# Risk Factor Detection with Methods from Explainable ML

**Natalie Packham**

joint work with

Marie Bernière (Humboldt University)

Quant Insights Conference

22 March 2023

Hochschule für
Wirtschaft und Recht Berlin
Berlin School of Economics and Law

# Introduction

# Overview

# Explainable ML

▶ Machine Learning (ML) techniques are now commonly used in finance applications to process large amounts of data.

▶ Ongoing challenges are missing transparency and missing interpretability: How do predictions and forecasts relate to the inputs?

▶ This will become more important with ongoing regulatory changes, e.g. (EC, 2021; EBA, 2021).

▶ Here: First results from a research project on explainable ML funded by IFAF for the next two years.



IFAF  Institut für angewandte Forschung Berlin

## Explainable ML

► Application in mind is **stress testing**.

► Classical setting: Factor model with **observable factors** (e.g. geographic regions, industries).

► Giving **latent factors** an interpretation extends range of stress scenarios.

► Concrete case: Use **Principal Component Analysis (PCA)** to determine latent factors from class factors and give them an interpretation.

► Idea goes back to work recent work on stress testing, (Packham and Woebbeking, 2019, 2023).

# Overview

## Linear factor model

▶ **Linear factor model**: Express vector of asset returns $(r_1, \ldots, r_p)$ as

$$r_i = \alpha_i + \beta_{i1}F_1 + \beta_{i2}F_2 + \cdots + \beta_{id}F_d + \varepsilon_i, \qquad i = 1, \ldots, p,$$

where

- $F_1, \ldots, F_d$: return in **common factors**,
- $\beta_{i1}, \beta_{i2}, \ldots, \beta_{id}$: **factor coefficients** or **factor weights**,[1]
- $\alpha_i$: constant,
- $\varepsilon_i$: **residual** or **idiosyncratic component**.

▶ Common assumption: residuals are uncorrelated.

▶ Number of factors small compared to number of securities, $d \ll p$.

---

[1]Sometimes called loadings. We will use the term "loadings" in a slightly different context.

# Linear factor model

▶ Factors $F_1, \ldots F_d$ **observable**, e.g. index returns of geographic regions and industries (MSCI GICS).

▶ Dependence structure of large portfolios expressed via covariances of common factors.

▶ Decompose $p \times p$ covariance matrix of returns $(r_1, \ldots, r_p)$ into

$$\Sigma \approx B\, \Omega\, B^T,$$

where

   – $B$: $p \times d$ matrix of factor coefficients,
   – $\Omega$: $d \times d$ covariance matrix of common factors, and
   – we ignore the variances of the residuals.

▶ Examples of factor models in credit risk management: Moody's KMV, CreditMetrics (by RiskMetrics), see e.g. Bluhm *et al.* (2003).

## Classical stress testing

▶ For "classical" stress testing method, see e.g. (Kupiec, 1998; Dowd, 2002; Packham and Woebbeking, 2019).

▶ Separate factors into **"core"** and **"peripheral" factors**.

▶ $\boldsymbol{F}_s$: $j < d$ **core factor returns** that are **stressed directly**.

▶ Remaining $d - j$ peripheral factor returns $\boldsymbol{F}_u$ indirectly affected by stress scenario.

▶ Under normal distribution assumption, optimal estimator of $\boldsymbol{F}_u | \boldsymbol{F}_s$ [2]:

$$\mathbb{E}(\boldsymbol{F}_u | \boldsymbol{F}_s) = \Sigma_{us} \Sigma_{ss}^{-1} \boldsymbol{F}_s,$$

where $\Sigma_{us}$ and $\Sigma_{ss}$ denote covariance and variance matrices of $\boldsymbol{F}_u$ and $\boldsymbol{F}_s$.

▶ See (Bonti *et al.*, 2006) for more advanced stress testing method.

---

[2]For simplicity, we assume the factor returns have expectation zero

## Stress testing with latent factors

▶ Goal here is to expand the universe of risk factors by **aggregating** existing factors into new factors.

▶ Examples: Global risk factor, European risk factor, cyclical industries, etc.

