

**UNIVERSIDADE DE SÃO PAULO
INSTITUTO DE FÍSICA DE SÃO CARLOS**

Pedro de Carvalho Braga Ilídio Silva

On bipartite decision forests

São Carlos

2024

Pedro de Carvalho Braga Ilídio Silva

On bipartite decision forests

Dissertation presented to the Graduate Program
in Physics at the Instituto de Física de São
Carlos da Universidade de São Paulo, to obtain
the degree of Master in Science.

Concentration area: Biomolecular Physics

Advisor: Prof. Dr. Otavio Henrique Thiemann

Coadvisor: Prof. Dr. Ricardo Cerri

Original version

São Carlos

2024

É possível elaborar a ficha catalográfica em LaTeX ou incluir a fornecida pela Biblioteca. Para tanto observe a programação contida nos arquivos USPSC-modelo.tex e fichacatalografica.tex e/ou gere o arquivo fichacatalografica.pdf.

A biblioteca da sua Unidade lhe fornecerá um arquivo PDF com a ficha catalográfica definitiva, que deverá ser salvo como fichacatalografica.pdf no diretório do seu projeto.

ACKNOWLEDGEMENTS

Primeira frase do agradecimento

Segunda frase

Outras frases

Última frase

This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Finance Code 001.

ABSTRACT

ILÍDIO, P. **On bipartite decision forests**. 2024. 95p. Dissertation (Master in Science) - Instituto de Física de São Carlos, Universidade de São Paulo, São Carlos, 2024.

The present study investigates decision forest algorithms for predicting interactions in bipartite networks. We concentrate on examples of such problems in the biological domain, such as drug-protein interactions, microRNA-gene interactions or long non-coding RNA-protein interactions. Notwithstanding, the proposed methods encompass the broad range of tasks satisfying i) the goal is to predict interactions between two entities; ii) the interacting pairs are composed of two different types of entities; and iii) each type of entity has its own set of input features. We refer to this paradigm as bipartite interaction learning or *bipartite learning*. Predicting interactions in such networks has fundamental challenges. For instance, the number of possible interactions is often very large in comparison to the number of known interactions. As a result, the data is frequently sparse, and negative annotations are unreliable. We explore a class of decision forest models specifically designed to address these challenges, that we broadly call *bipartite forests*. First, we demonstrate how these trees can be adapted to yield a $\log n$ speedup in training time. We also propose using weighted-neighbors approaches to determine each leaf's output, which resulted in improved generalization. Finally, we introduce semi-supervised impurity functions to bipartite forests. These functions result in trees that also consider clusters of instances in the feature space, rather than only their labels. This is shown to improve the forests' resilience to the missing annotations. Our models display highly-competitive performance across ten interaction prediction datasets. We believe the proposed methods can be a crucial step in developing effective and scalable machine learning models for interaction prediction. Further adaptations of these models could also impact other domains, such as recommendation systems, multilabel learning and weak-label learning.

Keywords: Decision forests. Interaction prediction. Bipartite learning. Positive-unlabeled learning.

RESUMO

ILÍDIO, P. **Florestas de decisão bipartidas**. 2024. 95p. Dissertação (Mestrado em Ciências) - Instituto de Física de São Carlos, Universidade de São Paulo, São Carlos, 2024.

O presente estudo investiga algoritmos de floresta de decisão para prever interações em redes bipartidas. Concentra-se em exemplos de tais problemas no domínio biológico, como interações fármaco-proteína, interações microRNA-gene ou interações entre moléculas de RNA longo não codificante e proteínas. No entanto, os métodos propostos abrangem uma ampla gama de tarefas de aprendizado, caracterizadas por i) o objetivo é prever interações entre duas entidades; ii) os pares de interação são compostos por dois tipos diferentes de entidades; e iii) cada tipo de entidade possui seu próprio conjunto de características de entrada. Refere-se a este paradigma como aprendizado de interações bipartidas, ou *aprendizado bipartido*. Prever interações em tais redes nos apresenta desafios fundamentais. Por exemplo, o número de interações possíveis é frequentemente muito superior ao número de interações conhecidas. Como resultado, os dados são muitas vezes esparsos, e as anotações negativas são incertas. Exploramos uma classe de florestas de decisão especificamente projetadas para enfrentar esses desafios, que chamamos de *florestas bipartidas* em geral. Primeiro, demonstramos como essas árvores podem ser adaptadas para obter uma melhora logarítmica no tempo de treinamento. Também propomos o uso de abordagens de vizinhos ponderados para determinar a saída de cada folha, resultando em melhora na capacidade de generalização dos modelos. Finalmente, introduzimos funções de impureza semi-supervisionadas para florestas bipartidas. Essas funções resultam em árvores cientes da densidade do espaço de características, em vez de apenas considerar os rótulos para o crescimento. Mostra-se que isso melhora a resiliência das florestas às anotações faltantes. Nossos modelos exibem desempenho altamente competitivo em dez conjuntos de dados de previsão de interação. Acreditamos que os métodos propostos podem ser um passo crucial no desenvolvimento de modelos de aprendizado de máquina eficazes e escaláveis para prever interações. Adaptações adicionais desses modelos também podem impactar domínios vizinhos, como sistemas de recomendação, aprendizado multi-rótulo e aprendizado de rótulos fracos.

Palavras-chave: Florestas de decisão. Predição de interações. Aprendizado bipartido. Aprendizado positivo-não rotulado.

CONTENTS

1	INTRODUCTION	11
1.1	Context and motivation	11
1.1.1	How to build bipartite models	13
1.1.2	Why decision trees?	15
1.2	Related work	16
1.3	Main objectives	17
1.4	Organization	17
2	DEVELOPMENT	19
2.1	Definitions	19
2.1.1	Mathematical notation	19
2.1.2	Problem statement	20
2.2	Data-centered adaptations of learning algorithms to bipartite interaction prediction	21
2.2.1	The standard global single-output adaptation	22
2.2.2	The standard local multi-output adaptation	22
2.3	Bipartite decision forests	23
2.3.1	Decision trees	24
2.3.2	Searching for the best split	26
2.3.3	Measuring the quality of a split	27
2.3.4	Bipartite decision trees	29
2.3.5	Bipartite global single-output trees	30
2.3.6	Prototype functions for bipartite trees	31
2.3.7	Asymptotic complexity analysis	32
2.3.8	Decision forests	34
2.3.9	Incorporating semi-supervision into decision trees	37
2.3.9.1	Semi-supervised learning	37
2.3.9.2	Semi-supervised decision trees	38
2.3.9.3	Unsupervised impurities	40
2.3.9.4	Heuristics for σ	41
2.4	Other estimator-centered strategies	42
2.4.1	Linear models	42
2.4.2	Neighborhood-Regularized Matrix Factorization	45
2.4.2.1	Traditional matrix factorization	45
2.4.2.2	Neighborhood regularization	48
2.5	Assessing the performance of bipartite models	52

2.5.1	Datasets	52
2.5.1.1	DPI-E, DPI-G, DPI-I, DPI-N	52
2.5.1.2	ERN and SRN	52
2.5.1.3	DAVIS	52
2.5.1.4	KIBA	53
2.5.1.5	mirTarBase	53
2.5.1.6	NPInter	54
2.5.2	There are multiple ways to measure model generality	54
2.5.3	Cross-validating in two dimensions	56
2.5.4	Prediction scoring metrics	56
2.5.4.1	Ideal descriptions of AUROC and AUPRC	60
2.5.4.2	AUPR and AUROC in terms of ranked decision values	62
2.5.4.3	AUROC is the normalized mean percentile ranks	64
2.5.5	General experimental settings	65
2.6	Experiments	65
2.6.1	What are the differences between AUROC and AUPR?	65
2.6.2	Are BGSO models faster than GMO models?	67
2.6.3	Which prototype should a GMO forest use?	68
2.6.4	Which adaptation strategy is the best for decision forests?	70
2.6.5	Can label imputation assist bipartite forests?	72
2.6.6	What is the best way of building semi-supervised forests?	74
2.6.7	Which forests are the best?	76
2.6.8	Can bipartite forests compete with other proposals?	77
3	CONCLUSION	91
3.1	Main findings	91
3.2	Main contributions	92
3.3	Future work	92
	REFERENCES	95

1 INTRODUCTION

As the starting point, we provide a brief overview of the concepts explored along this thesis. We begin by describing the learning paradigm of interest, and proceed by outlining its main challenges and current approaches. We then present the guiding research questions of our work, and describe how the remainder of this document will be organized.

1.1 Context and motivation

The effect a drug molecule has on our organism is tightly associated with the specific microscopic structures it physically and chemically interacts with (?). The regulatory roles of microRNAs are often mediated by their binding to a defined set of messenger RNAs. The activity of a transcription factor is determined by which genes it is capable of regulating (?). The success of a recommendation system depends on its ability to predict a user’s preference for a set of items in a catalog (?). All those scenarios share the same underlying structure: there are two distinct domains of objects that interact with one another, forming a heterogeneous *bipartite network* (?).

Frequently, great importance lies in computationally describing such relationships. For instance, developing new drug molecules is a costly and time-consuming process, requiring years of complex computer simulations, large screening assays, and extensive trials (). Being able to predict beforehand the interactions in a large dataset of drugs and protein targets holds a remarkable potential to accelerate drug development and reduce its costs, enabling the concentration of financial resources and human expertise to a set of most promising candidates. Furthermore, being able to characterize on a large scale the relationships between microRNAs and genes or other biomolecular entities can greatly aid the understanding of regulatory mechanisms of gene expression, a valuable step towards the discovery and development of new therapies.

Therefore, it is natural that numerous machine learning algorithms have been proposed to predict interactions in bipartite networks. The role of a learning algorithm in this setting is to receive a feature vector describing the heterogeneous pair, and characterize their interaction through a binary label. Additionally, we focus on cases where the input vector can be split into two: one vector representing each interacting entity.

Despite their importance, interaction prediction problems pose unique challenges to machine learning approaches, that we enumerate below.

1. **Heterogeneous data:** The input data is often composed of two very dissimilar types of objects, each with its own set of features. As a result, estimators may be required to simultaneously deal with descriptors of very different natures.

2. **High dimensionality:** Due to the intrinsic combinatorial nature of this type of problem, the interactions to be processed are usually very numerous, even if the number of training instances in each domain is not that expressive. Specifically, for each new object introduced to the training set, a new interaction could be considered for each training sample in the other domain, rendering the number of interactions to grow quadratically with the number of entities of each type.
3. **Knowledge sparsity:** As a consequence of the previous point, no amount of research efforts in characterizing new interactions can keep up with the rate at which possible relationships appear, resulting in a fundamental sparsity of confidently verified interactions.
4. **Lack of validated negative annotations:** In many cases, the absence of a relationship between two objects is much harder to validate than its presence. For instance, negative results in biochemical essays asserting molecular interactions can often be caused by a plethora of external experimental factors other than the actual lack of interaction under optimal conditions, unlike positive results that usually imply a successful assay (?, ?). Furthermore, since positive results are often much more informative, negative results in such interaction characterization settings are naturally less valued in the scientific community, and therefore less likely to be reported. These factors result in a general lack of high-quality negative data in interaction datasets, requiring special considerations from machine learning pipelines that are often overlooked in the literature.
5. **Impact on model evaluation:** Two rank-based metrics are mainly used to evaluate the performance of interaction prediction models: the area under the receiver operating characteristic curve (AUROC) and the area under the precision-recall curve (AUPR). These metrics are known to assess the performance of a model in different ways, and are often used in conjunction to provide a more comprehensive view. However, these differences are not fully clear, and neither is how they are affected by the missing annotations. This makes it difficult to interpret why different models are sometimes chosen, and to decide which metric to prioritize in each experiment. For instance, AUPR is often recommended for imbalanced scenarios (), but metrics similar to AUROC are suggested in interaction prediction contexts (?, ?).

These issues are an integral part of interaction prediction problems, that arise from how they are fundamentally defined around combinations of objects. However, these differences are often overlooked when developing machine learning algorithms. For instance, it is common to represent each dyad with the concatenated feature vectors of its components. Nevertheless, it can be challenging to collect representations for all dyads in this format, which would result in a very large matrix to be processed by the algorithm. Instead, the training set for a bipartite problem is more naturally represented as two separate design matrices rather than the usual single one. Ideally, we would want to avoid the combination of feature vectors as a data preprocessing

step, leaving the conjoint treatment of both design matrices as a matter of estimator design. This view defines a different machine learning paradigm, in which what we call "instances" or "samples" are not the dyads themselves, but rather each individual interacting object. The estimator receives both feature matrices and it is up to the algorithm to decide how to use them in conjunction to predict new interactions.

In conclusion, acknowledging such defining aspects and incorporating the bipartite nature of the problem are crucial steps in developing successful and scalable machine learning models for interaction prediction.

1.1.1 How to build bipartite models

There are two general strategies for developing machine learning algorithms for bipartite interaction prediction. We call them *data-centric* and *estimator-centric* approaches. The majority of methods proposed for bipartite learning is based on organizing the data to enable the use of traditional learning algorithms. These data-centric adaptations can be broadly classified into two categories ($?, ?, ?$), which we describe below.

- **Global single output:** The most common approach observed in the literature. The dataset is formatted so that each dyad is an instance. Each dyad has a single label associated with it, and the feature vector describing the dyad is the concatenation of the two feature vectors describing the interacting entities. Usual single-output learning algorithms can then be applied to the problem.
 - **Main advantage:** Enables the estimator to consider features from both instance domains simultaneously.
 - **Main disadvantage:** It can be computationally infeasible to build the feature vectors and train the estimator for all possible interactions.
- **Local multi-output:** This strategy separates the bipartite learning problem into two multi-output tasks. We first select one of the domains as input and the other as output. We train a multi-output estimator to receive the features of the input domain and predict its interactions with the output domain. Each instance in the output domain is considered an output to be modeled, so that the features of the output domain are completely disregarded in this step. The same procedure is then repeated with the roles of the domains reversed. Predicted interactions of these two models are fed into a new round of training that occurs in a similar fashion. A total of four models is necessary to predict interactions between instances that were both not present in the training set.
 - **Main advantage:** A much smaller number of input samples is considered, in comparison to the global single output approach. As a consequence, the models are usually faster to train.

- **Main disadvantage:** It requires training new models for each new batch of interactions to be predicted. This severely limits the application on online learning settings.

Other methods adapt the learning algorithm itself to take full advantage of the bipartite format. These strategies tend to be more computationally efficient than the global single output approach, since they consider the dataset in a more compact format. Notwithstanding, they are still able to consider both domains simultaneously, unlike the local multi-output procedure. There is, however, a limited number of learning algorithms that can be adapted in a more fundamental level. Notable examples of estimator-centered adaptations are listed below.

- **Linear regression:** Instance-instance similarities among each domain could be used to build a dyad-dyad similarity matrix. This matrix could then be used as a kernel matrix for training a linear regression model in a global single output fashion. The main problem with this idea is that a dyad-dyad kernel matrix would be often prohibitively large. Kron-RLS (?) solves this issue, building the hypothetical linear model we described without explicitly obtaining the dyad-dyad kernel matrix.
 - **Main advantage:** It is a remarkably efficient method.
 - **Main disadvantage:** Linear relationships frequently are too simplistic to capture the complexity of the underlying interactions.
- **Matrix factorization:** The adjacency matrix of the bipartite network is decomposed into two latent feature matrices, one for each domain. The product of these matrices approximates the original adjacency matrix, and can be used to predict new interactions between known instances. liu2016neighborhood introduces a method to extend this idea to predict interactions with new instances, by encouraging instance-instance vicinities to be transferred from the original to the latent feature space.
 - **Main advantage:** It is a naturally scalable technique, being designed from the ground up for interaction prediction tasks.
 - **Main disadvantage:** Interpretability is limited and a large number of hyperparameters are necessary.
- **Decision trees:** (?) proposes an adaptation to the split search procedure occurring at each decision tree node. The search is performed locally for each domain, similarly to the local multi-output approach, but the resulting model is a single decision tree.
 - **Main advantage:** Better scores are achieved in comparison to decision trees adapted with data-centric approaches (?).
 - **Main disadvantage:** No improvement in training complexity is observed relative to decision trees under the global single output adaptation.

1.1.2 Why decision trees?

In particular, decision trees are a very interesting option to explore. They show remarkable interpretability, flexibility to different tasks, and straightforward usage, besides demonstrating remarkable predictive power when combined into ensembles (?). We list below some of the main advantages of learning algorithms based on decision trees.

1. **Interpretability:** Decision trees and forests are well-known transparent models, with several methods being available to interpret their predictions (?). For instance, the frequency with which a feature is used by the tree can reflect its importance in the prediction process (?).
2. **No need for preprocessing:** In most cases, very little preprocessing is necessary to use decision trees. They can directly handle continuous or categorical features, and natively deal with missing feature values (?). They are also not sensible to feature scaling, so normalization is most often not necessary.
3. **Flexibility:** Decision trees can be used for both classification and regression tasks, and can handle either single or multi-output problems. They were successfully adapted to a wide range of scenarios including time series forecasting (), survival analysis (), and unstructured inputs in general (?).
4. **Small number of hyperparameters:** They have a small number of hyperparameters in comparison to other algorithms such as deep neural networks, which makes them easy to use and tune.
5. **Predictive power:** They can uncover complex non-linear relationships between features and labels, while being notably resistant to overfitting when combined into ensembles. They are widely recognized for their predictive performance especially on structured data (?).
6. **Fast inference:** Once trained, they can make predictions in logarithmic time complexity (?).

Hence, the main focus of this study is to explore the potential of decision forests for bipartite learning. We introduce three main methodology refinements, aimed at the high dimensionality, sparsity, and lack of information imposed by our learning settings.

1. **More scalable trees:** We experiment with different metrics to guide the tree growing procedure. Optimizing the algorithm for these metrics can offer improvements in training complexity, making bipartite trees more scalable.

2. **Weighted neighbors prototypes:** We explore using similarity scores between instances to calculate the output value of each leaf node. This enables shallower trees to be built while increasing their ability to generalize.
3. **Semi-supervised impurities:** We build trees that consider how the instances are grouped in the feature space, and not only their labels. This improves the model’s resilience to missing annotations.

1.2 Related work

REESCREVER E COMPLEMENTAR

(?) define the ideas of the GSO and LMO adaptations of traditional machine learning algorithms. In their study, (?) analyses the performance of Random Forests breiman2001random (RF) and Extremely Randomized Trees geurts2006extremely (ExtraTrees) adapted in both ways to be applied to DTI prediction, suggesting a greater potential to generalization of ExtraTrees in comparison to RF.

In ??, (?) presents the PBCT algorithm, an original adaptation of decision trees to bipartite scenarios, and demonstrates improved results over the previous, more superficial, adaptations. Although directly considering the bipartite data format, with separate design matrices for each domain, their proposal is still limited by the same asymptotic complexity of GSO-adapted trees. The authors later extend the PBCT algorithm to ensembles, building bipartite versions of Random Forests and ExtraTrees (?). Later still, they explore the use of Neighborhood-Regularized Logistic Matrix Factorization (NRLMF) liu2016neighborhood as a pre-training step to bipartite ExtraTrees, generating continuous pseudo-labels for negative interactions to be inputted to the forest estimator. They demonstrate significant improvements in predictive performance resulting from this procedure (?), reinforcing the hypothesized potential of semi-supervised techniques. Nevertheless, the authors refrain from experimenting with random forests, and their bipartite models are still subject to the same training complexity.

(?) and independently propose a way to consider semi-supervised assumptions on the level of a decision tree’s split search procedure rather than the whole model. They do so by penalizing the selection of split points that separate close training instances, so that instances within the same tree node, that tend to yield the same output, also will now tend to have similar feature vectors. Under the assumption that neighboring instances are more likely to have similar labels (the *smoothness* assumption), delegating part of the split’s influence to the feature space directly is expected to reduce the impact of missing or incorrect labels.

More recently, (?) applied a more computationally efficient variant of this idea to bipartite trees, using a semi-supervised evaluation not in every possible split but to choose between the two best supervised splits of each domain at each tree node. The authors, however, do not explore ensembles of bipartite trees or the more expensive semi-supervised evaluation of every

candidate split.

1.3 Main objectives

The present work will be centered at the following research questions:

1. Can bipartite trees be faster?

Investigate how can we exploit the bipartite nature of the problem to achieve a speedup in the training of decision trees.

2. Are semi-supervised techniques beneficial for bipartite forests?

Test the hypothesis that semi-supervised techniques can improve the resilience of bipartite trees to missing annotations.

3. How do AUROC and AUPR differ in their assessment of model performance?

Evaluate in detail the differences between the two most common threshold-agnostic metrics for binary classification tasks: the Area Under the Receiver Operating Characteristic Curve (AUROC) and the Area Under the Precision-Recall Curve (AUPR).

4. How do bipartite forests compare with proficient models in the field?

Compare the performance of bipartite forests with other state-of-the-art models for bipartite learning.

1.4 Organization

EESCREVER

2 DEVELOPMENT

2.1 Definitions

2.1.1 Mathematical notation

REESCREVER E PADRONIZAR

For any given matrix M , we denote by $M^{[ij]}$ its element on the i -th row and j -th column ($i, j \in \mathbb{N}^*$). Analogously, we represent by $M^{[i]}$ the vector containing M 's i -th row so that $M^{[i][j]} = M^{[ij]}$ and by $M^{[j]}$ the column vector $((M^\top)^{[j]})^\top$ referring to the j -th column of M so that $M^{[j][i]} = M^{[ij]}$. Defining the index notations as superscripts enables the subscripts to be used only as indentifiers, naming the matrix or vector as a whole and not each element. Indices are also always represented by a single letter, to dispense the use of separators between them.

We write the total number of M 's rows as $|M|_i$ and its number of columns as $|M|_j$. The total number of elements in M is written $|M| = |M|_j |M|_i$, not to be confused with the determinant of M .

We display filtered matrices or vectors by writing the condition as the index, optionally enclosed by parentheses when necessary or improving readability (Eq. 2.1.1).

$$M^{[(i<3)j]} \equiv \{M^{[kj]} \mid k < 3\}$$

When summing over all indices in a given dimension, we took the freedom of omitting the start and end positions (Section 2.1.1).

$$\sum_i M^{[ij]} \equiv \sum_{i=1}^{|M|_i} M^{[ij]}$$

We also made the choice of representing averages in a more concise way, optionally pondered by sets of w_1 and w_2 weights in each respective axis (Equation 2.1).

$$\begin{aligned} M^{\langle i \rangle [j]} &\equiv \frac{\sum_i w_1^{[i]} M^{[ij]}}{\sum_i w_1^{[i]}} \\ M^{[i] \langle j \rangle} &\equiv \frac{\sum_j w_2^{[j]} M^{[ij]}}{\sum_j w_2^{[j]}} \\ M^{\langle ij \rangle} &\equiv \frac{\sum_j \sum_i w_2^{[j]} w_1^{[i]} M^{[ij]}}{\sum_j \sum_i w_2^{[j]} w_1^{[i]}} \end{aligned} \tag{2.1}$$

The enclosing of indices within brackets also allows for the omission of parentheses when concomitantly using exponents, as exemplified by Eq. 2.2, and we additionally reserve ourselves the freedom of representing each index only once, which in the last two shown cases

requires preemptively defining that i and j respectively represent rows and columns. Notice that dispensing parentheses makes important the order in which the exponent and averaged indices (those within $\langle \cdot \rangle$) appear. The position of indices within $[\cdot]$ is however facultative, and is here chosen to be as close as M as possible in order to avoid confusion with $M^2 = MM$. Also notice that the indices within $[\cdot]$ will be the indices of the resulting matrix or vector.

$$\begin{aligned}
M^{[ij]2} &= ((M^{[ij]})^2)^{[ij]} \\
M^{\langle ij \rangle 2} &= (M^{\langle ij \rangle})^2 \\
M^{2 \langle ij \rangle} &= (M^{[ij]2})^{\langle ij \rangle} \\
M^{[i]2 \langle j \rangle} &= M^{[ij]2 \langle j \rangle} \\
M^{[j]2 \langle i \rangle} &= M^{[ij]2 \langle i \rangle}
\end{aligned} \tag{2.2}$$

2.1.2 Problem statement

The supervised machine learning applications focus on modeling a function $f: \mathbb{R}^{n_f} \rightarrow \mathbb{R}^{n_o}$ whose exact underlying mechanism is unknown or costly to implement. As a result, the only information available about such mapping is a set of inputs $\{x_i \in \mathbb{R}_f^n\}$ and their corresponding outputs $\{y_i \in \mathbb{R}_o^n\}$ of that given function. The goal is therefore to build an *in silico* model (or estimator) \tilde{f} that approximates f , yielding as similar as possible outputs for the same given input, even and especially for outputs not utilized in the process of building \tilde{f} .

The known input vectors are usually organized as rows of an X matrix so that $X^{[ij]} = x_i^{[j]}$, and we refer as *feature* or *attribute* to each specific horizontal position j of $x^{[j]}$, which corresponds to a column of X . Likewise, a Y matrix is built with their corresponding outputs ($Y^{[ik]} = y_i^{[k]}$). Commonly referred to as "targets" in the context of regression learning, we here call the known outputs of the modeled process by *labels*, as in classification, even if real-valued, to avoid confusion when referring to the protein targets of a drug.

In the present setting, we concentrate on problems involving the interaction of two domains of instances (also called sample groups). As such, each sample domain forms a different X matrix, that we term X_1 and X_2 . Only inter-domain interactions are allowed, that is, instances are restricted from interacting with others in the same sample group, so that the interaction network constitutes an undirected bipartite graph. Each interacting pair of instances is also called a *dyad*.

The output, in our case, is any scalar piece of information describing the interaction between a given instance pair, such as the rating of a movie given by a user or a kinetic parameter of an enzyme-substrate reaction. The labels are then disposed in a $|X_1|_i$ by $|X_2|_i$ adjacency matrix Y (also called interaction matrix) so that the function to be modeled can now be formulated as mapping the pair's vector representations to the interaction label $f: (X_1^{[i \cdot]}, X_2^{[k \cdot]}) \mapsto Y^{[ik]}$.

Since each sample in X_1 refers to a *row* of Y and each sample in X_2 corresponds to a *column* of Y , we sometimes refer to the sample domains of X_1 and X_2 as *row samples* and

column samples, respectively. We call these datasets *bipartite*, to differentiate from the more common problems in which a single X matrix is utilized.

2.2 Data-centered adaptations of learning algorithms to bipartite interaction prediction

As previously discussed, bipartite interaction problems fundamentally differ from the usual machine learning paradigm, in which input data represents a single entity to be labeled. Interaction tasks are instead concerned with labeling a relationship between two entities, and as such, each prediction is made upon a pair of feature vectors, each vector being specific to each of the two sample domains.

Such subtlety is most often bypassed by reformulating a bipartite dataset into the traditional machine learning setting (?). These strategies can be encompassed under two general approaches, initially termed *global* and *local* by vert2008reconstruction and later adopted and extended by sschrynemackers2015. For the sake of clarity and generality, we further specify these categories by defining *global* and *local* as general properties of estimators, rather than specific training procedures:

- **Global** estimators are those aware of both instance domains during the training procedure (X_1 and X_2).
- **Local** estimators are those which only have access to feature information from one of the two instance domains during training (either X_1 or X_2). As such, they are often employed in compositions, combining the predictions of several local models to produce the final output.

Furthermore, to be consistent with (?, ?, ?, ?), we assume the following definitions in our current context of bipartite interactions:

- **Single-output** estimators are those which consider all labels, i.e. all $Y^{[ij]}$ elements, irrespectively of the column j or row i they are in. They are all regarded as a *single type of output*.
- **Multi-output** estimators, on the other hand, are those which consider each instance, each row or column of Y , as a separate task, for example by defining a composite loss function formed by the combination of losses over each row or column.

Notice that the label matrix Y can still be represented in two dimensions even if the model is single-output in this sense, contrary to the usual case where bidimensionality of Y is a defining characteristic of a multi-output problem.

Finally, the two most common ways of adapting traditional estimators to bipartite data are then named *global single-output* (GSO) and *local multi-output* (LMO), as proposed by pliakos2018,pliakos2019,pliakos2020. We further denote them standard, to clearly distinguish them from new adaptation proposals that could share the globality, locality, or outputness properties, but work in an entirely different way. Specific definitions and shortcomings of these procedures are now presented.

2.2.1 The standard global single-output adaptation

Arguably the most straightforward way to provide bipartite data into a conventional learning algorithm is by presenting concatenated pairs of row-column samples, labeled by each element of Y . The interacting dyad is now what we consider a sample, and each feature vector is the combination of the feature vectors of the two interacting instances. This is usually done by converting the two X matrices and the interaction matrix Y to a single design matrix we call X_{SGSO} and a column-vector of labels we refer to as Y_{SGSO} .

One way of doing so is by choosing indices as described by Equation 2.3, where $\begin{bmatrix} x_1 & x_2 \end{bmatrix}$ denotes concatenation and all i_1 and i_2 combinations are explored (see ??).

$$\begin{aligned} i_g &= (i_1 - 1)|X_2|_i + i_2 - 1 \\ X_{\text{SGSO}}^{[i_g \cdot]} &= \begin{bmatrix} X_1^{[i_1 \cdot]} & X_2^{[i_2 \cdot]} \end{bmatrix} \\ Y_{\text{SGSO}}^{[i_g 1]} &= Y^{[i_1 i_2]} \end{aligned} \tag{2.3}$$

To consider all possible dyads, X_{SGSO} would have $|X_1|_i |X_2|_i$ rows and $|X_1|_j + |X_2|_j$ columns. Thus, X_{SGSO} as defined by Equation 2.3 is highly redundant. As a result, naively dealing with such a large X_{SGSO} matrix is impeditive in many cases, both in terms of memory usage and computation time. Therefore, a commonly used workaround is to undersample the negative annotations, yielding a dataset with equal amounts of positive and negative interactions (). This strategy is justified by the negative-annotated dyads being usually far more numerous than positives and much more likely to be truly negative interactions than false negatives (). Nonetheless, we demonstrate in Section 2.6.4 that undersampling negatives is significantly detrimental if the main goal is to indicate wich interactions are most likely to be true.

