characteristics. $R^2_{\beta'F}$ measures how much of the variation in excess returns can be explained by exposure to common factors. $R^2_{\beta'\lambda}$ measures how much of excess returns is explained by the factor risk premia, $R^2_{\beta'(F+\lambda)}$ measures how much can be explained by all risk-based components. $R^2_\alpha$ measures how much in-sample variation of excess returns can be explained by mispricing. Panel A shows the results for all firms in our CRSP/Compustat sample, Panel B focuses on the 20% largest firms and Panel C shows the results for the 80% smallest firms. All $R^2$ measure are in percentage. The sample period is 1955 - 2021.

| | 1 Layer | | | | | 2 Layers | | | | | 3 Layers | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $K$ | $R^2_{\hat{y}}$ | $R^2_{\beta'F}$ | $R^2_{\beta'\lambda}$ | $R^2_{\beta'(F+\lambda)}$ | $R^2_\alpha$ | $R^2_{\hat{y}}$ | $R^2_{\beta'F}$ | $R^2_{\beta'\lambda}$ | $R^2_{\beta'(F+\lambda)}$ | $R^2_\alpha$ | $R^2_{\hat{y}}$ | $R^2_{\beta'F}$ | $R^2_{\beta'\lambda}$ | $R^2_{\beta'(F+\lambda)}$ | $R^2_\alpha$ |
| Panel A: All Firms | | | | | | | | | | | | | | | |
| 1 | 74.14 | 20.73 | 0.99 | 21.68 | 0.98 | 83.43 | 21.10 | 0.98 | 22.03 | 1.19 | 84.39 | 21.13 | 0.99 | 22.06 | 1.21 |
| 5 | 74.14 | 35.58 | 1.27 | 36.76 | 0.70 | 83.43 | 37.20 | 1.28 | 38.35 | 0.89 | 84.39 | 37.84 | 1.28 | 38.99 | 0.91 |
| 10 | 74.14 | 45.18 | 1.35 | 46.41 | 0.62 | 83.43 | 48.28 | 1.35 | 49.49 | 0.81 | 84.39 | 48.63 | 1.36 | 49.84 | 0.83 |
| Panel B: Large Firms | | | | | | | | | | | | | | | |
| 1 | 64.05 | 23.80 | 0.54 | 24.71 | 1.28 | 70.42 | 23.82 | 0.51 | 24.67 | 1.41 | 70.33 | 23.81 | 0.50 | 24.67 | 1.40 |
| 5 | 64.05 | 38.61 | 1.01 | 40.10 | 0.87 | 70.42 | 38.97 | 0.96 | 40.42 | 1.01 | 70.33 | 38.91 | 0.95 | 40.37 | 1.01 |
| 10 | 64.05 | 44.63 | 1.15 | 46.35 | 0.76 | 70.42 | 45.94 | 1.09 | 47.60 | 0.94 | 70.33 | 45.95 | 1.09 | 47.58 | 0.93 |
| Panel C: Small Firms | | | | | | | | | | | | | | | |
| 1 | 75.34 | 20.37 | 1.04 | 21.32 | 0.95 | 84.97 | 20.77 | 1.04 | 21.72 | 1.16 | 86.06 | 20.82 | 1.04 | 21.76 | 1.19 |
| 5 | 75.34 | 35.22 | 1.31 | 36.36 | 0.67 | 84.97 | 36.99 | 1.32 | 38.11 | 0.87 | 86.06 | 37.72 | 1.32 | 38.83 | 0.90 |
| 10 | 75.34 | 45.24 | 1.37 | 46.42 | 0.61 | 84.97 | 48.56 | 1.38 | 49.72 | 0.80 | 86.06 | 48.95 | 1.39 | 50.11 | 0.82 |