▶ Idea:
  – Use PCA on **observable factors** to determine aggregated (latent) factors.
  – Give these factors an interpretation.

# Overview

# Principal Component Analysis

▶ In $\mathbb{R}^n$, **PCA** refers to a particular rotation of the axes, driven by random variables or data.

▶ Key idea is to align random variables / data such that
  – first dimension captures maximal variance,
  – second dimension is orthogonal and captures second-most variance,
  – etc.

▶ **Principal components (PCs)** are the **eigenvectors** of covariance / correlation matrix.

▶ **Eigenvalues** express **amount of variance** captured by each PC.

## Principal Component Analysis

▶ See James *et al.* (2013), Section 10.2, for the following.

▶ Given $n \times d$ data set $\mathbf{X}$ that is **standardised**.

▶ **First principal component**: find **scores**

$$z_{i1} = \phi_{11} x_{i1} + \phi_{21} x_{i2} + \cdots + \phi_{d1} x_{id}, \quad i = 1, \ldots, n,$$

that have largest sample variance, subject to constraint $\sum_{j=1}^{p} \phi_{j1}^2 = 1$.

▶ In other words, **first PC vector**[3] solves optimisation problem

$$\max_{\phi_{11}, \ldots, \phi_{p1}} \left\{ \frac{1}{n} \sum_{i=1}^{n} \underbrace{\left( \sum_{j=1}^{p} \phi_{j1} x_{ij} \right)^2}_{=z_{i1}^2} \right\} \quad \text{subject to} \quad \sum_{j=1}^{p} \phi_{j1}^2 = 1.$$

▶ Second (and higher) PCs: linear combination of data uncorrelated with first PC(s) and with largest variance (subject to constraint).

---

[3]Called loading vector in (James *et al.*, 2013).

## Principal Component Analysis

▶ Compact notation (recall that $\mathbf{X}$ is standardised):

$$\mathbf{Z} = \mathbf{\Phi}^T \mathbf{X}$$

▶ PCs can be viewed as factors, giving factor model

$$\mathbf{X} = \mathbf{\Phi}\, \mathbf{Z}.$$

▶ $\mathbf{\Phi}$ are the eigenvectors of correlation matrix of $\mathbf{X}$.

▶ Example:

# Principal Component Analysis

▶ Example from Mathematica:

# Overview

# Overview

# Data

► Geographical factors: 16 regions and countries represented by MSCI indices

► Industry factors: 11 MSCI Global Industry Classification Standard (GICS) sector indices

► Daily data, split into Jan 1999-Dec 2019 (train) and Jan 2020-Feb 2023 (test)

► Data from Refinitiv Eikon

► Data split into **six groups**:
  - Europe (developed)
  - Asia-Pacific (developed)
  - N. America
  - Emerging Markets (Europe, M. East, Africa, Asia, Latin Am.)
  - Cyclical industries
  - Defensive industries

# Overview

# Loadings

- Giving PC's an interpretation: correlation between data and scores ($=$ projection of data to PC).

- Assume that data standardised.

- Using that the PC's are uncorrelated and have variances $\lambda_i$, $i = 1, \dots, d$:

$$\mathsf{Corr}(x_{\cdot j}, z_{\cdot i}) = \frac{\mathsf{Cov}(x_{\cdot j}, z_{\cdot i})}{\sqrt{\lambda_i}} = \frac{\mathbb{E}[\phi_{ji} z_{\cdot i} z_{\cdot i}]}{\sqrt{\lambda_i}} = \phi_{ji} \sqrt{\lambda_i}.$$

- In words: correlation of data and scores are just PCs scaled with PC standard deviation ("importance" of PC).

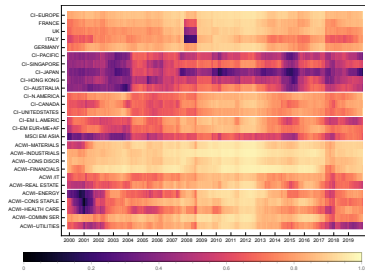- In-line with $\approx 50\%$ of the literature, we shall call these **loadings**.

# Loadings

▶ (Absolute) loadings, i.e., correlations of factor returns with first PCs:

# Loadings

- ▶ PCA at the end of each month on a rolling window of 250 days.