2.2.2 The standard local multi-output adaptation

The local approaches, in contrast to global methods, propose training different models on X_1 and X_2 , so that each estimator only has access to information regarding either row samples or column samples.

As such, multiple non-bipartite models need to be used in conjunction to predict interactions between new row samples and new column samples. The standard local multi-output

(SLMO) approach uses two traditional estimators for each axis (four in total), that we here refer to as *primary* and *secondary* estimators. In general, they must be multi-output estimators, each being able to output a bidimensional Y matrix.

The procedure for training the estimators in an SLMO setting consists of simply training the primary estimators. The primary rows estimator is trained on X_1 and Y , and the primary columns estimator on X_2 and Y^\top (function `TrainLocalModel`). The prediction step is however more complicated, involving the training of the secondary estimators on predictions of the primary models. The procedure is described in details by function `PredictLocalModel` and illustrated by `??`. The `combine` function used in line 15 of the function `PredictLocalModel` procedure can be arbitrarily chosen, and is usually defined as the simple element-wise average of both matrices (`combine` (Y_1, Y_2) = $\frac{1}{2}(Y_1 + Y_2)$).

Function `TrainLocalModel`(X, Y)

Input: The bipartite training dataset.

Output: A bipartite local model.

```

1 train(primaryRowsEstimator,  $X_1, Y$ );
2 train(primaryColumnsEstimator,  $X_2, Y^\top$ );
3 return primaryRowsEstimator, primaryColumnsEstimator

```

An unusual behaviour occurs if the secondary estimators are able to exploit inter-output dependencies. In this case, the outputs of the SLMO composition would depend on the amount of test data and specific combinations of test instances provided. This is usually not the case in traditional learning problems, and complicates the evaluation procedure. The behaviour forces us to consider the size and sampling method of the batches of test data. For online applications, the ideal strategy would be to retrain the secondary models for each batch of new instances, using the predictions from the primary models over all the previously received instances (labeled and unlabeled). As a result, the application of SLMO models is effectively restricted to offline learning contexts.

2.3 Bipartite decision forests

Some learning algorithms can be adapted in a deeper level to the bipartite context. These algorithms are able to directly receive the two X matrices as input, without needing to combine them as a preprocessing step for training. We term them *learner-centered* adaptations, in contrast with the data-centered strategies discussed in Section 2.2.

This section explores the main focus of this work: learner-centered adaptations of decision forest models. Section 2.4 will then present two other learner-centered strategies, based on linear models and matrix factorization respectively. A general description of top-down decision trees is now presented, as a theoretical foundation for the upcoming formal definition of bipartite decision trees (Section 2.3.4).

Function PredictLocalModel(primary models, X_{new})

Input: The trained primary models and the unseen sample matrices X_{new} for both axes.

Output: Y_{pred} predictions for each interaction provided.

```

1  $Y_{\text{new rows}} \leftarrow \text{predict}(\text{primaryRowsEstimator}, X_{1\text{new}})$ 
2  $Y_{\text{new cols}} \leftarrow \text{predict}(\text{primaryColumnsEstimator}, X_{2\text{new}})$ 
3 if Secondary estimators consider label dependencies then
    // Concatenate known rows and columns labels to the
    // primary predictions
4    $Y_{\text{new cols}} \leftarrow \begin{bmatrix} Y^{\top} \\ Y_{\text{new cols}} \end{bmatrix};$ 
5    $Y_{\text{new rows}} \leftarrow \begin{bmatrix} Y \\ Y_{\text{new rows}} \end{bmatrix};$ 
    // Otherwise, if label columns are considered
    // independently, this step is not necessary
6 end

7  $\text{train}(\text{secondaryRowsEstimator}, X_1, Y_{\text{new cols}}^{\top});$ 
8  $\text{train}(\text{secondaryColumnsEstimator}, X_2, Y_{\text{new rows}}^{\top});$ 

9  $Y_{\text{pred rows}} \leftarrow \text{predict}(\text{secondaryRowsEstimator}, X_{1\text{new}});$ 
10  $Y_{\text{pred cols}} \leftarrow \text{predict}(\text{secondaryColumnsEstimator}, X_{2\text{new}});$ 

11 if Secondary estimators consider label dependencies then
    // Skip predictions not referring to  $X_{1\text{new}}$  and  $X_{2\text{new}}$ 
12    $Y_{\text{pred rows}} \leftarrow Y_{\text{pred rows}}^{[j > |X_1|_i]};$ 
13    $Y_{\text{pred cols}} \leftarrow Y_{\text{pred cols}}^{[j > |X_2|_i]};$ 
14 end

15 return  $\text{combine}(Y_{\text{pred rows}}, Y_{\text{pred cols}}^{\top})$ 

```

2.3.1 Decision trees

Let's consider the scenario where a single protein of interest is selected, and we receive the task of determining which drug molecules will likely affect its physical structure or catalytic function. We wish to find a systematic procedure to decide whether a given drug molecule x_i will interact with our protein or not. To develop such a procedure, consider we have at our disposal a set of n_f known drug molecules, whose degree of interaction with our protein of interest was previously experimentally determined. We can then describe a drug molecule x_i in general by how similar it is to each of our n_f known molecules, organizing this information as a vector $x_i = [x_i^{[1]} \ x_i^{[2]} \ x_i^{[3]} \ \dots \ x_i^{[n_f]}]$ so that $x_i^{[j]}$ represents the similarity score between the drug x_i and the j -th of our n_f known drugs.

The hypothetical decision procedure we intend to determine could then be structured

as a path with consecutive bifurcations. We always start at the same place, and, at each bifurcation, a question is asked about the drug x_i in hands. The questions are in a standard format, exemplified by "Is x_i more than 60% similar to the 3rd known drug?", or $x_i^{[j]} > t$, for a general known drug j and similarity threshold t . The answer to the question in each bifurcation determines which of the two possible paths we should follow. No cycles are allowed in the path, and eventually, all routes reach final locations instead of bifurcations. Each final location contains a fixed value that will be returned as the final decision about the drug x_i 's effect on our protein of interest.

Such a decision procedure, structured as a binary tree path, is what is commonly referred to as a *decision tree* (DT) model, illustrated by ???. In this context, each bifurcation then represents a tree *node*, and the final locations are called *leaves*. The value outputted by each leaf is termed the leaf's *prototype*.

The main challenge, however, lies in the building process of such models: in how to determine the rules that define each fork and how we define the stopping criteria for a final decision to be yielded.

To build the decision tree, we are given a training set composed of the two bidimensional matrices X and Y . The dataset represents $|X|_i = |Y|_i$ instances, each described by $|X|_j$ numeric features and labeled by $|Y|_j$ labels. In the previous example, $|X|_i$ would be the number of drug molecules in the dataset $|X|_j$ would be the number of reference drugs (used to obtain the similarities), and $|Y|_j$ would be the number of proteins of interest.

Consider now executing the prediction process of a decision tree for each training instance, going through the branched path. Each bifurcation would divide the training instances between those who answer the question affirmatively and those who answer negatively. Eventually, each leaf will contain a partition of the training instances. From the resulting partitions, we can evaluate the decision tree: a "good" decision tree would be one in which the prototype value of each leaf is a good estimate of the labels in the leaf's partition. The training algorithm then approaches the problem from the other direction: what would be a good tree according to our training data? A greedy top-down procedure is usually followed, that we introduce below.

Formally, each rule is encoded by an index f representing a feature column and a threshold value t . Each rule then represents a split of the training dataset in two partitions, named *left* and *right*, as defined by Equation 2.4.

$$\begin{aligned}
 Y_{\text{left}} &= \{Y^{[k]} \vee k \mid X^{[kf]} \leq t\} \\
 Y_{\text{right}} &= \{Y^{[k]} \vee k \mid X^{[kf]} > t\} \\
 X_{\text{left}} &= \{X^{[kf]} \vee k \mid X^{[kf]} \leq t\} \\
 X_{\text{right}} &= \{X^{[kf]} \vee k \mid X^{[kf]} > t\}
 \end{aligned} \tag{2.4}$$

To build a decision tree, we start with the whole training set (the *root node*). A procedure FindSplit is executed to yield a good splitting rule. The training data is then partitioned

according to this rule, as defined in Equation 2.4. Two descendant nodes are created, each receiving one of the partitions. For each new node, we have two options: i) apply FindSplit to its data partition and continue recursively; or ii) set the node as a leaf, taking record of the partition it received.

The tree-building algorithms can thus be described by four key components:

1. **The FindSplit procedure**, which determines a set of candidate splitting rules and selects the best one according to a quality metric.
2. **The split quality metric** used by FindSplit to evaluate the candidate rules.
3. **The stopping criteria**, that determine when a node should be set as a leaf.
4. **The prototype function**, that determines the output value of a leaf.

function BuildTree describes in detail how these components come together to build a decision tree. function Predict then formally details how predictions are made given a model built by function BuildTree and a new data instance. These procedures encompass a wide range of decision tree algorithms, including the popular CART (?), ID3 (?), and C4.5 (?), as well as the bipartite trees we explore in the present work (?).

The following sections provide further considerations on specific components.

Function Predict(RootNode, x): Compute a Decision Tree's prediction.

Input: A new interaction sample to be evaluated and the root node of a Decision Tree.

Output: The Decision Tree's predicted value for the given sample attributes.

```

1 Node ← RootNode;
2 while Node is not a leaf do
3   if  $x^{[Node.feature]} > Node.threshold$  then
4     Node ← Node.childRight
5   else
6     Node ← Node.childLeft
7   end
8 end
9 return prototype(Node.X, Node.Y)
```

2.3.2 Searching for the best split

We explore two main strategies for the FindSplit procedure: the exhaustive search and the randomized search. The exhaustive search is the most common. The randomized search is a faster but less accurate alternative, intended to be used in the context of tree ensembles (??). We describe each procedure in the present section.

Function BuildTree(X, Y): Recursively build a Decision Tree

Input: The training data for the current node.
Output: Current node, with all information of subsequent splits.

```

1  $Q^*, f^*, t^* \leftarrow \text{FindSplit}(X, Y, \tilde{n}_f)$ ;
   // Many stopping criteria are possible
2 if DecideToStop( $Q^*, f^*, t^*, X, Y$ ) then
3   return NodeObject {
4     isLeaf  $\leftarrow$  True
5      $X_{\text{leaf}} \leftarrow X$ 
6      $Y_{\text{leaf}} \leftarrow Y$ 
7   };
8 else
9   Get  $X_l, Y_l, X_r, Y_r$  from  $f^*$  and  $t^*$  (Eq. 2.4);
10  return NodeObject {
11    isLeaf  $\leftarrow$  False
12    childLeft  $\leftarrow$  BuildTree( $X_l, Y_l$ )
13    childRight  $\leftarrow$  BuildTree( $X_r, Y_r$ )
14    feature  $\leftarrow f^*$ 
15    threshold  $\leftarrow t^*$ 
16  };
17 end

```

The most common approach to the FindSplit procedure is to consider all possible partitions of the given input data. For a feature column $X^{[f]}$, this exhaustive evaluation of partitions can be done by sorting $X^{[f]}$ and considering a threshold t between each two consecutive values in it. Notice that any threshold value between the same two consecutive $X_{\text{sorted}}^{[f]}$ elements will result in the exact same partitioning of the training set (Equation 2.4). The common practice is thus to take the averages between each two neighboring feature values. The same procedure is then repeated for each feature column, and the overall best t^* and corresponding feature index f^* are selected. The exhaustive split search procedure is detailed by the function FindSplitBest.

The randomized search is an alternative that avoids considering all partitioning options and greatly reduces the amount of operations performed. It consists of drawing a random threshold t between the minimum and maximum value of each feature, thus evaluating only $|X|_j$ splits when choosing the best (line 13). Although degrading the performance of a single tree, this procedure is an interesting option when building tree ensembles (??), being the core idea behind the extremely randomized trees (ERT) algorithm (?). Ensembles of decision trees will be discussed in ??.

2.3.3 Measuring the quality of a split

A split quality criterion must be defined so we can compare and select the best splitting rules at each node. The quality ΔI of a split is commonly framed as the decrease of an impurity

Function FindSplitBest(X, Y)

Input: A partition of the training data in a given node.

Output: The highest quality score Q^* found among all splits evaluated, with its corresponding feature column f^* and threshold value t^* .

```

1 Initialize  $S_r$  and  $S_l$  as a  $|Y|_j$ -sized vectors;
2  $Q^*, f^*, t^* \leftarrow 0$ ;
3 Draw  $\tilde{n}_f$  columns (features) of  $X$  without replacement;
4 foreach feature index  $f$  of the  $\tilde{n}_f$  drawn features do
5      $n_l \leftarrow 0$ ;                                     // Holds  $|Y_l|_i$ 
6      $S_r \leftarrow \sum_i Y^{[ij]}$ ;                       // Holds  $\sum_i Y_r^{[ij]}$ 
7      $S_l \leftarrow 0$ ;                                   // Holds  $\sum_i Y_l^{[ij]}$ 
8     Get the permutation  $P$  that sorts  $X^{[f]}$ ;
9     Apply  $P$  to  $Y$ 's and  $X^{[f]}$ 's rows:
10     $Y_P, X_P \leftarrow P(Y), P(X^{[f]})$ ;
11    foreach row index  $\hat{i}$  of  $Y_P$  do
12         $n_l \leftarrow n_l + 1$ ;
13        foreach column index  $\hat{j}$  of  $Y_P$  do
14             $S_r^{[\hat{j}]} \leftarrow S_r^{[\hat{j}]} - Y_P^{[\hat{i}\hat{j}]}$ ;
15             $S_l^{[\hat{j}]} \leftarrow S_l^{[\hat{j}]} + Y_P^{[\hat{i}\hat{j}]}$ ;
16        end
17        Use  $S_l, S_r$  and  $n_l$  to calculate  $Q$  (Eq. ??). Notice that other node-specific
            constants might be needed;
18        if  $Q > Q^*$  then
19             $Q^* \leftarrow Q$ ;
20             $f^* \leftarrow f$ ;
21             $t^* \leftarrow \frac{1}{2}(X_P^{[\hat{i}]} + X_P^{[\hat{i}+1]})$ ;
22        end
23    end
24 end
25 return  $Q^*, f^*, t^*$ ;
  
```

metric calculated over the partitions of training labels (Equation 2.5). This decrease is taken for the combined impurities of the generated children nodes (Equation 2.4). All impurities are multiplied by the size of each partition relative to the total number of training samples ($|Y_{\text{root}}|$), restricting the effect of nodes with less data that could introduce spurious variations of impurity. Notice that $|Y_{\text{node}}| = |Y_{\text{left}}| + |Y_{\text{right}}|$.

$$\Delta I(Y, t, f) = \frac{|Y_{\text{node}}|}{|Y_{\text{root}}|} I(Y_{\text{node}}) - \frac{|Y_{\text{left}}|}{|Y_{\text{root}}|} I(Y_{\text{left}}) - \frac{|Y_{\text{right}}|}{|Y_{\text{root}}|} I(Y_{\text{right}}) \quad (2.5)$$

Several metrics can be chosen as the impurity function $I(\cdot)$, such as the Gini impurity, the Shannon entropy or the Poisson loss (). In this study we utilize the variance of each output

Function FindSplit_{random}(X, Y)**Input:** A partition of the training data in a given node.**Output:** The highest quality score Q^* found among all splits evaluated, with its corresponding feature column f^* and threshold value t^* .

```

1  $Q^*, f^*, t^* \leftarrow 0$ ;
2 Draw  $\tilde{n}_f$  columns (features) of  $X$  without replacement;
3 foreach feature index  $f$  of the  $\tilde{n}_f$  drawn features do
4   Find  $\min(X^{[f]})$  and  $\max(X^{[f]})$ ;
5   Draw a random threshold value  $t \in \mathbb{R}$  so that  $\min(X^{[f]}) < t < \max(X^{[f]})$ ;
6   Calculate  $Q$  for the drawn  $t$  (Eq. ??); //  $O(|Y|)$ 
7   if  $Q > Q^*$  then
8      $Q^* \leftarrow Q$ ;
9      $f^* \leftarrow f$ ;
10     $t^* \leftarrow \frac{1}{2}(X_P^{[i]} + X_P^{[i+1]})$ ;
11  end
12 end
13 return  $Q^*, f^*, t^*$  1FindSplitrandom( $I$ ) FindSplitrandom

```

column, averaged over all outputs (Equation 2.6).

$$I_{\text{MSE}}(Y) = (Y^{[ij]} - Y^{\langle i \rangle [j]})^2 \langle ij \rangle = Y^{2 \langle ij \rangle} - Y^{\langle i \rangle 2 \langle j \rangle} \quad (2.6)$$

Most commonly, the prototype function returns the column averages of a leaf's partition of the training labels: $Y_{\text{leaf}}^{\langle i \rangle [j]}$. In this case, the column variances correspond to the *mean squared error* (MSE) for the training data in the node, as if the node holding $Y_{\text{partition}}$ were to become a leaf.

Also notice that I_{MSE} is equivalent to the Gini impurity if Y contains only binary values. That can be shown by noticing that $Y_{\text{bin}}^{[ij]2} = Y_{\text{bin}}^{[ij]}$ for binary labels, so $Y_{\text{bin}}^{2 \langle i \rangle [j]} = Y_{\text{bin}}^{\langle i \rangle [j]} \equiv p^{[j]}$, which yields Equation 2.7.

$$I_{\text{MSE}}(Y_{\text{bin}}) = Y_{\text{bin}}^{2 \langle ij \rangle} - Y_{\text{bin}}^{\langle i \rangle 2 \langle j \rangle} = (p^{[j]} - p^{\langle j \rangle 2 \langle j \rangle}) \langle j \rangle = [p^{[j]}(1 - p^{[j]})] \langle j \rangle = I_{\text{Gini}}(Y) \quad (2.7)$$

2.3.4 Bipartite decision trees

In this section, we describe how the decision tree algorithm can be adapted to directly operate on bipartite interaction data.

Essentially, the adaptation process consists of defining bipartite versions of each of the four key components of the decision tree algorithm (Section 2.3.1): the FindSplit procedure, the split quality metric, the stopping criteria, and the prototype function. From these, the most central is the FindSplit procedure adaptation.

For bipartite trees, as proposed by (?), we perform the split search locally in each domain, in a similar fashion to the SLMO adaptation (Section 2.2.2). Specifically, let FindSplit

trad be a procedure for finding a split threshold in a traditional multioutput decision tree (see Section 2.3.2). A bipartite tree applies $\text{FindSplit}_{\text{trad}}$ twice at each node: once over X_1 and Y and once over X_2 and Y^\top . This results in one split being chosen for each axis. Finally, the best overall split is selected between the vertical and horizontal splits. In essence, a bipartite decision tree algorithm uses a FindSplit procedure that composes two traditional $\text{FindSplit}_{\text{trad}}$ procedures, as detailed by function $\text{FindBipartiteSplit}$.

Function $\text{FindBipartiteSplit}(X, Y)$

Input: A partition of the bipartite training data in a given node. X encodes one design matrix for each axis, X_1 and X_2 .

Output: The highest quality score Q^* found among all splits evaluated in both row and column directions, with its corresponding feature column f^* and threshold value t^* .

```

1 if adapterStrategy is GSO (Section 2.3.5) then
    | // Build  $Y$  proxies  $\tilde{Y}_1$  and  $\tilde{Y}_2$  (Equation 2.10)
2   |  $\tilde{Y}_1 \leftarrow Y^{[ \cdot ](j)}$ ;
3   |  $\tilde{Y}_2 \leftarrow Y^{(i)[ \cdot ]}$ ;
4 else
    | // Using GMO strategy, no proxies are used
5   |  $\tilde{Y}_1 \leftarrow Y$ ;
6   |  $\tilde{Y}_2 \leftarrow Y^\top$ ;
7 end

    // Generate a split in each axis. Get each split's
    // position, feature and quality score
8  $Q_r^*, f_r^*, t_r^* \leftarrow \text{FindSplit}(X_1, \tilde{Y}_1)$ ;
9  $Q_c^*, f_c^*, t_c^* \leftarrow \text{FindSplit}(X_2, \tilde{Y}_2)$ ;
10 if  $Q_r^* > Q_c^*$  then
11 |   return  $Q_r^*, f_r^*, t_r^*$ 
12 else
    | //  $f_c^*$  value lets clear its  $X_2$  ownership
13 |   return  $Q_c^*, f_c^*, t_c^*$ 
14 end
```

2.3.5 Bipartite global single-output trees

This section shows how we can grow bipartite trees more efficiently, making tree estimators more scalable to large bipartite datasets.

The original bipartite tree proposed by (?) uses the impurity metric I_{MSE} (Equation 2.6) to evaluate splits, so that the authors classify their technique as a *Global MultiOutput* (GMO) estimator. *Multi-output* because the impurity evaluates the average variance of each label column (or row), as if each column (or row) were a different output. *Global* because each tree utilizes features from both feature domains for training (see Section 2.2). Their GMO trees, however,

are shown to have the same algorithmic complexity as traditional decision trees trained with the SGSO adaptation. We demonstrate that faster training algorithms can be developed if we assume a single-output format under the context of bipartite trees.

The split rules as defined by Equation 2.4 are agnostic to the specific arrangement of the bipartite data. They can be seamlessly applicable to either the $(X_{\text{SGSO}}, Y_{\text{SGSO}})$ format employed by the SGSO adaptation (??) or directly to the X_1 , X_2 and Y matrices. Furthermore, the impurity at each node of an SGSO tree, in many cases can be translated to the bipartite format. Equation 2.8 describes such a translation of the MSE impurity presented by Equation 2.6.

$$I_{\text{MSE}}(Y_{\text{SGSO}}) = Y_{\text{SGSO}}^{2\langle ij \rangle} - Y_{\text{SGSO}}^{\langle i \rangle 2\langle j \rangle} = Y^{2\langle ij \rangle} - Y^{\langle ij \rangle 2} \quad (2.8)$$

If we then define

$$I_{\text{GMSE}}(Y) \equiv I_{\text{MSE}}(Y_{\text{SGSO}}) = Y^{2\langle ij \rangle} - Y^{\langle ij \rangle 2}, \quad (2.9)$$

the exact same SGSO-adapted decision tree, with its node structure and split rules, can be grown by applying the bipartite procedure (function `FindBipartiteSplit`) with the I_{GMSE} impurity instead of the original I_{MSE} . This again follows from the generality of the split rules defined in Equation 2.4, invariant under the SGSO data rearrangement (??).

Notice how the squared divergences in ?? is computed relative to the average of each Y column, while in 2.9 the inner average is computed over the whole Y matrix.

We can explore this property of the GSO impurity to iterate more efficiently over candidate splits. This is done by pre-computing averages of each row and column of Y_{node} . We can then iterate over one-dimensional \tilde{Y}_{node} proxies (Equation 2.10) instead of the bi-dimensional matrix when evaluating splits (?? of ??).

$$\begin{aligned} \tilde{Y}_1^{[i]} &= Y^{[i] \langle j \rangle} \\ \tilde{Y}_2^{[j]} &= Y^{\langle i \rangle [j]} \end{aligned} \quad (2.10)$$

The same is not possible for the GMO trees, as Equation 2.6 requires storing averages for each Y column, to be squared individually. function `FindBipartiteSplit` describes the adapted procedure in each case. The complexity improvements are demonstrated theoretically in Section 2.3.7 and empirically in Section 2.6.2.

2.3.6 Prototype functions for bipartite trees

When a leaf is reached during the prediction step of a decision tree, the `prototype` function is called to determine the output value to be returned (line 9 of function `Predict`). With traditional datasets, the `prototype` function most often returns the average label of the leaf's partition (Equation 2.11) (?).

$$\text{prototype}(Y_{\text{leaf}})^{[j]} = Y_{\text{leaf}}^{\langle i \rangle [j]} \quad (2.11)$$

As an extension, the most natural approach to use on single-output bipartite trees is the analogous average of the whole partition of the interaction matrix (Equation 2.12).

$$\text{prototype}(Y_{\text{leaf}}) = Y_{\text{SGSO, leaf}}^{(i)[1]} = Y_{\text{leaf}}^{(ij)} \quad (2.12)$$

Nevertheless, some considerations are possible when dealing with bipartite data, since there are cases in which one of the entities of the interaction being predicted is already known from the training set. As introduced by pliakos2018global, if a row or column instance is in the training set, we have the option of averaging only the column or row (respectively) of Y_{leaf} corresponding to its known outputs. Specifically, when predicting the interaction between a sample pair $x_{1, \text{new}}$ and $x_{2, \text{new}}$, we can set prototype as in Equation 2.13.

$$\text{prototype}(Y_{\text{leaf}}) = \begin{cases} Y_{\text{leaf}}^{[k]\langle j \rangle} & \text{if } \exists k \mid x_{1, \text{new}} = X_{1, \text{leaf}}^{[k \cdot]} \\ Y_{\text{leaf}}^{\langle i \rangle [k]} & \text{if } \exists k \mid x_{2, \text{new}} = X_{2, \text{leaf}}^{[k \cdot]} \\ Y_{\text{leaf}}^{\langle ij \rangle} & \text{otherwise.} \end{cases} \quad (2.13)$$

A drawback of this approach, especially when working with very imbalanced interaction matrices and sufficiently small leaf partitions, is a possibly greater susceptibility to random fluctuations, since the label averages in the prediction step are taken over a much smaller sample size (a single row or column of Y_{leaf} instead of the whole Y_{leaf}). Given we are predominantly working with similarity scores as sample attributes, we propose an intermediate approach: to weight the rows and columns of Y_{leaf} by the similarity values in the form $w_{s1}^{[i]} \equiv \text{similarity}(x_1^{[i]}, X_1^{[i \cdot]})$ between x_{new} and the training samples in the leaf node (Equation 2.14).

$$\text{prototype}(Y_{\text{leaf}}) = \frac{\sum_{i \in \text{leaf}} w_{s1}^{[i]} Y_{\text{leaf}}^{[i]\langle j \rangle}}{2 \sum_{i \in \text{leaf}} w_{s1}^{[i]}} + \frac{\sum_{j \in \text{leaf}} w_{s2}^{[j]} Y_{\text{leaf}}^{\langle i \rangle [j]}}{2 \sum_{j \in \text{leaf}} w_{s2}^{[j]}} \quad (2.14)$$

Since we are dealing with precomputed pairwise similarities, X_1 and X_2 are square matrices in which $X_a^{[i_1 i_2]} = \text{similarity}(X_a^{[i_1 \cdot]}, X_a^{[i_2 \cdot]})$. We explore three different cases:

1. $w_s = x^{[i]}$
2. $w_s = (x^{[i]})^2$
3. $w_s = e^{x^{[i]}}$

2.3.7 Asymptotic complexity analysis

In this section, we derive the asymptotic complexity of the proposed algorithms for growing bipartite trees.

We consider that the number of horizontal instances has the same magnitude as the number of vertical instances, so that we define $n \sim |X_1|_i \sim |X_2|_i$. Similarly, we assume $m \sim |X_1|_j \sim |X_2|_j$ and define $f(m)$ as the effective number of features. That is, the split search procedure samples $f(m)$ features to consider in each node. The sampling is most often done without replacement, so that $f(m) \in O(m)$. Common choices are $f(m) = \sqrt{m}$ or $f(m) = \log_2(m)$.

From the algorithm description, we can infer that the complexity of both function `FindSplitBest` and line 13 will be given by

$$\begin{aligned} \text{FindSplit}(n, f) &\in \Theta(fS(|Y|_i) + f|Y|) \\ &= \Theta(fS(n) + fn^2) \\ &= \Theta(fn^2) \end{aligned} \quad (2.15)$$

where $S(n)$ is the complexity of the chosen sorting algorithm of when operating on n values. When features are real-valued, the current most effective algorithms are $\Theta(n \log n)$, as, for example, Quick Sort or Merge Sort (). For integer features, the complexity is $\Theta(n)$, as in Counting Sort or Radix Sort ().

When applied on X_{SGSO} and Y_{SGSO} , we still have that $|Y| = n^2$ and $f(m) \in O(m)$.

$$\begin{aligned} \text{FindSplit}_{\text{SGSO}}(n, f) &\in \Theta(fS(|Y_{\text{SGSO}}|_i) + f|Y_{\text{SGSO}}|) \\ &= \Theta(fS(n^2) + fn^2) \\ &= \Theta(fS(n^2)) \end{aligned} \quad (2.16)$$

For BGSO, the complexity is given by

$$\begin{aligned} \text{FindSplit}_{\text{BGSO}}(n, f) &\in \Theta(n^2 + fS(|\tilde{Y}_1|_i) + f|\tilde{Y}_1|) \\ &= \Theta(n^2 + fS(n) + fn) \\ &= \Theta(n^2 + fS(n)) \end{aligned} \quad (2.17)$$

The complexity of `FindBipartiteSplit` will be the same as the non-bipartite `FindSplit` that it wraps, since it simply applies the wrapped function twice at each node.

