# Addressing the Rashomon Effect through ranking aggregation

**Claudia Sessa**[1], **PhD Candidate**

`csessa@liuc.it`

Joint work with X. Lu [2], E. Borgonovo [3], A. Cillo [1] and G.P. Crespi [1]

[1] LIUC University, 21053 Castellanza, Italy.
[2] SKEMA Business School-Université Côte d'Azur, 92150 Suresnes, France.
[3] Department of Decision Sciences, Bocconi University, Milan, Italy, MI 20136.

April 22, 2025

# Table of contents

Given a dataset and a prediction task, there may exist several prediction models that exhibit equally good performance (Breiman 2001b):

⟶ the **Rashomon Set** (Fisher, Rudin, and Dominici 2019)

⊕ Analysts can choose from a range of acceptable models to solve their prediction problem.

⊖ Models in the Rashomon Set may produce **different explanations** for the phenomenon at hand, because they are **built on alternative rationales**.

# Prediction Models and the Rashomon Set

- Let $\mathbf{X} = (X_1, \ldots, X_p) \in \mathcal{X} \subseteq \mathbb{R}^p$, and $Y \in \mathcal{Y} \subseteq \mathbb{R}$ be random variables on the probability space $(\Omega, \mathcal{B}(\Omega), \mathbb{P})$.

- The data generating process (DGP) is given by $g^* : \mathcal{X} \to \mathcal{Y}$.

- A dataset $\mathcal{D} = \{\mathsf{X}, \mathsf{Y}\} = \{(\mathbf{x}^i, y^i)\}_{i=1}^n$ consists of $n$ realizations of the DGP.

- Given a set of parameters $\Theta$, a **prediction model** is a function $m : \mathcal{X} \times \Theta \to \mathbb{R}$ that can be used to approximate $g^*$ and make predictions given observations of $\mathbf{X}$.

- Let $\mathcal{L}(m)$ denote the loss function used to assess the out-of-sample performance of a model $m$.

# Prediction Models and the Rashomon Set

- Let $\mathbf{X} = (X_1, \ldots, X_p) \in \mathcal{X} \subseteq \mathbb{R}^p$, and $Y \in \mathcal{Y} \subseteq \mathbb{R}$ be random variables on the probability space $(\Omega, \mathcal{B}(\Omega), \mathbb{P})$.

- The data generating process (DGP) is given by $g^* : \mathcal{X} \to \mathcal{Y}$.

- A dataset $\mathcal{D} = \{\mathsf{X}, \mathsf{Y}\} = \{(\mathbf{x}^i, y^i)\}_{i=1}^n$ consists of $n$ realizations of the DGP.

- Given a set of parameters $\Theta$, a **prediction model** is a function $m : \mathcal{X} \times \Theta \to \mathbb{R}$ that can be used to approximate $g^*$ and make predictions given observations of $\mathbf{X}$.

- Let $\mathcal{L}(m)$ denote the loss function used to assess the out-of-sample performance of a model $m$.

## Population $\epsilon$-Rashomon Set (Fisher, Rudin, and Dominici 2019)

For a given class of machine learning models $\mathcal{F}$, the **population $\epsilon$-Rashomon Set**, for a fixed parameter $\epsilon$, is defined as the subset of models $m \in \mathcal{F}$ whose expected loss deviates at most $\epsilon$ from the minimum loss attained, that is:

$$\mathcal{R}(\mathcal{F}, \epsilon) := \left\{ m \in \mathcal{F} : \mathbb{E}[\mathcal{L}(m)] \leq \min_{m' \in \mathcal{F}} \mathbb{E}[\mathcal{L}(m')] + \epsilon \right\}$$

## Variable Importance Measure

Let $\mathcal{S} = \{1, 2, \ldots, p\}$ be the set of indices of the feature vector $\mathbf{X} = (X_1, \ldots, X_p)$.

Then, an **importance measure** is a map $\Phi : \mathcal{S} \to \mathbb{R}_+^p$.

We call the image of $\Phi$ importance vector.

Two main goals of variable importance measures (Williamson et al. 2022):

$\rightarrow$ understanding the features on which a given algorithm relies,

$\rightarrow$ measuring the population-level predictive potential of features.