- ▶ A few more loadings plots:



- ▶ And a movie: Open Movie, Download movie

## Loadings

▶ Loadings of PCs through time (top: PC1, PC2; bottom: PC3, PC4):



Results 24

# Overview

# Interpretation of PCs

► Two questions:
  – How many PCs are relevant?
  – Which geographic region or industry group does PC explain?

► Literature: (Fenn *et al.*, 2011)

# Overview

# How many PCs are relevant?



Percentage variances over time

# How many PCs are relevant?



Cumulative Percentage variances over time
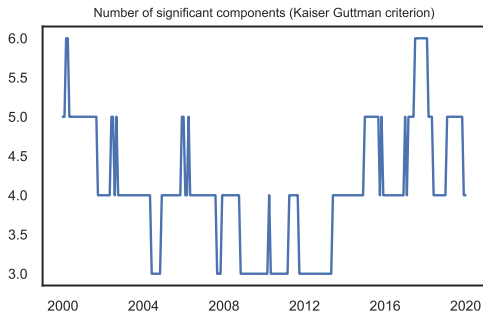
# Kaiser-Guttman criterion

▶ The **Kaiser-Guttman criterion** measures the number of significant PCs.

▶ The $i$-th PC is significant if its (normalised) eigenvalue $\lambda_i$ is greater than $1/d$, where $d$ is the number of eigenvalues.

▶ Idea: A PC that satisfies this criterion accounts for more than a fraction $1/d$ of the variance.

▶ See e.g. (Fenn *et al.*, 2011; Guttman, 1954).



Number of significant components (Kaiser Guttman criterion)
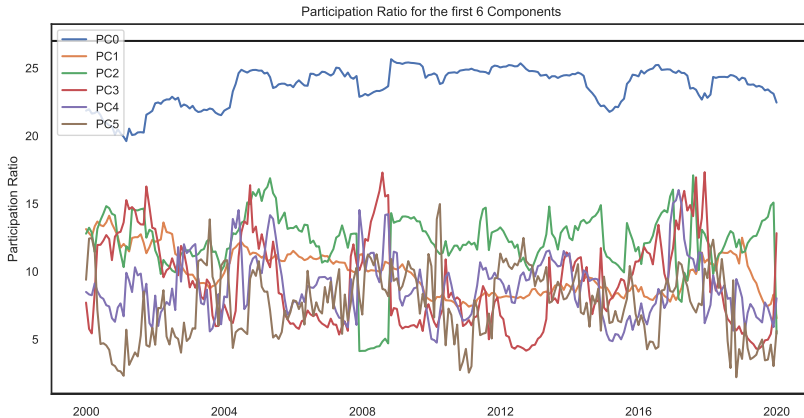
# Overview

## What drives changes in PCs?

► **Inverse participation ratio (IPR)** of $i$-th PC: (Fenn *et al.*, 2011; Guhr *et al.*, 1998):

$$I_k = \sum_{j=1}^{d} (\phi_{ji})^4.$$

► IPR measures number of assets participating in a PC:
  – eigenvector with equal contributions $\phi_{ji} = 1/\sqrt{d}$ has $I_k = 1/d$;
  – eigenvector with single contribution $\phi_{ji} = 1$ (others zero) has $I_k = 1$.

► **Participation ratio (PR)**: $1/I_k$

► Large PR: Many assets contribute

# Participation Ratio



Participation Ratio for the first 6 Components

# Overview

# Which geographic region or industry group does PC explain?

- Six groups:
  - Europe (developed)
  - Asia-Pacific (developed)
  - N. America
  - Emerging Markets
  - Cyclical industries
  - Defensive industries

- For a given PC and its PR, define the **PR group** as the group of **size PR** of indices with **highest loadings**.