To estimate the complexity of the whole tree-building process, let's consider the ideal scenario of a balanced decision tree. We also consider that the axis of a split alternates at each level, so that a matrix Y is separated in four equal sized pieces after two levels. This results in the following recurrence relationship:

$$T(n) = 4T\left(\frac{n}{2}\right) + \text{FindSplit}(n, f) \quad (2.18)$$

$T(n)$ is the time to build the tree from an n by n interaction matrix, and `FindSplit` (n, f) in this case is the time taken to select a split. The algorithm complexity of such recursive functions

then follows the master theorem (?). The theorem gives the time complexity of a function T obeying the recurrence relation $T(n) = aT(n/b) + F(n)$ and $c = \log_b a$:

$$T(n) \in \begin{cases} \Theta(n^c) & \text{if } F(n) \in O(n^{c-\epsilon}). \\ \Theta(n^c \log^{k+1} n) & \text{if } F(n) \in \Theta(n^c \log^k n). \\ \Theta(f(n)) & \text{if } F(n) \in \Omega(n^{c+\epsilon}) \text{ and } F(n) \text{ is regular.} \end{cases} \quad (2.19)$$

In Equation 2.19, ϵ is a positive infinitesimal constant and k is any non-negative integer. For the third case, we say a function $F(n)$ is regular if it satisfies the *regularity condition* $aF(n/b) \leq qF(n)$ for some constant $q < 1$ and all sufficiently large n (?).

For our tree algorithms, Equation 2.18 shows that $c = 2$ and $F(n)$ represents $\text{FindSplit}(n, f)$. The resulting complexities are presented by Table 1. We see that the choice of the sorting function only affects the SGSO-adapted trees. In the worst case, where $S(n) = n \log n$, the SGSO trees are $\log n$ times slower than the GMO trees. If $S(n) = n$ instead, both have the same complexity.

Table 1 also shows that BGSO trees are faster than GMO trees when considering a high number of features. Specifically, if $f(n) \in O(\log n)$, BGSO will be f times faster than GMO. If $f(n) \in \Omega(\log n)$, BGSO will be $\log n$ times faster than GMO. The only case in which both are equivalent is when $f \in O(1)$.

Strategy	Split search	Tree building	$f \propto n$
SGSO	$O(fnS(n))$	$O(fnS(n) \log n)$	$O(n^2 S(n) \log n)$
GMO	$O(fn^2)$	$O(fn^2 \log n)$	$O(n^3 \log n)$
BGSO	$O(n^2 + fS(n))$	$O(n^2(\log n + f))$	$O(n^3)$

Table 1 – Comparison between asymptotic time complexities of decision tree-building procedures. We assume $n \sim |X_1|_i \sim |X_2|_i$. Similarly, f represents the number of features to be considered for the split search in each node. $S(n)$ denotes the complexity of the sorting procedure, usually n for integral features or $n \log n$ for decimal values. The last column refers to the case where $f \sim n$. This scenario could arise, for instance, when X_a are pairwise similarities or kernel matrices. Source: made by the authors.

2.3.8 Decision forests

While single decision trees are valuable tools for comprehending the learning problem at hand, they often fall short in effectively modeling intricate relationships within the data, displaying limited generalization capabilities and high susceptibility to overfitting (). Their most significant impact on machine learning applications is observed when committees of such models are utilized, in which the final output values are, for example, the average predictions of all trees.

These compositions of estimators are usually referred to as *ensembles* (). Being studied in the context of machine learning since the 1970s, they reflect the intuitive idea that combining multiple opinions from a diverse set of experts frequently results in better decision taking.

In fact, it has been extensively demonstrated both empirically and theoretically that the predictive performance a group of learners always surpasses that of its individual components if and only if the individual estimators are sufficiently accurate and diverse (?). Importantly, requiring diversity means that the individual estimators must ideally commit errors on different instances for the composition to succeed (?).

dietterich2000ensemble, polikar2006ensemble describe three ways in which combining diverse estimators could benefit the ensemble's performance.

1. **Statistical:** Building each estimator can be seen as finding a hypothesis that explains the prediction problem as well as possible. If multiple different hypotheses are found, the likelihood that at least one of them is close to the true underlying function is increased. Furthermore, averaging multiple hypotheses can contribute to alleviate the influence of each hypothesis' variance, resulting in a more accurate approximation and reduced propensity to overfitting.
2. **Computational:** Finding the globally optimal decision tree or neural network is known to be an NP-complete problem (). Thus, finding and combining multiple approximate hypotheses, that may represent local maxima of the objective function, is often much cheaper than expending time on the search for a global solution.
3. **Representational:** Sometimes, a single estimator is not complex enough to represent the intricacies of the problem at hand. However, combination of models can be utilized to enhance the representational power of the ensemble, building a more general decision function as a combination of the simpler decision boundaries of each estimator.

Being diversity a key factor, mechanisms of introducing model heterogeneity are as important to consider as strategies for enhancing the base algorithm. In fact, many diversity-inducing techniques even degrade the performance of individual estimators, but still results in improvements for the ensemble as a whole (?). For example, one of the most straightforward ways to promote heterogeneity is to build each estimator on randomized subsets of the training data. In this and in many other cases, designing ensemble models is thus a matter of balancing diversity and individual strength.

This is the core idea of *decision forests* (the ensembles of decision trees). Even simple estimators such as decision trees, traditionally prone to overfitting, can be leveraged to compose powerful learning algorithms. We briefly describe some popular strategies for introducing heterogeneity in the decision tree growing procedure.

1. **Instance sampling:** A new version of the training set is built by randomly drawing instances from it, with or without replacement (selecting rows of X). If drawing with replacement, one can optionally draw the same number of instances as the original set, a procedure known as bootstrapping.
2. **Feature sampling:** A random subset of features is selected to be used (selecting columns of X). The number of features \tilde{n}_f selected in each node is commonly defined as a function of the total number of features $|X|_j$. Usual choices are $\tilde{n}_f = \lceil \sqrt{|X|_j} \rceil$ or $\tilde{n}_f = \lceil \log |X|_j \rceil$.
3. **Split threshold randomization:** Instead of searching for the best split threshold for each feature as in ??, a random threshold value is drawn between the minimum and maximum values of each feature column. This procedure is described in line 13.

Both instance and feature undersampling can be performed node-wise, occurring before the split search procedure of each node, or tree-wise, occurring once for each tree in the ensemble. Both can also be performed with or without replacement, but notice that sampling features with replacement only makes sense if randomization of split threshold is also employed. In such a case, a random candidate split threshold is selected for each repetition of a feature, whereas the same split threshold would be selected for all duplicates if the greedy approach was used, spending more time with no different result than if omitting the repeated features.

The concept of sampling instances with replacement to create a different training set for each estimator in the ensemble was first proposed by breiman1996bagging under the name of bootstrap aggregation or *bagging*.

The idea of selecting a subset of features was introduced independently by amit1997shape and tinkamho1998random. While amit1997shape explored feature sampling as a remedy for a shape recognition problem with impeditively large sample sets, tinkamho1998random was mainly focused on overfitting-prevention for decision forests.

It was breiman2001random who first combined the two ideas, proposing and greatly popularizing the Random Forest algorithm, with over 115 thousand citations according to Google Scholar as of september 2023 ().

geurts2006extremely later introduced the randomized split threshold concept, presenting the Extremely Randomized Trees algorithm (Extra-Trees). The algorithm shows competitive prediction scores and clear superiority in terms of training speed compared to bagging and Random Forests (?).

The two main ensemble-building strategies we explore in the current work can now be defined as follows:

- **Random Forests (RF):** Tree-wise instance sampling with replacement and node-wise feature undersampling without replacement are employed.

- **Extremely Randomized Trees (ERT):** Split threshold randomization and node-wise feature undersampling without replacement are employed. No resampling of instances is performed in the original proposal.

We refer the reader to [sagi2018ensemble](#) and [fawagreh2014random](#) for a more in-depth description of prominent decision forest strategies and previous work in the field. [amasyali2011comparison](#) provides an experimental comparison of the most popular ensemble methods. Regarding interaction problems, [schrynemackers2015classifying](#) explores the use of decision forests under the standard global multi-output and local single-output adaptations presented by Section 2.2.

The same tree-diversification and forest building techniques discussed in this section can also be applied to bipartite decision trees, with very small modifications regarding the data sampling procedures: instance sampling and feature sampling must now occur on both domains of the interaction dataset.

[\(?\)](#) explores these ideas, using bipartite global multi-output decision trees to build both Random Forests and Extra-Trees ensembles. We hereafter refer to these forests as Bipartite Random Forests (BRF) and Bipartite Extra-Trees (BXT), respectively. Their study suggests superior performance of BXT in comparison to SLMO- and SGSO-adapted forests, as well as in comparison to previously proposed algorithms [\(?\)](#). The authors later extend their work to include a label imputation step for predicting drug-target interactions [\(?\)](#). In that study, their model is built on a reconstructed version of the interaction matrix obtained through Neighborhood-Regularized Logistic Matrix Factorization (NRLMF, [liu2016neighborhood](#)).

2.3.9 Incorporating semi-supervision into decision trees

ADICIONAR MAIS REFERÊNCIAS E CONFERIR EXEMPLOS

The tree algorithms presented up until this point do not consider the positive-unlabeled (PU) property of interaction data. They work under the assumption of a *supervised* scenario, in which all annotations are reliable. In this section, we discuss adaptations to decision trees specifically designed to the challenges of PU data.

As several previous authors [\(? , ?\)](#), we argue that accounting for the partial availability of information could improve prediction performance of bipartite models. The strategies to deal with missing labels are the main concern of the *semi-supervised* learning paradigm [\(\)](#), which we briefly introduce.

2.3.9.1 Semi-supervised learning

Supervised problems are those we have been exploring, in which the goal is to build a model that can predict the label of a new instance based on its features. We have a matrix Y of target values and feature matrices X , and the model's objective is to learn the relationship between X and Y .

When the learning problem at hand does not involve predefined labels, the task is termed *unsupervised*. The most common example is clustering problems, where the goal is to group similar instances together based solely on their numerical attributes. In these problems, there is not a preconceived target to model, such as desired output values or classes to which each instance should belong. Only the X matrix is utilized, there is not an associated label matrix Y .

In summary, while supervised learning problems are concerned with the relationship between X and Y , unsupervised learning problems target the relationships between instances themselves. While supervised learning focuses on the graph of a function $f : X \rightarrow Y$, unsupervised learning focuses on the structure of the feature space.

Some problems, however, lie in the intersection of the supervised-unsupervised spectrum. They usually arise when the available annotations are incomplete or unreliable, so that both supervised and unsupervised objectives are of interest. Take for instance the case of an image gallery application capable of grouping all photos of a person into a folder with their name. While the names must be asked once to the user, all photos of each different person are already grouped beforehand. The algorithm is both concerned with labeling each photo with the correct name and with grouping similar photos together. Another example can occur when building a medical diagnosis database. Such a database could be a result of merging multiple datasets and considering a large collection of patients. Some amount of missing information is very common in these cases, and the algorithms must be able to deal with these vacancies. One possible strategy is to infer missing labels from the labels of similar instances, bringing the concept of clustering again into play.

In the examples, both label inference (that characterizes supervised problems) and clustering (that characterizes unsupervised problems) are intertwined as the model's objective. The goal is to simultaneously learn to correctly classify the instances with known labels while also considering similarities to infer the missing or uncertain annotations. The class of such hybrid learning tasks is called *semi-supervised learning* ().

Since missing labels are a defining characteristic of interaction prediction tasks, it is commonly suggested (?, ?) that applying semi-supervised concepts to our problems could significantly improve the performance of bipartite models.

For decision trees specifically, there are straightforward ways to incorporate semi-supervised assumptions, that we now describe.

2.3.9.2 Semi-supervised decision trees

Notice that the growing of decision trees unavoidably induces groupings of samples in the training set. Specifically, the structure of each tree represents a hierarchical clustering of the training samples, in which each tree node represents a partition of the training set composed by the training instances that reach that node. However, this clustering procedure is usually

performed under the objective of grouping instances with similar *labels*, which not necessarily means that instances in the same group will have similar *features*. This results from the definition of the impurity function I governing the split search procedure: I is usually chosen as to minimize the divergence of labels within each partition (see ??). As such, I commonly depends on Y_{node} alone. With that in mind, decision trees can be naturally adapted to unsupervised or semi-supervised tasks by redefining the impurity function to consider the feature matrix X_{node} instead of only the label matrix Y_{node} .

In a purely unsupervised context, an example would be to utilize the average of column variances in X_{node} as the impurity function (Equation 2.20). This would result in a tree that groups instances with similar features together, regardless of their labels. This idea is the main concept behind the CLUS algorithm (?).

$$I_{u, \text{MSE}}(X) = \text{MSE}(X) = \left(X^{[ij]} - X^{\langle i \rangle [j]} \right)^{\langle i \rangle 2 \langle j \rangle} \quad (2.20)$$

To address semi-supervised scenarios, we can consider both supervised and unsupervised objectives simultaneously. This can be done by using a linear combination of unsupervised and supervised impurities to guide the split selection, taking into account both the similarities between features and between labels to build a semi-supervised decision tree. Equation 2.21 defines such a hybrid impurity function.

$$I_{ss}(X_{\text{node}}, Y_{\text{node}}) = (1 - \sigma) \frac{I_u(X_{\text{node}})}{I_u(X_{\text{root}})} + \sigma \frac{I_s(Y_{\text{node}})}{I_s(Y_{\text{root}})} \quad (2.21)$$

We divide each term by the corresponding impurities on the root node, thus calculated over the entire training set. The reason is to avoid the influence of possible differences in scale between the values in X and in Y . It also compensates relative scale differences that could arise when choosing different functions for I_u and I_s . The parameter $\sigma \in [0, 1]$, that we call supervision balance, weights the relative importance given to each objective, with $\sigma = 0$ corresponding to a fully unsupervised tree and $\sigma = 1$ to a fully supervised tree. Strategies for adjusting σ are discussed in Section 2.3.9.4.

For semi-supervised tasks in general, where we have confirmed positive and confirmed negative annotations alongside missing entries, notice how I_s can only be calculated over the non-missing annotations, while I_u can always utilize both labeled and unlabeled instances. Nevertheless, please notice that all confirmed annotations are positive in our present scenario of positive-unlabeled learning. Thus, we must still consider the missing labels as zero entries for the calculation of I_s . Even so, we argue that decreasing the relative importance of the supervised objective can compensate the label uncertainty and improve generalization. The idea is to encourage the tree to appreciate the structures in the feature space. For example, it should select splits that isolate very compact groups of instances, even if the labels in a given group are not satisfactorily homogeneous.

A caveat of using a semi-supervised impurity as in Equation 2.21 is that trees could continue to find splits even if all instances in a given node have the same label. This is because

we could keep reducing the unsupervised impurity by further splitting even if the supervised impurity is already zero. However, the output value of each node is still calculated over Y_{node} alone, so all nodes descending from a homogeneous node would yield the same output. We avoid such redundant splits by forcefully stopping the split search for a node when $I_s(Y_{\text{node}}) = 0$.

Another important notice is that features must be normalized before training the tree: the impurity function I_u (Equation 2.20) is most often sensitive to the relative scale of the different columns of X . This is normally not a requirement for decision trees, since the split search procedure is performed separately for each feature and usually only depends on the order of the values, not on their specific magnitude.

2.3.9.3 Unsupervised impurities

We explore two different unsupervised impurity functions for the semi-supervised decision trees: the mean feature variance (Equation 2.20) and the mean pairwise distance between the samples (Equation 2.22). We describe and motivate them in this section.

The initial proposal of a semi-supervised impurity function (?) employed the mean variance of the feature matrix X as the unsupervised term (Equation 2.20), being concerned with traditional decision trees. The strategy, however, does not scale well with the number of features. This results from each node having to always consider the same number of features, even though the number of instances decrease. The scalability is especially a concern in the scenarios under study, in which the feature matrices are square similarity matrices.

Notice that, in the learning tasks under study, the X matrix already represents similarities between samples, and an additional metric such as the variance of each feature is not necessary to capture how close the samples are to each other. Exploring this property, we propose a more efficient unsupervised impurity function based on the average similarity between the samples in the tree node (equation 2.22).

$$I_{\text{MeanDistance}}(X) = \frac{1}{|\mathbf{s}_{\text{node}}|} \sum_{j \in \mathbf{s}_{\text{node}}} \sum_{i \in \mathbf{s}_{\text{node}}} (1 - X^{[ij]}) \quad (2.22)$$

\mathbf{s}_{node} denotes the set of indices representing the samples in the node. Notice that the number of features to consider equals the number of samples in the node: we always consider a square partition of X . As such, the number of operations required by I_u drops more steeply with respect to the node size in comparison to $I_{\text{MSE GMD}}$ or similar impurities, which is especially beneficial in our case of large X matrices. We can reduce even further the number of operations by considering only the upper or lower triangle of X , exploiting X 's symmetry.

We note that the subset of features is only used to calculate the impurity function, and not to select split points. The split search procedure is still performed considering all features as usual.

2.3.9.4 Heuristics for σ

Determining the ideal value of σ is not a trivial task. In fact, it is not even clear whether a single constant σ for the whole tree is enough or it should be adjusted for each node. We explore four different strategies for setting σ , that we present in this section.

alves2023semisupervised argues that the unsupervised impurity should be more important for nodes with a larger number of unlabeled instances, since the supervised information would be more likely to be unreliable. Thus, they propose updating σ according to the label density of each node, as defined by Equation 2.23.

$$\sigma = 0.1 + 0.9 \cdot Y_{\text{node}}^{\langle ij \rangle} \quad (2.23)$$

However, we must recall that the prototype value of each node is only dependent on its partition of the label matrix. Therefore, if Y_{node} is close to being homogeneous, a new split is unable to cause a big overall change in the outputs for the instances involved. Specifically, the most drastic prototype change possible occurs if we perfectly separate the instances with positive labels from the instances with negative labels. When the node partition is very imbalanced, even this ideal separation will only greatly affect the very few instances with the minority label. Thus, using the heuristic of Equation 2.23 prioritizes the unsupervised objective only in cases where no significant change in the output values is expected. Additionally, few new splits are possible in near-homogeneous nodes, and we are more likely to achieve label homogeneity (and thus stop the split search) right after we achieve the cases in which I_u is highly prioritized. In summary, the heuristic of Equation 2.23 is likely to undermine the effect of the unsupervised objective.

To test this hypothesis, we propose another heuristic for σ that prioritizes the unsupervised objective in the earliest stages of the tree growing process and the supervised objective in nodes closer to the leaves (Equation 2.24).

$$\sigma = 1 - \frac{|Y_{\text{node}}|}{|Y_{\text{root}}|} \quad (2.24)$$

This way, we take advantage of the unsupervised impurity prior to when the label partitions are already homogeneous. Essentially, we start by performing clustering in the feature space, and then gradually move towards separating instances based on their labels. This process should be similar to a two step procedure: we first identify large structures in the feature space (when a large number of instances is still available) and then apply the label clustering separately to each of these structures. The difference is that each of our semi-supervised trees represents a gradual transition between the two steps.

We also speculate that the observed benefit of a dynamic σ could emerge from the diversity it promotes in tree learners. As discussed in ??, the diversity of the individual estimators is a key factor for the success of ensemble models. In this sense, ensuring variety of the σ parameter could even be more important than estimating its "correct" value for each partition.

Therefore, we also evaluate the strategy of selecting a random σ at each node, drawn from a uniform distribution in the interval $[0, 1]$ (Equation 2.25).

$$\sigma \sim \mathcal{U}(0, 1) \quad (2.25)$$

2.4 Other estimator-centered strategies

This section presents two other methods for optimizing estimators for interaction prediction tasks. Like bipartite trees (Section 2.3), these methods focus on adapting the learning algorithm in a more fundamental level than data-centered approaches (Section 2.2). The first method is an adaptation of linear models to the bipartite setting, named RLS-Kron (?). It explores properties of the Kronecker product to build an estimator on a large global kernel matrix, without in fact having to calculate such matrix. The second method adapts matrix factorization approaches to take similarity side-features into account. Named Neighborhood-Regularized Logistic Matrix Factorization (NRLMF), this technique allows obtaining latent features for unseen instances, which is commonly not possible with traditional matrix factorization.

2.4.1 Linear models

One of the simplest approaches to learning problems in general is to assume a linear relationship between the input features and output labels. Formally, for the non-bipartite case, one assumes that the training output matrix Y can be approximated as follows:

$$\hat{Y} = XW \quad (2.26)$$

in which W is a matrix representing the set of parameters to be learned. To determine W , the mean squared error (MSE) is usually defined as the loss function to be minimized, to which we add an extra regularization term controlled by the hyperparameter α :

$$\mathcal{J} = \frac{1}{2} \|Y - \hat{Y}\|^2 + \frac{\alpha}{2} \|W\|^2 = \frac{1}{2} \|Y - XW\|^2 + \frac{\alpha}{2} \|W\|^2 \quad (2.27)$$

An analytical solution for W can be obtained by taking the derivative of \mathcal{J} with respect to W and setting it to zero:

$$\frac{\partial \mathcal{J}}{\partial W} = 0 = X^\top(XW - Y) + \alpha W \implies W = (X^\top X + \alpha I)^{-1} X^\top Y \quad (2.28)$$

There are scenarios, however, where the specific values of X are less interesting than the pairwise similarities between them. In those settings, while X may not be directly available, we do have access to similarity matrices S (also called *kernel* matrices) in which $S^{[ij]}$ designates a similarity score between $X^{[i]}$ and $X^{[j]}$. Rather than simply considering S in the same way we would treat X in linear regression, we could instead employ the *kernel trick* (?): replacing the XX^\top terms in the above equations by S . In this case, we are assuming that similarities $S^{[ij]}$ represent the internal product of the vectors $X^{[i]}$ and $X^{[j]}$ in some feature space, which is based

on the intuition that the internal product by itself can be regarded as a similarity metric. We also define W slightly differently, omitting the X^\top factor as $W = (S + \alpha I)^{-1}Y$, so that the final prediction is given by

$$\hat{Y} = SW = (S + \alpha I)^{-1}S \quad (2.29)$$

For the bipartite interaction prediction case, besides standard adaptations as described in Section 2.2, a unique formulation is presented by vanlaarhoven2011gaussian. Similar in concept to the standard global single output procedure (Section 2.2), the authors consider each interaction pair as a unitary instance. The authors then propose building a kernel matrix relating each pair of instances to another pair, and not each interacting entity to another of the same domain. If $S_1 \in \mathbb{R}^{n_1 \times n_1}$ and $S_2 \in \mathbb{R}^{n_2 \times n_2}$ are the intra-domain similarity matrices, the global kernel matrix S is defined as

$$S^{[(i_1 n_2 + i_2)(j_1 n_2 + j_2)]} = S_1^{[i_1 j_1]} S_2^{[i_2 j_2]} \quad (2.30)$$

or, more succinctly, as the Kronecker product of S_1 and S_2 :

$$S = S_1 \otimes S_2 \quad (2.31)$$

Each entry on S thus represents the similarity between the pair $X_1^{[i_1]}-X_2^{[j_1]}$ and another pair $X_1^{[i_2]}-X_2^{[j_2]}$ by the product of the similarities between $X_1^{[i_1]}$ and $X_1^{[j_1]}$ and between $X_2^{[i_2]}$ and $X_2^{[j_2]}$.

The bipartite linear regression is then framed on the vectorized Y , denoted $\text{vec}(Y)$, built by concatenating the columns of Y into a single $|Y|$ by 1 column vector. Purely for notation purposes, we organize the weight parameters as the vectorized version of a matrix W with the same dimensions of Y .

$$\text{vec}(\hat{Y}) = S \text{vec}(W) \approx \text{vec}(Y) \quad (2.32)$$

$$\text{vec}(W) = (S + \alpha I)^{-1} \text{vec}(Y) \quad (2.33)$$

As the reader may imagine, S gets prohibitively large for big datasets (it's a $n_1 n_2$ -sized square matrix!), both in terms of memory usage and the time needed to perform the matrix inversion. The authors, however, provide a clever way of circumventing this issue by decomposing each of the S_1 and S_2 kernel matrices separately and exploiting the properties of the Kronecker product.

Given that S_1 and S_2 are symmetric square matrices, it follows from the spectral theorem () that they can be decomposed as follows:

$$S_1 = U_1 \Lambda_1 U_1^\top$$

$$S_2 = U_2 \Lambda_2 U_2^\top$$

where, if λ_1 represents the vector of eigenvalues of S_1 , Λ_1 is the diagonal matrix of those eigenvalues ($\Lambda_1 = \text{diag}(\lambda_1)$), with U_1 columns representing their corresponding eigenvectors

$U^\top[i]$ for each $\lambda_1^{[i]}$. The symmetry of the similarity matrices also implies that U_1 and U_2 are orthogonal, i.e. $U_1^\top U_1 = U_1 U_1^\top = \mathbb{I}$ and $U_2^\top U_2 = U_2 U_2^\top = \mathbb{I}$, or, equivalently, $U_1^{-1} = U_1^\top$ and $U_2^{-1} = U_2^\top$. Utilizing the fact that $(AB) \otimes (CD) = (A \otimes C)(B \otimes D)$ (), the Kronecker product of S_1 and S_2 can be written as

$$S = S_1 \otimes S_2 = (U_1 \Lambda_1 U_1^\top) \otimes (U_2 \Lambda_2 U_2^\top) = (U_1 \otimes U_2)(\Lambda_1 \otimes \Lambda_2)(U_1 \otimes U_2)^\top = U \Lambda U^\top \quad (2.34)$$

in which we denote $U = U_1 \otimes U_2$ and $\Lambda = \Lambda_1 \otimes \Lambda_2$. W now becomes

$$\text{vec}(W) = (U \Lambda U^\top + \alpha \mathbb{I})^{-1} \text{vec}(Y)$$

Further exploring the orthogonality of U , $U U^\top = U \mathbb{I} U^\top = \mathbb{I}$, so that

$$\text{vec}(W) = U(\Lambda + \alpha \mathbb{I})^{-1} U^\top \text{vec}(Y)$$

The most crucial property of the Kronecker product for our application is its relationship with the vectorization operator $()$: $(A \otimes B) \text{vec}(C) = \text{vec}(B C A^\top)$, which allows us to write

$$\text{vec}(W) = U(\Lambda + \alpha \mathbb{I})^{-1} (U_1^\top \otimes U_2^\top) \text{vec}(Y) = U(\Lambda + \alpha \mathbb{I})^{-1} \text{vec}(U_2^\top Y U_1)$$

Since $(\Lambda + \alpha \mathbb{I})^{-1}$ is diagonal, its multiplication by the vector $\text{vec}(U_2^\top Y U_1)^\top$ can be expressed as a Hadamard product (element-wise multiplication, denoted by \odot) between two vectors. Acting element-wise, the Hadamard product is unaffected by vectorization, so that we can simplify the above expression by employing the matrix

$$(\alpha + \lambda_1 \otimes \lambda_2^\top)^{\circ-1[ij]} = \frac{1}{\alpha + \lambda_1^{[i]} \lambda_2^{[j]}} \quad (2.35)$$

In this context, $\lambda_1 \otimes \lambda_2^\top$ represents the n_1 by n_2 matrix resulting from the *outer product* of the vectors containing the eigenvalues of S_1 and S_2 , respectively, while $A^{\circ-1}$ denotes the Hadamard inverse, corresponding to the matrix formed by taking the reciprocal of each element in A . Thus,

$$\begin{aligned} \text{vec}(W) &= U \text{diag}(\Lambda + \alpha^{-1} \mathbb{I}) \odot \text{vec}(U_2^\top Y U_1) = \\ &= (U_1 \otimes U_2) \text{vec}((\alpha + \lambda_1 \otimes \lambda_2^\top)^{\circ-1} \odot (U_2^\top Y U_1)) = \\ &= \text{vec}(U_2 [(\alpha + \lambda_1 \otimes \lambda_2^\top)^{\circ-1} \odot (U_2^\top Y U_1)] U_1^\top) \end{aligned}$$

Which yields

$$W = U_2 [(\alpha + \lambda_1 \otimes \lambda_2^\top)^{\circ-1} \odot (U_2^\top Y U_1)] U_1^\top \quad (2.36)$$

Finally, predictions for a new group of instances in the test set are obtained as follows from the similarities with the train instances ($S_{1, \text{test}}^{[ij]}$ specifies the similarity between $X_{1, \text{test}}^{[i]}$ and $X_{1, \text{train}}^{[j]}$).

$$\text{vec}(\hat{Y}_{\text{test}}) = (S_{1, \text{test}} \otimes S_{2, \text{test}}) \text{vec}(W) = \text{vec}(S_{2, \text{test}} W S_{1, \text{test}}^\top) \quad (2.37)$$

which summarizes to

$$\hat{Y}_{\text{test}} = S_{2, \text{test}} U_2 [(\alpha + \lambda_1 \otimes \lambda_2^\top)^{\circ-1} \odot (U_2^\top Y U_1)] U_1^\top S_{1, \text{test}}^\top \quad (2.38)$$

2.4.2 Neighborhood-Regularized Matrix Factorization

This section describes the NRLMF (?) algorithm, which is a matrix factorization technique that incorporates similarity information between samples in the learning process. It enables latent features of new instances to be inferred from the features of their neighbors. The method is an extension of the Logistic Matrix Factorization (LMF) algorithm, which is designed for recommendation tasks (?). NRLMF was successfully applied to diverse bipartite interaction prediction problems, such as drug-target (?), long non-coding RNA-protein (?), and lncRNA-microRNA (?).

