

Automation of risk profiling and asset allocation processes with machine learning

Simone Gigante Gassó
Universitat Pompeu Fabra
05 - 07 - 2022

Table of contents

- Introduction
- Risk profile
- Clustering
- Clustering Results
- Asset allocation
- Reinforcement learning
- Results
- Conclusion

Introduction: What? Why? How?

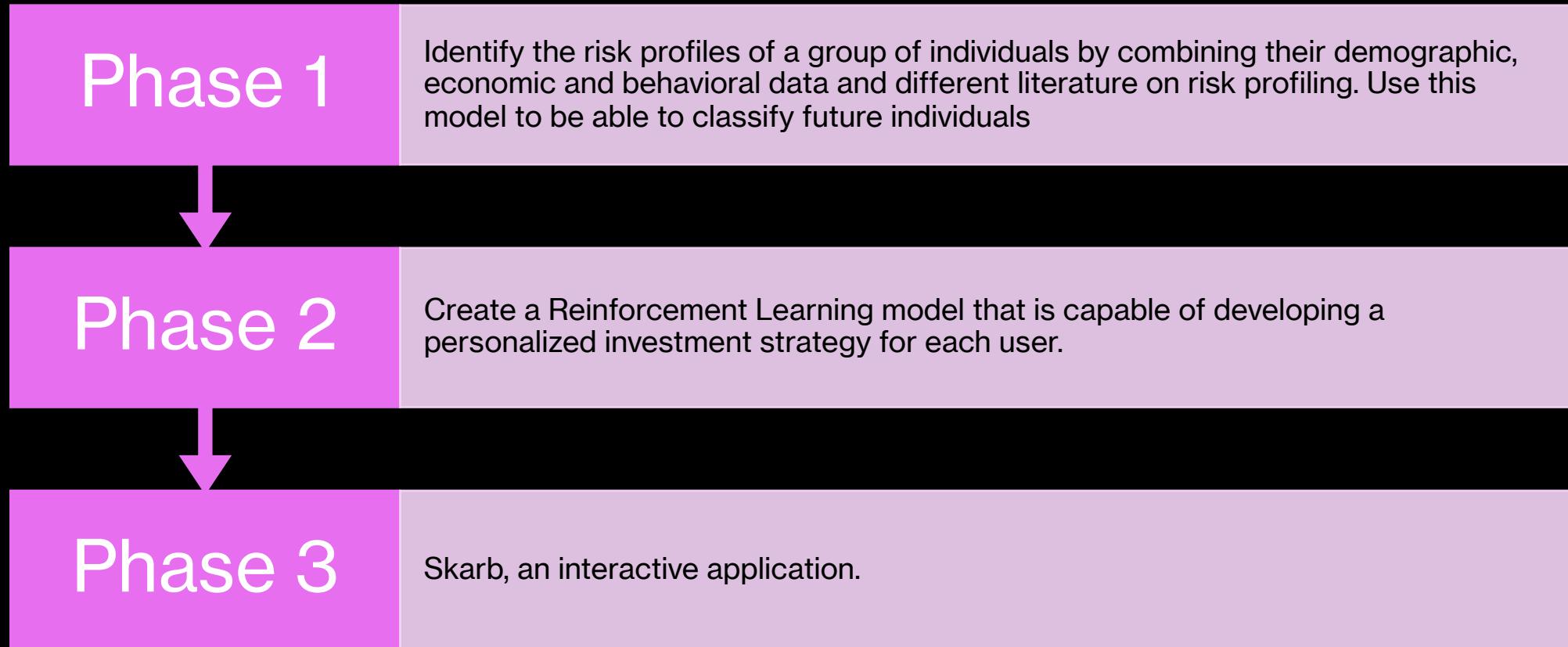
- The main objective of this thesis is to find a way to automate two of the most important processes in the wealth management sector: The creation of an individual's **risk profile**, and the development of a diversified investment strategy (**asset allocation**).



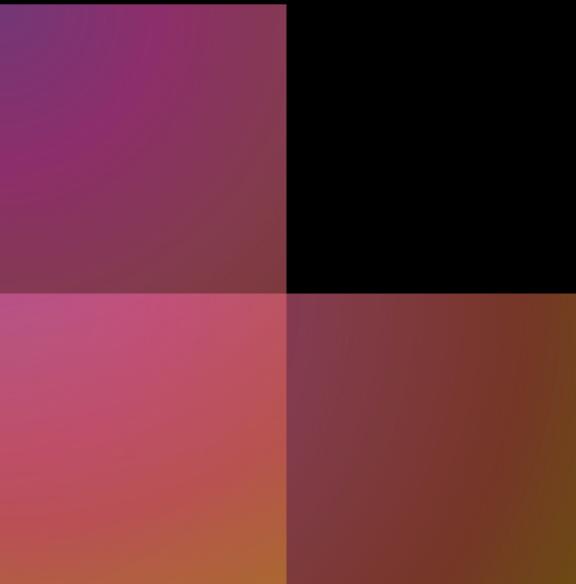
Introduction: What? Why? How?

- Many reasons:
 - There is no standard despite the fact that important institutions such as the EU or the US government insist on finding a generalized and proven form of financial advice in order to avoid catastrophes such as the 2008 crisis.
 - Apply machine learning to complex real-world problems, having to overcome various difficulties along the way.
 - Democratize wealth management
 - Personal interests

Introduction: What? Why? How?



Risk profile



Objective and subjective characteristics of an individual.



Risk aversion is the investor's willingness to take financial risk and the degree of psychological pain the investor experiences when faced with a financial loss.



Risk capacity is the objective capacity of an individual to assume a financial risk based on his economic circumstances.

Data

Encuesta Financiera de
las Familias (EFF 2017)

6,413 Spanish
households

+2500 questions

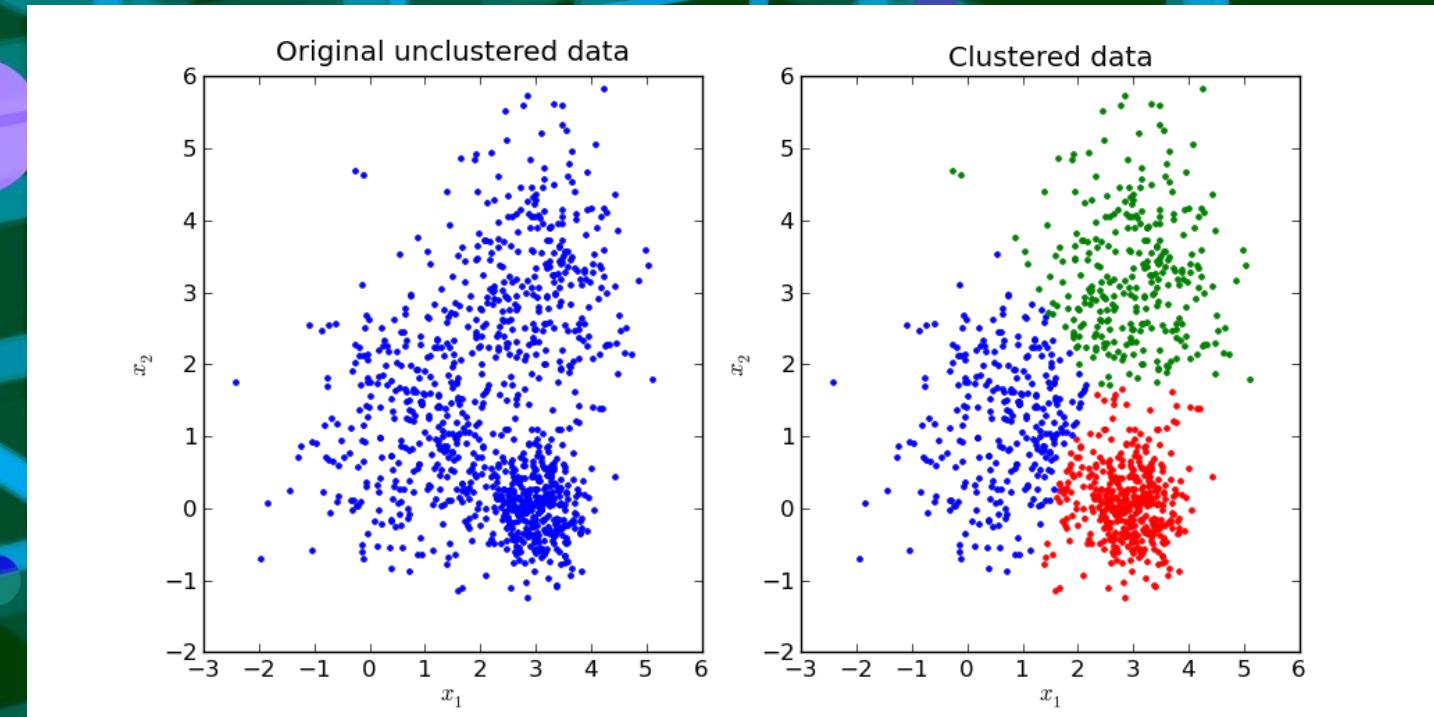
Data

Question
Age
Marital status
Maximum educational level achieved
Value of the main residence
Future expenses
Percentage you would spend in the next 12 months if you would win a lottery
Financial risk you are willing to take when you save or make an investment
Total amount allocated each year to pension plans
Total value of pension plans
Total monthly income
Total monthly household expenses
Total amount of money in bank accounts
Total amount owed on loans still to be paid
Total value of investments
Annual investment income