- Group explained / not explained by a particular PC:

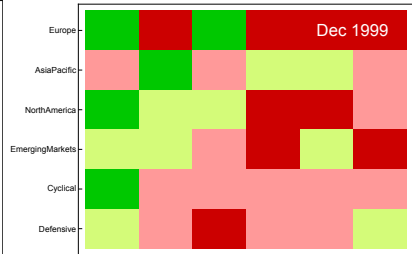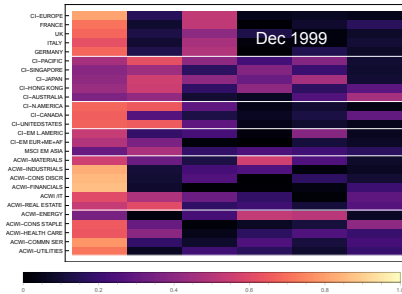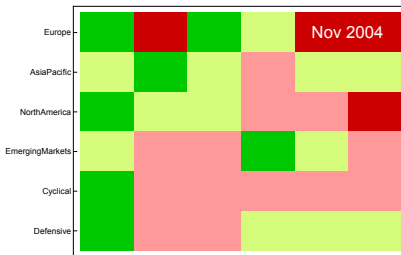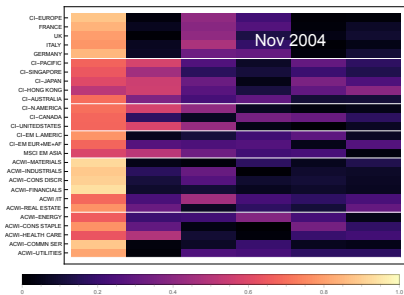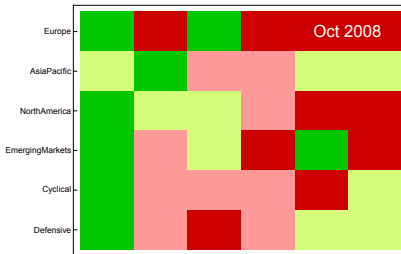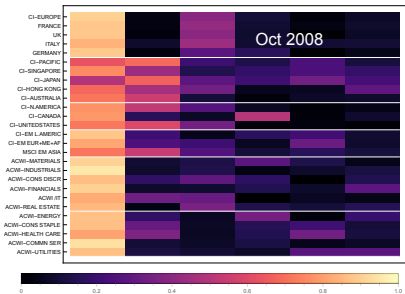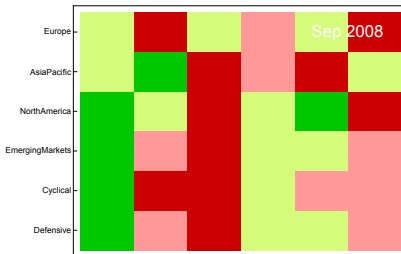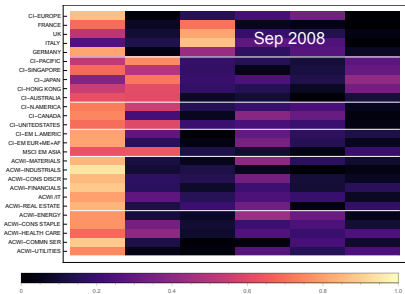| | |
|---|---|
| Strong In:<br>*All* indices in a group are in the PR group. | Strong Out:<br>*No* indices in a group are in the PR group. |
| Weak In:<br>*More than half* of indices in a group are in the PR group. | Weak Out:<br>*Half or less* of indices in a group are in the PR group. |

# PC interpretation



- ▶ First PC: global ex-AsiaPacific
- ▶ Second PC: AsiaPacific, EM, NA
- ▶ Third PC: Europe, NA

# PC interpretation

# PC interpretation

# PC interpretation

# PC interpretation

# Some observations

- Procedure selects into appropriate basket, but does not indicate strength of correlation (e.g. May 2012 vs. Oct 2017).

- For strength, consider eigenvalue.

- First PC is always a global risk factor, often ex-Asia-Pacific.

- Cyclical industries are always strong in global factor; defensive industries less strong.

- Second PC is Asia-Pacific factor, mostly with North America.

- Third PC is typically Europe with North America.

## Stress testing with aggregated risk factors

▶ Global stress scenario: adjust first PC or first two PCs, e.g. by choosing an explicit historical scenario or a historical realisation at a specific quantile.