2.4.2.1 Traditional matrix factorization

The idea behind matrix factorization is to find an approximation of the interaction matrix Y by decomposing it into two matrices U and V of lower dimensions, such that

$$\hat{Y} = UV^T \approx Y \quad (2.39)$$

The rows of U and V can be seen as learned representations in a new vector space (usually called *latent* space) of each sample in the row and column domains, so that $U^{[i]}$ represents $X_1^{[i]}$ and $V^{[j]}$ represents $X_2^{[j]}$. Notice that the number of latent features is constrained by Equation 2.39 to be the same for both instance domains: $|U|_j = |V|_j$, being an arbitrary hyperparameter to be defined by the user. Usually, the number of latent features is set to be much smaller than the number of original features ($|U|_j \ll |X_1|_j$ and $|V|_j \ll |X_2|_j$), requiring less computational labor and generating models less susceptible to overfitting.

The learning procedure of matrix factorization algorithms thus consists of obtaining such latent feature matrices U and V so to approximate Y as best as possible. The most common approach is to define a loss function that penalizes the difference between the predicted and the true values of Y and to employ gradient descent techniques to gradually change U and V in the direction that minimizes such loss.

As can be deduced from Equation 2.39, the dot product of each row of U and V approximates the corresponding element of Y :

$$\hat{Y}^{[ij]} = U^{[i]} \cdot V^{[j]} \quad (2.40)$$

Logistic matrix factorization (LMF) slightly redefines the problem by introducing one more step to obtain \hat{Y} from U and V (?): it assumes that the interaction matrix Y is the result of the logistic function applied to UV^T and not only UV^T anymore.

$$\hat{Y}^{[ij]} = \text{logistic}(U^{[i]} \cdot V^{[j]}) = \frac{\exp(U^{[i]} \cdot V^{[j]})}{1 + \exp(U^{[i]} \cdot V^{[j]})} \quad (2.41)$$

where $\mathbf{a} \cdot \mathbf{b}$ represents the dot product between the vectors \mathbf{a} and \mathbf{b} . When applied to a matrix, we assume that \log and \exp functions operate element-wise:

$$(\log M)^{[ij]} = \log M^{[ij]}$$

If \hat{Y} is interpreted as the probability of being positive as predicted by the model ($\hat{Y}^{[ij]} = P(Y^{[ij]} = 1)$ and $1 - \hat{Y}^{[ij]} = P(Y^{[ij]} = 0)$), the optimization objective is usually based on maximizing the joint probability of correctly guessing all interactions in Y :

$$P(\hat{Y} = Y) = \prod_{ij} |\hat{Y}^{[ij]} + Y^{[ij]} - 1| \quad (2.42)$$

in which $|a|$ represents the absolute value of a . A few modifications are then further made:

1. The logarithm of the objective function is taken to simplify the expression without affecting the optimization problem, since \log is a monotonic function;
2. Since positive interactions are usually less numerous but more important in matrix completion problems, a factor α is introduced to prioritize them, multiplying the terms corresponding to $Y^{[ij]} = 1$ in the objective function. It results as if *alpha* copies of each positive interaction are present in the training set;
3. To discourage overfitting and avoid U and V being arbitrarily large, quadratic regularization terms are added, penalizing the magnitude of the elements of U and V .
4. Similarity information between samples is incorporated by NRLMF,

To simplify notation, we define matrices J whose combined sum of all elements corresponds to the objective function \mathcal{J} :

$$\mathcal{J} = \sum J_{\text{labels}} + \sum J_{1, \text{regularization}} + \sum J_{2, \text{regularization}} + \sum J_{1, \text{neighborhood}} + \sum J_{2, \text{neighborhood}} \quad (2.43)$$

where by $\sum J$ we denote $\sum_i^{|J|_i} \sum_j^{|J|_j} J$. The sums must be performed individually due to the J matrices having different dimensions.

Applying the logarithm to Equation 2.42 we have our first term of \mathcal{J} :

$$J_{\text{labels}} = \log |\hat{Y} + Y - 1| = Y \odot \log \hat{Y} + (1 - Y) \odot \log(1 - \hat{Y})$$

where we separate the cases in which $Y^{[ij]} = 1$ and $Y^{[ij]} = 0$. \odot represents the Hadamard product (element-wise multiplication) between matrices:

$$(A \odot B)^{[ij]} = A^{[ij]} B^{[ij]}$$

Adding the aforementioned positive importance factor α and expanding \hat{Y} according to Equation 2.41 we have

$$\begin{aligned} J_{\text{labels}} &= \alpha Y \odot \{UV^\top - \log[1 + \exp(UV^\top)]\} - (1 - Y) \odot \log[1 + \exp(UV^\top)] = \\ &= \alpha Y \odot UV^\top + [(1 - \alpha)Y - 1] \odot \log[1 + \exp(UV^\top)] \end{aligned} \quad (2.44)$$

To discourage large values in U and V , we consider quadratic regularization terms weighted by hyperparameters λ_1 and λ_2 :

$$J_{1, \text{regularization}} = -\frac{\lambda_1}{2} U \odot U = -\frac{\lambda_1}{2} \mathbb{I} \odot (UU^\top) \quad (2.45)$$

$$J_{2, \text{regularization}} = -\frac{\lambda_2}{2} V \odot V = -\frac{\lambda_2}{2} \mathbb{I} \odot (VV^\top) \quad (2.46)$$

If the initial objective proposal can be interpreted as maximizing the probability of guessing all labels correctly given specific U and V ($\sum J_{\text{labels}} = \log P(\hat{Y} = Y \mid U, V)$), adding the regularization terms is equivalent to introduce prior assumptions about U and V distributions and define a slightly different objective: maximizing the posterior probability to obtain the current U and V given that $\hat{Y} = Y$. Applying Bayes' theorem and assuming $P(Y) = 1$ we have

$$P(U, V \mid \lambda_1, \lambda_2, Y) = P(Y \mid U, V) P(U \mid \lambda_1) P(V \mid \lambda_2) \quad (2.47)$$

Under the assumption that the values in U and V follow zero-centered spherical gaussian distributions with variances given by $\frac{1}{\lambda}$, that is, $U^{[ij]} \sim \mathcal{N}(0, \lambda_1^{-1} \mathbf{I})$ and $V^{[ij]} \sim \mathcal{N}(0, \lambda_2^{-1} \mathbf{I})$, we recover the regularized objective function of Equation 2.46 (?).

$$\begin{aligned} \log P(U, V \mid \lambda_1, \lambda_2, Y) &= \\ &= \log P(Y \mid U, V) + \sum \log(\exp(-\frac{\lambda_1}{2} U \odot U)) + \sum \log(\exp(-\frac{\lambda_2}{2} V \odot V)) = \\ &= \sum J_{\text{labels}} - \sum \frac{\lambda_1}{2} U \odot U - \sum \frac{\lambda_2}{2} V \odot V \end{aligned} \quad (2.48)$$

Therefore, if multiple values of U and V possibly generate the same $\hat{Y} = Y$, applying the regularization can be understood as not only finding one of such combinations but, between those U and V that satisfy $\hat{Y} = Y$, finding the U and V that are most likely to be randomly sampled. If $P(U, V \mid \lambda_1, \lambda_2, Y)$ continuously varies as a function of U and V , pooling U and V from a region of maximal $P(U, V \mid \lambda_1, \lambda_2, Y)$ should improve generality, arguably ensuring that stochastic deviations of U and V would still result in high $P(\hat{Y} = Y)$ and justifying the use of regularization as a way to avoid overfitting.

One may have noticed that the original feature matrices X_1 and X_2 were not considered in any regard when describing matrix factorization and detailing the objective functions. Born in the context of recommendation systems where the relationship labels are usually the only information available, matrix factorization algorithms in general encounter a significant issue when brought to our current scenario of bipartite interaction prediction: in its canonic formulation, they do not take sample-level features into account, often called *side information* or *side features* in the recommendation field (), possibly overlooking valuable data. As a consequence, they are unable to provide predictions for new samples that were not present in the training set, since no information about them is available to be inputted to the model. This is commonly regarded as the *cold-start problem* (). As such, matrix factorization usage is usually restricted to the task of matrix completion, in which the goal is to predict the missing values of a matrix

given the values of the remaining elements, without receiving completely new rows or columns during model evaluation ().

2.4.2.2 Neighborhood regularization

Targeting these issues, liu2016neighborhood, liu2017lpinrlmf, liu2020predicting proposes a modification to LMF that incorporates side information into the model, enabling predicting interactions for new samples. The core idea of their technique lies in adding one more term to the objective function, penalizing instances regarded as close when considering the original features but were separated by U and V , placed far from each other in the latent space. As such, the algorithm is called *Neighborhood-Regularized Logistic Matrix Factorization* (NRLMF). To precisely define it, let's consider similarity-weighted adjacency matrices A_1 and A_2 referring to each sample domain that specifies neighborhood relationships between samples. If $S_1^{[ij]}$ denotes a similarity score between $X_1^{[i]}$ and $X_1^{[j]}$, $A_1^{[ij]}$ is set to this similarity value if $X_1^{[j]}$ belongs to the neighborhood of $X_1^{[i]}$, denoted $N(X_1^{[i]})$, and 0 otherwise:

$$A_1^{[ij]} = \begin{cases} S_1^{[ij]} & \text{if } X_1^{[j]} \in N(X_1^{[i]}) \\ 0 & \text{otherwise} \end{cases} \quad (2.49)$$

Multiple options are available for the definition of neighborhoods, such as considering all samples within a certain radius of each other or only the k nearest neighbors of each sample. In this work, following the original proposal of NRLMF (?), we will consider the latter, defining A_1 and A_2 as the adjacency matrices of the k -nearest neighbors graphs of X_1 and X_2 , respectively. In other words, $N(X_1^{[i]})$ is the set formed by the k rows $X_1^{[j]}$ with the k highest s_{ij} . In the interaction prediction problems we analyze, similarities are precalculated so that the X matrices directly provide the distance metric over which the nearest neighbors are selected. That is, X_1 and X_2 themselves already constitute pairwise similarity matrices: $S_1^{[ij]} = X_1^{[ij]}$. In general, however, one may need to define a kernel matrix S as a preprocessing step, choosing a distance metric over the original features to be used in the nearest neighbors search such as the Euclidean distance or a radial basis function (?). In any case, notice that A is a function of X alone for NRLMF. While Y may also be considered in similar scenarios (as will be discussed ahead), A does not depend on U or V and can be built as a single pre-training step.

The loss term proposed by NRLMF is then given by the sum of the Euclidean distances in the latent space between samples in the same neighborhood, weighted by their similarities:

$$\mathcal{J}_{1, \text{neighborhood}} = - \sum_{ij} A_1^{[ij]} \|U^{[i]} - U^{[j]}\|^2 \quad (2.50)$$

$$\mathcal{J}_{2, \text{neighborhood}} = - \sum_{ij} A_2^{[ij]} \|V^{[i]} - V^{[j]}\|^2 \quad (2.51)$$

in which $\|\mathbf{v}\|$ represents the Euclidean norm of a vector \mathbf{v} . Concentrating on the row instances

and expanding the last definition we have

$$\begin{aligned} \sum_{ij} A_1^{[ij]} \|U^{[i]} - U^{[j]}\|^2 &= \sum_{ij} A^{[ij]} (U^{[i]} \cdot U^{[i]} + U^{[j]} \cdot U^{[j]} - 2U^{[i]} \cdot U^{[j]}) = \\ &= \sum_i \left(\sum_j A_1^{[ij]} \right) U^{[i]} \cdot U^{[i]} + \sum_j \left(\sum_i A_1^{[ij]} \right) U^{[j]} \cdot U^{[j]} \end{aligned}$$

The terms in which U appears with the same index ($U^{[i]} \cdot U^{[i]}$ and $U^{[j]} \cdot U^{[j]}$) can be rewritten to include both by multiplying them by the identity matrix \mathbb{I} . Essentially, we consider $\sum_i U^{[i]} \cdot U^{[i]} = \text{trace}(UU^\top) = \sum_{ij} (\mathbb{I} \odot UU^\top)^{[ij]}$.

$$\sum_i \left(\sum_j A_1^{[ij]} \right) U^{[i]} \cdot U^{[i]} = \sum_i \left(\sum_j A_1^{[ij]} \right) \sum_k \mathbb{I}^{[ik]} U^{[i]} \cdot U^{[k]} = \sum_{ij} \left(\sum_k A_1^{[ik]} \right) \mathbb{I}^{[ij]} U^{[i]} \cdot U^{[j]}$$

This allows us to write

$$\begin{aligned} \sum_{ij} A_1^{[ij]} \|U^{[i]} - U^{[j]}\|^2 &= \\ &= \sum_{ij} \left[\left(\sum_k A_1^{[ik]} + \sum_l A_1^{[lj]} \right) \mathbb{I}^{[ij]} - 2A_1^{[ij]} \right] U^{[i]} \cdot U^{[j]} = \\ &= \sum L_1 \odot (UU^\top) \end{aligned}$$

in which we define

$$L_1^{[ij]} = \left(\sum_k A_1^{[ik]} + \sum_l A_1^{[lj]} \right) \mathbb{I}^{[ij]} - 2A_1^{[ij]} = \left(\sum_k A_1^{[ik]} + A_1^{[kj]} \right) \mathbb{I}^{[ij]} - 2A_1^{[ij]} \quad (2.52)$$

Notice that taking A_1^\top instead of A_1 has no effect on the final result, since $\sum A_1 UU^\top = \sum A_1^\top UU^\top$. We could then work with a symmetrized version of A_1 from the start:

$$\tilde{A}_1 = A_1 + A_1^\top$$

yielding

$$L_1 = \left(\sum_k \tilde{A}_1^{[ik]} \right) \mathbb{I}^{[ij]} - \tilde{A}_1^{[ij]}$$

We can see that the first term of L_1 acts in a similar way to the quadratic regularization terms presented by Equation 2.46, multiplying the main diagonal of UU^\top and thus penalizing the model for latent vectors with large Euclidean norms (the diagonal of UU^\top holds the squared norms $U^{[i]} \cdot U^{[i]}$). The amount of regularization is however pondered by the weighted number of neighbors of each sample in this case: $\sum_k A_1^{[ik]}$ represents the sum of similarities of $X_1^{[i]}$ with all its neighbors, also called the *degree* of a sample. This results in samples with larger and more compact neighborhoods being more heavily penalized for having large norms in the latent space. The second term of Equation 2.52, on the other hand, rewards the model for placing close

neighbors colinear to each other in the latent space, summing over $S^{[ij]}U^{[i]} \cdot U^{[j]}$ terms between each sample and its neighbors ($A_1^{[ij]}$ is 0 if $U^{[i]}$ and $U^{[j]}$ are not neighbors).

Finally, the neighborhood regularization terms are written as

$$J_{1, \text{neighborhood}} = -\frac{\beta_1}{2} L_1 \odot (UU^\top) \quad (2.53)$$

$$J_{2, \text{neighborhood}} = -\frac{\beta_2}{2} L_2 \odot (VV^\top) \quad (2.54)$$

Combining the matrix terms as in Equation 2.43, NRLMF's objective function is given by

$$\begin{aligned} \mathcal{J} = & \sum \alpha Y \odot UV^\top + [(1 - \alpha)Y - 1] \odot \log [1 + \exp(UV^\top)] \\ & - \sum \frac{1}{2} (\lambda_1 \mathbb{I} + \beta_1 L_1) \odot (UU^\top) \\ & - \sum \frac{1}{2} (\lambda_2 \mathbb{I} + \beta_2 L_2) \odot (VV^\top) \end{aligned} \quad (2.55)$$

and the derivatives of the objective function with respect to U and V to be used in the gradient descent procedure are given by

$$G_U = \frac{\partial \mathcal{J}}{\partial U} = \{[(1 - \alpha)Y - 1] \odot \hat{Y} + \alpha Y\}V - (\lambda_1 \mathbb{I} + \beta_1 L_1)U \quad (2.56)$$

$$G_V = \frac{\partial \mathcal{J}}{\partial V} = \{[(1 - \alpha)Y - 1] \odot \hat{Y} + \alpha Y\}^\top U - (\lambda_2 \mathbb{I} + \beta_2 L_2)V \quad (2.57)$$

The training procedure of NRLMF is presented by function TrainNRLMF, and consists of alternated updates on U and V in the gradient's direction until certain stop criteria are satisfied. Common choices for stopping conditions are a maximum number of iterations or a minimum change in the objective function between iterations.

Since faster convergence is reported by the original authors (?), we follow johnsonlogistic,liu2016neighborhooda by implementing the AdaGrad procedure (?), in which the length of each gradient step is divided by the square-root sum of squared previous steps:

$$U_{t+1} = U_t + \frac{\eta G_{U,t}}{\sqrt{\sum_{t'=0}^t G_{U,t'}^2}} \quad (2.58)$$

where $G_{U,t}$ is the partial derivative $\frac{\partial \mathcal{J}}{\partial U}$ of the objective function with respect to U at step t , and η is the user-defined learning rate. The same is done for V .

The main importance of NRLMF however lies in the inference phase. As mentioned, matrix factorization methods are not designed to deal with new input samples, that are not present in the training set. Specifically, traditional matrix factorization is incapable of generating latent vectors for the unseen samples to be used for label prediction. An idea that may seem natural at first glance is to delay training to the arrival of new instances, including them in the training set with zero-only labels before performing the optimization. But even then, using an

Function TrainNRLMF($Y, S_1, S_2, \alpha, \lambda_1, \lambda_2, \beta_1, \beta_2, \eta$): Train an NRLMF model.

Input: Y , the training labels matrix to be approximated.
Input: S_1, S_2 , Similarity matrices among instances of each axis.
Input: α , the positive importance factor.
Input: λ_1, λ_2 , quadratic regularization factors.
Input: β_1, β_2 , neighborhood regularization factors.
Input: η , the learning rate.
Output: U, V , the resulting latent feature matrices.

- 1 Optionally precompute constant factors of the gradient (Equations 2.56 and 2.57), such as $(\lambda_2 \mathbb{I} + \beta_2 L_2)$, $[(1 - \alpha)Y - 1]$ or αY ;
- 2 Initialize U and V with normally-distributed random values;
- 3 $T_U, T_V \leftarrow \mathbf{0}$; // Initialize gradient accumulators
- 4 **while** *Stop conditions are not met* **do**
 - // Update U
 - 5 Obtain G_U through Equation 2.56;
 - 6 $T_U \leftarrow T_U + G_U^2$;
 - 7 $U \leftarrow U + \eta \frac{G_U}{\sqrt{T_U}}$;
 - // Update V
 - 8 Obtain G_V through Equation 2.57;
 - 9 $T_V \leftarrow T_V + G_V^2$;
 - 10 $V \leftarrow V + \eta \frac{G_V}{\sqrt{T_V}}$;
- 11 **end**
- 12 **return** U, V

objective function based only on Y as is traditionally done, no new information is brought by the new instances and adding new zeroed rows or columns to Y will mainly introduce noise to the training data and likely degrade the model's predictive performance.

NRLMF, however, leverages proximity information encoded by X to remarkably enable determining latent feature vectors for completely new instances. The neighborhood regularization terms in the objective function now reveal their full importance: they support proximity as a transferable property between the original and the latent spaces. By encouraging that neighbors in X_1 and X_2 remain close in U and V , we can infer latent features of new instances based on their neighborhood.

Consider the test similarity matrices $S_{1, \text{test}}$ and $S_{2, \text{test}}$ respectively derived from $X_{1, \text{test}}$ and $X_{2, \text{test}}$, relating the new instances to the known training samples. For instance, $S_{1, \text{test}}^{[ij]}$ represents the similarity between $X_{1, \text{test}}^{[i]}$ and $X_{1, \text{train}}^{[j]}$. If $A_{1, \text{test}}$, like before in Equation 2.49, accordingly restricts the similarity matrix to the neighborhood of each sample,

$$A_{1, \text{test}}^{[ij]} = \begin{cases} S_{1, \text{test}}^{[ij]} & \text{if } X_{1, \text{train}}^{[j]} \in N(X_{1, \text{test}}^{[i]}) \\ 0 & \text{otherwise} \end{cases} \quad (2.59)$$

the latent feature vector of a new instance is simply estimated as the weighted average of its neighbors' latent representations:

$$U_{\text{test}}^{[i]} = \frac{A_{1, \text{test}}^{[i]} U_{\text{train}}}{\sum A_{1, \text{test}}^{[i]}} \quad (2.60)$$

the analogous being held for V , so that new predictions are made as usual with Equation 2.41, where U_{test} , U_{train} , V_{test} and V_{train} can be used in accordance with the prediction task under study (see Section 2.5.3 for details on the different prediction scenarios):

$$\hat{Y}_{\text{TT}} = \frac{\exp(U_{\text{test}} V_{\text{test}}^{\top})}{1 + \exp(U_{\text{test}} V_{\text{test}}^{\top})} \quad \hat{Y}_{\text{TL}} = \frac{\exp(U_{\text{test}} V_{\text{train}}^{\top})}{1 + \exp(U_{\text{test}} V_{\text{train}}^{\top})} \quad \hat{Y}_{\text{LT}} = \frac{\exp(U_{\text{train}} V_{\text{test}}^{\top})}{1 + \exp(U_{\text{train}} V_{\text{test}}^{\top})} \quad (2.61)$$

2.5 Assessing the performance of bipartite models

In this section, we discuss the evaluation procedure of estimators in bipartite learning settings. The PU property of such datasets and the presence of two instance domains pose special considerations on how we assess the performance of our models. We start by presenting the datasets used in this study, followed by a description of the evaluation metrics and the cross-validation strategy we utilize.

2.5.1 Datasets

COMPLEMENTAR

We gathered ten publicly available interaction datasets to evaluate the performance of the proposed models. Quantitative information about each of them is presented by Table 2, and more detailed descriptions are provided in this section.

2.5.1.1 DPI-E, DPI-G, DPI-I, DPI-N

These datasets comprise drug-protein interactions for four distinct classes of proteins: enzymes, GPCRs, ion channels, and nuclear receptors, respectively. Drug similarities were computed using the SIMCOMP metric, while protein similarities were computed as normalized scores of Smith-Waterman pairwise alignments (?).

2.5.1.2 ERN and SRN

The datasets represent interactions between genes and transcription factors in *E. coli* and *S. cerevisiae*, respectively. Features for genes and transcription factors are initially composed of experimentally measured expression levels and, in SRN, gene motif features (?, ?). We compute the RBF kernel of such values to obtain the final similarity matrices.

2.5.1.3 DAVIS

The DAVIS dataset contains experimentally measured drug-kinase dissociation constants (?). The dataset was binarized by considering interactions with dissociation constants

Dataset	Type of interaction	Y shape	Density
DPI-E	Drug-enzyme	664×445	0.9902%
DPI-G	Drug-GPCR	95×223	2.997%
DPI-G	Drug-GPCR	95×223	2.997%
DPI-I	Drug-ion channel	204×210	3.445%
DPI-N	Drug-nuclear receptor	26×54	6.410%
ERN	Gene-transcription factor	1164×154	1.837%
SRN	Gene-transcription factor	1821×113	1.780%
DAVIS	Inhibitor-kinase	68×442	5.011%
KIBA	Inhibitor-kinase	2111×229	19.74%
NPInter	lncRNA-protein	586×446	18.12%
mirTarBase	miRNA-mRNA	1873×415	7.065%

Table 2 – Summary of the datasets used in this study. The similarity scores for mirTarBase and NPInter were obtained from the raw sequences as their normalized Smith-Waterman alignment scores. Original references are: DPI (?), ERN (?), SRN (?), DAVIS (?), KIBA (?), NPInter (?), mirTarBase (?). DAVIS and KIBA are provided by huang2020deeppurpose. DPI datasets, SRN, and ERN are provided by schrynemackers2015. Source: made by the authors.

$\leq 30nM$ as the positive ones, as suggested by (?). Drug similarities were computed using the Extended Connectivity Fingerprints (ECFP4) (?, ?) while protein similarities were taken as the normalized Smith-Waterman score (?, ?).

2.5.1.4 KIBA

The KIBA dataset was initially built by ??(?) and contains experimentally verified affinity scores between kinase and kinase inhibitors.

(?) further processed the dataset by removing all drugs and targets with less than 10 observations. In alignment with (?, ?), we consider positive interactions as those with \log_{10} KIBA-scores ≤ 3.0 to reframe the task as binary classification.

The utilized version of the dataset with corresponding amino acid sequences and SMILES representations were provided by (?). From them, we generated the protein similarity matrix using the same procedure employed in the preprocessing of NPInter proteins. The drug similarities were computed similarly to how (?) processed the DAVIS dataset, using the Tanimoto distances of ECFP4 fingerprints (?, ?). The Python library `rdkit` (?) was used to this calculation.

2.5.1.5 mirTarBase

The mirTarBase dataset contains experimentally validated microRNA-messengerRNA interactions. MicroRNA sequences were obtained from miRBase (?) while transcript sequences were obtained from GENCODE (?). The longest transcript for each gene was selected and the 3' UTR exonic sequences were recovered from the genome and annotation files provided by GENCODE. The similarity matrices were then built from the normalized Smith-Waterman (?)

alignment scores among microRNAs and among the genes' 3' UTRs. The alignments were performed using the BLASTN substitution matrix and no gap penalty, with the help of the Biopython package (?).

Each miRNA was required to have at least 10 interactions in the dataset, and each gene was required to have at least 100 interactions.

2.5.1.6 NPInter

Interactions between long non-coding RNAs (lncRNA) and proteins were recovered from NPInter (?). The lncRNA sequences were obtained from NONCODE (?) and the protein sequences were obtained from UniProt (?). The similarity matrices were built from the normalized Smith-Waterman (?) alignment scores among lncRNAs and among the proteins. Similarly to the preprocessing of mirTarBase, we utilized the Biopython package (?) to perform the alignments, using the BLASTN and BLOSUM62 substitution matrices for the lncRNA and protein alignments, respectively, and no gap penalty in both cases.

Each lncRNA was required to interact with 50 proteins or more to be incorporated in the dataset, and each protein was required to have at least 2 interactions.

2.5.2 There are multiple ways to measure model generality

To evaluate machine learning models, the standard procedure consists of separating a subset of data samples not to be used in the training process. These samples are subsequently inputted to the trained model and its known labels are compared to the model's predictions in order to estimate the algorithm performance. The held-out samples are collectively called the *test set* while the ones used for model building are called the *training set*.

Having two distinct sample groups in bipartite learning settings makes the concept of a held-out set more nuanced. We list five reasonable test configurations we encounter in the literature. Given two samples x_1 and x_2 from each respective domain, a "test set" could refer to:

1. **Test x_1 , test x_2 (TT):** both x_1 and x_2 are not present in the training set, the model has never seen either of them before;
2. **Learned x_1 , test x_2 (LT):** x_1 is part of the training feature matrix $X_{1, \text{train}}$, but x_2 was never seen before;
3. **Test x_1 , learned x_2 (TL):** x_2 is part of the training feature matrix $X_{2, \text{train}}$, but x_1 was never seen before;
4. **Learned x_1 , learned x_2 , masked label (LL-M):** Some of the positive annotations in the training set are randomly masked (replacing 1 by 0). The model is trained on the masked dataset and evaluated based on the predictions for the masked positives and the negative annotations.

5. **Learned x_1 , learned x_2 , unknown label (LL-U):** Models trained with the SGSO adaptation (Section 2.2) can exclude specific labels from the training set. All instances are still used for training, but some combinations of them are not presented to the model. This corresponds to applying standard validation procedures directly on the X_{SGSO} and Y_{SGSO} matrices (Equation 2.3).

Accordingly, we sometimes call the training set the *LL set*. LL, TL, LT, and TT were already presented with similar names by (?) (?) and (?). The authors call them $L_r \times L_c$, $T_r \times L_c$, and so on. (?, ?) denotes LL-U, LT, TL and TT as C1, C2, C3 and C4, respectively. (?) uses the same notation as (?) but uses C4 to designate LL-M instead of LL-U. The LL-U test set is most frequently used when building SGSO models (). (?) also reports TL and LT scores. Using LL-M is more common for matrix factorization methods ().

Note that a model that performs well on TT is likely to also perform reasonably well on the other test sets, but the opposite is not necessarily true. For instance, performing well on the TL or LT sets does not mean a model has learned information about both domains. As an example, consider a drug-protein scenario where the model is tested with new drugs and known proteins. Given a specific known protein, a model might correctly infer interactions based on its knowledge that similar drugs have interacted with that protein, or that a specific drug trait correlates well with the protein’s interactome. No knowledge about the protein’s features is required for the model to make such predictions. Therefore, there is no guarantee that the model will also perform well for new proteins.

Similarly, an effective model according to the LL-U or LL-M settings will not necessarily perform well on any of TT, LT or TL sets. The reason is: an estimator would still be able to perform significantly well on LL-U or LL-M using only the information from the interaction matrix (as matrix factorization algorithms do). There is no guarantee that the model will be a representation of the underlying function: $(\mathbf{x}_1, \mathbf{x}_2) \mapsto y$. As such, there is no guarantee that new instances will be correctly classified, since the only information available about them would be their feature vectors.

On the other hand, the TT test sets are the most informative in terms of generalization. They require the model to fully rely on both input vectors to predict labels, providing a more reliable representation of how attributes from both objects come together to determine the interaction. TT is therefore especially useful when one is interested in gaining insight about the underlying process that generates the interactions. However, they are also the most challenging to predict, and the model’s performance on them is expected to be lower than on the other test sets.

In summary, TT should be preferred when the main interest is modeling how input features determine the interaction. TL and LT are useful if generalization is only interesting for one of the domains. LL-U and LL-M are useful when the main goal is to obtain predictions for

a defined set of instances, regardless of modeling the function $(\mathbf{x}_1, \mathbf{x}_2) \mapsto y$ or not.

For the learning problems under study, explainability is an important factor to be explored in the future. We thus report our results in terms of TT, LT and TL test sets. Additionally, the learning tasks that we approach represent very diverse phenomena (Section 2.5.1), which makes the specific LT or TL sets not comparable between datasets. We thus average both LT and TL results together to report a single LT+TL score for each dataset.