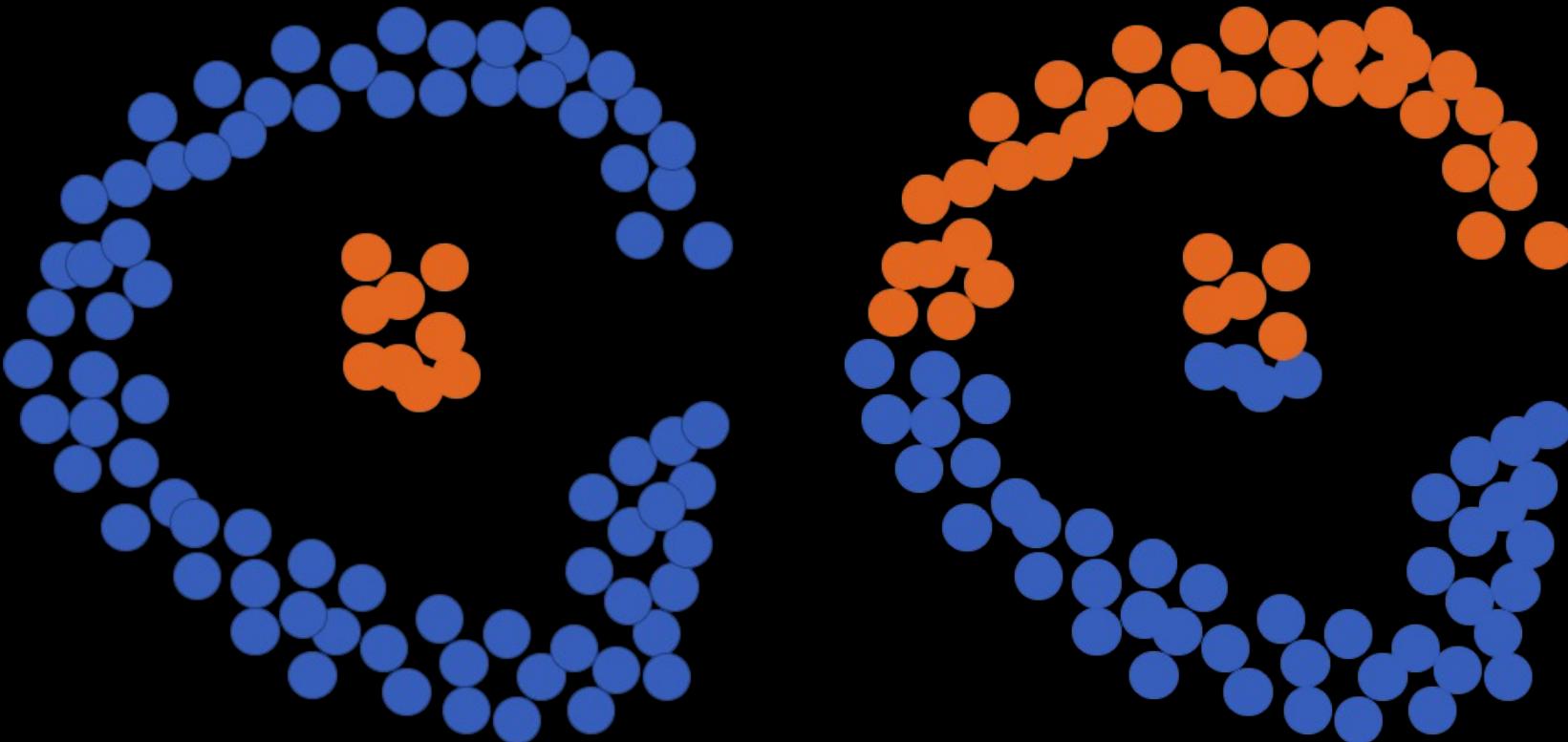
Clustering

Clustering is the task of identifying similar instances and assigning them to groups of similar instances.

Is an unsupervised task.



Clustering



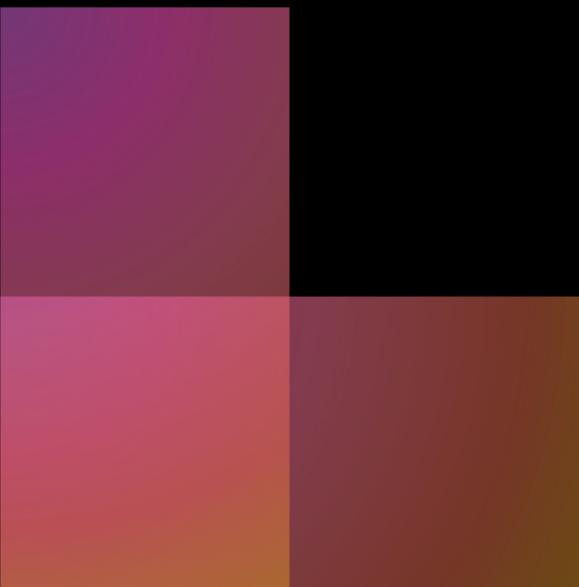
Evaluation metrics

- There are many different metrics (Silhouette score, Calinski Harabasz, Davies Bouldin, S_Dbw index etc...)
- *Understanding and enhancement of internal clustering validation measures* (Liu et al. 2013) shows why S_Dbw is the best option to compare different clustering techniques.

$$Dens_bw(c) = \frac{1}{c \cdot (c - 1)} \sum_{i=1}^c \left(\sum_{\substack{j=1 \\ i \neq j}}^c \frac{density(u_{ij})}{\max\{density(v_i), density(v_j)\}} \right) \quad (1)$$

$$Scat(c) = \frac{1}{c} \sum_{i=1}^c \|\sigma(v_i)\| / \|\sigma(S)\|$$

Data visualization: t-SNE



t-SNE (t-distributed stochastic neighbor embedding) algorithm allow interpreting n-dimensional data for $n \geq 3$ in two-dimensional graphs.

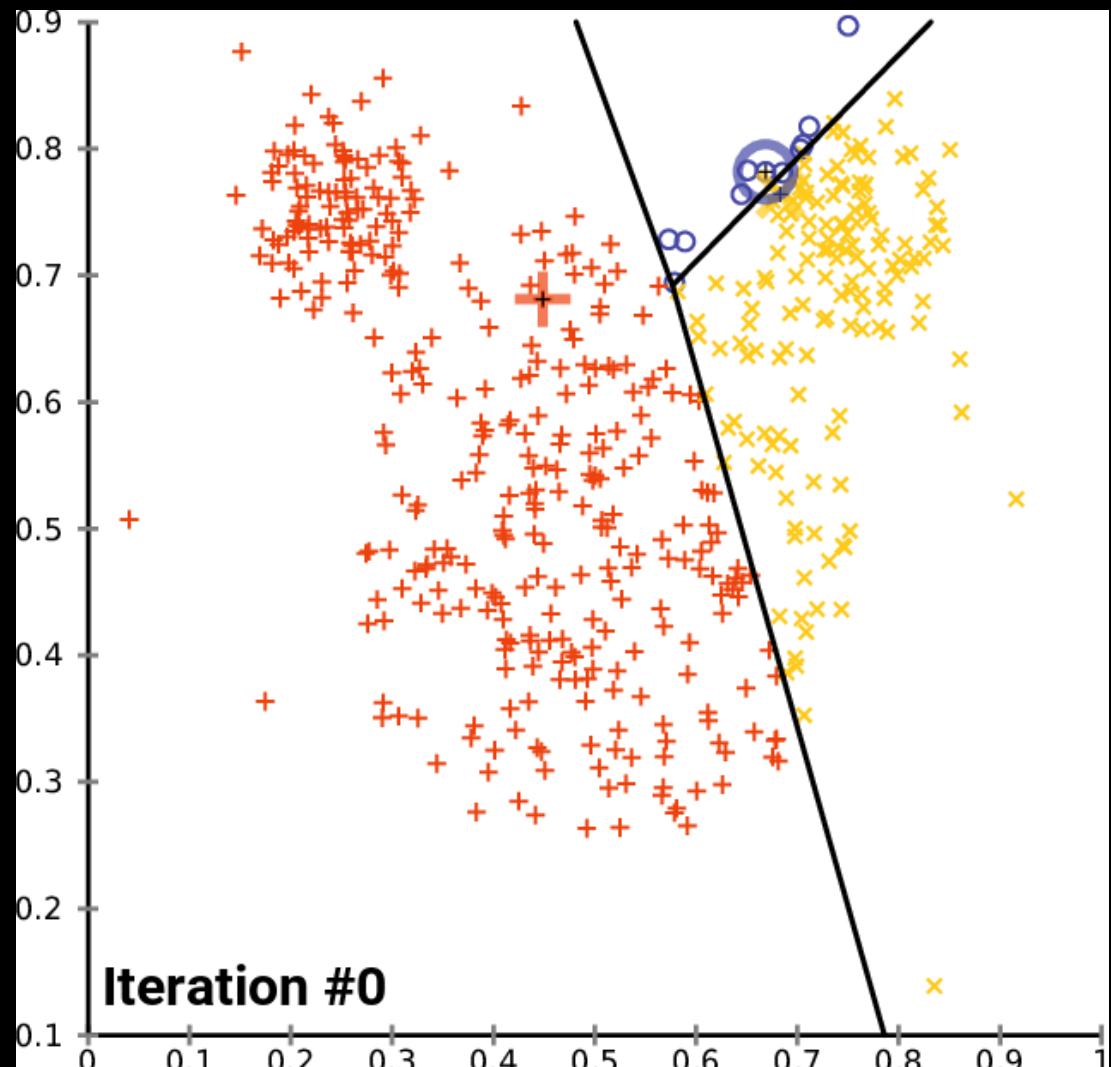
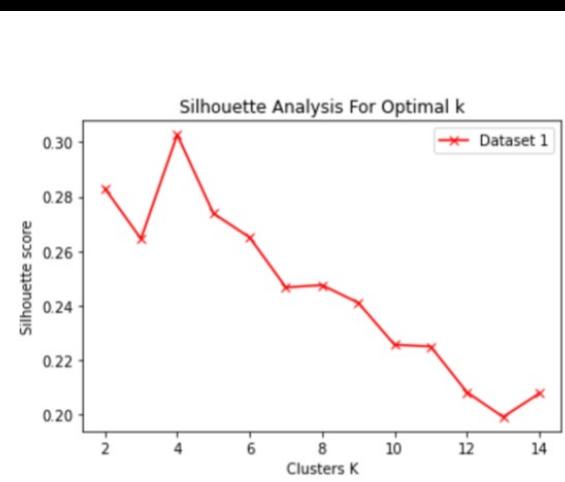
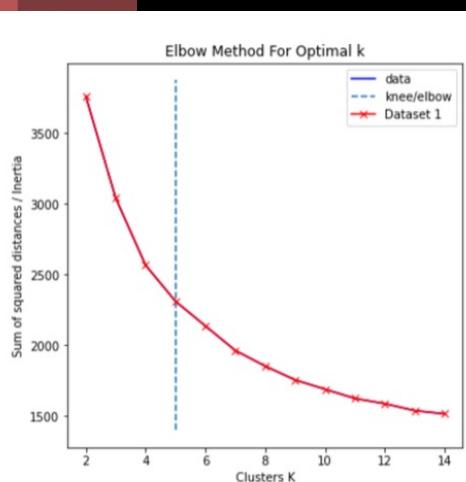
1 - It calculates the unscaled similarities between the points