The NN displays overfitting performance for $R^2_{\hat{y}}$, but the local-PCA to extract the estimated factor components and mispricing separately doesn´t (column $R^2_{\beta'(F+\lambda)}$)
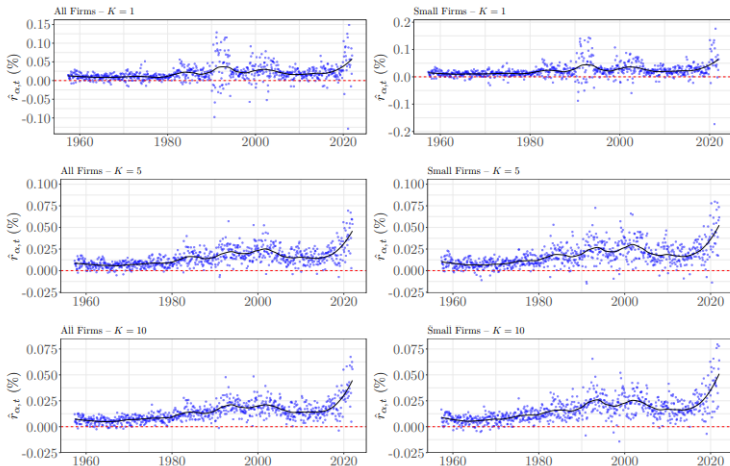
# Dynamics of the Pricing Error

- To investigate the time-variation, they analyse the dynamics of the pricing error, using a "denoised" version of the mispricing portfolio return:

$$\hat{r}_{\alpha,t} = \frac{1}{N_t} \hat{G}'_{\alpha,t-1} \hat{y}_t$$

- This quantity can be interpreted as an estimate of the squared pricing error often used in the examination of factor models

- $\hat{r}_{\alpha,t}$ cannot be interpreted as an excess return to a traded portfolio (because $\hat{y}_t$ are not the returns of traded assets), but it is a good measurement of the returns' magnitude because the idiosyncratic components have been removed

Figure 5.1: Evolution of Pricing Error over Time

This figures shows estimates of the average squared pricing error computed as $\frac{1}{N_t}\hat{\mathbf{G}}_{\alpha,t-1}(\mathbf{x})'\hat{\mathbf{y}}_t$ for all firms and $K=1$, $K=5$ and $K=10$ for the full sample (blue dots). We also present a local regression smoothing curve as an estimate of the local average (black line). The red dashed horizontal line is at zero.



The pricing error is correlated with volatility (the dot-com episode, the 2008/2009 financial crises and "covid-alpha"), with $\hat{a}=1.3502$ and $\hat{b}=0.11839$, significant at 1%
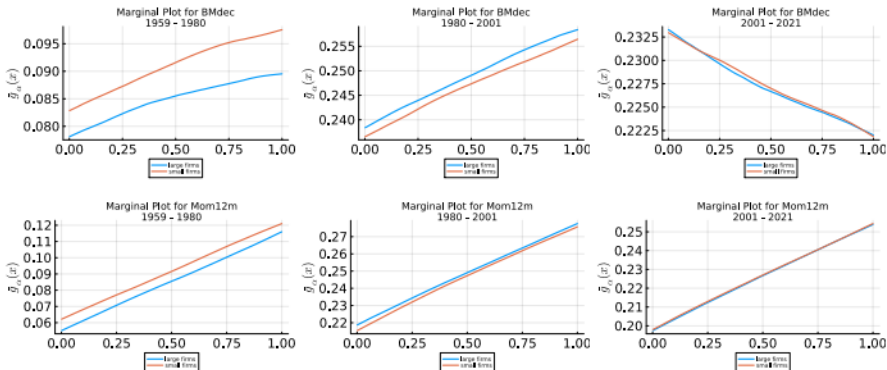
# Time-varying Functions

- To capture time-varying mispricing and risk premium dynamics, the authors use time-invariant functions, while the characteristics are time-varying

- The estimated marginal mispricing and the risk premium functions with respect to a single characteristic are:

$$h_j(z) := \frac{1}{|\mathcal{S}|} \sum_{t \in \mathcal{S}} g_t(z; x_{-j} = 0.5)$$

- where $g_t(z; x_{-j} = 0.5)$ denotes the function (of mispricing or risk-premium) marginally evaluated at the j-th characteristic, holding remaining characteristics at the mean-level 0.5 aggregated over a period S.

Figure 5.2: Time Variation in the Mispricing Function

The estimated marginal mispricing function $h_j(z)$ with respect to the book-to-market (top three panels) and the 12-month-momentum (bottom three panels). The estimated functions are separately estimated using data for large firms and small firms, and aggregated in three periods: 1959-1980, 1980-2001, and 2001-2021. We use ten factors ($K = 10$) for this analysis.