2.5.3 Cross-validating in two dimensions

In traditional learning tasks, k -fold cross-validation consists in equally and randomly dividing both X and Y together in k non-overlapping partitions (or folds). The model is then evaluated k times, in each round selecting a fold as the test set and the remaining ones together as the training set.

In the bipartite interaction setting, fold division can be done in each of the two axis of the interaction matrix, corresponding to each of the two X_a sample domains. Each of the k_1 folds of X_1 can be combined with one of the k_2 folds of X_2 , resulting in a corresponding partition of Y . Each bidimensional fold can be used as the TT set in a CV round, which yields the corresponding four LL, LT, TL and TT sets for each round. A k_1 by k_2 bidimensional CV then has a total of $k_1 k_2$ folds.

We can optionally enforce independent test sets, so that each dyad composes only one bidimensional fold. This can be done by selecting $k = k_1 = k_2$ and pairing each X_1 partition with a single X_2 partition, yielding a total of k folds, not k^2 as before. We refer to this procedure as *diagonal* cross-validation.

In our settings, we average the results of all folds of each dataset before computing test statistics. Therefore, the number of folds and the independence between test folds will not affect the final significance estimates. We thus use all k^2 folds in our experiments. Diagonal cross-validation is still used, but for the nested validation used for parameter tuning in some experiments.

An important note is that the similarity matrices must be thoughtfully handled when splitting the datasets: we must not include similarities with the test samples in the training set.

2.5.4 Prediction scoring metrics

This section is concerned with defining the two metrics used throughout this work to evaluate the predictive performance of an estimator and enable comparison between them.

Consider a test set of N interaction labels to be inferred by a classifier (we are not concerned with the shape of Y in this section, and $N = |Y|$). Let the classifier's predictions then be represented by a matrix \hat{Y} of the same shape as Y , with $\hat{Y}^{[ij]}$ being the predicted value

for the ground-truth label $Y^{[ij]}$. Since Y and \hat{Y} are both binary matrices, there are four possible outcomes when a prediction is made, traditionally quantified () as follows:

- **True Positives (TP):** the number of positive labels correctly predicted, where both the predicted and actual labels are 1.

$$TP = \sum_{i,j} \mathbb{I}(Y^{[ij]} = 1 \text{ and } \hat{Y}^{[ij]} = 1) \quad (2.62)$$

- **True Negatives (TN):** the number of negative labels correctly predicted, where both the predicted and actual labels are 0.

$$TN = \sum_{i,j} \mathbb{I}(Y^{[ij]} = 0 \text{ and } \hat{Y}^{[ij]} = 0) \quad (2.63)$$

- **False Positives (FP):** the number of instances where the predicted label is positive (1), but the actual label is negative (0).

$$FP = \sum_{i,j} \mathbb{I}(Y^{[ij]} = 0 \text{ and } \hat{Y}^{[ij]} = 1) \quad (2.64)$$

- **False Negatives (FN):** the number of instances where the predicted label is negative (0), but the actual label is positive (1).

$$FN = \sum_{i,j} \mathbb{I}(Y^{[ij]} = 1 \text{ and } \hat{Y}^{[ij]} = 0) \quad (2.65)$$

where $\mathbb{I}(A)$ is the indicator function that equals 1 if statement A is true and 0 otherwise:

$$\mathbb{I} = \begin{cases} 1 & \text{if } A \\ 0 & \text{otherwise} \end{cases} \quad (2.66)$$

Notice that the sum of TP, TN, FP and FN is equal to the total number of instances T , and we also define $P = TP + FN$, the total number of *a priori* positive labels in the test set, and $N = TN + FP$, the total number of *a priori* negative labels. The total number of predicted positives is termed $PP = TP + FP$ and the remaining predicted negatives are called $PN = TN + FN$.

Those quantities can be organized in the so-called *confusion matrix*, illustrated by ???. Naturally, one wants their estimator to maximize TP and TN while minimizing FP and FN.

Most commonly, we do not use these metrics directly, but instead normalize them to the interval $[0, 1]$ in numerous ways. This enables score comparisons across datasets with different numbers of samples and different densities of positive annotations. Below we list the most common normalized scoring metrics for binary classification problems, from which the metrics used in this work are derived.

- **True positive rate (TPR) or recall:** the ratio of correctly predicted positive labels to the total number of positive labels.

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} \quad (2.67)$$

- **True negative rate (TNR):** the ratio of correctly predicted negative labels to the total number of negative labels.

$$TNR = \frac{TN}{N} = \frac{TN}{TN + FP} \quad (2.68)$$

- **False positive rate (FPR):** the ratio of incorrectly predicted positive labels to the total number of negative labels.

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - TNR \quad (2.69)$$

- **False negative rate (FNR):** the ratio of incorrectly predicted negative labels to the total number of positive labels.

$$FNR = \frac{FN}{P} = \frac{FN}{FN + TP} = 1 - TPR \quad (2.70)$$

- **Precision:** the ratio of correctly predicted positive labels to the total number of predicted positive labels.

$$Pr = \frac{TP}{PP} = \frac{TP}{TP + FP} \quad (2.71)$$

These metrics are better visualized in ??.

An optimal binary classifier will thus present high TPR, TNR and precision while minimizing FPR and FNR.

Notice that each of these metrics still allows trivial solutions: if a classifier outputs positive labels for all instances irrespectively of the input features, it will achieve perfect TPR, FNR and Pr, but its TNR and FPR will be null. On the other hand, if a negative label is outputted every time, the opposite will happen. Thus, we will consider pairs of these metrics simultaneously, such as TPR and TNR. The selected pair must include at least three of TP, TN, FP, or FN to ensure the whole confusion matrix is taken into account.

Minimizing all the four components of the confusion matrix is not always possible. Most often, the learning algorithms are subject to a tradeoff between the ability to correctly infer positive annotations and the ability to correctly infer negative annotations, that we express as a balance between TPR and TNR. When evaluating models, we must be aware of the relative importances being assigned to each of these tendencies. This choice is highly application-specific. For instance, in the process of diagnosing medical conditions, TPR is often prioritized (), favoring the identification of all true cases at the expense of some misleading positive results. In this

case, the cost of missing a positive result is usually far greater than that of a false positive. In other scenarios such as spam email filtering (), TNR might be favored, minimizing the number of legitimate emails marked as spam even if some spam emails go undetected.

Considering this balance across a variety of learning tasks is a challenging factor to be taken into account, especially in cases where the estimators can be easily adjusted to favor each class. Many estimators, such as those employed in the present study, do not directly output a binary label but instead provide us with a continuous decision value (such as a probability of interaction) that additionally depends on a threshold parameter to establish the final predicted classes. Formally, the predicted labels \hat{Y} are obtained from a threshold t applied to the decision values \tilde{Y} :

$$\hat{Y}^{[ij]} = \mathbb{I}(\tilde{Y}^{[ij]} > t) \quad (2.72)$$

The selection of t directly affects the propensity to positive or negative outputs, so that a single model can yield multiple different results with varying levels of TPR and TNR depending on the chosen thresholds. A common practice is then to consider TPR and TNR for all t , avoiding the influence of threshold selection on the estimator comparisons.

Notice that, since a finite set of outputs is considered for evaluation (the test set), a finite set of thresholds will cover all possible classification results of a model. These results can be easily displayed in a two-dimensional plot, using a point for each considered threshold so that its corresponding TPR and TNR are each indicated by an axis. Conventionally, the TPR is plotted in the y axis while the FPR values (representing the TNR, $FPR = 1 - TNR$) are presented as x coordinates, which results in the traditional *receiver operating characteristic* (ROC) curve, exemplified by ???. An ideal threshold of an ideal estimator would then be close to the top-left corner of the plot, where TPR and TNR are both 1. On the other hand, consider a completely random classifier outputting uniformly random values $\hat{Y}_{\text{proba}}^{[ij]}$ in the 0-1 interval. We would have $\hat{Y}^{[ij]} = \mathbb{I}(\hat{Y}_{\text{proba}}^{[ij]} > t)$ for a given threshold t . This implies that the number of correctly guessed positive labels is $1 - t$ (the probability of yielding 1) times the total number of positive labels: $TP = (1 - t)(TP + FN)$, which results in $TPR = 1 - t$ from the definition. Similarly, $FPR = TPR = 1 - t$, so that the ROC curve of a random classifier is a diagonal line from the bottom-left to the top-right corners.

To summarize a classifier's performance across all thresholds, the area under the ROC curve (AUC ROC) is often employed, with values ranging from 0.5 to 1, where 1 represents a perfect classifier and 0.5 represents a random classifier. Although theoretically possible, values below 0.5 would signify the opposite label is being consistently predicted, most likely indicating misconfiguration of the estimator. If $AUROC < 0.5$, the result can be easily converted to a value greater than 0.5 by simply inverting the predicted labels (turning 0s into 1s and vice-versa).

Despite considerably common, the use of ROC curves requires additional considerations when dealing with heavily imbalanced classification datasets (where some classes are greatly overrepresented) (?, ?, ?). For instance, consider the case where $NEG \gg POS$. Since the de-

nominator of TNR is far greater than the denominator of TPR, a change in TN (for example, missing one more negative label) will have a much smaller impact on TNR than a change in TP (missing a positive label) will have on TPR. Specifically, an increase of k in TN will cause an increase in TNR $\frac{NEG}{POS}$ times greater than the increase in TPR caused by the same k increase in TP.

Given that ROC equally considers both metrics, classifiers that are more sensitive to positive labels will arguably be favored over those prioritizing negative outputs. By the same logic, ROC will also be more lenient towards false positives (that decrease TNR) rather than false negatives (that decrease TPR) ().

To address this issue, it is commonly suggested the usage of precision-recall (PR) curves instead (?, ?), where the precision is plotted as the vertical coordinate while the *recall* (another name for TPR) is represented horizontally. We can see the FP term in ?? as the proxy for the true negatives (FP = N - TN). Notice that FP in the definition of precision (??) is not divided by the total number of negative labels, as TN is in TNR. Thus, it is argued that AUPR is less likely than AUROC to prioritize the minority class in imbalanced scenarios (?, ?).

We explore these claims in further detail in the following section, formally defining the ROC and PR curves in terms of ideal label probability distributions.

2.5.4.1 Ideal descriptions of AUROC and AUPRC

Consider a general estimator outputting a decision value $s \in \mathbb{R}$ for each input instance. The final class to be assigned is still to be defined by a threshold s^* so that $\hat{Y}^{[ij]} = \mathbb{I}(s^{[ij]} > s^*)$. For the ideal case of an infinite number of test samples, the possible scoring results of the estimator would be fully determined by the two theoretical distributions of s given the true label $Y^{[ij]}$, i.e. the probability density functions $P(s \mid Y^{[ij]} = 1)$ and $P(s \mid Y^{[ij]} = 0)$. Similar to hand2009measuring, we thus define the probability density functions f_k and their corresponding cumulative distribution functions F_k for each of the two classes $k \in \{0, 1\}$:

$$f_0(s) = P(s \mid Y^{[ij]} = 0) \quad (2.73)$$

$$f_1(s) = P(s \mid Y^{[ij]} = 1) \quad (2.74)$$

$$F_0(s) = \int_{-\infty}^s f_0(u) du \quad (2.75)$$

$$F_1(s) = \int_{-\infty}^s f_1(u) du \quad (2.76)$$

We further define $p = \frac{P}{T} = P(Y_{\text{test}}^{[ij]} = 1)$ and $n = \frac{N}{T} = P(Y_{\text{test}}^{[ij]} = 0)$, the fractions of positive and negative labels in the test set, respectively. We can then express the expected values of the

confusion matrix as functions of a given threshold s^* for the decision value s :

$$\text{TP}(s^*) = p(1 - F_1(s^*)) \quad (2.77)$$

$$\text{TN}(s^*) = nF_0(s^*) \quad (2.78)$$

$$\text{FP}(s^*) = n(1 - F_0(s^*)) \quad (2.79)$$

$$\text{FN}(s^*) = pF_1(s^*) \quad (2.80)$$

which, in turn, yields

$$\text{TPR}(s^*) = 1 - F_1(s^*) \quad (2.81)$$

$$\text{TNR}(s^*) = F_0(s^*) \quad (2.82)$$

$$\text{FPR}(s^*) = 1 - F_0(s^*) \quad (2.83)$$

$$\text{FNR}(s^*) = F_1(s^*) \quad (2.84)$$

$$\text{Pr}(s^*) = \frac{p(1 - F_1(s^*))}{p(1 - F_1(s^*)) + n(1 - F_0(s^*))} \quad (2.85)$$

The area under the ROC curve can now be expressed as

$$A_{\text{ROC}} = \int_0^1 \text{TPR}(\text{FPR}) d\text{FPR} = \int_{-\infty}^{\infty} \text{TPR}(s) \frac{d\text{FPR}(s)}{ds} ds = \int_{-\infty}^{\infty} (1 - F_1(s)) f_0(s) ds \quad (2.86)$$

The integration limits arise from the fact that FPR is maximal when all instances are classified as positives and minimal when all instances are classified as negatives, which respectively corresponds to $s^* \rightarrow -\infty$ and $s^* \rightarrow \infty$. This formulation of AUROC by ?? leads to the most common intuition behind the metric. The first factor $(1 - F(s^*))$, represents the fraction of positive instances that are ranked higher than the threshold s^* , alternatively expressed as $\int_{s^*}^{\infty} f_1(s) ds$. f_0 , as defined by ??, represents the probability of finding a negative instance within s^* and $s^* + ds$. Hence, the product of both factors represents the joint probability of having a negative instance between s^* and $s^* + ds$ while also finding a positive instance with $s > s^*$. After integration over all possible thresholds, we conclude that the AUROC score represents the overall probability of randomly selecting a positive instance ranked higher than a randomly selected negative instance ().

The baseline score 0.5 can be derived as follows. A random classifier is defined as an estimator incapable of distinguishing between the true label distributions of each class. That is, a classifier is a random classifier if and only if $f_0(s) = f_1(s) \forall s$. As a consequence,

$$A_{\text{ROC}} = \int_{-\infty}^{\infty} (1 - F_0(s)) f_0(s) ds = - \int_{-\infty}^{\infty} (1 - F_0) d(F_0) = \left[\frac{(1 - F_0)^2}{2} \right]_{-\infty}^{\infty} = \frac{1}{2} \quad (2.87)$$

?? also shows a characteristic of AUROC discussed in the previous section: the AUROC score is independent of the relative prevalence of each class in the test set. There is no influence of p or n and only the dependency on the learned decision value distributions.

The area under the PR curve (AUPRC) is defined as

$$\begin{aligned} A_{\text{PRC}} &= \int_0^1 \text{Pr}(\text{TPR}) \, d\text{TPR} = \int_{-\infty}^{\infty} \text{Pr}(s) \frac{d\text{TPR}(s)}{ds} \, ds = \\ &= \int_{-\infty}^{\infty} \text{Pr}(s) f_1(s) \, ds = \int_{-\infty}^{\infty} \frac{p(1 - F_1(s)) f_1(s) \, ds}{p(1 - F_1(s)) + n(1 - F_0(s))} \end{aligned} \quad (2.88)$$

From Equation 2.88, AUPRC can be interpreted as the average precision weighted by the distribution of positive instances. Analogously, it corresponds to collecting the decision values attributed to each positive instance, and then calculating the average precision considering only these values as classification thresholds. Furthermore, unlike AUROC, the true label cumulative distributions (F_0 and F_1) appear each weighted by their respective class prevalences (p and n).

2.5.4.2 AUPR and AUROC in terms of ranked decision values

When computing AUROC or AUPR for a given estimator on a test set, the values of the decision function s are usually not directly considered. Instead, the decision values outputted for each test instance are used to rank them from lowest to highest, and from these ranked test labels the curves and respective areas are obtained. The specific values of s are thus indifferent to the scoring process, as long as the ranking is preserved. If the percentile rank of each test instance is denoted by $r \in [0, 1]$ and each test instance is associated with a decision value, there is a one-to-one monotonic correspondence between r and s in the limit of an infinite number of test samples. Formally, we have

$$r(s) = pF_1(s) + nF_0(s) \quad (2.89)$$

$$dr = [pf_1(s) + nf_0(s)] \, ds \quad (2.90)$$

We also take the liberty to represent $g(r) = g(s(r))$, so that

$$F_1(r) = \int_0^r \frac{f_1(r)}{pf_1(r) + nf_0(r)} \, dr \quad (2.91)$$

$$\frac{dF_1(r)}{dr} = \frac{f_1(r)}{pf_1(r) + nf_0(r)} \quad (2.92)$$

AUPRC as defined in Equation 2.88 can now be written as

$$\begin{aligned} A_{\text{PRC}} &= p \int_0^1 \frac{1 - F_1(r)}{(p + n) - [pF_1(r) + nF_0(r)]} \frac{f_1(r)}{[pf_1(s) + nf_0(s)]} \, dr = \\ &= p \int_0^1 \frac{[1 - F_1(r)]}{[1 - r]} \frac{dF_1(r)}{dr} \, dr = -\frac{p}{2} \int_0^1 \frac{1}{1 - r} \frac{d[1 - F_1(r)]^2}{dr} \, dr = \\ &= -\frac{p}{2} \left| \frac{[1 - F_1(r)]^2}{1 - r} \right|_0^1 + \frac{p}{2} \int_0^1 \left[\frac{1 - F_1(r)}{1 - r} \right]^2 \, dr = \\ &= \frac{p}{2} \left\{ 1 + \int_0^1 \left[\frac{1 - F_1(r)}{1 - r} \right]^2 \, dr \right\} = \frac{p}{2} \left\{ 1 + \int_0^1 [\text{Pr}(r)]^2 \, dr \right\} \end{aligned} \quad (2.93)$$

in which the upper boundary term is determined by using L'Hôpital's rule and noticing that $F_1(r = 1) = 1$ and $\frac{dF_1(r)}{dr} \leq \frac{1}{p}$.

$$\lim_{r \rightarrow 1} \frac{[1 - F_1(r)]^2}{1 - r} = 2 \cdot \lim_{r \rightarrow 1} [1 - F_1(r)] \frac{dF_1(r)}{dr} = 0$$

Equation 2.93 reveals that AUPRC is closely related to the average squared precision across all ranks.

We can also express the AUROC in terms of the percentile ranks:

$$\begin{aligned} A_{\text{ROC}} &= \frac{1}{n} \int_0^1 (1 - F_1(r)) \left(\frac{nf_0(r)}{nf_0(r) + pf_1(r)} \right) dr = \\ &= \frac{p}{n} \int_0^1 (1 - F_1(r)) \left(1 - \frac{pf_1(r)}{nf_0(r) + pf_1(r)} \right) dr = \\ &= \frac{1}{n} \int_0^1 (1 - F_1(r)) \left(1 - p \frac{dF_1(r)}{dr} \right) dr = \\ &= \frac{1}{n} \int_0^1 (1 - F_1(r)) dr - \frac{p}{n} \int_0^1 (1 - F_1(r)) d(F_1(r)) = \\ &= \frac{1}{n} \int_0^1 (1 - F_1(r)) dr + \frac{p}{n} \left[\frac{(1 - F_1(r))^2}{2} \right]_0^1 = \\ &= \frac{1}{n} \left\{ \int_0^1 (1 - F_1(r)) dr - \frac{p}{2} \right\} \end{aligned}$$

Equations 2.93 and 2.94 put AUPR and AUROC in a similar format, better delineating the differences between the two metrics. Consider expressing both now in terms of the precision.

$$A_{\text{ROC}} = \frac{1}{n} \left\{ \frac{1}{p} \int_0^1 (1 - r) \text{Pr}(r) dr - \frac{p}{2} \right\} \quad (2.94)$$

$$A_{\text{PR}} = \frac{1}{2} \left\{ p + \frac{1}{p} \int_0^1 \text{Pr}(r)^2 dr \right\} \quad (2.95)$$

Ignoring constant terms and factors, both metrics are centered on integrating the precision over all possible ranks, each rank representing a classification threshold. The crucial difference is that AUPR equally considers all precision values, while AUROC weights each $\text{Pr}(r)$ value by the the reversed ranks.

The precision metric is normalized by the number of instances it considers, so precision values in different r have comparable magnitude. With this in mind, notice that the precision values calculated for $r \approx 1$ are obtained from a very small number of labeled instances. As a result, each label considered when $r \approx 1$ has a large influence on the $\text{Pr}(r)$ value, and precisions will tend to have larger variances when $r \approx 1$. Therefore, the AUPR metric is more influenced by the label value of the highest ranked instances, since it equally considers all precision values.

Additionally, the quadratic exponent of AUPR's integrand emphasizes higher precision values overall, independently of the number of samples on which they were calculated ($1 - r$).

This would amplify the effect of high $Pr(r)$ values for $r \approx 1$ resulting from stochastic label variations. Hence, we suggest that the AUPR metric is more prone to noise in the labels of highly ranked samples, prioritizing models that strictly maintain high precision for a smaller selection of highest ranks.

On the other side, AUROC compensates this effect by weighting each $Pr(r)$ value by $1-r$, so that precision values calculated for lower ranks are prioritized. The prioritized precision values ($r \approx 0$) are calculated from a larger number of labels, having lower variances and being more robust to label noise. Each label considered when $r \approx 0$ has a smaller influence on the $Pr(r)$ value relative to the other labels being considered. However, AUROC assigns a larger weight to these precision values, so the influence of smaller ranks also tend to increase.

This characteristic of AUROC might be more suitable for model comparison under the Positive-Unlabeled (PU) assumption, even though AUPR is usually recommended for imbalanced scenarios (?, ?, ?). In PU datasets, it is common for negative-labeled instances to rank highly due to the possibility of them being unannotated positives. These highly-ranked negatives would have a larger influence on the AUPR score in comparison to the AUROC score. In fact, we show in the next section that AUROC is closely related to the Mean Percentile Rank metric, which has been suggested for PU learning contexts of interaction prediction and recommendation systems (?, ?). Section 2.6.1 presents further discussion on the specific usecases of AUPR and AUROC, besides a numerical analysis of their dependence on the percentile ranks.

2.5.4.3 AUROC is the normalized mean percentile ranks

From Equation 2.94, we can express AUROC as

$$\begin{aligned} A_{\text{ROC}} &= \frac{1}{n} \left\{ \int_0^1 (1 - F_1(r)) dr - \frac{p}{2} \right\} = \\ &= \frac{1}{n} \left\{ 1 - \int_0^1 F_1(r) dr - \frac{p}{2} \right\} = \frac{1}{n} \left\{ 1 - \int_0^1 \frac{d(r)}{dr} F_1(r) dr - \frac{p}{2} \right\} = \\ &= \frac{1}{n} \left\{ 1 - [r F_1(r)]_0^1 + \int_0^1 r \frac{dF_1(r)}{dr} dr - \frac{p}{2} \right\} = \frac{1}{n} \left[\int_0^1 r dF_1 - \frac{p}{2} \right] \quad (2.96) \end{aligned}$$

Equation 2.96 offers another perspective on AUC ROC. The term $\int_0^1 r dF_1(r)$ represents the expected percentile rank of the positive samples:

$$\int_0^1 r dF_1(r) = E[r \mid y = 1] = \text{MPR} \quad (2.97)$$

This quantity is sometimes referred to as the *mean percentile ranking* (MPR) in the previous literature (?, ?), being proposed in contexts of recommendation systems (?, ?) and bipartite interaction prediction (?, ?, ?).

Consider now the maximum and minimum values of MPR, achieved, respectively, for

the ideal $f_1(r)$ distributions:

$$f_{1,\max}(r) = \frac{1}{p} \mathbb{I}(r > n) \quad (2.98)$$

$$f_{0,\max}(r) = \frac{1}{n} \mathbb{I}(r < n) \quad (2.99)$$

$$f_{1,\min}(r) = \frac{1}{p} \mathbb{I}(r < p) \quad (2.100)$$

$$f_{0,\min}(r) = \frac{1}{n} \mathbb{I}(r > p) \quad (2.101)$$

Applying these definitions to determine $\frac{dF_1(r)}{dr}$ and using the results in Equation 2.97, we obtain

$$\text{MPR}_{\max} = \frac{1}{p} \int_n^1 r \, dr = \frac{1 - n^2}{2p} = \frac{1 + n}{2} \quad (2.102)$$

$$\text{MPR}_{\min} = \frac{1}{p} \int_0^p r \, dr = \frac{p}{2} \quad (2.103)$$

from which is straightforward to show that

$$A_{\text{ROC}} = \frac{\text{MPR} - \text{MPR}_{\min}}{\text{MPR}_{\max} - \text{MPR}_{\min}} \quad (2.104)$$

Therefore, the AUROC can also be interpreted as the normalized MPR. To the best of our knowledge, this relationship is not clearly shown in previous explorations.

This result corroborates the argument that the AUROC could be preferable for PU learning scenarios, since the closely-related MPR is a known metric specifically recommended for this context (?,?,). Furthermore, the normalization enables the comparison across different tasks, ensuring that the scores on each dataset are always in the same range. Hence, the AUROC should be preferred over the MPR for the majority of cases.

2.5.5 General experimental settings

ESCREVER

2.6 Experiments

PADRONIZAR NOMES DE ESTIMADORES

This section describes the experiments performed in this work, designed to explore the proposed research questions (??), assess the effectiveness of the developed methods, and to validate predictions made in the theoretical analyses.

2.6.1 What are the differences between AUROC and AUPR?

Key findings:

- AUPR prioritizes a smaller number of highest-ranked interactions, while AUROC considers a larger number of both highest and lowest ranks.

- AUPR should be used when the goal is to select a small number of most-likely interactions. AUROC should be used i) when both likely-positive and likely-negative interactions are important; or ii) when interested in a large fraction of the predictions.
- AUROC could be also preferable for PU learning.

This experiment was designed to evaluate the hypotheses raised in Section 2.5.4.2. Our main objective is to elucidate how much importance each metric assigns to each percentile rank. To do so, we perform a Monte Carlo simulation. First, we select R equally spaced rank values between 0 and 1 (excluding 0 and 1). For each rank, we generate N random binary values to be used as labels. Therefore, each of the N iterations of the simulation will produce a random set of R binary labels, corresponding to each percentile rank. In each iteration, we use the R labels and R percentile ranks to calculate AUROC and AUPR, resulting in N values for each metric. Finally, for each of the R ranks, we calculate the point biserial correlation (?) between the N random labels and the N AUROC and AUPR values. The result is displayed in Figure 1.

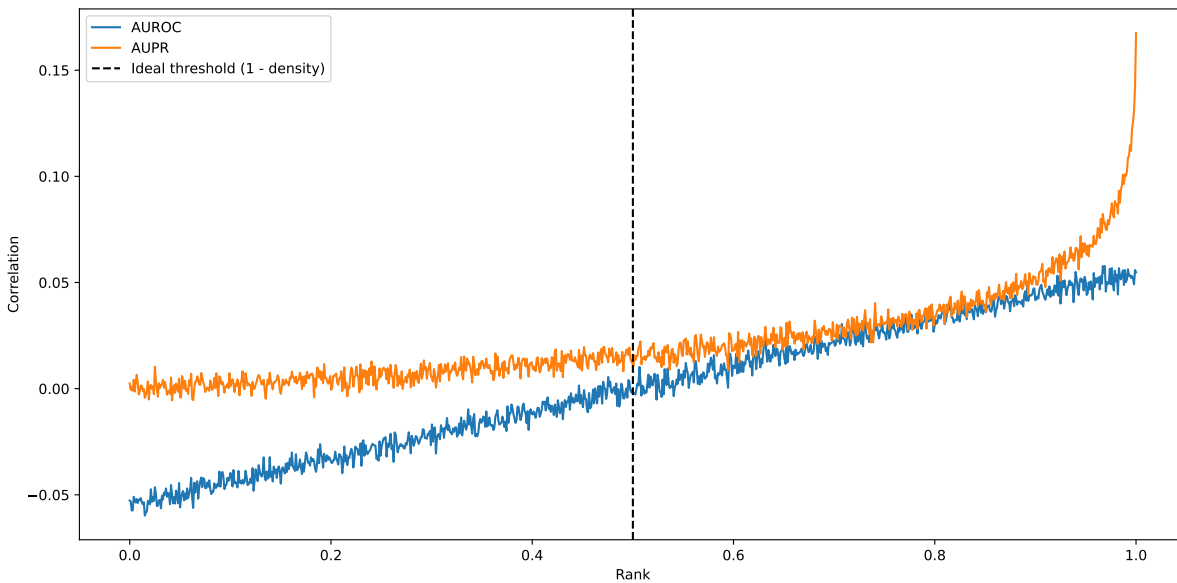


Figure 1 – Point biserial correlation between the binary labels at each percentile rank and the AUROC and AUPR scores. The correlation was calculated between the ranks and the scores for $N = 10^5$ random sets of $R = 10^3$ binary labels. The results show that AUPR is very sensitive to a small group of highest ranks, while AUROC has a more distributed weighting profile. For AUROC, both higher and lower ranks have a higher impact than ranks around 0.5, contributing in opposite directions to the final metric. The results confirm the theoretical analyses from Section 2.5.4.2. Source: made by the authors.

The plot shows that AUPR is highly sensitive to a small set of highest ranks, while ranks closer to 0 have little influence on the metric. On the other hand, the correlations with the AUROC score are more distributed across ranks. Also, for AUROC, both higher and lower ranks have a higher impact, but in opposite directions. This behaviour is in agreement with the

theoretical analyses from Section 2.5.4.2. The correlation with AUROC is also symmetrical around the median ($r = 0.5$), linearly increasing from around -0.05 at $r = 0$ to around 0.05 in $r = 1$. This demonstrates the class symmetry of the AUROC metric described in Section 2.5.4.2, meaning that negative and positive labels could be swapped without affecting the final results.