2 – The points are randomly projected to a lower dimension and the previous process is repeated, but the second time the unscaled similarities are calculated on a t-distribution.

3 - Compares the two distance matrixes until the second one is similar to the first one.

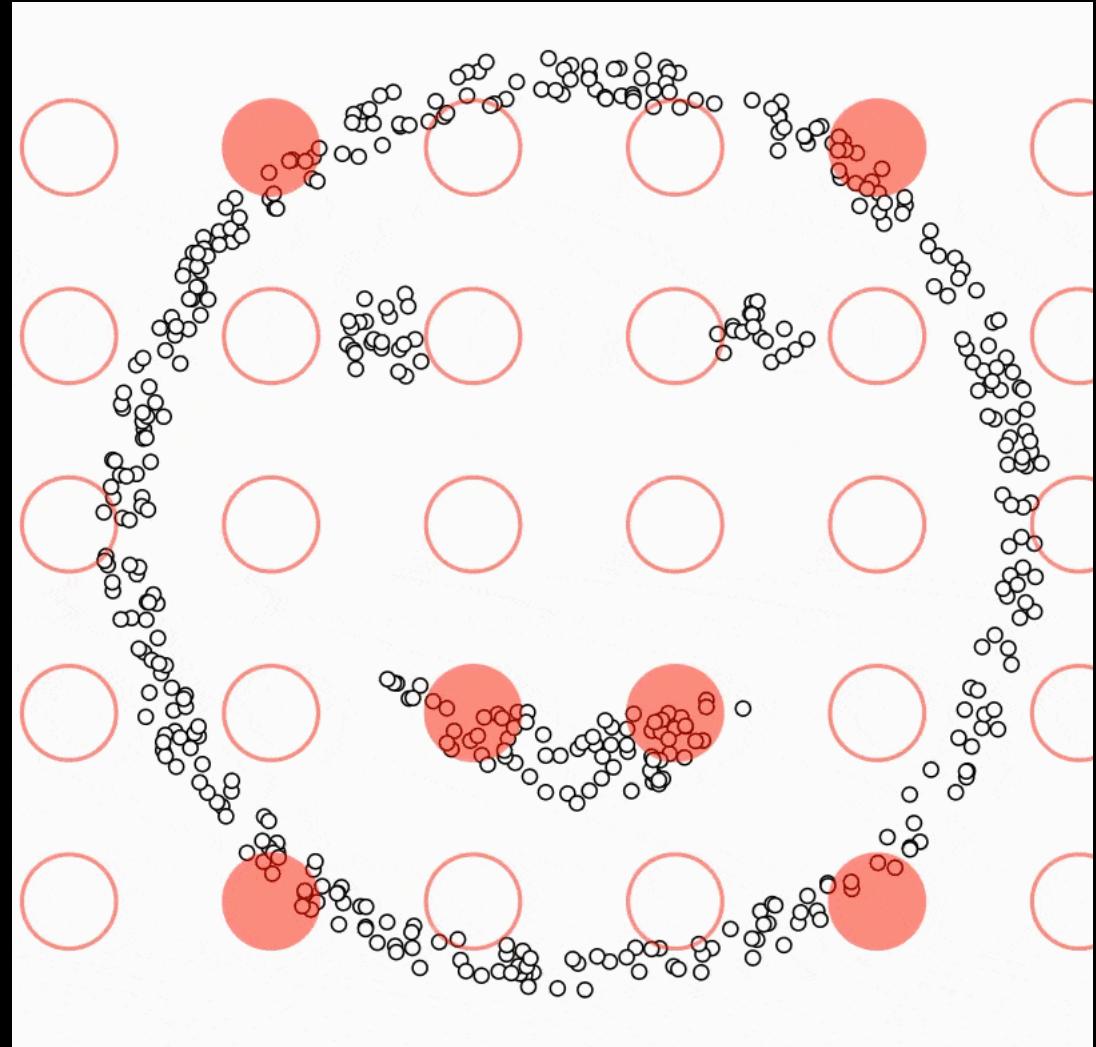
K-means

- Centroid based method
- It needs the number of clusters as a parameter
- Best results with $k=4$ and $k=5$



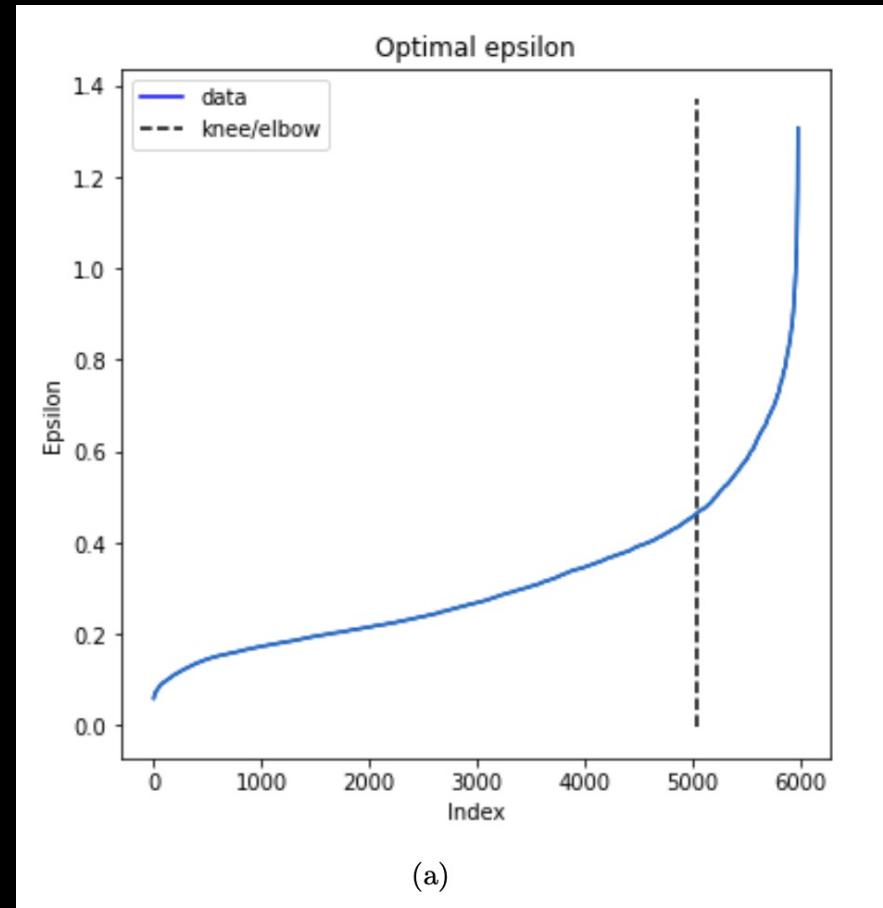
DBSCAN

- DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density based algorithm.
- **Core instances** and ϵ – neighborhoods.
- It includes anomalies.
- The minimum number of neighbour points n and the maximum distance ϵ have to be defined by the user.