## Variable Importance Measure

Let $\mathcal{S} = \{1, 2, \ldots, p\}$ be the set of indices of the feature vector $\mathbf{X} = (X_1, \ldots, X_p)$.

Then, an **importance measure** is a map $\Phi : \mathcal{S} \to \mathbb{R}_+^p$.

We call the image of $\Phi$ importance vector.

Two main goals of variable importance measures (Williamson et al. 2022):

$\to$ understanding the features on which a given algorithm relies,

$\to$ measuring the population-level predictive potential of features.

Several indicators fall under this definition, for instance:

- **Permutation Importance** (Breiman 2001a):

$$\nu_j^{(m)} = \mathbb{E}[\mathcal{L}(Y, m(\mathbf{X}_\mathbf{j}'))] - \mathbb{E}[\mathcal{L}(Y, m(\mathbf{X}))]$$

where $X_j'$ represents an independent copy of $X_j$, and $\mathbf{X}_\mathbf{j}' = (X_1, \ldots, X_{j-1}, X_j', X_{j+1}, X_p)$ is the random vector where the feature $X_j$ is substituted with an independent copy.

- **Global sensitivity measures** (Borgonovo, Hazen, and Plischke 2016):

$$\xi^\zeta(Y, X_j) := \mathbb{E}_{X_j}[\zeta(\mathbb{P}_Y, \mathbb{P}_{Y|X_j})]$$

where $\mathbb{P}_Y$ and $\mathbb{P}_{Y|X_j}$ denote the marginal distribution of $Y$ and the conditional distribution of $Y$ given $X_j$, respectively, and $\zeta(\cdot, \cdot)$ is a separation measurement quantifying the discrepancy between two probability measures on the reference probability space.

Few Variable Importance approaches have been proposed that use the entire Rashomon Set:

- **Model Class Reliance** (Fisher, Rudin, and Dominici (2019))

- **Variable Importance Cloud** (Dong and Rudin (2020))

- **Method Agnostic Model Class Reliance Range** (Gunasekaran, Mistry, and Chen (2024))

- **Rashomon Importance Distribution** (Donnelly et al. (2023))

.

**1**

Clarify the connection between the coefficients of linear models in the Rashomon Set and the Permutation importance measure, exploring its relation to total indices.

**2**

Propose an actionable methodological framework that exploits the whole Rashomon Set to obtain a unified explanation for the data.

**1**

**Clarify the connection between the coefficients of linear models in the Rashomon Set and the Permutation importance measure, exploring its relation to total indices.**

**2**

Propose an actionable methodological framework that exploits the whole Rashomon Set to obtain a unified explanation for the data.

- Let $\mathcal{D} = \{X, Y\}$ be a dataset. We call $m^* \in \mathcal{F}_l$ the best linear regression model trained on this dataset, characterized by the *p*-dimensional vector of coefficients $\beta^*$.

- Let $\mathcal{L}^*$ be the estimated mean squared loss from this model, that is $\mathcal{L}^* = \mathcal{L}(m^*) = \frac{1}{n} \sum_{i=1}^{n} (y^i - \mathbf{x}^i \beta^*)^2$.

- Let $\mathcal{D} = \{X, Y\}$ be a dataset. We call $m^* \in \mathcal{F}_l$ the best linear regression model trained on this dataset, characterized by the $p$-dimensional vector of coefficients $\beta^*$.

- Let $\mathcal{L}^*$ be the estimated mean squared loss from this model, that is $\mathcal{L}^* = \mathcal{L}(m^*) = \frac{1}{n} \sum_{i=1}^{n} (y^i - \mathbf{x}^i \beta^*)^2$.