Overall, these results demonstrate that using AUPR is more suited when one is interested in selecting a restricted number of top-ranked instances from a pool of predictions, such as in recommendation systems or drug discovery tasks. Conversely, AUROC should be favored when the goal is to rank a large batch of interactions. Examples would be modeling genetic interactions in a genome-wide fashion or building interaction databases.

If one intends to select a small number of negative instances instead, a possible strategy would be to swap the binary labels and use AUPR. However, one should prioritize AUROC over AUPR if the goal is to select both negative and positive predictions with the highest confidences. This conclusion results from AUPR disregarding ranks close to 0, so models selected with AUPR are not guaranteed to yield the most confident true negatives.

We also argue that AUROC could be preferable for comparing models under the PU assumption, at least in purely theoretical settings. For PU datasets, we naturally expect some negative-labeled instances to be very highly ranked since they could be, in fact, unannotated positives. AUPR would more strictly penalize such predictions, favoring models that consider the labeling mechanism itself rather than only the underlying interaction mechanism. AUPR could thus undermine the model's potential to discover new interactions, possibly failing to gauge the generalization capabilities of algorithms in a PU context. The results from Section 2.6.8 seem to corroborate this hypothesis, with AUROC ranking estimators more consistently than AUPR across different levels of label noise.

2.6.2 Are BGSO models faster than GMO models?

Key findings:

- The GMO complexity is measured to be $O(n^3 \log n)$, while the BGSO complexity is measured to be $O(n^3)$.
- The empirical results match the theoretical expectations from Section 2.3.7.
- The difference in complexity between GMO and GSO is less pronounced for the ExtraTrees, but the GSO models still present significantly lower complexity than the GMO models.

This experiment empirically measures the training time complexity of the tree models under study. We artificially generate a series of bipartite datasets by filling three n by n matrices with pseudo-random values, representing the two X matrices and the Y matrix on each interac-

tion. Values were taken uniformly from the interval $[0, 1]$ for the feature matrices and from the interval $[0, 100]$ for the target matrix.

We then train the GMO and the optimized GSO versions of a single bipartite decision tree (BDT) and a single bipartite ExtraTree (BXT) on each of the generated datasets, measuring their training duration in seconds. The results are shown in Figure 2. From the least squares linear regression on the log-log plot, we see that the estimated training time complexities closely follow the theoretical expectations developed under Section 2.3.7, with slopes referring to the GSO models (predicted to be $O(n^3)$) approaching 3 while the GMO models (predicted to be $O(n^3 \log(n))$) produce slope between 3 and 4.

Statistical testing further shows that the empirical time complexity of the proposed GSO algorithms are indeed significantly lower than that of their GMO counterparts (see the caption for Figure 2).

We also notice slightly lower slopes for the ExtraTrees in comparison with the BDTs. In fact, they seem to be closer to the theoretical expectations for the GSO models. This is expected, since the randomized split search performs far less operations than the exhaustive search (Section 2.3.2). Notwithstanding, both algorithms have the same order of complexity. We then suggest that much larger datasets would be required to observe the asymptotic behavior of the ExtraTrees. In spite of that, the empirical complexity of `bdt_gso` is still observed to be highly significantly lower than that of `bxt_gmo`, validating once more the prediction that `bdt_gso` should present faster training times than `bxt_gmo` on sufficiently large datasets.

2.6.3 Which prototype should a GMO forest use?

Key findings:

- The `square` strategy is the best for both BXT and BRF models, except for the LT+TL average precision score, for which the fully-grown trees are the best.
- The weighted-neighbors strategies seem to improve generalization.

In this work we propose a different prototyping strategy to determine the output value of each leaf in a GMO decision tree, taking the similarity matrices of our use cases into consideration (Section 2.3.6). In this experiment we compare such strategies, building BXT and BRF models for every option. The minimum rows per leaf and minimum columns per leaf were both set to 5, ensuring that at least 5 samples of each domain are considered when calculating the prototype values. To observe the effect of this early-stopping criterion by itself, we also include forests of fully-grown trees in the comparison. The compared models are described below.

- **gmosa:** proposed by `pliakos2018global`, the output of each leaf is the average of the labels of the learned samples that reach that leaf (Equation 2.11).

- **uniform:** also proposed by pliakos2018global, the only difference from the `gmosa` strategy occurs in TL or LT test settings, when the average is taken only among the labels of the known sample of the pair being predicted (Equation 2.13).
- **precomputed:** the labels in each leaf are weighted by the similarities between the learned samples and the pair being predicted (Equation 2.14).
- **square:** the labels in each leaf are weighted by the squared similarities between the learned samples and the pair being predicted (Equation 2.14).
- **softmax:** the labels in each leaf are weighted by the exponential of the similarities between the learned samples and the pair being predicted (Equation 2.14).
- **full:** the trees are grown until a single interaction remains in each leaf.

The results for BRF (Figure 3) and BXT (Figure 4) were similar. In all cases except LT+TL average precision, using the square of the similarities to weight the labels in each leaf resulted in the best scores. For the LT+TL average precision score, growing the decision tree to its maximal size was the best strategy, followed by `uniform`, the original proposal by pliakos2018global of averaging only the outputs of the learned samples (known from the training set) in each leaf. With the exception of the LT+TL AUROC metric for BXT, the mentioned winning models were statistically distinguished from all the remaining estimators. For LT+TL AUROC, the superiority of `square` could not be attested when compared to the `uniform` strategy.

The fully-grown versions of both forest algorithms are shown to be especially advantageous when considering the LT+TL average precision. This suggests that building trees to their maximum depth is the best strategy for learning tasks in which

1. one of the domains is fixed, with the final goal being to model how new instances will bind to this known set of entities;
2. the goal is to select a small number of top-ranked interactions (see Section 2.6.1).

On the other hand, weighted averages seem to improve the forest’s generalization ability, as they seem to perform best under the AUROC metric and TT contexts.

We propose that this distinction results from the ability of a fully-grown tree to independently consider the labels of each learned instance when calculating the prototype. However, the methods using weighted averages invariably mix the labels of a pool of neighbors in each leaf. This hypothesis is supported by the fact that the second-best model regarding LT+TL average precision is the `uniform` strategy, which also uses the labels of individual learned samples to generate the predictions. The hypothesis alone, however, do not explain the superiority of the

deeper trees over the `uniform` weights. In this case, the larger tree depth is likely beneficial through i) an increase in the predictive power of each individual tree, and/or ii) an increase in tree diversity, both of which would improve the ensemble’s performance as discussed in ???. We let to future work the more specific investigation of these effects.

On the other hand, `full` and `uniform` show notable inferior performance in the TT evaluation settings, suggesting that these strategies have inferior ability to consider completely new interacting pairs.

We select the squared weighting strategy and the fully grown trees to be further investigated in the downstream analyses.

2.6.4 Which adaptation strategy is the best for decision forests?

Key findings:

- SLMO is the best strategy for BRF models on LT+TL sets.
- GMO is the best strategy for completely new dyads.
- Undersampling of negative annotations is beneficial for GSO BXT in terms of AU-ROC, but should be avoided if the goal is to select the highest ranked interactions.

We now compare each of the described approaches for adapting decision forests to bipartite data, including data-centered adaptations (Section 2.2) and the natively bipartite forests (Section 2.3.4 and Section 2.3.5). We briefly describe below the suffixes in the model names of this section, indicating the employed bipartite adaptations.

- **lmo**: implements the standard local multi-output approach (SLMO; Section 2.2.2) by training four separate multioutput models, two for each domain. First explored by schrynemackers2015classifying.
- **lso**: also implements the SLMO approach (Section 2.2.2), but each multioutput model is instead a composition of several local single-output (LSO) models, i.e. one model is trained for each row or column of the interaction matrix. This setting is similar to the early proposal by bleakley2009supervised, but employing decision forests as the base algorithm.
- **gmosa**: a global multi-output forest with single-label averaging (GMOSA; ??), as explored by pliakos2019network under the name eBICT.
- **gmo**: a global multi-output (GMO; ??) forest with minimum leaf dimensions of 5 by 5, implementing our proposed squared-similarities weighting for the prototype function (Section 2.3.6). Apart from the new prototype, this model is based on the original GMO

trees proposed by pliakov2018global. Despite their original results suggesting advantage over GMOSA, to the best of our knowledge, this is the first time the GMO trees are employed in building decision forests.

- **gso**: a bipartite global single-output (BGSO; ??) forest, as initially explored by schrynemackers2015classifying but now implementing our proposed algorithm with improved computational complexity.
- **sgso_us**: implements the standard global single-output (SGSO; Section 2.2.1) adaptation, additionally employing undersampling of the non-interacting pairs to yield a balanced training set (Section 2.2.1).

All the global models were built with 100 trees. For SLMO, each of the four forests used 50 trees, while for LSO 50 trees were used for each row or column of the interaction matrix.

In the LT+TL AP evaluation setting, the pattern observed in Section 2.6.3 again emerges: the label averaging strategy employed by GMO performs considerably worse in comparison to forests that separately consider each known instance. In particular, the SLMO and SLSO adaptations yield the clear best BRF models in terms of LT+TL AP, and these adaptations interpret each instance as a separate output to be predicted. Among the local adaptations, SLMO significantly outperforms SLSO. SLSO treats each training instance as a completely independent task, building a separate forest for each row and column of the interaction matrix. As such, the previous result shows that this complete independence is not desirable for the learning problems under study, and exploring label correlations between instances of the same domain is beneficial (as demonstrated by SLMO). This result could partially be a consequence of the very sparse nature of our problems: if the interaction information of each instance is limited, it becomes advantageous to aggregate information from other instances with correlated interactions. We then speculate that the advantage of SLMO over SLSO could become less prominent once the number of known interactions per instance increases and more data are available for training.

Still considering LT+TL AP but focusing on BXT models instead, we notice that SLMO loses the advantage to the fully-grown bipartite trees GMO and GSO. A possible explanation comes from the fact that local approaches yield shallower trees, and much more randomized trees could be required to achieve comparable performance. In more detail, first notice that the individual performance of each tree in a BXT forest is lower than that of a tree in a BRF. Thus, a BXT ensemble requires a larger number of trees to reach satisfactory performance. Additionally, if the trees are shallow, they tend to be less representative of the training data (each tree node brings a little more information on the dataset). Even more trees then should be required to compensate the randomness of each individual. Finally, building a forest locally as in SLMO or SLSO results that each tree is trained on a much smaller number of samples in comparison to considering each dyad as a separate instance. Therefore, local approaches generate smaller and

less representative trees. We then suggest that the advantage of SLMO could also manifest for BXT estimators if the number of trees in the ensemble were to be increased.

For both BXT and BRF in all the TT settings, the GMO model significantly outperforms the other estimators. GMO is also the best model in the LT+TL AUROC setting when comparing BXTs and the second best for BRFs. This reinforces the findings of ?? suggesting that larger leaves and label weighting seems to improve generalization to unseen instances. We also highlight SLMO as a prominent strategy for BRFs under TT AUROC and TT AP, being the second best model in both cases.

The SGSO US strategy performed consistently worse than the others under the AP metric. Conversely, it was significantly the second best BXT model in both LT+TL and TT AUROC settings, while surpassing BRF GMOSA also under both AUROC test sets. This shows that undersampling is a viable technique when considering the overall ranking of interactions. However, it is detrimental when the goal is to select the most likely interactions. That is, SGSO US allows a larger number of false positives in the highest ranked positions.

This behaviour is expected. Notice that SGSO US models are trained in a balanced dataset, while the other models are trained in the original dataset where negative annotations are much more frequent. With more training examples, the other models tend to be considerably better in correctly classifying negative annotations. This suggests that AUPR tends to favor models of high specificity, that are strictly avoid predicting false positives.

We select LMO, GMO, GSO, and GMOSA to be further analysed in the next section.

2.6.5 Can label imputation assist bipartite forests?

Key findings:

- Imputing positive annotations with NRLMF improves the performance of bipartite forests in the majority of cases.
- For LT+TL AP, the best BRF is BRF LMO and the best BXT is BXT GMOSA NRLMF.
- GMO NRLMF is the overall best model for TT.

It was previously suggested to employ logistic matrix factorization to create a dense representation of the interaction matrix. Using this representation as the training data for a BXT forest could improve the performance on DTI datasets (?). To test this hypothesis, we compare the bipartite forests cross-validation scores with and without the interaction matrix reconstruction step. As done by (?), the reconstruction step was performed using neighborhood-regularized logistic matrix factorization (NRLMF)(?).

We performed a randomized search to select hyperparameters for the NRLMF algo-

rithm. 100 different combinations of hyperparameters were evaluated in terms of their resulting mean squared error in a nested bipartite 5-fold diagonal cross-validation. The best combination of parameters according to the inner CV loop was then used to reconstruct the interaction matrix of each outer CV fold. Only then the resulting matrices were used as the training data for the bipartite forests. Note that a single forest was built per outer CV fold, so that the NRLMF hyperparameter search was performed independently from the downstream forest performance. The hyperparameters `lambda_rows`, `lambda_cols`, `alpha_cols`, `alpha_rows`, and `learning_rate` were all independently sampled from a log-uniform distribution bounded by $\frac{1}{4}$ and 2. The number of latent vector components was set to be equal for both axes, and chosen between 50 and 100. The number of neighbors was randomly selected as 3, 5 or 10 in each iteration, and the maximum number of optimization steps was always set to 100. The parameter c , as in the original paper, was set to 5. This parameter specifies the emphasis on positive annotations, so each positive annotation contributes 5 times more to the loss function in our settings, as if 5 copies of each positive interaction were present in the training set. See Section 2.6.4 for descriptions on the bipartite forests being compared.

The results for BRF and BXT are shown by figures 7 and 8, respectively. The `nrlmf` suffix to a model indicates that it was built on the output of NRLMF.

The interaction matrix reconstruction step is shown to be especially beneficial in terms of AUROC scores, improving this score for almost all forest algorithms investigated. The only exceptions were `bxt_gmo` and `brf_gmo` in LT+TL ROC AUC, where the improvement is still observed but not statistically significant.

Even without the help from NRLMF, the `bxt_gmo` and `brf_gmo` models employing our squared similarities output weighting (sections Section 2.3.6 and Section 2.6.3) are placed second in their corresponding TT ROC AUC results, both significantly outperforming all NRLMF-combined models except for their own versions, `bxt_gmo__nrlmf` and `brf_gmo__nrlmf`. `brf_gmo` was also the top performing BRF model in terms of TT average precision, significantly surpassing all but `brf_gmo__nrlmf` and `brf_lmo__nrlmf`. Similarly, `bxt_gmo` was the top performing BXT model under the same metric, significantly surpassing all estimators but `brf_gmo__nrlmf`, `brf_lmo__nrlmf` and BXT GMOSA NRLMF.

Regarding TT AP, NRLMF still provides a significant improvement for all models except `brf_gmo`, `bxt_gmo` and `bxt_lmo`. While both BRF GMO and BXT GMO resulted in the highest average ranks for TT AP, the top position could not be statistically resolved between the BRFs using GMO, GMO NRLMF, and LMO NRLMF, nor between the BXTs using GMO, GMO NRLMF, LMO NRLMF, and GMOSA NRLMF.

In the LT+TL AP setting for BXT ensembles, NRLMF is shown to significantly degrade the performance of `bxt_gmo`. Therefore, the combination of the squared neighbors prototype and NRLMF likely results in more high-rank false positives under LT+TL, especially since such

degradation is not observed under LT+TL AUROC.

For BRF models under LT+TL AP, the advantage of employing NRLMF is not clear. GSO was the only adaptation to significantly benefit from NRLMF in this setting, while GMO and LMO were significantly impaired by it. In fact, the BRF LMO model without label matrix reconstruction significantly outperformed all the other random forests for this evaluation setting.

In summary, interaction matrix reconstruction by NRLMF seems to consistently improve BRF and BXT results in terms of AUROC and also the results of BXT regarding LT+TL AP. On TT AP, while less evidence is found, the results still seem to point in the same direction of a beneficial NRLMF transformation step. On the other hand, the comparisons of BRF under LT+TL average precision evaluation indicate the opposite conclusion, disfavoring usage of NRLMF especially for `brf_lmo` and `brf_gmo`. We thus discourage the application of NRLMF with BRF in scenarios where i) one of the entities of the dyad is always known to the model, such as in drug repositioning; and ii) the main interest is to select a small set of most-likely interactions.

2.6.6 What is the best way of building semi-supervised forests?

Key findings:

- The MD unsupervised impurity is the best option for LT+TL.
- AD is the most promising strategy for TT.
- The best strategy for determining the supervision balance seems to be either size or fixed, but further investigation is needed.

In this experiment, we compare the performances of semi-supervised bipartite forests Section 2.3.9. All forests in this section are based on the BGSO BXT, as described in Section 2.6.4. We evaluate all combinations of the three strategies for calculating the unsupervised impurity and the four strategies for determining the supervision balance at each node.

Strategies for determining the unsupervised impurity:

- **mse**: it refers to forests using the mean squared error as the unsupervised impurity (??).
- **md**: corresponds forests employing the mean distance unsupervised impurity (Equation 2.22).
- **ad**: also uses the mean squared error (??), however, the semi-supervised impurity is only calculated twice on each tree node. The search for the best feature and split point in each instance domain uses the supervised impurity only. Then, the semi-supervised impurity is only used to evaluate the best split in each domain. Finally, the two resulting values for the semi-supervised impurity are used to select between the horizontal and vertical split.

Strategies for determining the supervision balance (σ):

- **fixed:** the supervision balance (σ) is fixed at 0.5 for mse and md and at 0 for ad.
- **density:** σ is determined by the density of positive annotations in each tree node (Equation 2.23).
- **size:** σ is determined by the total number of interactions current in the tree node (Equation 2.24).
- **random:** σ is set to a random value between 0 and 1 drawn at each tree node (Equation 2.25).

The comparison results are displayed by Figure 9. Regarding LT+TL AUROC, the 0% ILR and 50% ILR do not result in statistically significant differences in performance. For ILR=70%, md random and md size significantly outperform the remaining models, while under ILR=90%, md random, md size, and md density are shown to surpass the others. These results indicate that the MD unsupervised impurity is the most suited to the LT+TL AUROC evaluation setting. Furthermore, we see that the random strategy to select the supervision amount is among the best in both ILR=70% and ILR=90%, suggesting that increasing diversity among the trees could be the main mechanism behind the improvement of the semi-supervised models, rather than necessarily guessing the best σ at each tree node. Since md size is also among the best, it is not clear if the advantage of md size over the other models is due to a better choice of σ or to a more stochastic nature of the σ s it selects. As a future investigation, we suggest using completely random values for the unsupervised impurity, to assess the possibility of tree-diversity being the main factor behind the observed improvements.

As for TT AUROC, only ILR=90% yielded significant comparisons. In this case, the mse unsupervised impurity with size, random, and density σ selection were the best strategies, surpassing the remaining. MSE thus seems the better option in terms of AUROC for scenarios with very scarce information and completely unknown instances.

Under LT+TL AP, md fixed significantly outperformed the other models for ILR=50% and ILR=70%. For ILR=0%, md fixed is also the first place, but could not be statistically resolved from the second place ad size. AD size was also the second best model for ILR=0%, 50%, 70%, statistically outperforming all but md fixed in ILR=70%.

With respect to TT AP, the AD strategies prevail, being the best four models for ILR=50% and ILR=70%, and being among the five best models for ILR=0%. The presence of AD random among the best models again suggests that tree-diversity is an important factor for performance improvement. It seems that AD models perform overall better according to AP in comparison to AUROC. This is consistent with the discussion in Section 2.6.1: AP tends to prioritize mod-

els more tightly related to the observed labels, assuming less risk of false positives. AD is the strategy that is less influenced by the unsupervised impurity function.

In conclusion, the best unsupervised impurity under LT+TL seems to be the mean distance, and AD seems to be the best for AP TT. The best strategy for determining σ seems to be either size or fixed, but more experiments are needed to confirm this finding.

2.6.7 Which forests are the best?

Key findings:

- Semi-supervised impurities are beneficial in terms of AUROC, especially when more annotations are missing.
- Label imputation with NRLMF is better than semi-supervised impurities.
- BXT are superior to BRF in all settings.

In this experiment, we compare the performances of various models from the previous sections. For more information on each model see the corresponding section in the list bellow.

- **BXT-GSO, BRF-LMO:** Bipartite forests without label imputation. See Section 2.6.4 for more information.
- **MD SIZE, MD FIXED, AD SIZE, AD FIXED:** Bipartite forests employing semi-supervised impurities. See Section 2.6.6.
- **Models with the NRLMF suffix:** Bipartite forests using NRLMF to impute positive annotations. See Section 2.6.5.

The comparisons are presented by Figure 10. The results reveal a clear superiority of forests employing the NRLMF as a label imputation strategy, in comparison to those using semi-supervised impurities. The four models employing NRLMF were the four highest ranked estimators in almost all evaluation settings, the two exceptions being LT+TL AP with ILR=0% and ILR=50%. In the first exception (LT+TL AP ILR=0%), BRF GMO NRLMF and BXT GMO NRLMF are the two worst-performing models, whereas BXT GMOSA NRLMF and BXT GSO NRLMF are the first and second best, respectively. In the second exception (LT+TL AP ILR=50%), the two best models are the same, and BRF GMO NRLMF is still one the worst performers. However, BXT GMO NRLMF jumps to the third best position.

Furthermore, under the other metrics (LT+TL AUROC, TT AUROC, and TT AP), BXT GMO NRLMF was the highest ranked estimator in almost all cases, the only exception being TT AP IRL=0%, where it was only behind BXT GMOSA NRLMF. Notwithstanding, the comparison between BXT GMO NRLMF and BXT GMOSA NRLMF was still not statistically

significant in this setting. Similarly, no statistical difference is found between these models in the LT+TL AUROC 0% setting, where they also occupy the first positions. The leadership of BXT GMO NRLMF is also not significant in comparison to BXT GSO NRLMF under TT AP 50% and TT AP 70%. For all remaining cases where BXT GMO NRLMF was the best model, it was statistically significantly better than all other estimators (TT AP 90, LT+TL AUROC 50, LT+TL AUROC 70, LT+TL AUROC 90, and all TT AUROC).

Another observation is that some semi-supervised impurities significantly improve predictive performance relative to the original BXT GSO, especially when more annotations are missing. This conclusion is based on the fact that, in all settings with $ILR \neq 0$ except LT+TL AP, the original BXT GSO is among the three worst performers. For all AUROC settings with $ILR \neq 0$, it was the lowest ranked model. On the other hand, md fixed is noted to significantly surpass BXT GSO in 10 of the 16 evaluation settings, the exceptions being LT+TL AUROC 0%, TT AUROC 0%, TT AUROC 90%, TT AP 0%, TT AP 70% and TT AP 90%. Therefore, although not as effective as the NRLMF reconstruction technique, using semi-supervised trees seems indeed beneficial when label information is scarce.

When comparing the BXT against the BRF models, we notice that BXT GMO NRLMF significantly outperformed BRF GMO NRLMF in all settings. The other random forest, BRF LMO, was also significantly outperformed by BXT GMO NRLMF in all cases but LT+TL AP 0%, where the opposite was observed. Similar results hold for the other BXT models as well: they significantly surpassed BRF LMO in the vast majority of test configurations, with the only exception being TT AUROC 0%, in which the comparison between BXT GMOSA NRLMF and BRF LMO was not significant. These results suggest that BXT models could offer significant advantages over BRF models in the context of DTI prediction. This conclusion is especially relevant given that the BXT training algorithm is considerably faster than the procedure for building BRFs, as discussed in Section 2.3.7.

2.6.8 Can bipartite forests compete with other proposals?

Key findings:

- Forests with NRLMF surpass NRLMF alone.
- BXT GMO NRLMF is the best overall performer, with the exception of LT+TL AP.
- For LT+TL AP, BXT GSO NRLMF and BXT GMOSA NRLMF are the best models overall.

In this section we compare the two most promising methods we developed (BXT GMO NRLMF and BXT GSO NRLMF) with several prominent models from the literature. The algorithms being considered in this section are listed below, and their scoring results are shown by Figure 11.

- **BXT GSO NRLMF, BXT GMOSA NRLMF, BXT-GMO-NRLMF:** Ensembles of randomized trees using the NRLMF model to impute missing annotations. The GSO model implements our optimized method for growing the trees. GMOSA uses the global multi output strategy developed by (?). Both GSO and GMOSA grows the trees to their maximum depth. GMO uses the weighted-neighbors prototype we developed in Section 2.3.6, and enforces at least 5 samples of each domain in each leaf. See Section 2.6.3 and Section 2.6.5 for more information.
- **NRLMF:** Neighborhood-Regularized Logistic Matrix Factorization, as proposed by (?) and described in Section 2.4.2.2. See Section 2.6.5 for the hyperparameters we utilized.
- **Kron-RLS:** Kronecker Regularized Least Squares, as proposed by (?) and described in ???. Each of the two input kernel matrices was taken as a linear combination of the similarity matrix and the gaussian interaction profiles (?) (??). The weight of the similarity kernel in this combination (the α parameter), was selected between the values {0.0, 0.1, 0.25, 0.5, 0.75, 0.9, 1.0}. The selection was performed in each fold by an internal 5 by 5 bipartite diagonal CV procedure (Section 2.5.3).
- **LMO RLS:** SLMO adaptation (Section 2.2.2) of Regularized Least Squares, as proposed by (?) as RLS-avg. Both the primary and secondary models were multi-output kernel Ridge regressors (). As Kron-RLS, the input features were linear combinations of similarity matrices and gaussian interactions profiles. The α parameter was selected in the same way as in Kron-RLS.
- **BLMNII RLS:** Bipartite local models (?, ?) with neighbor-based interaction-profile inferring (?). They use the SLMO strategy (Section 2.2.2), employing a weighted-neighbors technique (??) for the primary estimators and kernel Ridge regression () models as secondary estimators. As Kron-RLS, the input features were linear combinations of similarity matrices and gaussian interactions profiles (?). The α parameter was selected in the same way as in Kron-RLS.
- **BLMNII SVM:** The same as BLMNII RLS, but using support vector machines () as secondary estimators instead of regularized least squares.
- **DTHybrid:** Method proposed by (?) that combines the similarity kernels with network features of each domain calculated from the interaction matrix. A weighted-neighbors approach is then used to predict new interactions.
- **MLP:** Multi-layer perceptron () model adapted with the SGSO strategy (Section 2.2.1). We performed random undersampling of negative annotations so that the model was trained on equal number of negative and positive labels. The architecture was selected at each fold between four options: 5 hidden layers of 100 neurons; 10 hidden layers of 50 neurons; 5 hidden layers with [200, 100, 100, 100, 50] neurons; and 6 hidden layers

of [1024, 512, 256, 128, 64, 32] neurons. The architecture selection in each fold was performed by an internal 5 by 5 bipartite diagonal CV procedure (Section 2.5.3). The activation function was always the rectified linear unit, and ADAM (?) was the chosen optimizer. Other parameters were kept as the defaults from the scikit-learn library (?).

BXT GMO NRLMF was the best ranked model in 11 out of the 16 test settings analysed. The exceptions were the four LT+TL AP configurations and TT AP 0%. In TT AP 0%, BXT GMO NRLMF was only behind lmo rls, and among the three models that could not be statistically distinguished from lmo rls (the others being kron rls bxt and gmosa nrlmf). In LT+TL AUROC 0%, the comparison between BXT GMO NRLMF and the second best model, bxt gmosa nrlmf, was also not statistically significant. In the remaining 10 cases where BXT GMO NRLMF prevails, it was significantly better than all other learning algorithms analysed.

Regarding LT+TL AP, in ILR=50 and ILR=70 bxt gso nrlmf and bxt gmosa nrlmf significantly outperform the other models. In ILR=0, bxt gmosa nrlmf was the best model and bxt gso nrlmf was the second best, but neither could be statistically distinguished from the third best model, kron rls. In ILR=90, bxt gso nrlmf was the best model and bxt gmosa nrlmf was the second best, but neither could be statistically distinguished from the third best model, kron rls. Under ILR=90, bxt gso nrlmf and bxt gmo nrlmf are shown to significantly outperform the others.

We notice that the highest ranked bipartite forests in each evaluation setting are always able to significantly surpass NRLMF alone. This demonstrates that the forests are able to capture different patterns than the NRLMF model, even though they are trained on NRLMF's outputs. If this result were to not hold, one could argue that the forests might be merely approximating the predictive function learned by NRLMF, rather than expanding on the information extracted in the matrix factorization step.

We also note that the competitiveness of the kron rls model for TT AP 0% and LT+TL AP 0% is remarkable, since kron rls has notably low training times in comparison to the other models. However, the algorithm seems to not perform as well when a larger number of positive annotations is missing. This can be seen when comparing those same test metrics under higher ILR: TT AP 70, TT AP 90, LT+TL AP 70 LT+TL AP 90. Furthermore, Kron-RLS seems to be prioritized by AP in comparison to AUROC. This is especially suggested when observing the four TT AUROC, in which Kron-RLS performs poorly in comparison to the other cases. According to the discussion in Section 2.6.1, this could mean that Kron-RLS is mostly effective when considering a restricted number of top predictions.

Conversely, the performance of BXT GMO NRLMF for the LT+TL AP setting seems to improve relative to the other models as the number of missing annotations increases. A possible explanation comes from BXT GMO having higher generalization capabilities, as demonstrated in the TT AUROC setting. Notice that the LT+TL setting could possibly still benefit from some

form of overfitting, since it considers instances that are present in the training set (Section 2.5.3). However, being able to generalize becomes more important as the number of missing annotations increases, since the model has to rely on less information to make predictions. And that could be why we see the BXT GMO NRLMF model being more valued in the LT+TL AP setting for higher ILR. Additionally, we argue that the AP metric is less sensitive than AUROC to this ability of discovering missing positives Section 2.6.1, which explains why AUROC does not show the same pattern as AP in the LT+TL setting.

