



HIDDEN-RBP: Harnessing Intelligent Data Discovery to Explore Networks of RNA-Binding Proteins

Master's in Bioinformatics Thesis

Iresha De Silva

Supervisors:

Xavier Bofill De Ros and Thomas Bataillon

Spring 2025

Aarhus University

ABSTRACT

Understanding RNA-binding proteins (RBPs) is crucial for elucidating the regulatory mechanisms that control post-transcriptional gene expression and for identifying molecular pathways implicated in development and disease. In this study, I trained and compared supervised and unsupervised machine learning models to predict functional interactions involving RBPs using a combination of curated features and CRISPR-Cas9 derived phenotypic screening data. Two supervised models, Lasso and Random Forest, demonstrated strong classification performance on known interaction datasets, with the Random Forest model achieving 84% accuracy and a ROC-AUC of 0.90. In parallel, a Bayesian Latent Variable Model was employed to infer interaction likelihoods in an unsupervised manner, capturing hidden structure in the data by modeling uncertainty and leveraging prior knowledge.

The Bayesian model effectively identified latent signals associated with phenotypic similarity, particularly those derived from CRISPR screens, and generalized these patterns to unlabeled gene interaction pairs. Among 25,000 unverified pairs, 12 were identified with latent scores overlapping those of known interactions, including BRD4_TWIST2 and COX4I1_COX6B2, which show potential functional relevance based on existing literature. While posterior parameters convergence, referring to the stabilization of key latent variables or interaction strength estimates, was limited in some cases, particularly at lower training iterations or when key biological signals or features were missing, the findings nonetheless underline the value of probabilistic models in uncovering biologically meaningful patterns in sparse and noisy datasets.

These results highlight the critical role of CRISPR-Cas9 data in modeling RBP-related functions and suggest that incorporating additional features from RBPs datasets and structural, sequence, and pathway databases could further improve prediction accuracy. This integrative approach offers a framework for prioritizing novel RBP interactions and supports their downstream validation through experimental studies.

ACKNOWLEDGEMENTS

First and foremost, I would like to express my sincere gratitude to everyone who has supported me throughout the completion of this thesis.

I am deeply grateful to my supervisors, Asst. Professor Xavier Bofill De Ros and Professor Thomas Bataillon, for your invaluable guidance, continuous support, and patience during my research. Your immense knowledge and experience have been influential in shaping this thesis.

Your enthusiasm for the research is truly inspiring.

I would also like to thank the members of the XBR Lab Group and my DDSA mentor Myky Tran for the insightful discussions and help through the thesis for assistance in navigating and working with the dataset. Your contributions have significantly improved the quality of my work.

Special thanks to the Bioinformatics Research Center (BiRC) and the Department of Molecular Biology at Aarhus University for providing a stimulating environment. I am particularly thankful to my fellow Master's in Bioinformatics students for their moral support and the fun times we shared.

I would like to extend my heartfelt thanks to my family for their unwavering support and encouragement. To my husband and parents, thank you for your endless love and for believing in me.

ABBREVIATIONS

AUC	Area Under the Curve
Cas9	CRISPR-associated protein 9
CLIP-seq	Crosslinking and Immunoprecipitation sequencing
CRISPR	Clustered regularly interspaced palindromic repeats
DepMap	Dependency Map
DL	Deep learning
eCLIP	enhanced crosslinking and immunoprecipitation
ENCODE	Enhanced Crosslinking and Immunoprecipitation
ESS	Effective sample size
FN	False negative
FP	False positive
FPR	False Positive Rate
GI	Gene interaction
GO	Gene ontology
HITS-CLIP	High-throughput sequencing of RNA isolated by crosslinking immunoprecipitation
iCLIP	individual-nucleotide resolution cross-linking and immunoprecipitation
Lasso	Least Absolute Shrinkage and Selection Operator for classification
lincRNAs	long intergenic noncoding RNA
lncRNAs	long non-coding RNAs
MCMC	Markov Chain Monte Carlo
miRNAs	MicroRNAs
ML	Machine learning
mRNA	Messenger RNA
n_eff	Effective sample size
NIR	No information rate
OOB	Out-of-Bag
PAM	protospacer adjacent motif
POSTAR	Post-transcriptional Regulation Database