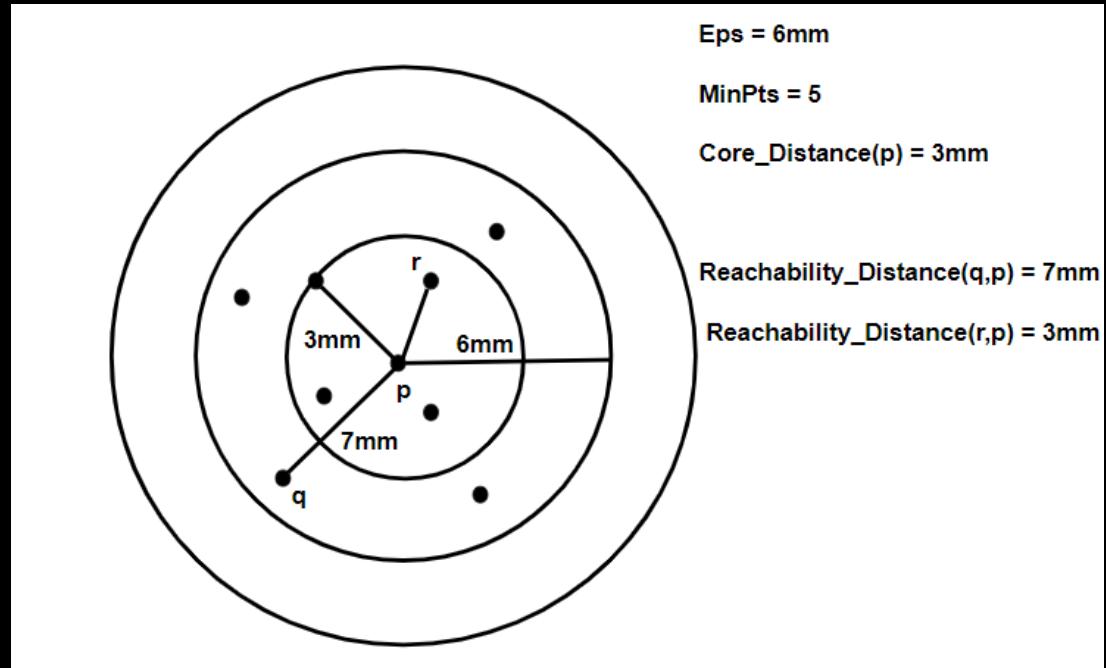
DBSCAN: Optimization

- Schubert et al. (2017) suggest $n = 2 * \text{dimensions}$, for big datasets (>30 dimensions) else $n = \text{dimensions}$. In this case $n = 15$ dimensions.
- To find the parameter ε , Rahmah and Sitanggang (2016)[19] Schubert et al. (2017) and Sander et al. (1998) suggest that it can be found by calculating the distance to the nearest n points for each point, sorting and plotting the results and then identifying the point where the change is most accentuated (similar to the elbow method) -> $\varepsilon = 0.466$
- KNeighbors classifier to assign anomalies.



OPTICS

- OPTICS (Ordering Points To Identify the Clustering Structure) is the improved version of DBSCAN.
- It introduces two new parameters: **core distance** (smallest distance ϵ to a point in its ϵ – neighborhood) and **reachability distance** (maximum value between the distance of two points and the core distance of p).
- We have to specify the same parameters than in DBSCAN as Sklearn calculates the rest. The same parameters have been used.

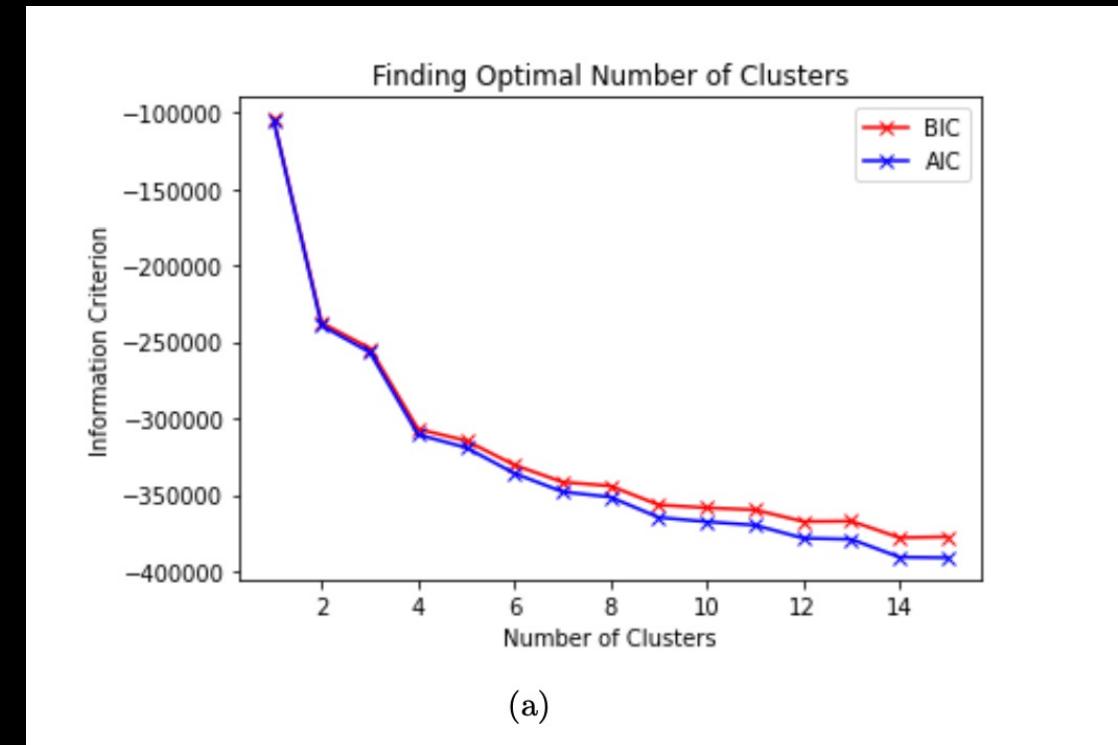


BIRCH

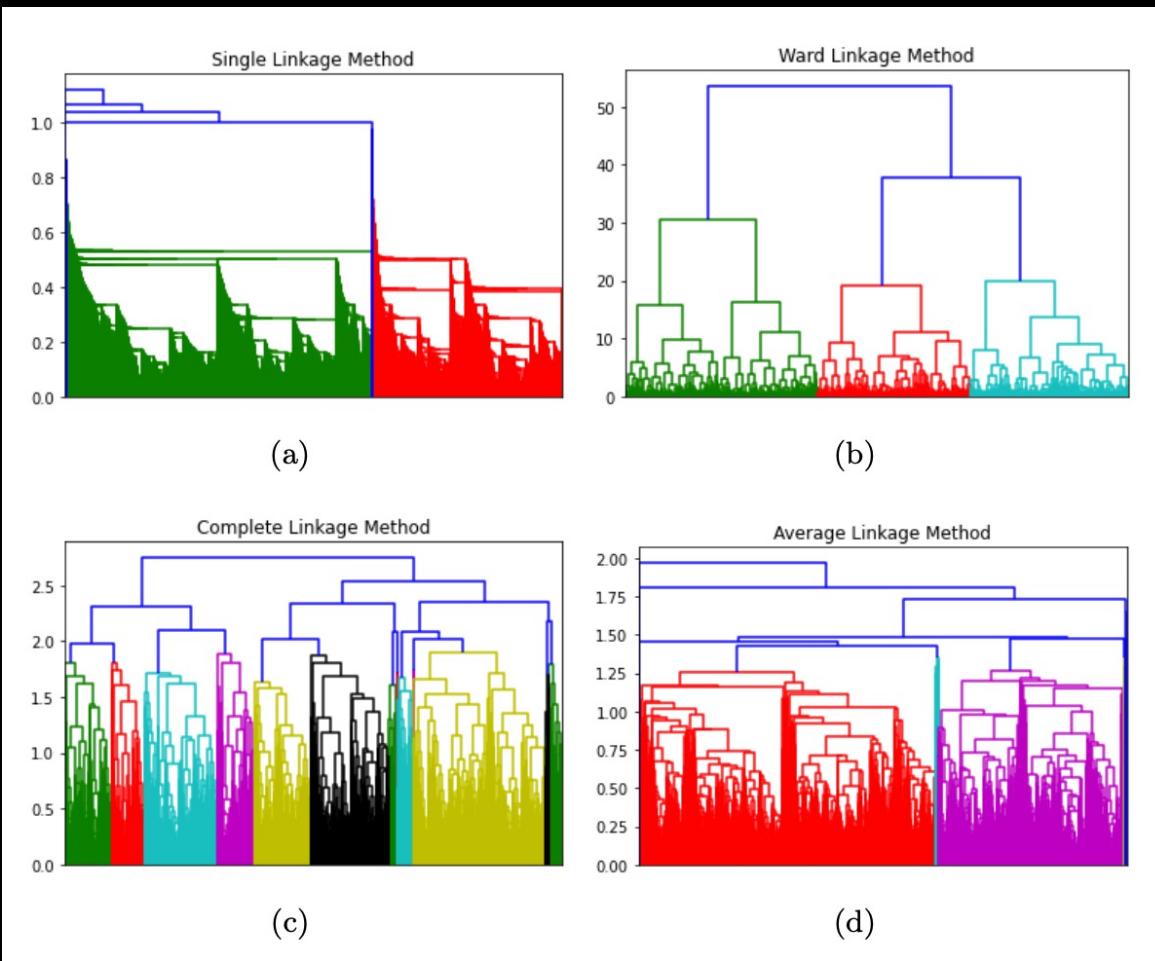
- The BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies) algorithm converts data into a tree data structure using the tree leaves as centroids.
- The final tree can be then the input for another clustering algorithm.
- The number of clusters has to be declared, in this case the minimum S_Dbw score (that is the best one) it's reached with $k = 8$.

Gaussian Mixtures

- As its name indicates, clustering by Gaussian mixtures is applied assuming that the data are the result of different Gaussian distributions.
- They need the number of clusters k .
- Find the number of clusters k by finding the k value that minimizes the value of a theoretical information criterion like the Akaike information criterion (AIC) or the Bayesian information criterion (BIC). In this case tryed $k = 12$ and $k = 14$.