- The **Rashomon Set of Linear Regression models** (Semenova, Rudin, and Parr 2022) is the set of all models $m \in \mathcal{F}_l$ characterized by a vector of coefficients $\beta$ that satisfies the relation:

$$\mathbb{E}[(Y - X\beta)^2] \leq (1 + \epsilon)\mathcal{L}^*.$$

- When the features are <u>uncorrelated</u>, the $\beta$'s characterizing all models in the Rashomon Set form an ellipsoid in $\mathbb{R}^p$ centered at $\beta^*$ (Dong and Rudin 2020), satisfying the following inequality

$$(\beta - \beta^*) \frac{X^\top X}{\epsilon \mathcal{L}^*} (\beta - \beta^*) \leq 1.$$

# New Analytical Insights for Linear Models

- Let $\beta_j$ b the coefficient associated with $X_j$, $\sigma_j$ the standard deviation of $X_j$ and $\sigma_Y$ the standard deviation of $Y$.

- Let $\rho_j = \frac{Cov(Y, X_j)}{\sigma_Y \sigma_j}$ and $SRC_j = \frac{\beta_j \sigma_j}{\sigma_Y}$ be the **Pearson linear regression coefficient** and the **Standardized Regression Coefficient** between $Y$ and $X_j$.

- Let $\tau_j = \sigma_Y^2 - \mathbb{E}[Var[Y|X_j]] = \frac{1}{2}\mathbb{E}\left[\left(g^*(\mathbf{X_j'}) - g^*(\mathbf{X})\right)^2\right]$ denote the **Total Index** of $X_j$ (Homma and Saltelli 1996).

- Let $\nu_j^{(m)} = \mathbb{E}\left[\left(Y - m(\mathbf{X_j'}; \beta)\right)^2\right] - \mathbb{E}\left[(Y - m(\mathbf{X}; \beta))^2\right]$ denote the the **Permutation Importance** of feature $X_j$ for model $m$ under a quadratic loss.

## New Analytical Insights for Linear Models

- Let $\beta_j$ b the coefficient associated with $X_j$, $\sigma_j$ the standard deviation of $X_j$ and $\sigma_Y$ the standard deviation of $Y$.

- Let $\rho_j = \frac{Cov(Y, X_j)}{\sigma_Y \sigma_j}$ and $SRC_j = \frac{\beta_j \sigma_j}{\sigma_Y}$ be the **Pearson linear regression coefficient** and the **Standardized Regression Coefficient** between $Y$ and $X_j$.

- Let $\tau_j = \sigma_Y^2 - \mathbb{E}[Var[Y|X_j]] = \frac{1}{2}\mathbb{E}\left[\left(g^*(\mathbf{X'_j}) - g^*(\mathbf{X})\right)^2\right]$ denote the **Total Index** of $X_j$ (Homma and Saltelli 1996).

- Let $\nu_j^{(m)} = \mathbb{E}\left[\left(Y - m(\mathbf{X'_j}; \beta)\right)^2\right] - \mathbb{E}\left[(Y - m(\mathbf{X}; \beta))^2\right]$ denote the the **Permutation Importance** of feature $X_j$ for model $m$ under a quadratic loss.

### Proposition 1

Let $m$ be a linear model with coefficients $\beta$. If $m$ is a perfect predictor, and features are uncorrelated, then for a feature $X_j$

$$SRC_j^2 = \rho_j^2 = \frac{\tau_j}{\sigma_Y^2} \implies \nu_j^{(m)} = 2\tau_j = 2(\beta_j)^2 \sigma_j^2$$

## Proposition 1

Let $m$ be a linear model with coefficients $\beta$. If $m$ is a perfect predictor, and features are uncorrelated, then for a feature $X_j$

$$SRC_j^2 = \rho_j^2 = \frac{\tau_j}{\sigma_Y^2} \quad \implies \quad \nu_j^{(m)} = 2\tau_j = 2(\beta_j)^2 \sigma_j^2$$

$\longrightarrow$ If we fix $\epsilon$ so that models in the Rashomon Set exhibit a negligible squared error, then we expect $\nu_j^{(m)} \approx 2(\beta_j)^2 \sigma_j^2$.

$\longrightarrow$ As the model performance worsens, the approximation does not hold and the Permutation importance $\nu_j^{(m)}$ may significantly deviate.

## Proposition 2

Let $\mathcal{R}(\mathcal{F}_l, \epsilon)$ be the Rashomon Set of linear models defined by the parameter $\epsilon$, with $\epsilon$ small enough so that the squared error of the models within the Rashomon Set is negligible.