▶ Asia-Pacific scenario: adjust second PC

▶ European scenario: adjust first and third PC

▶ North America scenario: adjust first and second PC

▶ Scenario "global economy more (less) connected": choose historical scenario where first PC's loadings are high (low)

# Overview

# Conclusion

► Factor models are used in various finance applications e.g. to estimate high-dimensional covariance matrices or in stress testing.

► Principal component analysis on a multivariate data set yields a factor model with latent factors.

► This is considered an unsupervised learning method.

► We attempt to give PCs on a data set consisting of risk factors (geographic regions and industries) an interpretation.

► Possible applications:
  – increase range of stress test scenarios
  – further decrease number of factors required for robust covariance matrix estimation

## Outlook

▶ Possibly of interest: Alternative methods find relevant factors across a number of PCs (e.g. (Mao, 2005; Masaeli *et al.*, 2010; Enki *et al.*, 2013; Chang *et al.*, 2016).

▶ Possibly use Varimax instead of PCA (Kaiser, 1958). Varimax attempts to find axes with few large loadings and many near-zero loadings.

▶ Non-linear relationships: Kernel-PCA, Autoencoder.

# References I

Bluhm, C., L. Overbeck, and C. Wagner. *An Introduction to Credit Risk Modeling*. Chapman & Hall/CRC, London, 2003.

Bonti, G., M. Kalkbrener, C. Lotz, and G. Stahl. Credit risk concentrations under stress. *Journal of Credit Risk*, 2(3):115–136, 2006.

Chang, X., F. Nie, Y. Yang, C. Zhang, and H. Huang. Convex sparse pca for unsupervised feature learning. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 11(1):1–16, 2016.

Dowd, K. *Measuring market risk*. Wiley, 2002.

EBA. EBA discussion paper on Machine Learning for IRB Models. European Banking Authority, EBA / DP / 2021 /04, November 2021.

EC. Laying down harmonised rules on Artifical Iintelligence (Artificial Intelligence Act) and amending certain Union legislative acts. European Commission, April 2021. https://eur-lex.europa.eu/legal-content/EN/TXT/HTML/?uri=CELEX:52021PC0206.

Enki, D. G., N. T. Trendafilov, and I. T. Jolliffe. A clustering approach to interpretable principal components. *Journal of Applied Statistics*, 40(3):583–599, 2013.

Fenn, D. J., M. A. Porter, S. Williams, M. McDonald, N. F. Johnson, and N. S. Jones. Temporal evolution of financial-market correlations. *Physical review E*, 84(2):026109, 2011.

Guhr, T., A. Müller-Groeling, and H. A. Weidenmüller. Random-matrix theories in quantum physics: common concepts. *Physics Reports*, 299(4-6):189–425, 1998.

# References II

Guttman, L. Some necessary conditions for common-factor analysis. *Psychometrika*, 19(2):149–161, 1954.

James, G., D. Witten, T. Hastie, and R. Tibshirani. *An introduction to statistical learning*, volume 112. Springer, 2013.

Kaiser, H. F. The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23(3):187–200, 1958.

Kupiec, P. Stress testing in a Value at Risk framework. *Journal of Derivatives*, 6:7–24, 1998.

Mao, K. Identifying critical variables of principal components for unsupervised feature selection. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 35(2):339–344, 2005.

Masaeli, M., Y. Yan, Y. Cui, G. Fung, and J. G. Dy. Convex principal feature selection. In *Proceedings of the 2010 SIAM international conference on data mining*, pages 619–628. SIAM, 2010.

Packham, N. and C. F. Woebbeking. A factor-model approach for correlation scenarios and correlation stress testing. *Journal of Banking & Finance*, 101:92–103, 2019.

Packham, N. and F. Woebbeking. Correlation scenarios and correlation stress testing. *Journal of Economic Behavior & Organization*, 205:55–67, 2023.

**Thank you!**

**Prof. Dr. Natalie Packham**
Professor of Mathematics and Statistics
Berlin School of Economics and Law
Badensche Str. 52
10825 Berlin
natalie.packham@hwr-berlin.de

Hochschule für
Wirtschaft und Recht Berlin
Berlin School of Economics and Law