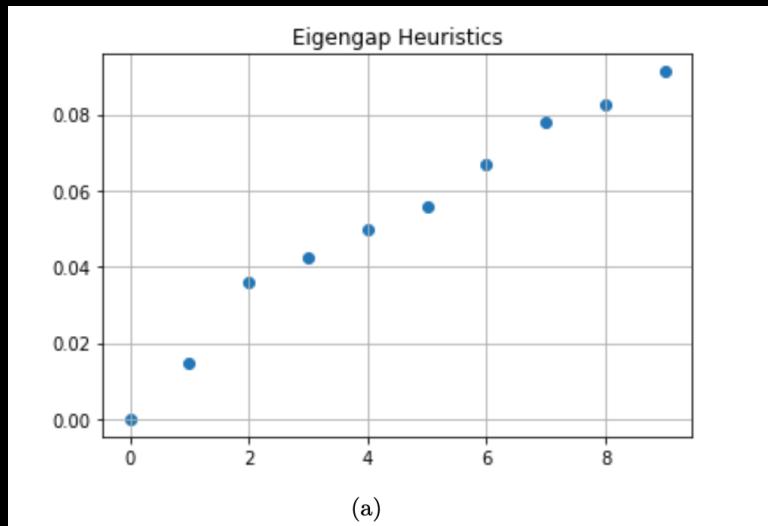
Agglomerative clustering



- Agglomerative clustering is a kind of hierarchical clustering with a **bottom-up** approach.
- There are different kinds of linkage criteria that determines the merge strategy.

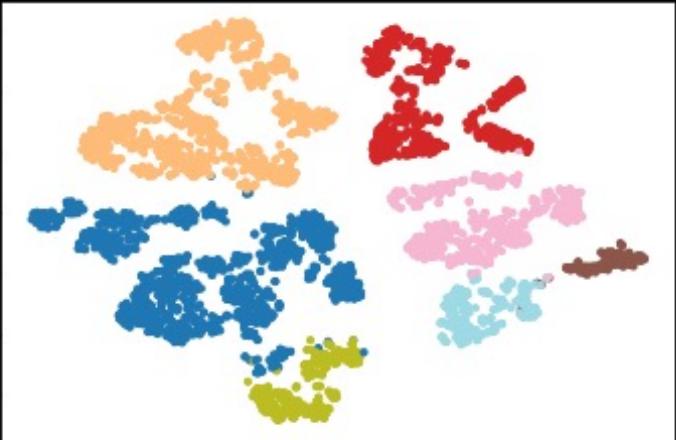
Spectral clustering

- It is very popular because it is capable of extracting very complex structures.
- It calculates the adjacency matrix of the data and then it extracts the eigenvalues. Once it has the eigenvalues, it uses them to reduce the dimensionality of the data and then it applies some other clustering algorithm, in this case k-means.
- In this case the best result is achieved applying k-nearest neighbors for calculating the adjacency matrix.
- To find the number of clusters k it has been used the Eigengap Heuristic technique (Von Luxburg (2007))
- . Finally $k = 6$ and $k = 7$ have been tested.

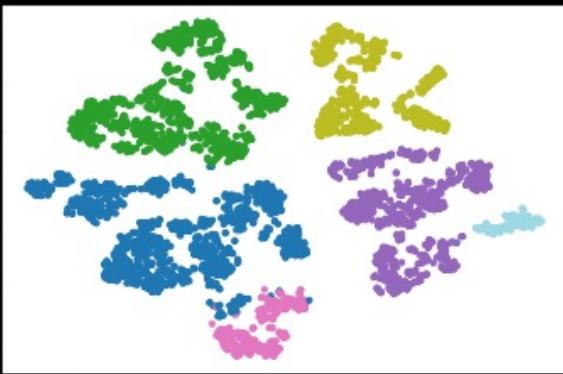


Results

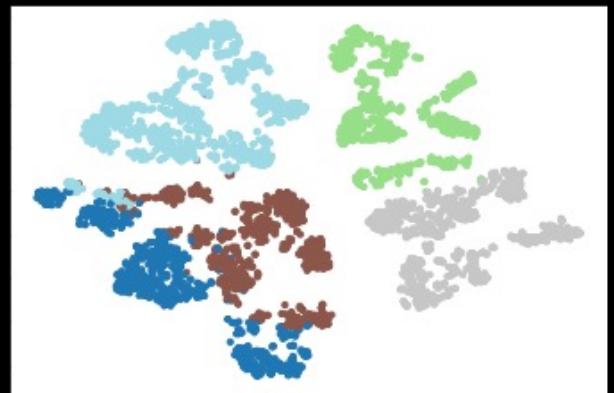
Spectral clustering k = 7



Spectral clustering k = 6



K-means k = 5



2

1

3

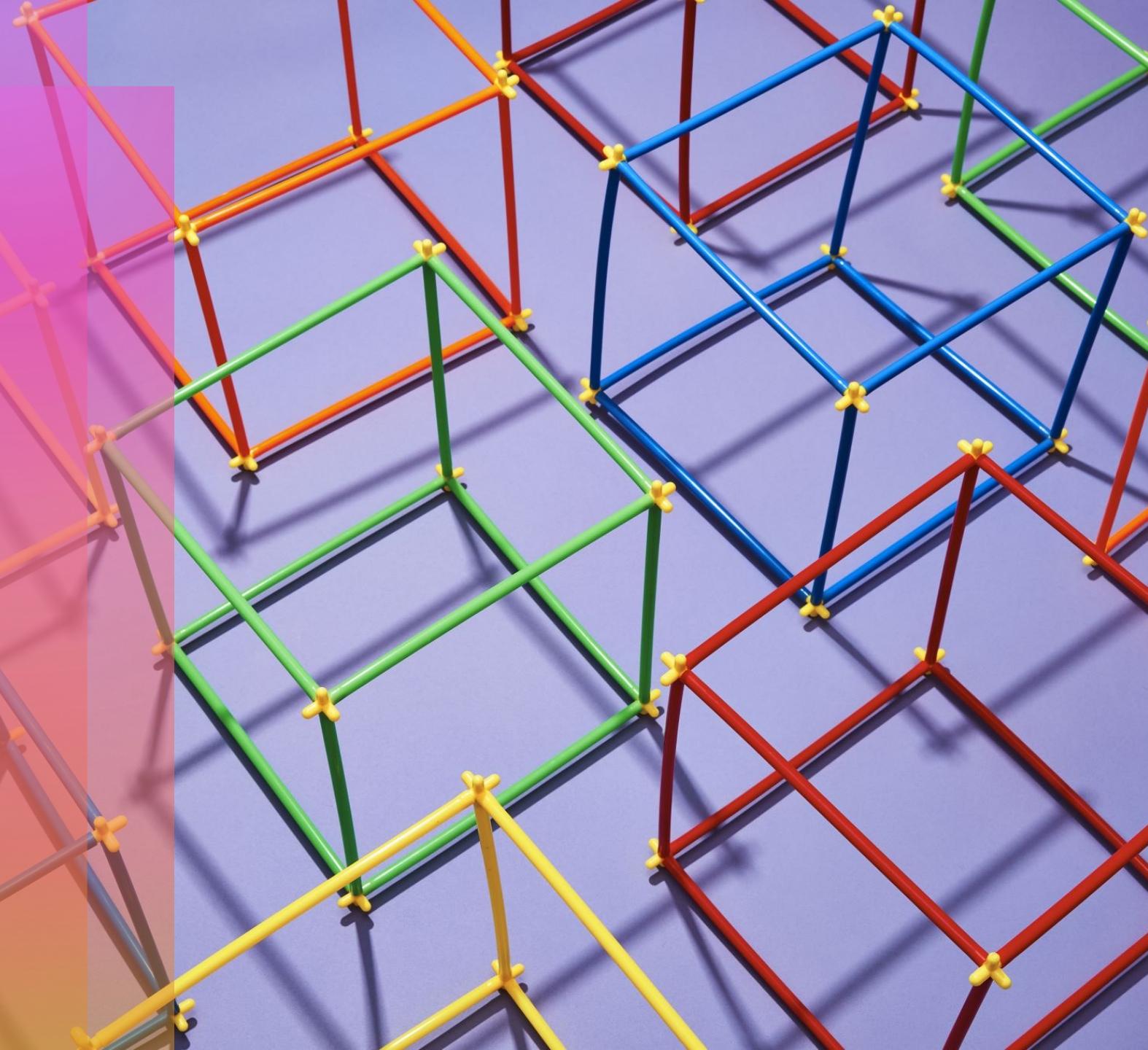
Results

Cluster	Risk level	Score
0	Medium - Low	25 + Z-score
1	Medium - High	60 + Z-score
2	Medium	45 + Z-score
3	High	80 + Z-score
4	Low	10 + Z-score
5	Medium - Low	25 + Z-score
6	Low	10 + Z-score

Asset allocation

- Strategy based on the risk level of the financial instruments and the client's risk profile.
- In this case just focused on securities.