Then, under the assumption of uncorrelated features, the coefficients associated to a given feature $X_j$ by models in $\mathcal{R}(\epsilon)$ lie within the range $\left[\beta_j^-, \beta_j^+\right]$ where

$$\beta_j^- = \beta_j^* - \delta \quad \text{and} \quad \beta_j^+ = \beta_j^* + \delta, \quad \text{with} \quad \delta = \sqrt{\frac{\epsilon L^*}{e_j^\top (X^\top X)^{-1} e_j}} ((X^\top X)^{-1} e_j)_j,$$

$\beta_j^*$ representing the coefficient of the best model, and $e_j$ the unit vector in the direction $j$.

Moreover, the Permutation Importance $\nu_j$ of $X_j$, lies within the range $[\nu_j^-(\epsilon), \nu_j^+(\epsilon)]$, where

$$\nu_j^-(\epsilon) = 2(\max\left\{0, sign\{\beta_j^- \cdot \beta_j^+\} \cdot \min\{|\beta_j^-|, |\beta_j^+|\}\right\})^2 \sigma_j^2 \quad,$$

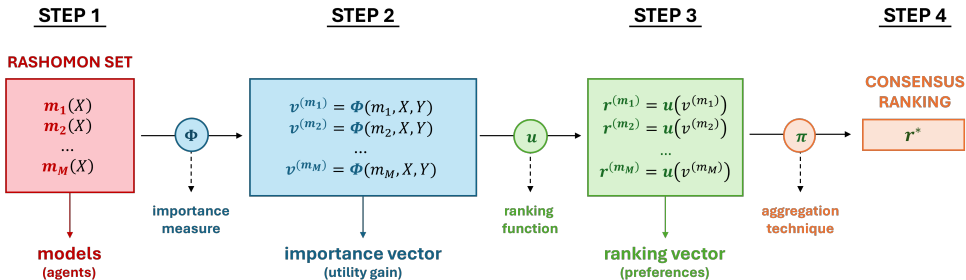$$\nu_j^+(\epsilon) = 2(\max\{|\beta_j^-|, |\beta_j^+|\})^2 \sigma_j^2 \quad.$$

**1**

Clarify the connection between the coefficients of linear models in the Rashomon Set and the Permutation importance measure, exploring its relation to total indices.

**2**

**Propose an actionable methodological framework that exploits the whole Rashomon Set to obtain a unified explanation for the data.**

The proposed framework consists of 4 steps:

## Computing the Rashomon Set (in practice)

Relatively few methods exist to compute or approximate a Rashomon Set for specific model classes:

- linear models (Dong and Rudin 2020)
- sparse decision trees (Xin et al. 2022)
- generalized additive models (Zhong et al. 2024)
- rule lists (Mata, Kanamori, and Arimura 2022)

## Computing the Rashomon Set (in practice)

Relatively few methods exist to compute or approximate a Rashomon Set for specific model classes:

- linear models (Dong and Rudin 2020)
- sparse decision trees (Xin et al. 2022)
- generalized additive models (Zhong et al. 2024)
- rule lists (Mata, Kanamori, and Arimura 2022)

### Empirical Rashomon Set

Let $\mathcal{D}_{train}$ and $\mathcal{D}_{test}$ be the training dataset and the test dataset, respectively, obtained by partitioning the dataset $\mathcal{D}$.

Let $\mathcal{F}$ be the set of machine learning models trained by the analyst on the dataset $\mathcal{D}_{train}$.

The **empirical Rashomon Set** is the subset of $\mathcal{F}$ containing those models whose accuracy on the test dataset, computed in terms of some measures of goodness of fit $\mathcal{A}$ (e.g. the $R^2$), exceeds a threshold $\alpha$:

$$\tilde{\mathcal{R}}(\alpha) = \{m \in \mathcal{F} : \mathcal{A}\left(Y_{test}, m(\mathbf{X}_{test})\right) \geq \alpha\}$$

- For each model *m* in the empirical Rashomon Set, obtain the corresponding importance vector.

- Sort the features in descending order of importance and assign each feature its corresponding position in this order to obtain a **ranking vector**.

- For each model $m$ in the empirical Rashomon Set, obtain the corresponding importance vector.