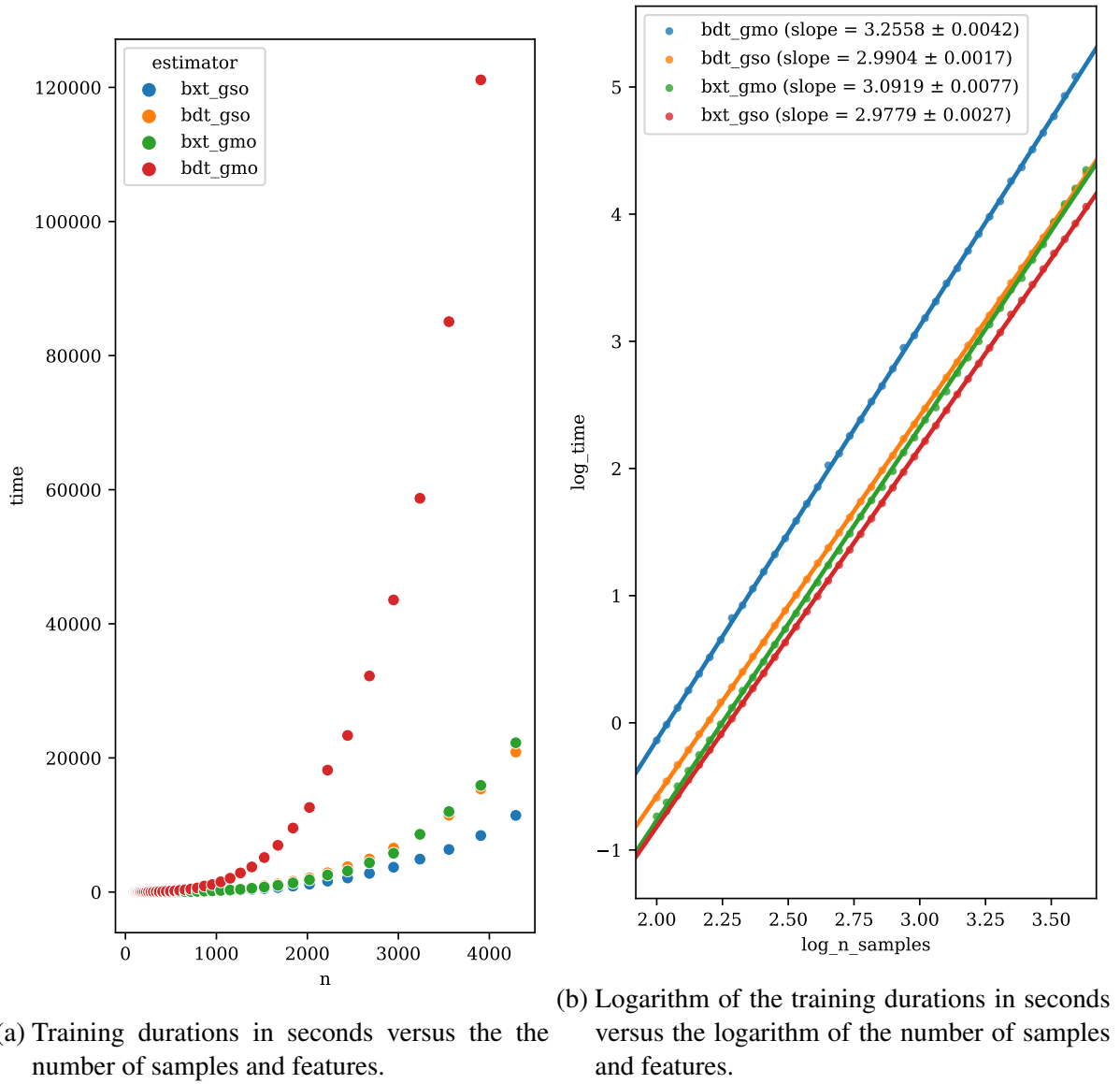


Figure 2 – Empirical time complexity estimation of the proposed bipartite global single-output (BGSO) and the global multi-output (GMO) (?) algorithms. Bipartite versions of both extremely randomized trees (?) (BXT) and greedy decision trees (?) (BDT) were built under the GSO and GMO scheme and trained over artificial datasets of varying numbers of samples (as described in Section 2.3.7). Applying least-squares linear regression to the logarithm of the values in (a), as shown in (b), the empirical slopes and respective standard deviations are obtained. Independent two-sample t-tests comparing the slope estimates reveal that the time complexity of `bdt_gso` is highly significantly lower than `bdt_gmo` (p-value $< 10^{-64}$) and even `bxt_gmo` (p-value $< 10^{-20}$), and also that `bxt_gso` significantly exhibits lower complexity than `bxt_gmo` (p-value $< 10^{-22}$). Those values corroborate the theoretical estimates from Section 2.3.7. Source: made by the authors.

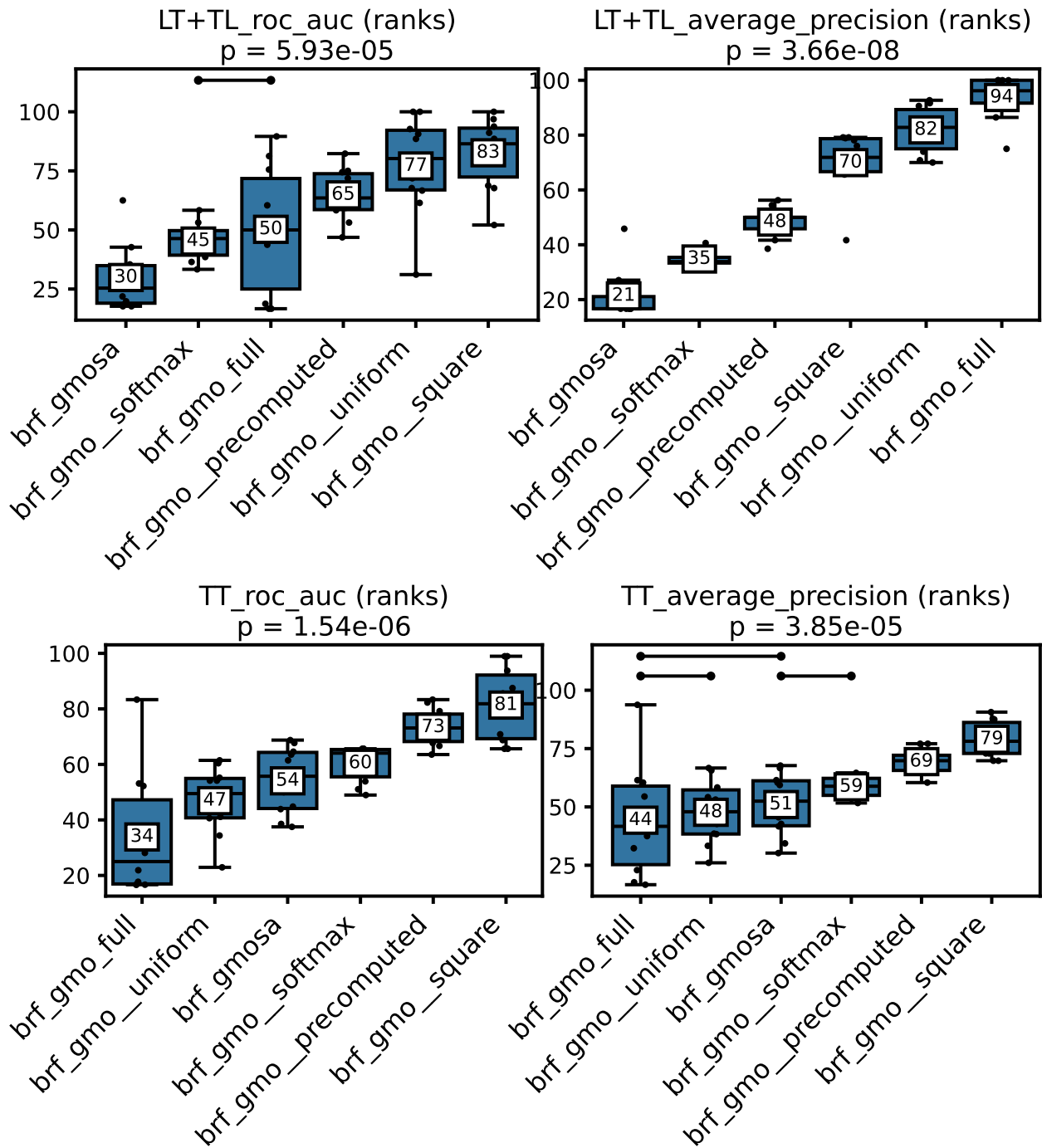


Figure 3 – Percentile rankings of prediction scores of bipartite random forests for different prediction weighting strategies, across all 10 datasets. The omnibus p-value obtained through a Friedman test is indicated below each plot’s title. Significant results in each case lead us to perform pairwise Wilcoxon rank-sum tests, whose non-significant results are indicated by crossbars above each plot, that is, statistically undistinguished boxes are connected by the crossbars. The pairwise test results were corrected by the Benjamini-Hochberg procedure among each subfigure, considering that all possible estimator pairs were evaluated even if their comparison is not indicated in the plot. See Section 2.6.2 for further descriptions of each estimator. Source: made by the authors.

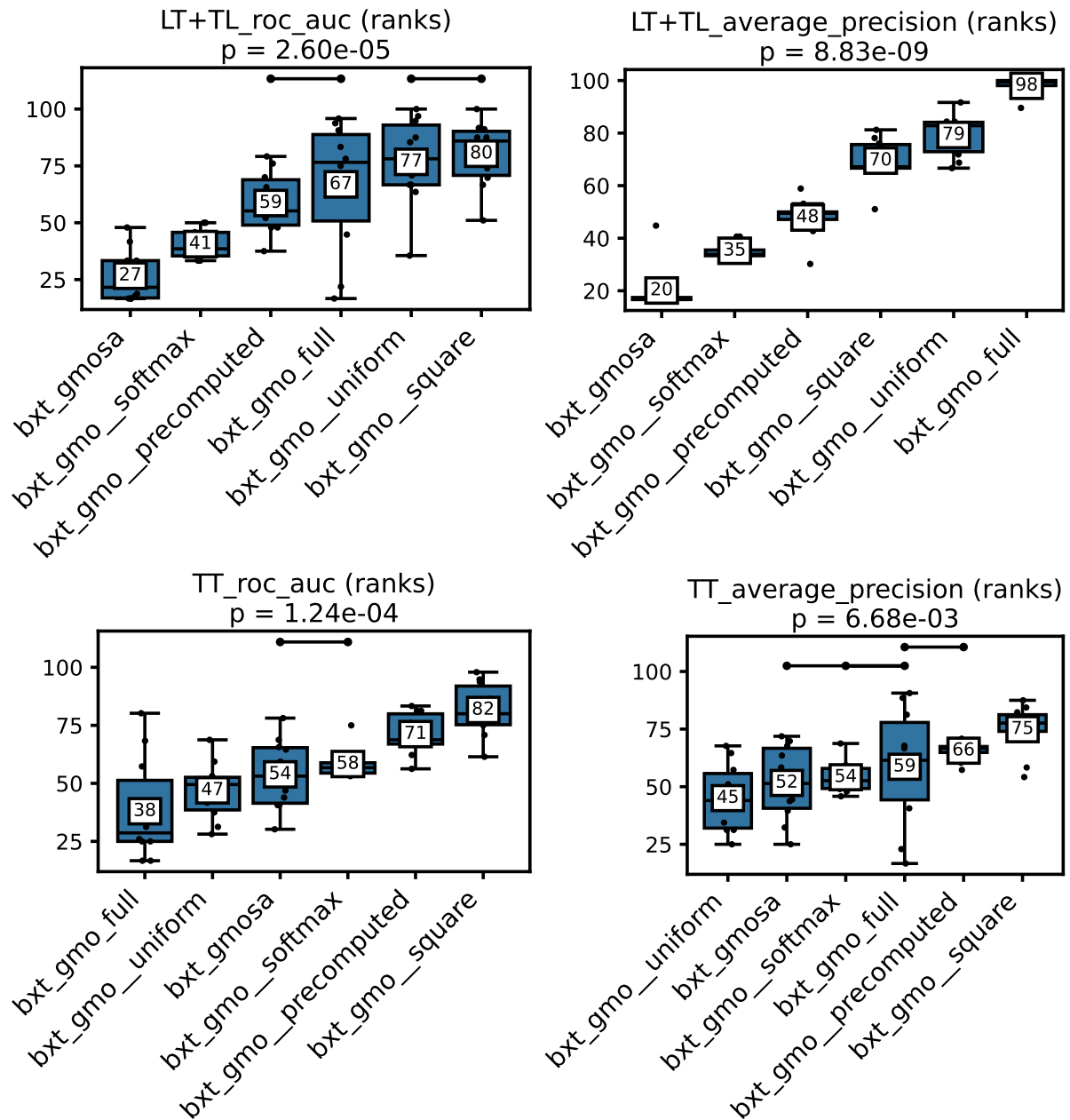


Figure 4 – Percentile rankings of prediction scores of bipartite extremely randomized trees for different prediction weighting strategies, across all 10 datasets. The omnibus p-value obtained through a Friedman test is indicated below each plot’s title. Significant results in each case lead us to perform pairwise Wilcoxon rank-sum tests, whose non-significant results are indicated by crossbars above each plot, that is, statistically undistinguished boxes are connected by the crossbars. The pairwise test results were corrected by the Benjamini-Hochberg procedure among each subfigure, considering that all possible estimator pairs were evaluated even if their comparison is not indicated in the plot. See Section 2.6.2 for further descriptions of each estimator. Source: made by the authors.

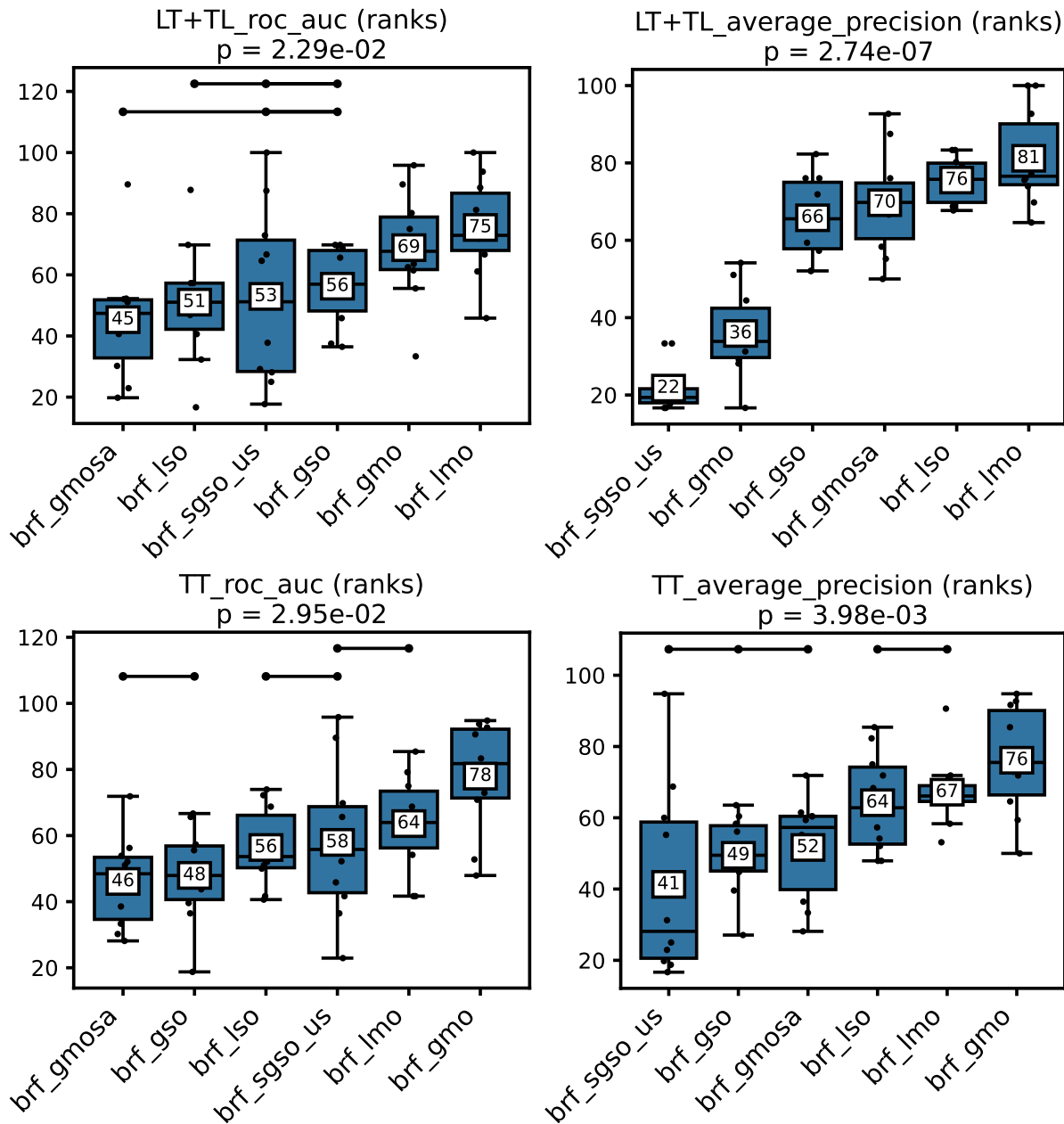


Figure 5 – Percentile rankings of prediction scores of random forest models under different adaptation strategies to bipartite interaction data, across all 10 datasets. The omnibus p-value obtained through a Friedman test is indicated below each plot's title. Significant results in each case lead us to perform pairwise Wilcoxon rank-sum tests, whose non-significant results are indicated by crossbars above each plot, that is, statistically undistinguished boxes are connected by the crossbars. The pairwise test results were corrected by the Benjamini-Hochberg procedure among each subfigure, considering that all possible estimator pairs were evaluated even if their comparison is not indicated in the plot. See Section 2.6.4 for further descriptions of each estimator. Source: made by the authors.

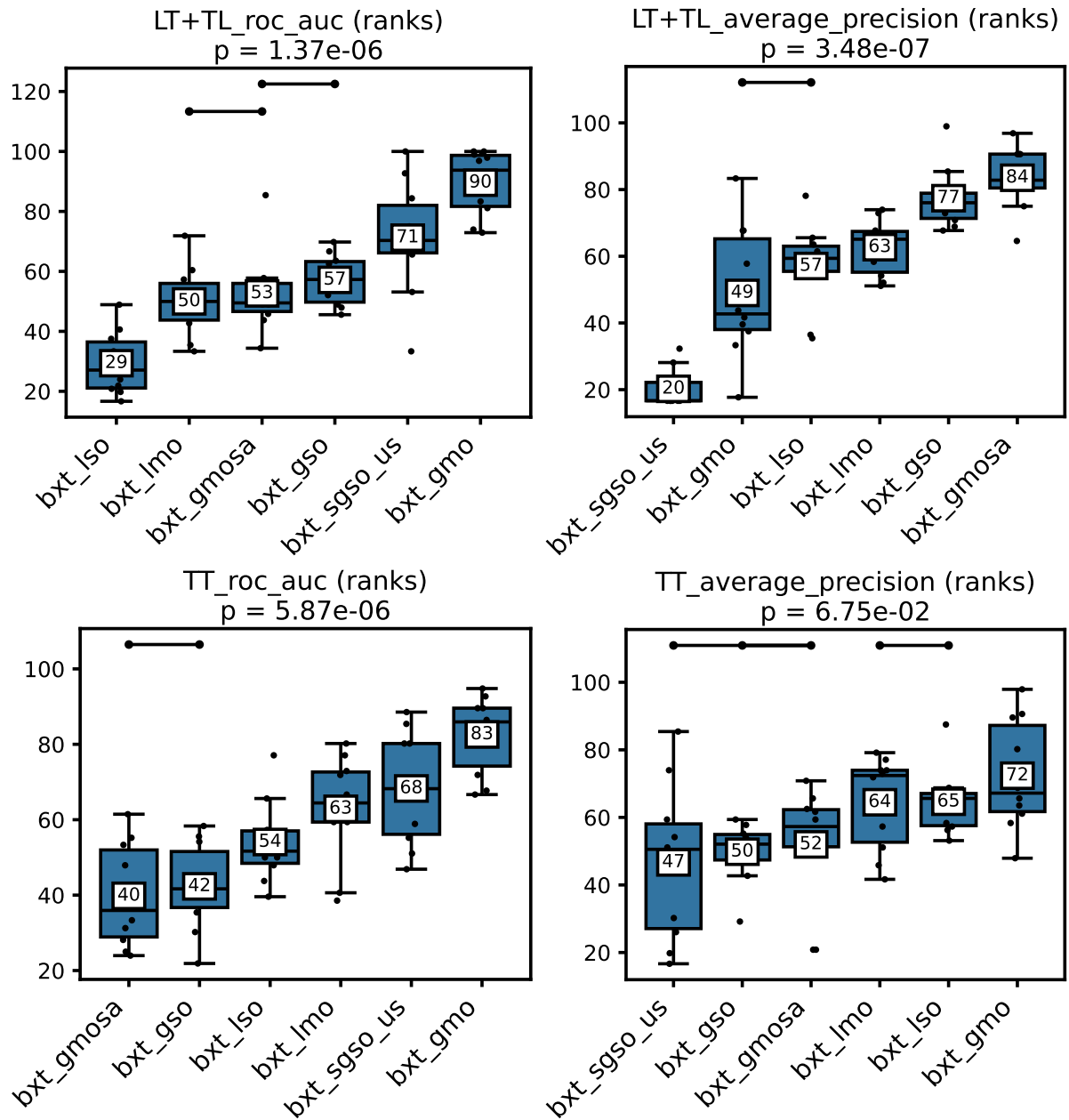


Figure 6 – Percentile rankings of prediction scores of extremely randomized tree ensembles under different adaptation strategies to bipartite interaction data, across all 10 datasets. The omnibus p-value obtained through a Friedman test is indicated below each plot's title. Significant results in each case lead us to perform pairwise Wilcoxon rank-sum tests, whose non-significant results are indicated by crossbars above each plot, that is, statistically undistinguished boxes are connected by the crossbars. The pairwise test results were corrected by the Benjamini-Hochberg procedure among each subfigure, considering that all possible estimator pairs were evaluated even if their comparison is not indicated in the plot. See Section 2.6.4 for further descriptions of each estimator. Source: made by the authors.

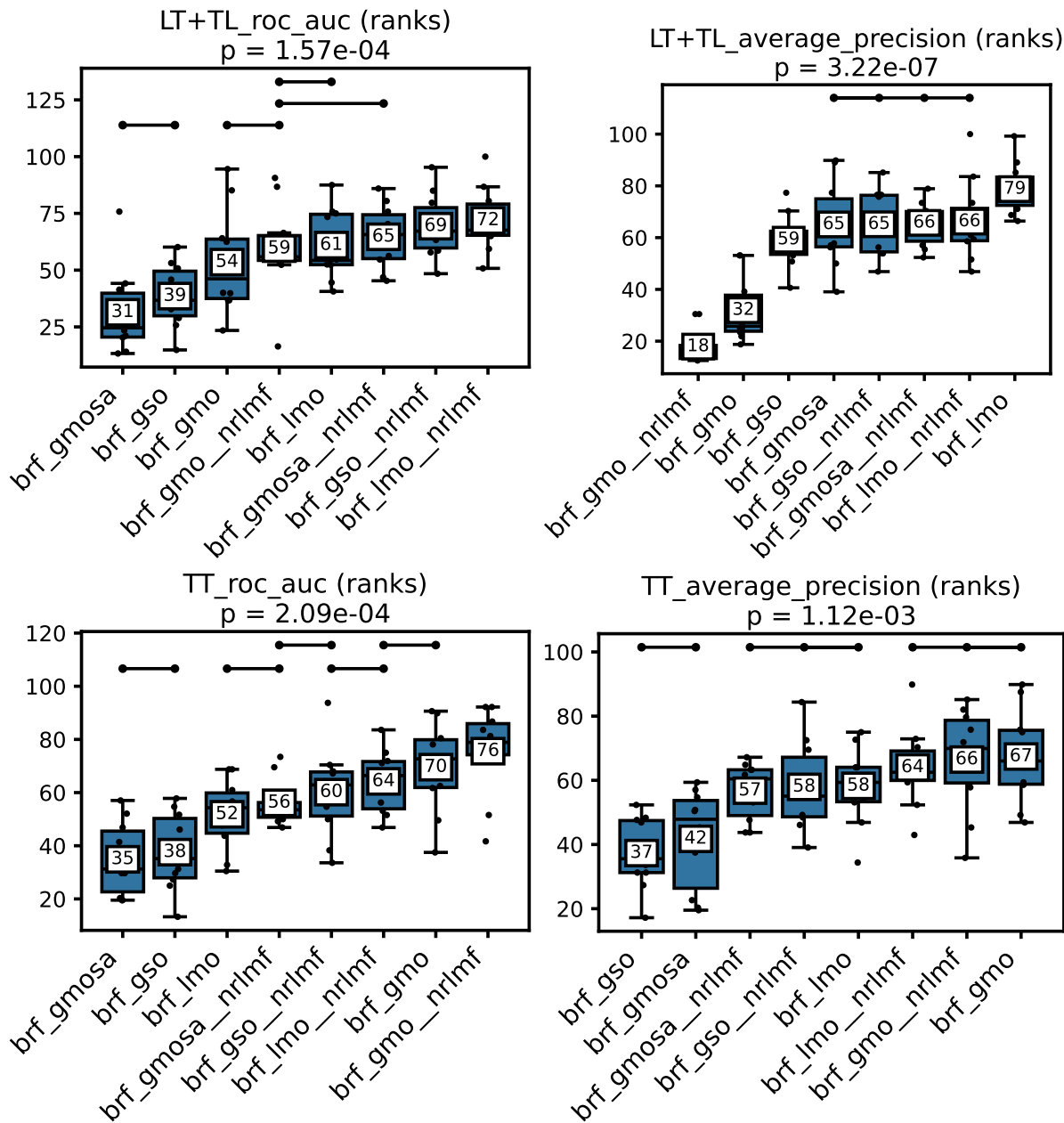


Figure 7 – Percentile rankings of prediction scores of extremely randomized tree ensembles under different adaptation strategies to bipartite interaction data, across all 10 datasets. The omnibus p-value obtained through a Friedman test is indicated below each plot's title. Significant results in each case lead us to perform pairwise Wilcoxon rank-sum tests, whose non-significant results are indicated by crossbars above each plot, that is, statistically undistinguished boxes are connected by the crossbars. The pairwise test results were corrected by the Benjamini-Hochberg procedure among each subfigure, considering that all possible estimator pairs were evaluated even if their comparison is not indicated in the plot. See Section 2.6.4 for further descriptions of each estimator. Source: made by the authors.

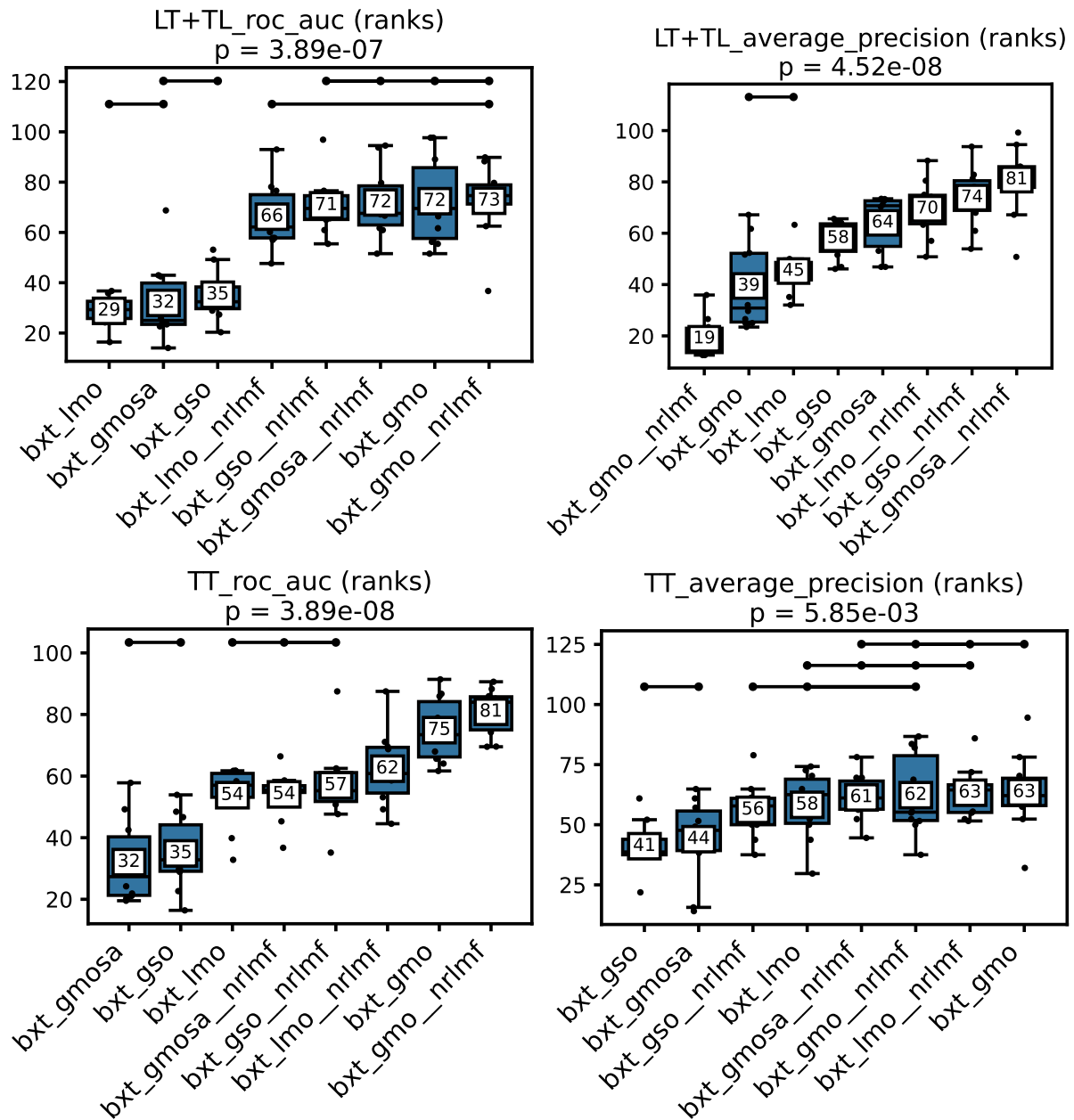


Figure 8 – Percentile rankings of prediction scores of extremely randomized tree ensembles under different adaptation strategies to bipartite interaction data, across all 10 datasets. The omnibus p-value obtained through a Friedman test is indicated below each plot’s title. Significant results in each case lead us to perform pairwise Wilcoxon rank-sum tests, whose non-significant results are indicated by crossbars above each plot, that is, statistically undistinguished boxes are connected by the crossbars. The pairwise test results were corrected by the Benjamini-Hochberg procedure among each subfigure, considering that all possible estimator pairs were evaluated even if their comparison is not indicated in the plot. See Section 2.6.4 for further descriptions of each estimator. Source: made by the authors.