The plot demonstrates some degrees of time-varyingness when S changes.

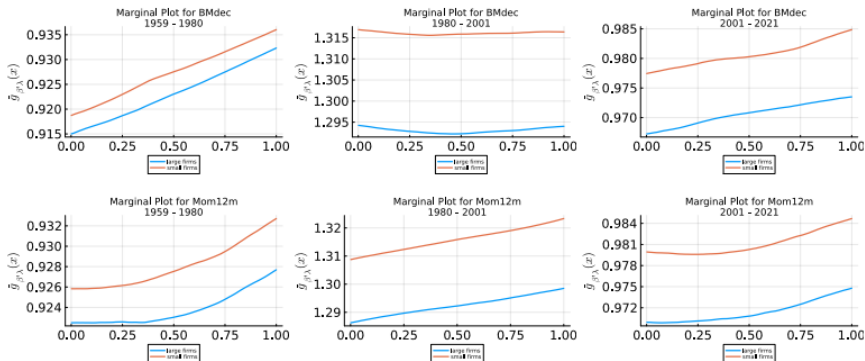## Figure 5.3: Time Variation in the Risk Premium Function

The estimated marginal risk premium function $h_j(z)$ with respect to the book-to-market (top three panels) and the 12-month-momentum (bottom three panels). The estimated functions are separately estimated using data for large firms and small firms, and aggregated in three periods: 1960-1980, 1980-2001, and 2001-2021. We use ten factors ($K = 10$) for this analysis.



The plot demonstrates some degrees of time-varyingness when S changes.

## Figure 5.4: Predictive $R^2$ for Different Re-Training Frequencies

The vertical axis is the predictive $R^2$ for predicting individual stock returns using mispricing function (left panel) and risk premium function (right panel). The functions are trained using either one-, two- or three- layer neural networks. The horizontal axis is the $\log(\tau)$, where $\tau$ is the re-training interval. The larger $\tau$ is, the less frequently the functions are trained. We use $K = 10$ factors in this analysis.



To study the effect of time-varying functions on the predictability, they retrain NN (for estimating $g_\alpha$ and $g_{\beta'\lambda}$) every $\mathcal{T}$ months with the frequencies $\mathcal{T} \in \{1, 12, ..., 360\}$. Both plots show a downward trend, indicating the predictive power decay as intervals increase

# Out-of-sample decomposition

- Computation the out-of-sample predictive $R^2$ as:
$$R^2 = 1 - \frac{\sum_{t \in \mathcal{T}} \sum_i (y_{i,t} - prediction_{i,t})^2}{\sum_{t \in \mathcal{T}} \sum_i y_{i,t}^2}$$

- Structural decomposition of the quantities learned from NN predictions: (i) expected return $\hat{y}_{t+1} := \hat{m}_t(x_{new})$, (ii) mispricing $\hat{g}_{\alpha,t}(x_{new})$, (iii) risk premium $\hat{g}_{riskP,t}(x_{new})$, (iv) alpha plus risk premium $\hat{g}_{\alpha+\beta'\lambda,t}(x_{new})$.

- Where $x_{new}$ refers to firm-level out-of-sample characteristics. Each of the four functions are fitted as a separate NN function using their in-sample estimates as the "data" and plug-in $x_{new}$ for predictions.

- They also use model averaging (MA) of previously estimated functions to make a one-period-ahead prediction

Table II: Out-of-Sample Decomposition - Expected Returns

This table shows empirical estimates for the out-of-sample decomposition of realized returns (equation (2.7)). $R_{\hat{y}}^2$ measures the quality of the out-of-sample fit from the period-by-period DNN regressions of excess returns onto characteristics. $R_{\beta'\lambda}^2$ measures how much of excess returns is explained by the factor risk premia, $R_{\beta'(F+\lambda)}^2$ measures how much can be explained by all risk-based components. $R_{\alpha}^2$ measures how much out-of-sample variation of excess returns can be explained by mispricing. All $R^2$ measure are in percentage. The sample period is 1960 - 2021.