- Sort the features in descending order of importance and assign each feature its corresponding position in this order to obtain a **ranking vector**.

- **Ranking aggregation**: combining multiple rankings into a single, representative *consensus ranking* that synthesize information from input rankings.

- Several ranking aggregation techniques are available in the literature:
  - **Heuristic techniques** (e.g. *CombMIN*, *CombMAX*, *CombSUM* (Fox and Shaw 1994), *Borda counts* (Borda 1953), *Condorcet method* (De Condorcet et al. 1785))
  - **Optimization-based techniques** (e.g. minimizing Kendall's tau (Kendall 1948) or Spearman's footrule similarity (Spearman 1987))
  - **Distribution-based techniques** (e.g. *MC4*, *MCT* (DeConde et al. 2006), *Robust Rank Aggregation* (Kolde et al. 2012))

## Simulated data

## Real data

### Hooker test case

- 10 features, linear DGP
- Permutation Importance
- Two empirical Rashomon Sets of Linear Regressions based on the $R^2$ with respectively $\alpha = 0.94$ and $\alpha = 0.89$

### Friedman test case

- 6 features (1 irrelevant)
- 2-Wasserstein Sensitivity Index
- Classification task
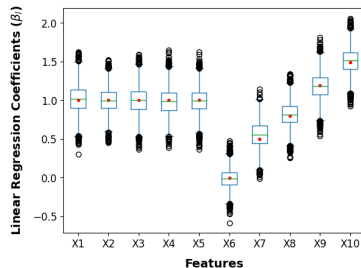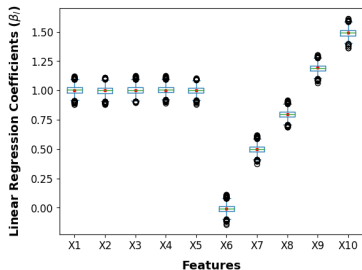- Empirical Rashomon Set based on Accuracy with $\alpha = 0.8$

### Boston Housing

- 13 features
- Permutation Importance
- Empirical Rashomon Set based on the $R^2$ with $\alpha = 0.8$

### Procedure

1. Compute the empirical Rashomon Set
2. Choose one importance measure and compute the importance vector and importance ranking vector
3. Aggregate rankings to compute the consensus ranking using different aggregation techniques
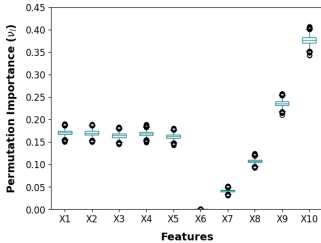
- $X_j \sim \text{Uniform}[0, 1], \quad \text{for } j = 1, ..., 10$

- $Y = X_1 + X_2 + X_3 + X_4 + X_5 + 0X_6 + 0.5X_7 + 0.8X_8 + 1.2X_9 + 1.5X_{10} + \eta, \quad \text{with } \eta \sim N(0, 0.1^2)$

- Training test: 10000 observations

- Test set: 5000 observations

- Two empirical Rashomon Sets of Linear Regressions based on the $R^2$, with $\alpha = 0.94$ and $\alpha = 0.89$:
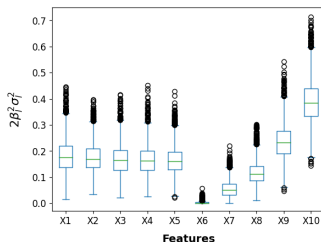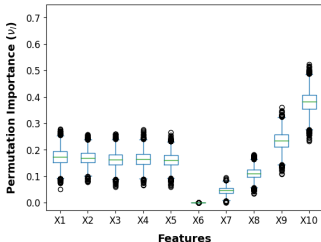
$\alpha = 0.94$

$\alpha = 0.89$

- $X_j \sim \text{Uniform}[0, 1]$ for $j = 1, \ldots, 6$

- $Y = \mathbb{1}[10 \sin(\pi X_1 X_2) + 20(X_3 - 0.5)^2 + 10X_4 + 5X_5 + \eta \geq 15], \quad \text{with } \eta \sim Normal(0, 1^2)$