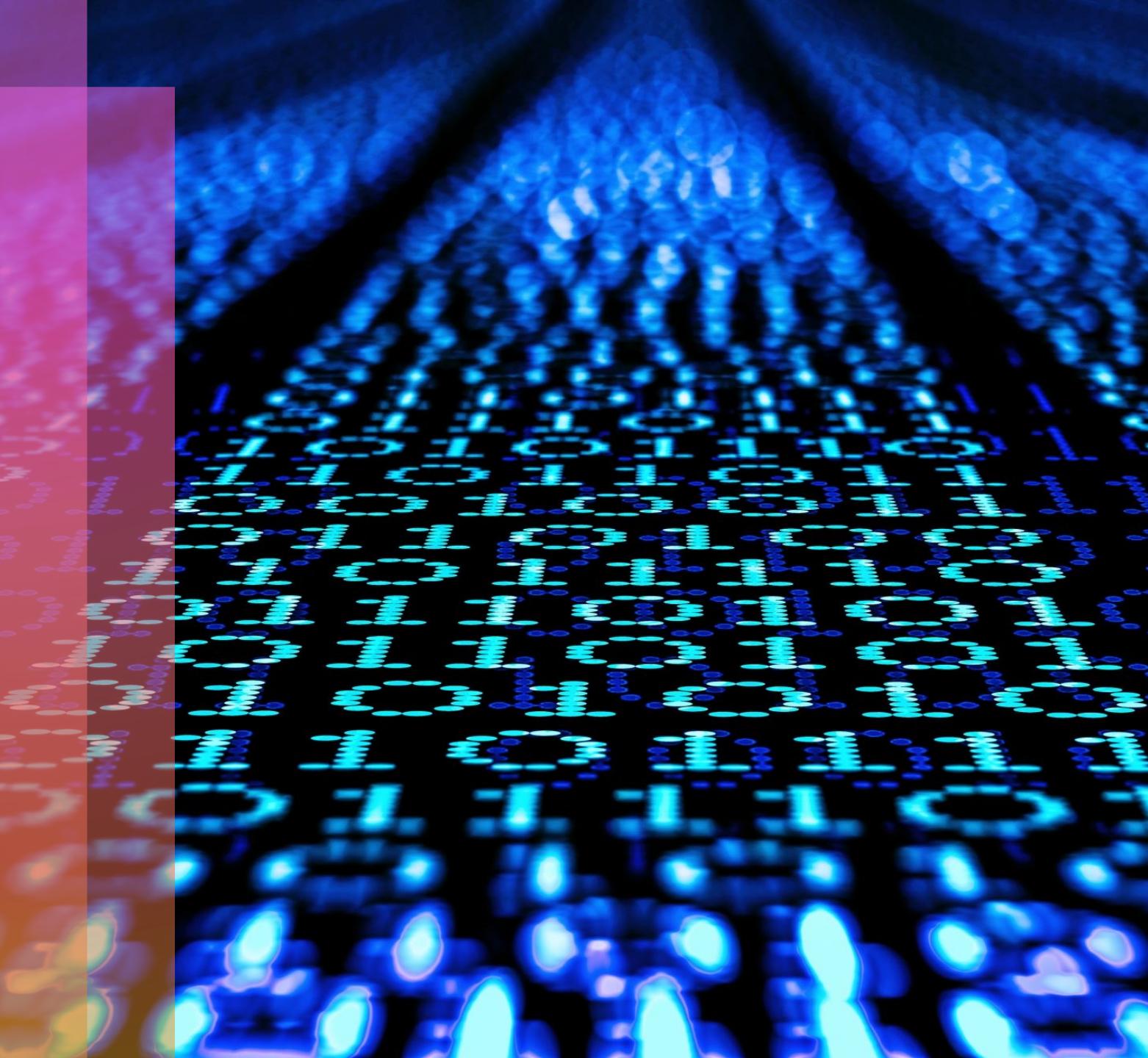


Reinforcement Learning in finance

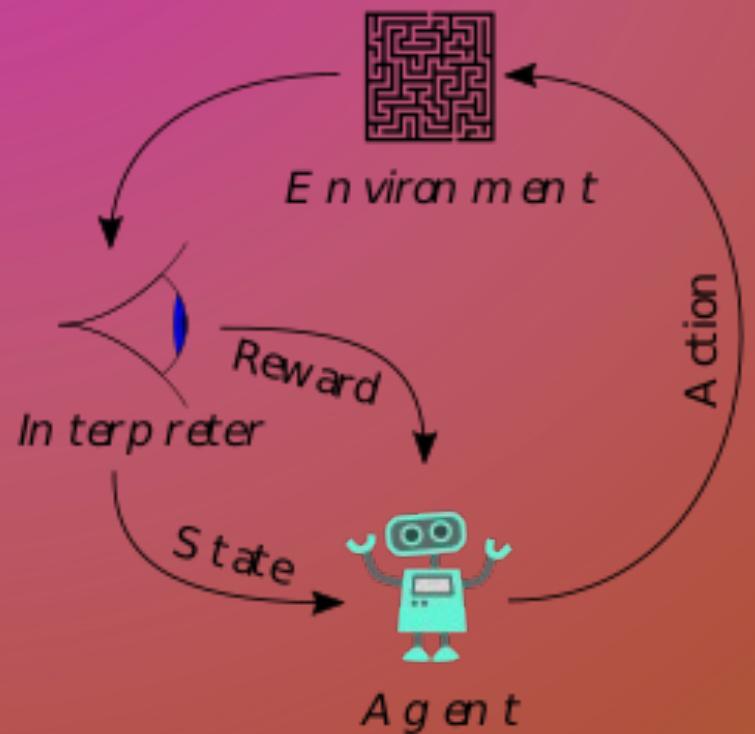
- Reinforcement Learning is a branch of machine learning in which the model learn how to act in a certain environment by trial and error, learning the strategies that result in a higher reward.
- Mostly focused on creating a model that provides a single strategy that is the most successful from a returns point of view.

Data

- 28 stocks from the S&P500 of companies that has a record in the markets of at least 10 years.
- US bond to 30 years and the Russell 2000 index.



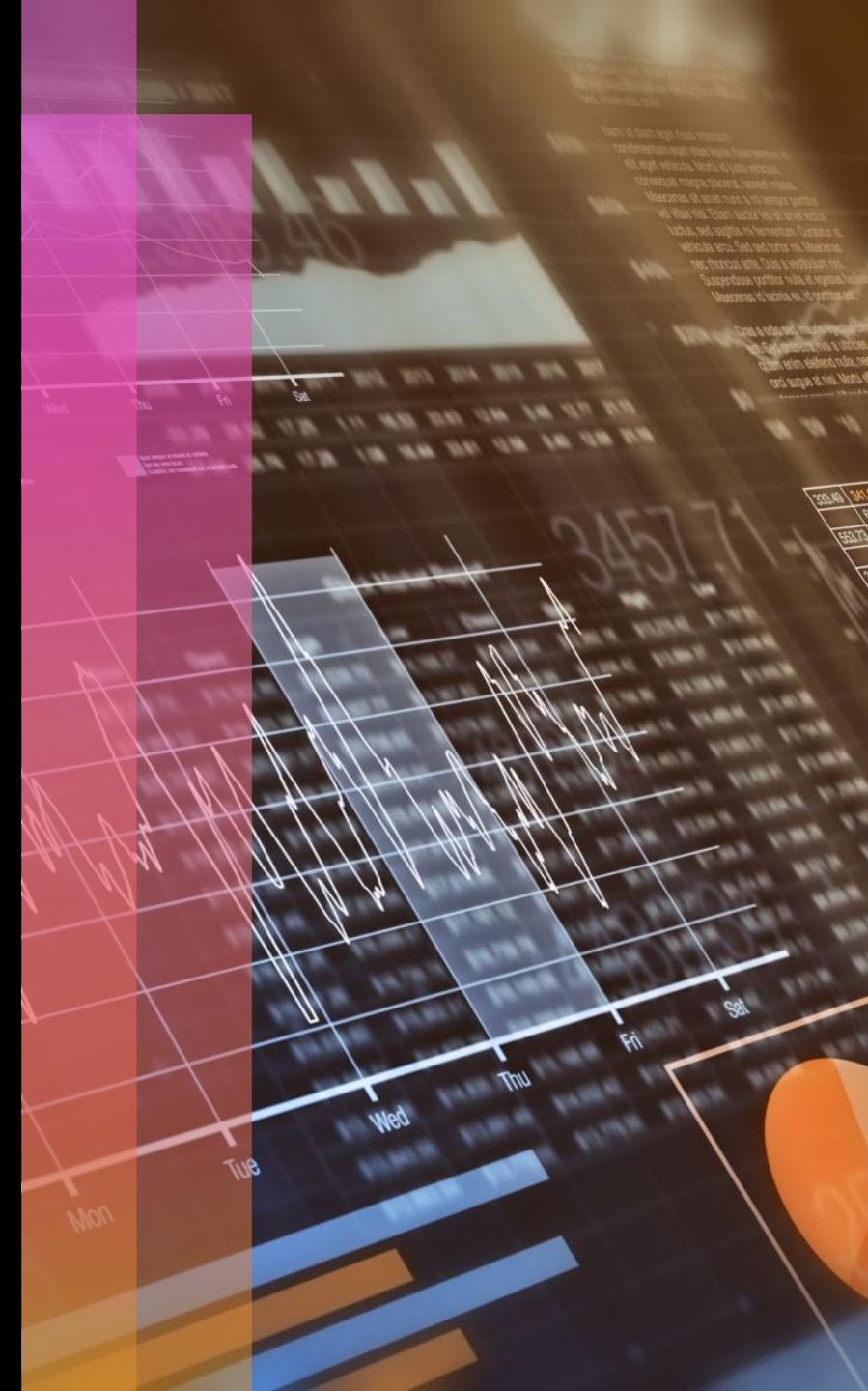
Environment



- In Reinforcement Learning, the environment is the representation of our problem and the scenario with which the agent will act, from which it will receive the data and to which it will apply certain actions.
- It has 5 core parts: initialization, step and reset functions, and state, reward and action.