| K | 1 Layer | | | | 2 Layers | | | | 3 Layers | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $R_{\hat{y}}^2$ | $R_{\beta'\lambda}^2$ | $R_{\alpha}^2$ | $R_{\beta'\lambda+\alpha}^2$ | $R_{\hat{y}}^2$ | $R_{\beta'\lambda}^2$ | $R_{\alpha}^2$ | $R_{\beta'\lambda+\alpha}^2$ | $R_{\hat{y}}^2$ | $R_{\beta'\lambda}^2$ | $R_{\alpha}^2$ | $R_{\beta'\lambda+\alpha}^2$ |
| Panel A: All Firms (period-by-period) | | | | | | | | | | | | |
| 1 | ≪ 0 | 0.25 | 0.08 | 0.30 | ≪ 0 | 0.32 | 0.10 | 0.33 | ≪ 0 | 0.44 | 0.11 | 0.43 |
| 5 | ≪ 0 | 0.32 | 0.01 | 0.31 | ≪ 0 | 0.40 | 0.02 | 0.36 | ≪ 0 | 0.48 | 0.04 | 0.44 |
| 10 | ≪ 0 | 0.32 | 0.02 | 0.27 | ≪ 0 | 0.38 | 0.05 | 0.33 | ≪ 0 | 0.49 | 0.04 | 0.43 |
| Panel B: Large Firms (period-by-period) | | | | | | | | | | | | |
| 1 | ≪ 0 | 0.51 | 0.32 | 0.48 | ≪ 0 | 0.59 | 0.39 | 0.64 | ≪ 0 | 0.90 | 0.42 | 0.83 |
| 5 | ≪ 0 | 0.56 | 0.21 | 0.50 | ≪ 0 | 0.67 | 0.20 | 0.65 | ≪ 0 | 0.97 | 0.29 | 0.82 |
| 10 | ≪ 0 | 0.61 | 0.13 | 0.44 | ≪ 0 | 0.67 | 0.24 | 0.59 | ≪ 0 | 0.96 | 0.22 | 0.80 |
| Panel C: Small Firms (period-by-period) | | | | | | | | | | | | |
| 1 | ≪ 0 | 0.23 | 0.06 | 0.28 | ≪ 0 | 0.30 | 0.07 | 0.31 | ≪ 0 | 0.40 | 0.09 | 0.40 |
| 5 | ≪ 0 | 0.30 | 0.00 | 0.29 | ≪ 0 | 0.38 | 0.01 | 0.33 | ≪ 0 | 0.45 | 0.02 | 0.41 |
| 10 | ≪ 0 | 0.29 | 0.01 | 0.25 | ≪ 0 | 0.35 | 0.03 | 0.31 | ≪ 0 | 0.45 | 0.03 | 0.41 |
| Panel D: All Firms (MA) | | | | | | | | | | | | |
| 1 | – | 0.60 | 0.08 | 0.60 | – | 0.63 | 0.09 | 0.65 | – | 0.64 | 0.09 | 0.66 |
| 5 | – | 0.61 | 0.06 | 0.61 | – | 0.63 | 0.06 | 0.65 | – | 0.65 | 0.06 | 0.66 |
| 10 | – | 0.60 | 0.05 | 0.60 | – | 0.63 | 0.05 | 0.65 | – | 0.65 | 0.06 | 0.66 |
| Panel E: Large Firms (MA) | | | | | | | | | | | | |
| 1 | – | 1.16 | 0.40 | 1.14 | – | 1.25 | 0.44 | 1.24 | – | 1.32 | 0.43 | 1.29 |
| 5 | – | 1.20 | 0.19 | 1.15 | – | 1.27 | 0.15 | 1.24 | – | 1.32 | 0.17 | 1.29 |
| 10 | – | 1.18 | 0.14 | 1.13 | – | 1.27 | 0.13 | 1.23 | – | 1.32 | 0.15 | 1.29 |
| Panel F: Small Firms (MA) | | | | | | | | | | | | |
| 1 | – | 0.56 | 0.06 | 0.56 | – | 0.58 | 0.06 | 0.60 | – | 0.59 | 0.07 | 0.61 |
| 5 | – | 0.56 | 0.05 | 0.57 | – | 0.58 | 0.05 | 0.60 | – | 0.60 | 0.06 | 0.61 |
| 10 | – | 0.55 | 0.04 | 0.56 | – | 0.59 | 0.05 | 0.60 | – | 0.60 | 0.05 | 0.61 |