- Training test: 10000 observations

- Test set: 5000 observations

- Epirical Rashomon Set made of 24 classification models whose accuracy is greater than or equal to 80%
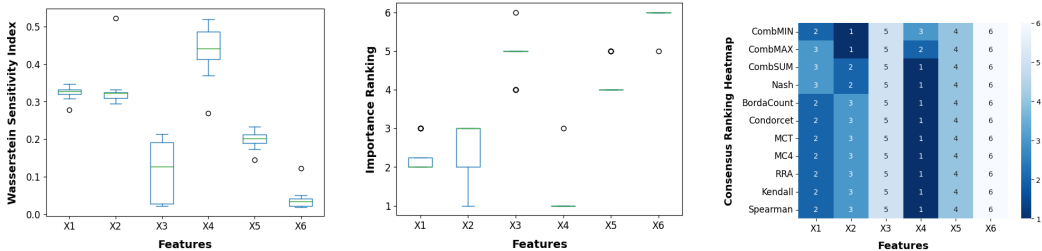
- Importance measure:

## 2-Wasserstein Sensitivity Index (Borgonovo, Figalli, et al. 2024)

Let $Y \in \mathbb{R}$ and let $Q_Y(u)$ and $Q_{Y|X_j}(u)$ represent the $u^{th}$ quantile of the distribution of $Y$ and of $Y$ given $X_j$, respectively. Then, the **2-Wasserstein Sensitivity Index** is defined as:

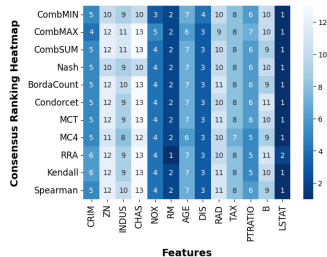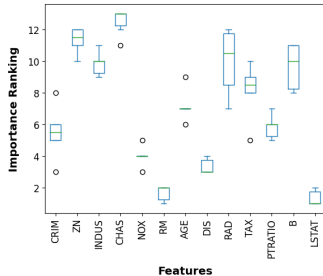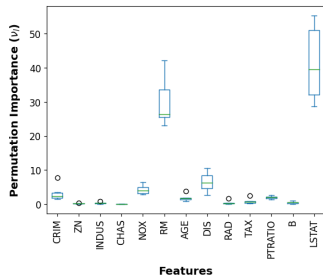$$\iota(Y, X_j) = (2\mathbb{V}[Y])^{-1} \mathbb{E}\left[\int_0^1 \left(Q_Y(u) - Q_{Y|X_j}(u)\right)^2 du\right]$$

- $X_j \sim$ Uniform$[0, 1]$ for $j = 1, \ldots, 6$

- $Y = \mathbb{1}[10 \sin(\pi X_1 X_2) + 20(X_3 - 0.5)^2 + 10X_4 + 5X_5 + \eta \geq 15]$, with $\eta \sim Normal(0, 1^2)$

- Training test: 10000 observations

- Test set: 5000 observations

- Epirical Rashomon Set made of 24 classification models whose accuracy is greater than or equal to 80%

- 506 observations of 13 features capturing information about houses and their surrounding neighborhoods

- Output variable: **MEDV**, the median value of owner-occupied homes

- Train-test split: 80% train, 20% test

- Epirical Rashomon Set made of 6 tree-based models:
  a gradient boosting ($R^2 = 92\%$), an XGBoost ($R^2 = 91\%$), a bagging regressor ($R^2 = 91\%$), a light gradient boosting machine (LGBM) ($R^2 = 89\%$), an extra-trees regressor ($R^2 = 88\%$) and a random forest ($R^2 = 88\%$)

## Contributions

- We established a theoretical connection between linear models' coefficients and Permutation Importance that facilitates interpretation

- We introduced a framework that uses the Rashomon Set to derive a unique feature importance ranking

- Our approach offers an intuitive means for practitioners to interpret model outputs in real-world contexts, where decision-makers may lack advanced technical expertise

## Limitations

- The choice of the rank aggregation technique is crucial to achieve a reliable consensus

- We do not address the issue of explanation multiplicity

# **Thank You!**

---

**Claudia Sessa**

PhD Candidate at LIUC University

`csessa@liuc.it`

Borda, Jean C de (1953). **"Memoire sur les Elections au Scrutin, 1781".** In: *Histoire de l'Academie Royale des Sciences, Paris* 99.