RBDs	RNA binding domains
RBPs	RNA binding proteins
RBPDB	RNA-Binding Protein Database
RBS	Ribosome binding sites
RF	Random forest
RIP-seq	RNA Immunoprecipitation sequencing
RNAi	RNA interference
RNPs	Ribonucleoprotein
ROC	Receiver Operating Characteristic
RRM	RNA recognition motifs
scaRNA	small Cajal body-specific RNA
sd	Standard deviation
se	Standard error
sgRNA	single-guide RNA
shRNA	small hairpin RNA
siRNA	Small interfering RNA
snoRNAs	small nucleolar RNA
TN	True negative
TP	True positive
TPR	True Positive Rate
X_{bin}	Binary variable
X_{cont}	Continuous variable
Z_i	Latent score

CONTENTS

ABSTRACT	ii
ACKNOWLEDGEMENTS	iii
ABBREVIATIONS	iv
CHAPTER 1: INTRODUCTION	1
1.1. Background	1
1.2. Research Questions.....	3
1.3. Research Objectives.....	4
1.4. Scope and Delimitations	4
1.5. Thesis Outline	5
CHAPTER 2: RELATED WORKS	6
2.1. Inference Methods for RNA-Binding Proteins (RBPs)	6
2.2. Bayesian Models for Complex Biological Systems.....	7
2.3. Machine Learning Models for Gene and Protein Interaction Inference	7
CHAPTER 3: BACKGROUND THEORY	10
3.1. RBPs in Gene Regulation	10
3.1.1. RNA Target Recognition and Complex Assembly Scaffolds.....	10
3.1.2. Role of RBPs in Development and Differentiation	10
3.1.3. RNA-Binding Proteins in Disease	11
3.1.4. Research Techniques for Studying RBPs	11
3.1.5. Addressing the Challenges in RBPs Identification	12
3.2. Cancer Dependency Map and CRISPR-Cas9 Screening	12
3.2.1. CRISPR-Cas9 Screening	13
3.2.2. Dependency Score	13
3.3. Machine Learning	14
3.3.1. Least Absolute Shrinkage and Selection Operator for classification (Lasso).....	14
3.3.1.1. Model Evaluation.....	15
3.3.2. Random Forest Model for Classification.....	15
3.3.3. Evaluating Classification Performance.....	16
3.3.4. Bayesian Latent Variable Model	17
3.3.4.1. Model Overview for Binary and Continuous Features	19

3.3.4.2.	Markov Chain Monte Carlo (MCMC) Sampling.....	20
CHAPTER 4: METHODS	22
4.1.	Data Curation and Preprocessing.....	22
4.1.1.	Achilles CRISPR-Cas 9 data	22
4.1.2.	Sanger CRISPR-Cas9 data.....	23
4.1.3.	ShRNA Screen data	23
4.1.4.	Additional Data.....	23
4.1.5.	Benchmarking Data	26
4.1.6.	Defining True and Non-Interacting Gene Pairs	26
4.2.	Train, Validation and Test Split.....	27
4.3.	Supervised Learning Model.....	27
4.3.1.	Supervised Learning Technique – Lasso	28
4.3.2.	Supervised Learning Technique – Random Forest.....	28
4.3.3.	Unsupervised Learning Model – Bayesian Latent Variable Model.....	28
4.3.3.1.	Benchmarking Model.....	29
4.3.3.2.	Prediction for Unknown Set.....	31
CHAPTER 5: RESULTS AND DISCUSSION	32
5.1.	Brief Overview of the Data.....	32
5.1.1.	Sample Overview of Co-Dependency Data.....	32
5.1.2.	Feature Illustration	33
5.1.2.1.	Overview of Gene-Gene Correlations Across Datasets.....	33
5.1.2.2.	Integration of Additional Biological Features and Benchmarking Resources.....	35
5.2.	Supervised Learning Lasso	38
5.2.1.	Model Training	38
5.2.2.	Model Performance.....	39
5.3.	Supervised Learning Model – Random Forest	41
5.3.1.	Model Training	41
5.3.2.	Model Performance.....	42
5.4.	Bayesian Latent Variable Model	43
5.4.1.	Model Training on Benchmarking Data	43
5.4.2.	Model Performance.....	46

5.4.3.	Model Training on Unknown Data	48
5.4.4.	Comparison Between Benchmarking and Inferred Latent Scores	50
5.5.	Discussion and Future Perspectives.....	52
CHAPTER 6: CONCLUSION.....		54
REFERENCES.....		55
APPENDIX A		62
APPENDIX B		65