$R_{\hat{y}}^2$ indicates bad predictability, due to low temporal dependence of returns. But $R_{\alpha+\beta'\lambda}^2$ is high, indicating greater predictive accuracy using risk premium and mispricing

# Out-of-Sample Factors

- True out-of-sample risk factors are unknown, however, they construct an out-of-sample proxy:

$$\hat{f}_{t+1} := \frac{1}{N} \hat{G}'_{\beta,t} y_{t+1}$$

- And to examine how well these factors explain the out-of-sample returns, they look at the total $R^2$:

$$R^2 = 1 - \frac{\sum_{t \in \mathcal{T}} \sum_i (y_{i,t} - \hat{G}'_{\beta,t} y_{t+1})^2}{\sum_{t \in \mathcal{T}} \sum_i y_{i,t}^2}$$

- For this analysis, the out-of-sample period starts in 1968, which is later than in the previous section, due to data availability which start in 1963. Five years were required for the estimation of parameter so out-of-sample results are from 1968 onward.

## Table III: Out-of-Sample $R^2_{\text{total}}(\%)$ Comparison

This table presents total R square $R^2_{\text{total}}$ of the out-of-sample factors explaining the returns. The out-of-sample period is from 1968-07-31 to 2021-12-31. We compare among models: the proposed $1 \sim 3$ layer structural neural networks; the Fama-French and Carhart models (FF). Here FF-1 refers to the CAPM; FF-3 refers to the Fama-French-three-factor model (Fama and French, 1992); FF-4 refers to the Carhart-four-factor model (Carhart, 1997); FF-5 refers to the Fama-French-five-factor model (Fama and French (2015), and FF-6 refers to the FF-5 factors plus the momentum factor.

|          | Number of Factors |       |       |       |       |       |
|----------|-------|-------|-------|-------|-------|-------|
| Model    | 1     | 3     | 4     | 5     | 6     | 10    |
| 1 layer  | 12.99 | 14.93 | 15.52 | 15.86 | 16.17 | 17.40 |
| 2 layers | 12.65 | 14.49 | 15.11 | 15.45 | 15.76 | 16.87 |
| 3 layers | 12.57 | 14.39 | 14.95 | 15.32 | 15.62 | 16.74 |
| FF       | 9.77  | 9.39  | 9.78  | 5.40  | 2.91  | -     |

The out-of-sample factors constructed by the structural NN, in all scenarios, explain higher percentage of out-of-sample total variations than the models of "observable factors" (Fama-French and Carhart models).

# Characteristic Importance

- They use the measure "characteristic relative importance" (CRI), defined using the marginal function with respect to the j-th characteristic, $h_j(\cdot)$.

- The identities of most important characteristics are different between the mispricing and the risk-premium functions.

- For mispricing, the two most important characteristics are "12-monthmomentum" and the "off-season-momentum", while for the risk premium are "beta" and "beta-FP".

- Results are intuitive: momentum and "beta" are insightful to respectively explain the mispricing and the risk premium.

- Few variables are relatively important to the risk function, while mispricing component seems to be explained by many more characteristics. This is consistent with economic models and intuition.

- Finally, the CRI also varies over time for both functions

# Conteúdo

# Conclusion

- The authors develop the theoretical justification for using machine learning estimators to predict and explain returns in the cross section

- If the regression is run period-by-period the regression approximates "$\alpha + \beta' F + \beta' \lambda$" and if it is pooled over time it converges to "$\alpha + \beta' \lambda$".

- They provide the asymptotic theory for neural network estimators, but the general approach is suitable for a generic machine learning method (random forests would generate similar results).

- The main empirical application is cross-sectional asset pricing of U.S. equities.

# Figure 6.1: Variable Importance in Mispricing Function

This figure shows characteristic relative importance measure (CRI) for the mispricing function. The left panel depicts the average measure of importance for the time period from 1960 to 1980, the middle panel represents the period from 1980 to 2001 and the left panel represents the period from 2001 through 2021. We use ten factors ($K = 10$) for this analysis.

# Figure 6.2: Variable Importance in Risk Premium Function

This figure shows characteristic relative importance measure (CRI) for the risk premium function. The left panel depicts the average measure of importance for the time period from 1960 to 1980, the middle panel represents the period from 1980 to 2001 and the left panel represents the period from 2001 through 2021. We use ten factors ($K = 10$) for this analysis.