Borgonovo, Emanuele, Alessio Figalli, et al. (2024). **"Global sensitivity analysis via optimal transport".** In: *Management Science*.

Borgonovo, Emanuele, Gordon B Hazen, and Elmar Plischke (2016). **"A common rationale for global sensitivity measures and their estimation".** In: *Risk Analysis* 36.10, pp. 1871–1895.

Breiman, Leo (2001a). **"Random Forests".** In: *Machine Learning* 45.1, pp. 5–32.

— (Aug. 2001b). **"Statistical Modeling: The Two Cultures (with comments and a rejoinder by the author)".** In: *Statistical Science* 16.3. ISSN: 0883-4237.

De Condorcet, Nicolas et al. (1785). ***Essai sur l'application de l'analyse à la probabilité des décisions rendues à la pluralité des voix***. Limprimerie royale.

DeConde, Robert P et al. (2006). **"Combining results of microarray experiments: a rank aggregation approach".** In: *Statistical applications in genetics and molecular biology* 5.1.

Dong, Jiayun and Cynthia Rudin (2020). **"Exploring the cloud of variable importance for the set of all good models".** In: *Nature Machine Intelligence* 2.12, pp. 810–824.

Donnelly, Jon et al. (2023). **"The rashomon importance distribution: Getting rid of unstable, single model-based variable importance".** In: *Advances in Neural Information Processing Systems* 36, pp. 6267–6279.

Fisher, Aaron, Cynthia Rudin, and Francesca Dominici (2019). **"All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously".** In: *Journal of Machine Learning Research* 20.177, pp. 1–81.

Fox, Edward and Joseph Shaw (1994). **"Combination of multiple searches".** In: *NIST special publication SP*, pp. 243–243.

Gunasekaran, Abirami, Pritesh Mistry, and Minsi Chen (2024). **"Which Explanation Should be Selected: A Method Agnostic Model Class Reliance Explanation for Model and Explanation Multiplicity".** In: *SN Computer Science* 5.5, pp. 1–20.

Homma, Toshimitsu and Andrea Saltelli (Apr. 1996). **"Importance measures in global sensitivity analysis of nonlinear models".** In: *Reliability Engineering & System Safety* 52.1, pp. 1–17. ISSN: 0951-8320.

Hooker, Giles, Lucas Mentch, and Siyu Zhou (2021). **"Unrestricted permutation forces extrapolation: variable importance requires at least one more model, or there is no free variable importance".** In: *Statistics and Computing* 31, pp. 1–16.

Kendall, Maurice George (1948). *Rank correlation methods*. Griffin.

Kolde, Raivo et al. (2012). **"Robust rank aggregation for gene list integration and meta-analysis".** In: *Bioinformatics* 28.4, pp. 573–580.

Mata, Kota, Kentaro Kanamori, and Hiroki Arimura (2022). **"Computing the Collection of Good Models for Rule Lists".** In: *Proc. the 18th International Conference on Machine Learning and Data Mining (MLDM 2022)*.

Semenova, Lesia, Cynthia Rudin, and Ronald Parr (2022). **"On the existence of simpler machine learning models".** In: *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pp. 1827–1858.

Spearman, C. (1987). **"The Proof and Measurement of Association between Two Things".** In: *The American Journal of Psychology* 100.3/4, pp. 441–471. ISSN: 00029556.

Williamson, Brian D. et al. (Jan. 2022). **"A General Framework for Inference on Algorithm-Agnostic Variable Importance".** In: *Journal of the American Statistical Association* 118.543, pp. 1645–1658. ISSN: 1537-274X.

Xin, Rui et al. (2022). **"Exploring the whole rashomon set of sparse decision trees".** In: *Advances in neural information processing systems* 35, pp. 14071–14084.

Zhong, Chudi et al. (2024). **"Exploring and Interacting with the Set of Good Sparse Generalized Additive Models".** In: *Advances in Neural Information Processing Systems* 36.