Environment: The state

- The state is what the agent analyzes and uses to take an action.
- Includes the covariance matrix of the returns of the financial instruments in the last year (252 days), technical metrics ('macd', 'boll_ub', 'boll_lb', 'rsi_30', 'cci_30', 'dx_30', 'close_30_sma', 'close_60_sma') of each instrument for each day, the turbulence of the markets, the maximum level of turbulence supported by the user and the current weights of each instrument in the portfolio.
- It's initialized during the initialization process



Environment: Step

- The function receives an action as a parameter.
- It normalizes the action and then advance by n days, being n the period for rebalancing the portfolio's weights. (During all the experiment n = 1)
- Next, the scaled covariance matrix of the returns of the last year until the day in question is taken and updated in the state.
- The new updated weights are also added to the state.
- Finally the reward function is calculated

Environment: Reward function

- The reward function is the function that returns feedback to the agent to indicate how well it is doing the required task.
- It is not easy to design a good and suitable reward function, since any imprecision can result in the agent learning the wrong strategies.
- Many experiments
- How to punish losses with a threshold? Volatility? Ratios?

A vertical strip of handwritten mathematical notes and diagrams, likely from a notebook, showing various calculations and diagrams related to mathematics and physics.

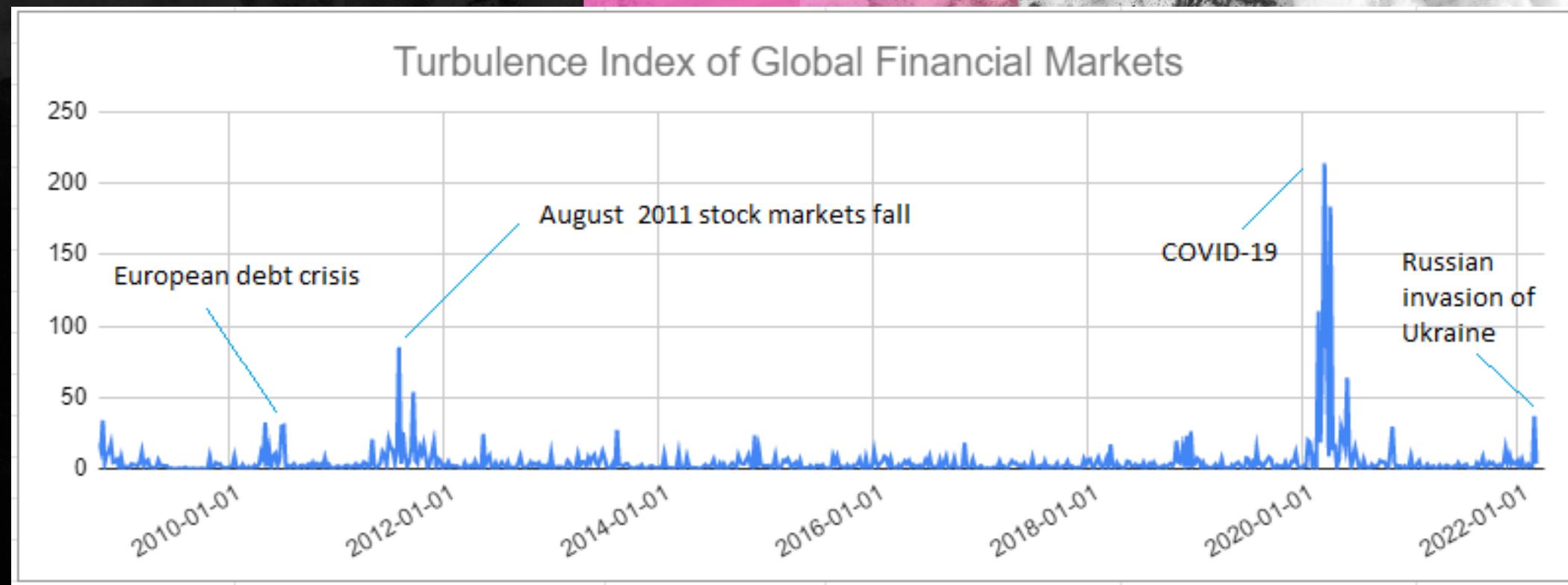
Some of the visible text and symbols include:

- At the top right, there is a diagram of a rectangle divided into four quadrants, with arrows indicating a flow or transformation process.
- Below the diagram, there is a series of equations:
 - $\sum x^2 = 9478 - 0.00$
 - $\sqrt{2456.56} - 0.00$
 - $\frac{B^2}{2} = 4.31447$
 - $\frac{a^2 + b^2}{2} = x^2$
 - $c(x, y) \begin{cases} xy = c \\ cx - cy = 0 \\ 2\pi = c \end{cases}$
 - $\frac{2x+3y}{c} + \frac{d^2 + 3^2}{c} + x^2 = 384$
 - $x = 14! + n^{3v}$
 - $\sum N^{30} - x - L$
 - $\beta = 9 + x^2$
- On the left side, there is a circular diagram with a central point labeled 'c' and several radii extending outwards.
- At the bottom, there is a diagram of a circle with a radius labeled 'r' and a central angle labeled 'a'.

Turbulence

- Financial turbulence measures the statistical unusualness of a set of returns given their historical pattern of behavior. (Chow et al. (1999), Kritzman and Li (2010))
- Coincides with the Mahalanobis distance
- In periods of great crises the turbulence values skyrocketed.
- Turbulence persists, and therefore once it has risen, it remains high for a while, which is useful for predicting moments of low profitability.