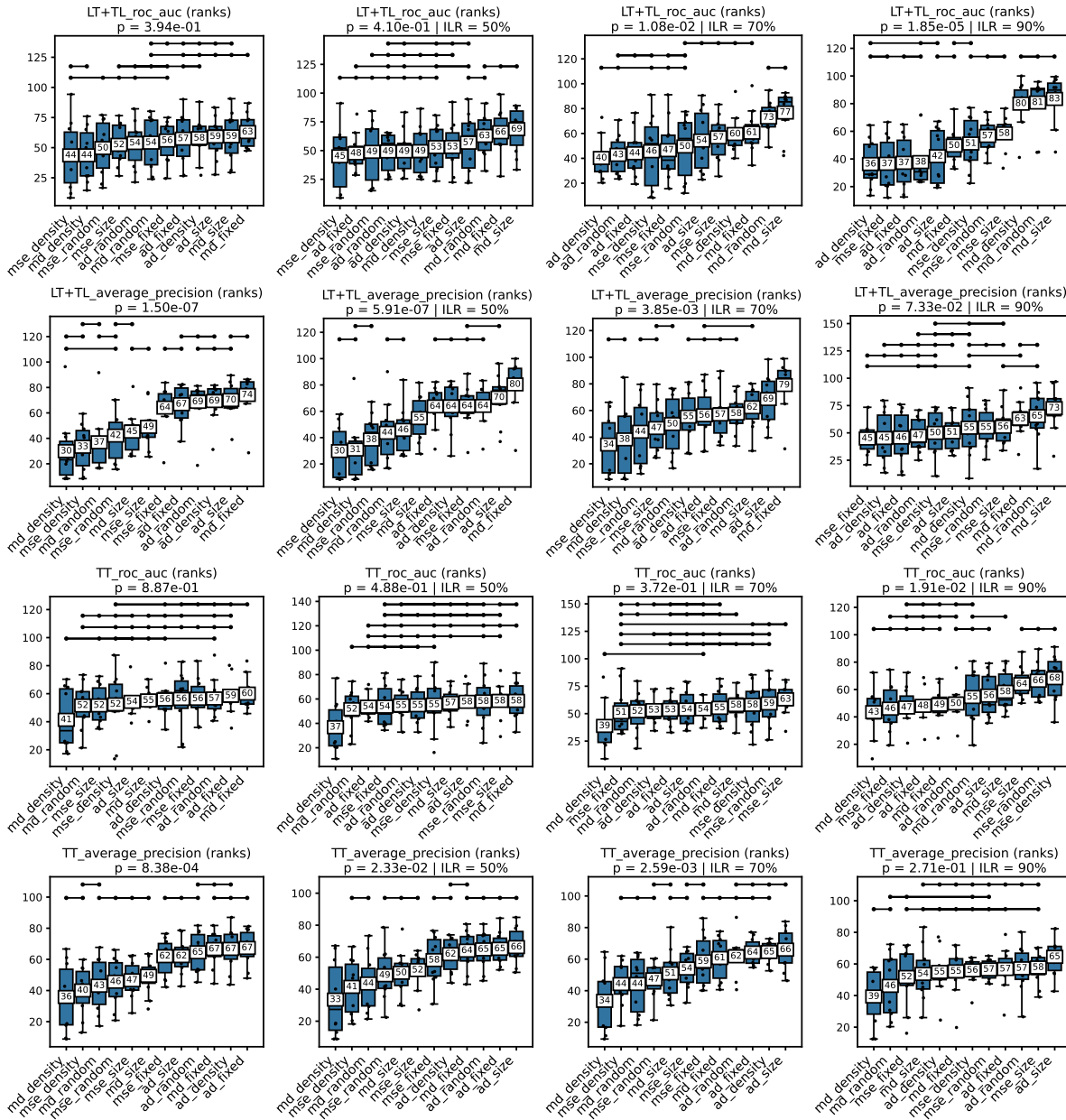


Figure 9 – Percentile rankings of prediction scores of extremely randomized tree ensembles under different adaptation strategies to bipartite interaction data, across all 10 datasets. Source: made by the authors.

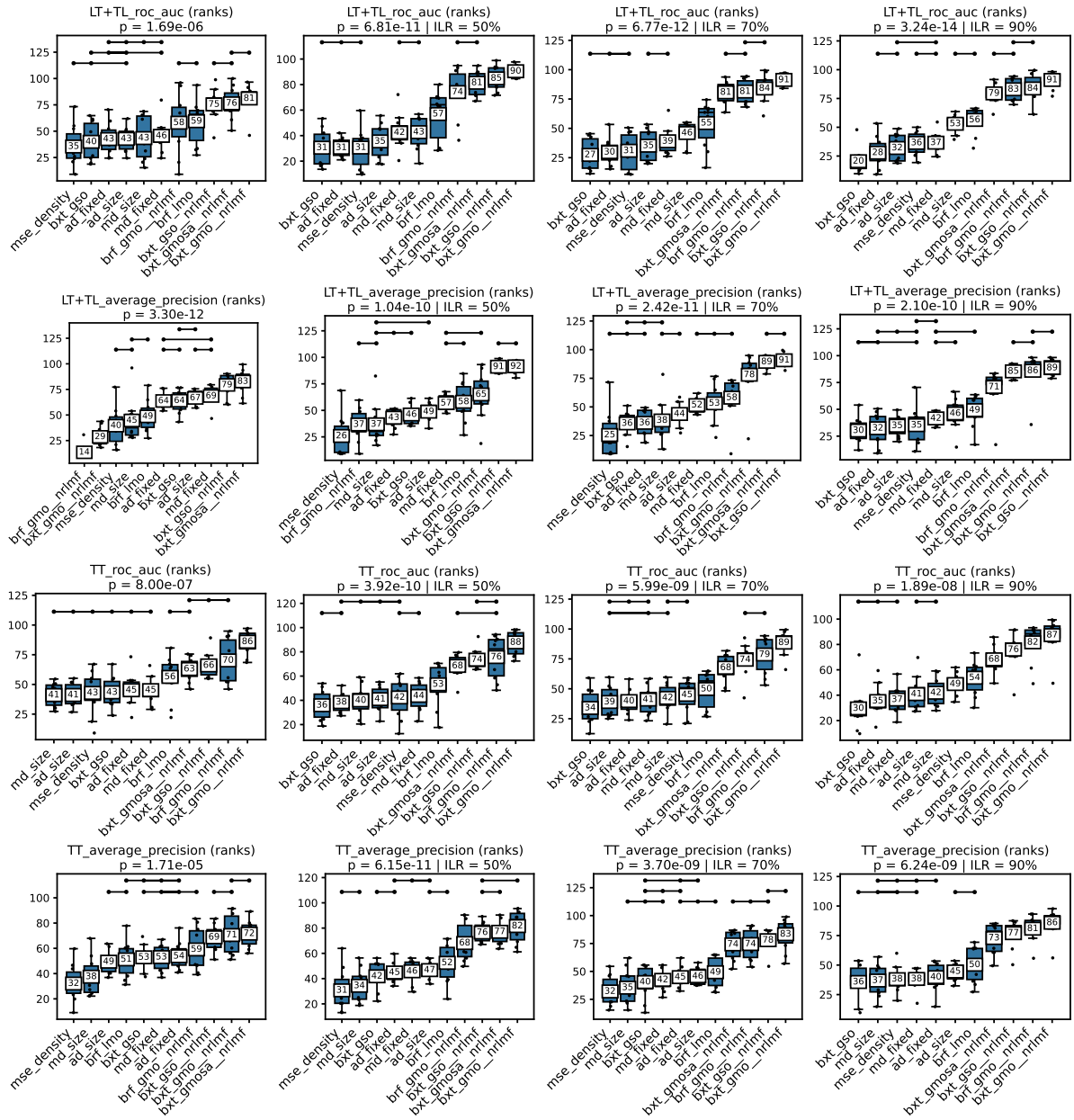


Figure 10 – Percentile rankings of prediction scores of extremely randomized tree ensembles under different adaptation strategies to bipartite interaction data, across all 10 datasets. Source: made by the authors.

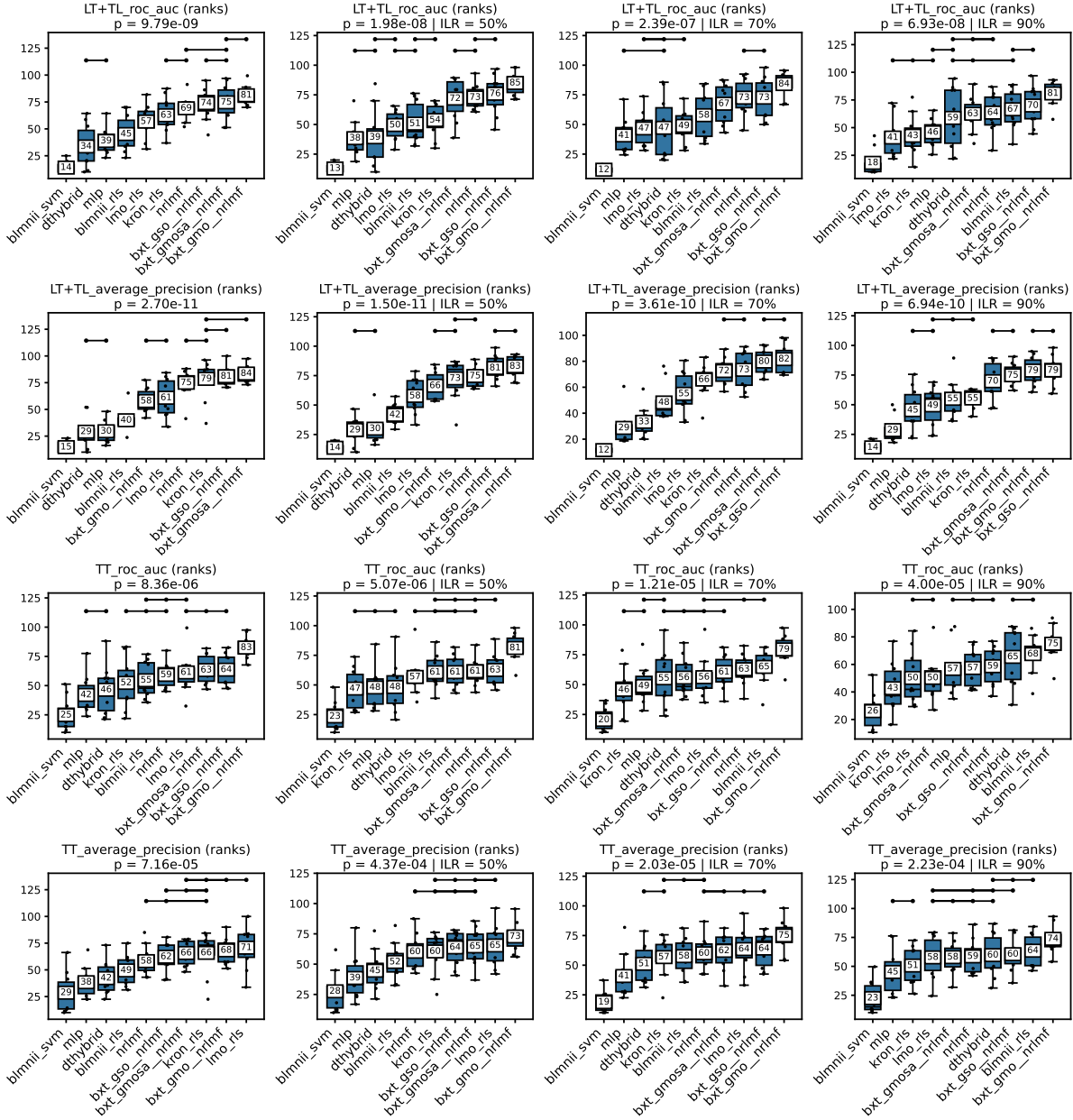


Figure 11 – Percentile rankings of prediction scores of extremely randomized tree ensembles under different adaptation strategies to bipartite interaction data, across all 10 datasets. Source: made by the authors.

3 CONCLUSION

In this chapter, we summarize the main findings of this work and discuss future research directions.

3.1 Main findings

We now list our final answers to the research questions proposed in ??.

1. Can bipartite trees be faster?

Yes, we developed significantly faster-growing bipartite trees by exploiting the bipartite nature of the problem.

We demonstrate both theoretically (Section 2.3.7) and empirically (Section 2.6.2) that our proposed bipartite tree algorithm achieves a $\log n$ improvement in training time complexity relative to its predecessors, with comparable predictive performance (Section 2.6.8).

2. Are semi-supervised techniques beneficial for bipartite forests?

Yes. Overall, we show that both semi-supervised impurities and label imputation by matrix factorization improve the predictions of bipartite forests.

For most scenarios, using our proposed semi-supervised impurity improves the predictive performance relative to the corresponding supervised model (Section 2.6.7). However, further investigation is needed to evaluate the benefits of other impurity functions. Additional experiments would also be interesting to investigate the effect of the methods that define the supervision balance parameter.

In any case, the most significant improvements were achieved by reconstructing the interaction matrix. As proposed by (?), we employ neighborhood-regularized matrix factorization (?) (NRLMF) to impute positive annotations prior to training the forests, and show that this approach consistently improves the scores in almost all settings (Section 2.6.5 and Section 2.6.7).

3. How do AUROC and AUPR differ in their assessment of model performance?

AUPR prioritizes a smaller number of highest-ranked interactions, while AUROC considers a larger number of both highest and lowest ranks.

We show that AUPR is very sensitive to which interactions from the test set the estimator selects as the most likely to occur. Also, AUPR mostly disregards the interactions selected as unlikely by the estimator. On the other hand, AUROC equally considers both the highest and lowest ranked interactions, and is less influenced by the extremal ranks in comparison to AUPR (Section 2.6.1).

Therefore, we argue that AUROC should be generally favored for model comparison under positive-unlabeled settings. One should favor AUPR instead when the main goal is to select a small number of most likely interactions.

4. **How do bipartite forests compare with proficient models in the field?**

The bipartite forests showed the best results in our experimental comparisons.

Notably, the BXT GMO NRLMF model stood out in most test settings analysed, especially when predicting interactions between unknown instances (Section 2.6.8). This model employs NRLMF label imputation and our proposed prototype function, demonstrating the effectiveness of these techniques in improving generalization.

3.2 **Main contributions**

ESCREVER

3.3 **Future work**

The present study opened several paths for future research. We list some of the most promising ones below.

1. **Adding noise to the impurity**

It is possible that using completely random values for the unsupervised impurity could still result in improved performance for the semi-supervised forests. This is because the diversity introduced by the noisy impurity to the individual trees could be enough, in theory, to improve the performance of the ensemble.

Hence, we propose adding noise to the supervised impurity and compare the results with supervised and semisupervised trees. The experiment could bring insights into the importance of different impurities in comparison with adding diversity to the trees.

2. **There is a lot of unsupervised impurities to be explored**

The present work explored two options of unsupervised impurity functions, that are utilized for training semi-supervised decision trees. However, unsupervised decision trees are a vast topic to explore. Many other approaches could be brought to the bipartite setting, with potential benefits for the predictive performance and scalability of the forests. For instance, several works use a single feature column to determine an unsupervised impurity (). Their reasoning is that the most significant change in impurity will result from the feature being used to select the split points. This feature is previously ordered by the algorithm. Therefore, other feature columns can be disregarded.

3. **Semi-supervised impurities for all forest models**

Only the bipartite global single-output ExtraTrees were trained with semi-supervised impurities. However, other models could benefit from the same approach, maybe even to a greater extent. Future investigations could explore the use of semi-supervised impurities with Random Forests and global multi-output forests. Another promising strategy would be to combine such impurities with label imputation by matrix factorization, since imputation by itself proved to be effective.

4. Other label imputation approaches

As previous studies, we used neighborhood-regularized matrix factorization (NRLMF) to impute positive annotations before training bipartite forests. The results were promising, but many other methods for label imputation could be explored.

The neighborhood regularization terms are not strictly necessary if the labels are used to train a second model. Their main importance is in inferring labels for new instances, which could be done solely by the downstream model in this case. Therefore, it is possible that simpler methods achieve similar results in less time (logistic matrix factorization (?), for example).

Additionally, there have been further developments of NRLMF that could be explored. NRLMF- β (?) adds beta-distribution-rescoring to the outputs of NRLMF. Dual network integrated logistic matrix factorization (?, ?) uses a concept similar to NRLMF, but merging the similarity matrices with the interaction matrix before factorization.

Other than matrix factorization, potentially any estimator could be utilized, even the same forest algorithms used for the final inference. This would be a form of self-learning, where the model is trained with the labels it has inferred (). In a similar fashion, co-training could be used, where models are trained in separate folds and then used to infer labels for the other models ().

5. Other forest algorithms

There is a variety of tree-based algorithms that could be adapted to the bipartite setting but were not explored in this work. Some of them are: Gradient Boosting Machines (), Rotation Forests (), Fuzzy Forests (), Oblique Random Forests () and Deep forests ().

6. Other fields

The present study concentrated on scenarios in which both instance domains have respective feature matrices. However, bipartite forests could also be extended to contexts where only one of the domains has side-features (dyadic prediction, multi-label learning, weak-label learning) or even to scenarios where both domains are featureless (collaborative filtering, recommendation systems). In these cases, features can be extracted from the interaction matrix itself, for instance, as gaussian interaction profiles (?).

For multi-output scenarios, bipartite forests can learn to partition the label matrix both horizontally and vertically, learning patterns that are specific for certain groups of outputs. This property is explored by (?) for single trees and hierarchical multi-label learning.

Weak-label learning could also be approached in a similar way. The difference from traditional multi-label learning is that negative annotations are not validated, similarly to the PU learning (?) setting of the present work. As such, the semi-supervised techniques we developed could be especially advantageous for weak-label learning.

7. Explainable artificial intelligence

Decision Forests have remarkable interpretability properties that were not explored in the present work. Future studies could exploit explanation techniques (?, ?, ?) to provide new insights into the nature of each learning task. For instance, in microRNA-gene interaction prediction, it would be interesting to pinpoint specific positions or sequence motifs that are relevant for the interaction.

COMPLEMENTAR

REFERENCES

- 1 BAGHERIAN, M. *et al.* Machine learning approaches and databases for prediction of drug–target interaction: A survey paper. **Briefings in Bioinformatics**, Oxford University Press (OUP), v. 22, n. 1, p. 247–269, jan. 2020.
- 2 CHEN, R. *et al.* Machine Learning for Drug-Target Interaction Prediction. **Molecules**, MDPI AG, v. 23, n. 9, p. 2208, ago. 2018.
- 3 FAITH, J. J. *et al.* Large-scale mapping and validation of Escherichia coli transcriptional regulation from a compendium of expression profiles. **PLoS biology**, Public Library of Science San Francisco, USA, v. 5, n. 1, p. e8, 2007.
- 4 LÜ, L. *et al.* Recommender systems. **Physics reports**, Elsevier, v. 519, n. 1, p. 1–49, 2012.
- 5 ASRATIAN, A. S.; DENLEY, T. M.; HÄGGKVIST, R. **Bipartite Graphs and Their Applications**. Cambridge, USA: Cambridge university press, 1998. v. 131.
- 6 PAHIKKALA, T. *et al.* Toward more realistic drug-target interaction predictions. **Briefings in Bioinformatics**, v. 16, n. 2, p. 325–337, 2015.
- 7 JOHNSON, C. C. Logistic Matrix Factorization for Implicit Feedback Data.
- 8 SCHRYNEMACKERS, M. *et al.* Classifying pairs with trees for supervised biological network inference. **Molecular BioSystems**, Royal Society of Chemistry, v. 11, n. 8, p. 2116–2125, 2015.
- 9 PLIAKOS, K.; GEURTS, P.; VENS, C. Global multi-output decision trees for interaction prediction. **Machine Learning**, Springer, v. 107, p. 1257–1281, 2018.
- 10 PLIAKOS, K.; VENS, C. Drug-target interaction prediction with tree-ensemble learning and output space reconstruction. **BMC bioinformatics**, Springer, v. 21, p. 1–11, 2020.
- 11 LAARHOVEN, T. V.; NABUURS, S. B.; MARCHIORI, E. Gaussian interaction profile kernels for predicting drug–target interaction. **Bioinformatics**, Oxford University Press, v. 27, n. 21, p. 3036–3043, 2011.
- 12 BREIMAN, L. Random forests. **Machine learning**, Springer, v. 45, p. 5–32, 2001.
- 13 BREIMAN, L. *et al.* **Classification and Regression Trees**. First edition. New York, USA: Routledge, 1984. ISBN 978-1-315-13947-0.
- 14 PLIAKOS, K.; VENS, C. Network inference with ensembles of bi-clustering trees. **BMC bioinformatics**, Springer, v. 20, p. 1–12, 2019.
- 15 LEVATÍĆ, J. *et al.* Semi-supervised classification trees. **Journal of Intelligent Information Systems**, Springer, v. 49, p. 461–486, 2017.
- 16 ALVES, A.; ILIDIO, P.; CERRI, R. Semi-Supervised Hybrid Predictive Bi-Clustering Trees for Drug-Target Interaction Prediction. *In: Proceedings of the 38th ACM/SIGAPP Symposium on Applied Computing*. Tallinn, Estonia: Association for Computing Machinery (ACM), 2023. p. 1163–1170.

- 17 VERT, J.-P. **Reconstruction of Biological Networks by Supervised Machine Learning Approaches**. [S.l.: s.n.]: arXiv, 2008.
- 18 GEURTS, P.; ERNST, D.; WEHENKEL, L. Extremely randomized trees. **Machine learning**, Springer, v. 63, p. 3–42, 2006.
- 19 CORMEN, T. H. *et al.* **Introduction to Algorithms**. Fourth edition. Cambridge, USA: MIT press, 2022. ISBN 978-0-262-04630-5.
- 20 POLIKAR, R. Ensemble based systems in decision making. **IEEE Circuits and Systems Magazine**, v. 6, n. 3, p. 21–45, 2006. ISSN 1531-636X.
- 21 LIU, H. *et al.* LPI-NRLMF: lncRNA-protein interaction prediction by neighborhood regularized logistic matrix factorization. **Oncotarget**, Impact Journals, LLC, v. 8, n. 61, p. 103975, 2017.
- 22 HE, T. *et al.* SimBoost: A read-across approach for predicting drug–target binding affinities using gradient boosting machines. **Journal of cheminformatics**, BioMed Central, v. 9, n. 1, p. 1–14, 2017.
- 23 LIU, Y. *et al.* Neighborhood Regularized Logistic Matrix Factorization for Drug-Target Interaction Prediction. **PLOS Computational Biology**, Public Library of Science (PLoS), v. 12, n. 2, p. e1004760, fev. 2016.
- 24 MURPHY, K. P. **Machine Learning: A Probabilistic Perspective**. Cambridge, MA: MIT Press, 2012. (Adaptive Computation and Machine Learning Series). ISBN 978-0-262-01802-9.
- 25 DUCHI, J.; HAZAN, E.; SINGER, Y. Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. **Journal of Machine Learning Research**, v. 12, n. 61, p. 2121–2159, 2011. ISSN 1533-7928.
- 26 TENG, X. *et al.* NPInter v4. 0: An integrated database of ncRNA interactions. **Nucleic acids research**, Oxford University Press, v. 48, n. D1, p. D160–D165, 2020.
- 27 HUANG, H.-Y. *et al.* miRTarBase update 2022: An informative resource for experimentally validated miRNA–target interactions. **Nucleic acids research**, Oxford University Press, v. 50, n. D1, p. D222–D230, 2022.
- 28 HUANG, K. *et al.* DeepPurpose: A deep learning library for drug–target interaction prediction. **Bioinformatics**, Oxford University Press, v. 36, n. 22-23, p. 5545–5547, 2020.
- 29 FRANKISH, A. *et al.* GENCODE 2021. **Nucleic acids research**, Oxford University Press, v. 49, n. D1, p. D916–D923, 2021.
- 30 COCK, P. J. *et al.* Biopython: Freely available Python tools for computational molecular biology and bioinformatics. **Bioinformatics**, Oxford University Press, v. 25, n. 11, p. 1422, 2009.
- 31 WU, T. *et al.* NPInter: The noncoding RNAs and protein related biomacromolecules interaction database. **Nucleic acids research**, Oxford University Press, v. 34, n. suppl_1, p. D150–D152, 2006.
- 32 HUANG, K. *et al.* MolTrans: Molecular Interaction Transformer for drug–target interaction prediction. **Bioinformatics**, Oxford University Press (OUP), v. 37, n. 6, p. 830–836, out. 2020.

- 33 HE, H.; GARCIA, E. A. Learning from Imbalanced Data. **IEEE Transactions on Knowledge and Data Engineering**, v. 21, n. 9, p. 1263–1284, set. 2009. ISSN 1558-2191.
- 34 FERNÁNDEZ, A. *et al.* **Learning from Imbalanced Data Sets**. Cham, Switzerland: Springer International Publishing, 2018. ISBN 978-3-319-98073-7 978-3-319-98074-4.
- 35 OZENNE, B.; SUBTIL, F.; Maucourt-Boulch, D. The precision–recall curve overcame the optimism of the receiver operating characteristic curve in rare diseases. **Journal of Clinical Epidemiology**, v. 68, n. 8, p. 855–859, ago. 2015. ISSN 0895-4356.
- 36 HU, Y.; KOREN, Y.; VOLINSKY, C. Collaborative Filtering for Implicit Feedback Datasets. *In: 2008 Eighth IEEE International Conference on Data Mining*. Pisa, Italy: IEEE Press, 2008. p. 263–272. ISSN 2374-8486.
- 37 EZZAT, A. *et al.* Computational prediction of drug–target interactions using chemogenomic approaches: An empirical survey. **Briefings in Bioinformatics**, v. 20, n. 4, p. 1337–1357, jul. 2019. ISSN 1477-4054.
- 38 HAO, M.; BRYANT, S. H.; WANG, Y. Open-source chemogenomic data-driven algorithms for predicting drug–target interactions. **Briefings in Bioinformatics**, v. 20, n. 4, p. 1465–1474, jul. 2019. ISSN 1477-4054.
- 39 KORNBROT, D. Point Biserial Correlation. *In: Wiley StatsRef: Statistics Reference Online*. Hatfield, UK: John Wiley & Sons, Ltd, 2014. ISBN 978-1-118-44511-2.
- 40 YAMANISHI, Y. *et al.* Prediction of drug–target interaction networks from the integration of chemical and genomic spaces. **Bioinformatics**, Oxford University Press, v. 24, n. 13, p. i232–i240, 2008.
- 41 BLEAKLEY, K.; YAMANISHI, Y. Supervised prediction of drug–target interactions using bipartite local models. **Bioinformatics**, v. 25, n. 18, p. 2397–2403, jul. 2009. ISSN 1367-4803.
- 42 MEI, J.-P. *et al.* Drug–target interaction prediction by learning from local information and neighbors. **Bioinformatics**, Oxford University Press, v. 29, n. 2, p. 238–245, 2013.
- 43 BAN, T.; OHUE, M.; AKIYAMA, Y. NRLMF_{upbeta}: Beta-distribution-rescored neighborhood regularized logistic matrix factorization for improving the performance of drug–target interaction prediction. **Biochemistry and Biophysics Reports**, Elsevier BV, v. 18, p. 100615, jul. 2019.
- 44 LI, Y.; LI, J.; BIAN, N. DNILMF-LDA: Prediction of lncRNA-disease associations by dual-network integrated logistic matrix factorization and Bayesian optimization. **Genes**, MDPI, v. 10, n. 8, p. 608, 2019.
- 45 BEKKER, J.; DAVIS, J. Learning from positive and unlabeled data: A survey. **Machine Learning**, v. 109, n. 4, p. 719–760, abr. 2020. ISSN 1573-0565.