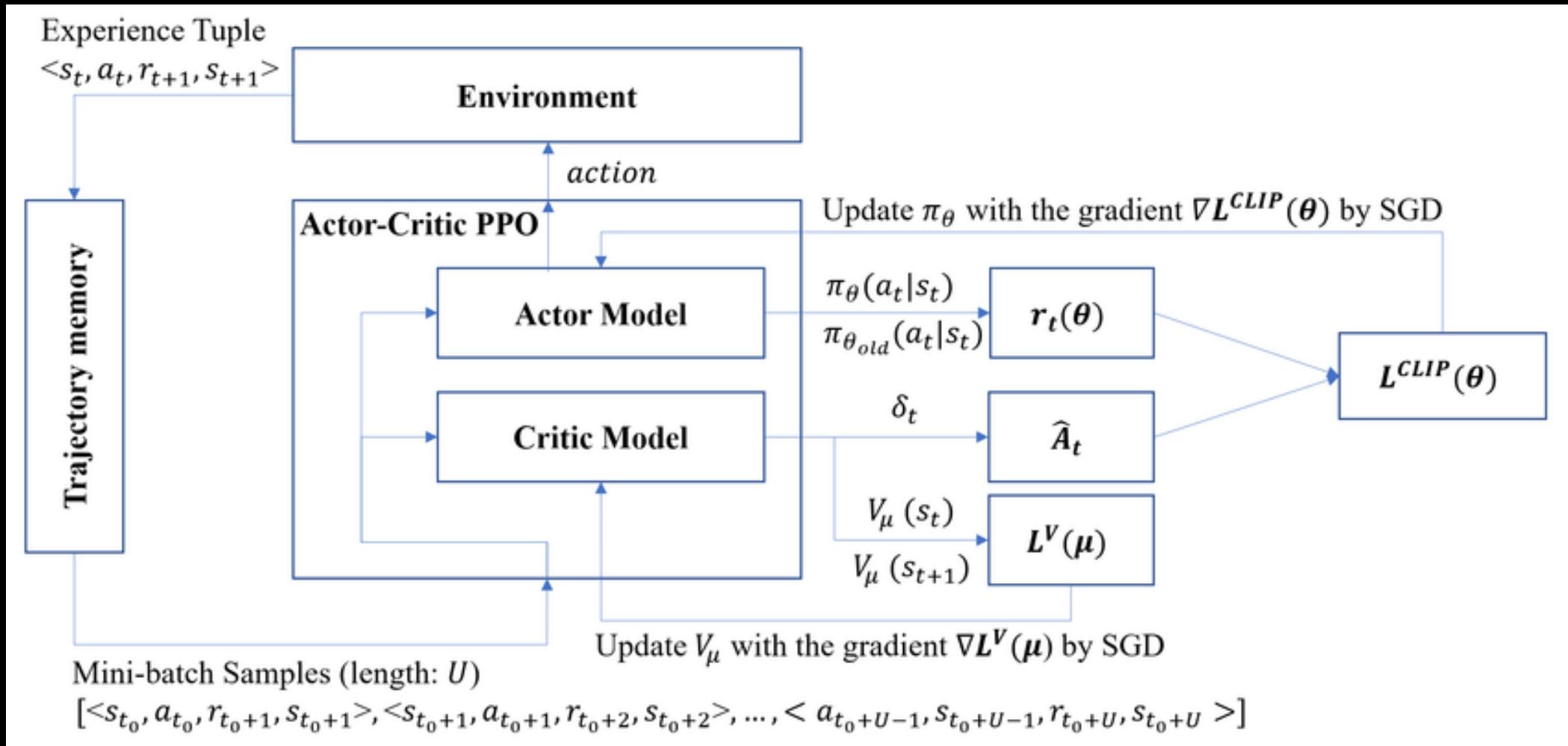
Turbulence



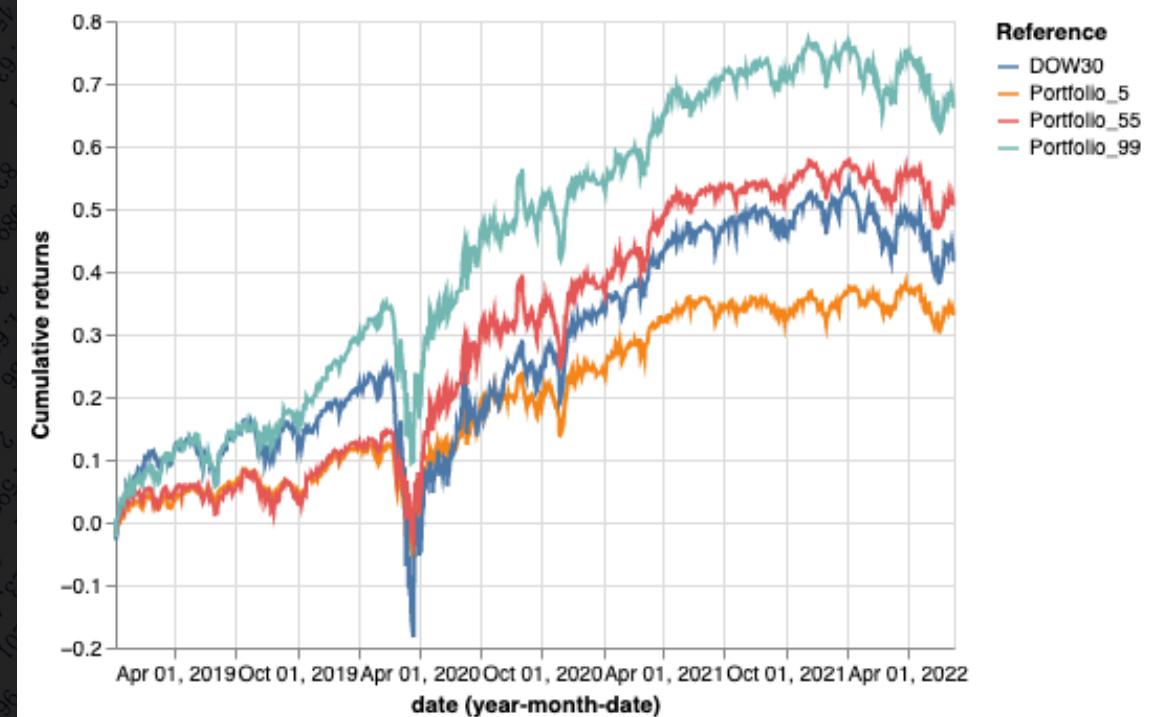
Reward function with turbulence

$$r = P - (P * |c - t| * I)$$

PPO



Results



Metrics	Score:5	Score:55	Score:99	DOW30
An. return	9%	14%	18%	9%
Cum. return	34%	57%	80%	38%
An. volatility	13%	17%	20%	23%
Sharpe	0.68	0.83	0.93	0.52
Calmar	0.51	0.82	0.78	0.26
Sortino	0.98	1.21	1.34	0.72
Stability	0.9	0.9	0.92	0.7
Max drawdown	-17%	-17%	-23%	-0.37%

Results



Metrics	Score:5	Score:55	Score:99	DOW30
An. return	-1.8%	-3%	-5.1%	-6.3%
Cum. return	-1.8%	-3%	-5.1%	-6.3%
An. volatility	9%	11%	16%	15%
Sharpe	-0.15	-0.25	-0.24	-0.33
Calmar	-0.25	-0.33	-0.36	-0.41
Sortino	-0.2	-0.33	-0.33	-0.44
Stability	0	0.04	0.15	0.19
Max drawdown	-7%	-10%	-15%	-15%

Skarb: The graphical interface

The image shows a screenshot of a web browser window titled "localhost". The left side of the screen displays a survey form with various input fields and dropdown menus. The right side features a welcoming message and a rocket launch illustration.

Welcome to Skarb!

Start investing in a smarter way

Take the survey to start

Tell us something about you

How old are you?
18

What is your marital status?
Single

What is the highest educational level you have reached?
Never went to school

Lets talk about money 💰

What is your total monthly income?
0

What are the total monthly expenses of your household?
0

In the future your expenses will

LAUNCH

A cartoon character in a purple shirt and tie holds a circular button labeled "LAUNCH". To the right is a purple rocket ship launching from a base of purple smoke, with green leaves floating around it.

Conclusions



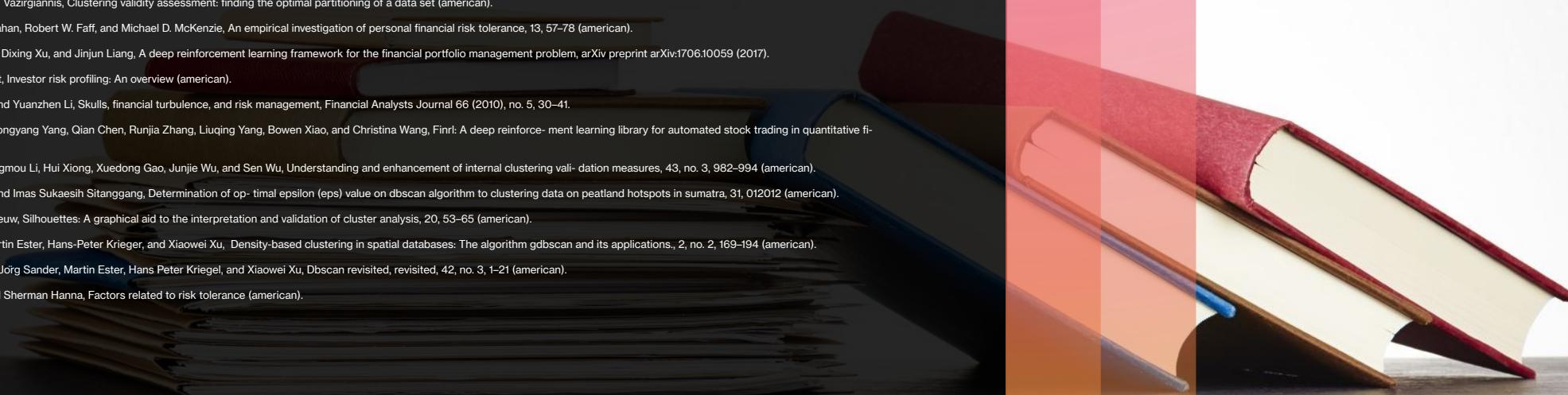
Github repository

- github.com/simoncraf/skarb



Bibliography

- [1] sklearn.cluster.optics.
- [2] Mihael Ankerst, Markus M. Breunig, Hans-Peter Kriegel, and Jörg Sander, Optics, 28, no. 2, 49–60 (american).
- [3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba, Openai gym, arXiv preprint arXiv:1606.01540 (2016).
- [4] T. Calinski and J. Harabasz, A dendrite method for cluster analysis, 3, no. 1, 1–27 (american).
- [5] Ricardo J. G. B. Campello, Peer Kröger, Jörg Sander, and Arthur Zimek, Density-based clustering, 10, no. 2 (american).
- [6] George Chow, Eric Jacquier, Mark Kritzman, and Kenneth Lowry, Optimal portfolios in good times and bad, Financial Analysts Journal 55 (1999), no. 3, 65–73.
- [7] D. Comaniciu and P. Meer, Mean shift: a robust approach toward feature space analysis, 24, no. 5, 603–619 (american).
- [8] David L. Davies and Donald W. Bouldin, A cluster separation measure, PAMI, no. 2, 224–227 (american).
- [9] John E. Grable, Financial risk tolerance: A psychometric review (american).
- [10] Mao Guan and Xiao-Yang Liu, Explainable deep reinforcement learning for portfolio management: An empirical approach (american).
- [11] Aurélien Géron, Hands-on machine learning with scikit-learn, keras, and tensorflow: Concepts, tools, and techniques to build intelligent systems, 2 ed., O'Reilly Media (american).
- [12] M. Halkidi and M. Vazirgiannis, Clustering validity assessment: finding the optimal partitioning of a data set (american).
- [13] Terrence A. Hallahan, Robert W. Faff, and Michael D. McKenzie, An empirical investigation of personal financial risk tolerance, 13, 57–78 (american).
- [14] Zhengyao Jiang, Dixin Xu, and Jinjun Liang, A deep reinforcement learning framework for the financial portfolio management problem, arXiv preprint arXiv:1706.10059 (2017).
- [15] Joachim Klement, Investor risk profiling: An overview (american).
- [16] Mark Kritzman and Yuanzhen Li, Skulls, financial turbulence, and risk management, Financial Analysts Journal 66 (2010), no. 5, 30–41.
- [17] Xiao-Yang Liu, Hongyang Yang, Qian Chen, Runjia Zhang, Liqing Yang, Bowen Xiao, and Christina Wang, Finrl: A deep reinforcement learning library for automated stock trading in quantitative finance (american).
- [18] Yanchi Liu, Zhongmou Li, Hui Xiong, Xuedong Gao, Junjie Wu, and Sen Wu, Understanding and enhancement of internal clustering validation measures, 43, no. 3, 982–994 (american).
- [19] Nadia Rahmah and Imas Sukaesih Sitanggang, Determination of optimal epsilon (eps) value on dbscan algorithm to clustering data on peatland hotspots in sumatra, 31, 012012 (american).
- [20] Peter J. Rousseeuw, Silhouettes: A graphical aid to the interpretation and validation of cluster analysis, 20, 53–65 (american).
- [21] Jörg Sander, Martin Ester, Hans-Peter Kriegel, and Xiaowei Xu, Density-based clustering in spatial databases: The algorithm gdbcscan and its applications, 2, no. 2, 169–194 (american).
- [22] Erich Schubert, Jörg Sander, Martin Ester, Hans Peter Kriegel, and Xiaowei Xu, Dbscan revisited, revisited, 42, no. 3, 1–21 (american).
- [23] Jaimie Sung and Sherman Hanna, Factors related to risk tolerance (american).



Thank you

Simone Gigante Gassó
simone.gigante01@estudiant.upf.edu
Universitat Pompeu Fabra
05 – 06 - 2022