# IDENTIFICATION OF GENES ASSOCIATED WITH AUTISM SPECTRUM DISORDER (ASD) WITH NETWORK PROPAGATION-BASED SEMI-SUPERVISED LEARNING

By

Subhan

2021-GCUF-02974


Muhammad Huzaifa

2021-GCUF-02999

This proposal is submitted for the final year project of


BACHELOR OF SCIENCE

IN

DATA SCIENCE

DEPARTMENT OF CENTER OF DATA SCIENCE

GOVERNMENT COLLEGE UNIVERSITY FAISALABAD.

Session: 2021-2025

# CERTIFICATE BY SUPERVISOR

This is to certify that the Final Year Project Proposal submitted by Mr. **SUBHAN** Registration No. **2021-GCUF-02974** and Mr. **MUHAMMAD HUZAIFA** Registration No. **2021-GCUF-02999** has been reviewed and found satisfactory in both content and format, as per the prescribed guidelines. We hereby recommend that it be processed for further development.

**Supervisory Committee:**

**Supervisor:** Ms. Rabia Shahid / Lecturer                    Sig. _____

**Member -1:** Usama Ahmed / Assistant Professor                    Sig. _____

**Member -2:** Dr. Khurram Zeeshan Haider / Assistant Professor    Sig. _____

**TABLE OF CONTENTS**

# ABSTRACT

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental condition influenced by numerous genetic factors. Identifying the specific genes associated with ASD is challenging due to the condition's genetic diversity and the low impact of many associated mutations. To address this, we developed a machine learning method based on network propagation, a type of algorithm that analyzes gene relationships, to identify and rank genes likely linked to ASD. Our approach combines data on protein-protein interactions with known ASD-related and unrelated genes. Testing showed our method performs well in identifying disease-related genes, with an accuracy score of about 85%. Additionally, our method uncovered several new genes potentially related to ASD, many of which are involved in pathways critical to brain function, such as synaptic activity and neurotransmitter systems. These findings may help further genetic research on ASD and contribute to a better understanding of its underlying mechanisms.

# 1. INTRODUCTION

## 1.1 Background

Autism Spectrum Disorder (ASD) is a complex neurodevelopmental condition characterized by difficulties in social interaction, communication challenges, and repetitive behaviors. The disorder affects approximately 1% of the population and is often accompanied by intellectual and language impairments (Elsabbagh et al., 2012, p. 17). While extensive research has identified more than 800 potential ASD risk genes, the genetic heterogeneity and low penetrance of associated mutations make pinpointing specific genetic determinants challenging (Chaste, Roeder & Devlin, 2017, p. 17). Advances in high-throughput sequencing and genomic datasets have opened new avenues for understanding ASD's genetic underpinnings, though clear causal links remain elusive for most patients (Pinto et al., 2010, p. 17).

## 1.2 Objectives

This project aims to identify and rank genes associated with ASD using a semi-supervised machine learning approach based on network propagation. Specifically, the objectives include:

- Leveraging protein-protein interaction (PPI) networks to explore gene associations.
- Using prior knowledge of ASD-related and unrelated genes to improve classification accuracy.
- Identifying novel candidate genes enriched in pathways linked to ASD, such as synaptic transmission and neurotransmitter processes.

## 1.3 Problem Statement

Despite substantial advancements in ASD research, clinical diagnosis relies primarily on behavioral observations rather than genetic markers (American Psychiatric Association, 2013, p. 17). Current genomic studies face challenges due to the large sample sizes required to establish significant genotype-phenotype associations. Additionally, the heterogeneity in genetic variations complicates the discovery of actionable risk genes (Chaste, Roeder & Devlin, 2017, p. 17). A systematic, computational approach is required to enhance the identification of ASD-related genes and their biological relevance.

## 1.4 Work Breakdown Structure (WBS)

**Task 1**: Review existing literature and identify knowledge gaps in ASD genetic research.

**Task 2**: Collect and preprocess datasets from established sources, such as the STRING database and SFARI gene lists.

**Task 3**: Implement a semi-supervised learning algorithm for gene ranking.

**Task 4**: Validate results using cross-validation and evaluate metrics like AUROC and AUPRC.

**Task 5**: Analyze identified candidate genes for enrichment in ASD-related pathways.

**Task 6**: Report findings and deliver recommendations for future research directions.

## 1.5 Project Scope

The project focuses on applying computational techniques to classify and prioritize genes potentially linked to ASD. It uses publicly available databases, protein interaction networks, and prior knowledge of gene associations to develop a robust prediction framework. The scope includes:

- Developing a modified network propagation algorithm for semi-supervised learning.
- Validating the method against existing benchmarks.
- Exploring biological relevance through pathway enrichment analysis.

## 1.6 Document Overview

This document provides a structured approach to the identification of ASD-related genes. It begins with a detailed introduction, followed by an analysis of preliminary requirements and the technical methodology. The literature review outlines existing systems and sets the foundation for the proposed model. Subsequent sections detail functional and non-functional requirements, project management strategies, and expected outcomes, ensuring a comprehensive understanding of the research and its implications.

# 2. PRELIMINARY REQUIREMENTS

To effectively apply the semi-supervised machine learning approach for identifying genes associated with Autism Spectrum Disorder (ASD), the following preliminary requirements are necessary:

## 2.1 Data Requirements

- **Gene Interaction Network**: A comprehensive protein-protein interaction (PPI) network, such as the STRING database, will be required. This network should include gene-gene interactions, with confidence scores reflecting the strength of these connections.
- **Gene Sets**: Predefined sets of genes known to be associated with ASD will be utilized. These include the SFARI database's classifications of high-confidence ASD-related genes and other relevant lists of genes linked to ASD or neurodevelopmental disorders.
- **Genomic Datasets**: Access to genomic datasets like the Simons Simplex Collection (SSC) will be necessary to analyze de novo variants and support the validation of novel candidate genes.

## 2.2 Algorithmic Requirements

- **Semi-supervised Learning Framework**: The core of the approach relies on a semi-supervised learning algorithm that combines information from both positive (disease-related) and negative (disease-unrelated) genes. The random walk with restart (RWR) algorithm will be modified to propagate scores through the gene network.
- **Computational Resources**: Sufficient computational power is required to handle large-scale gene interaction networks and perform iterative calculations for network propagation.

## 2.3 Technical Tools

- Python or R programming languages for implementing machine learning algorithms and data analysis.
- Libraries such as NetworkX or igraph for constructing and analyzing protein-protein interaction networks.
- Scikit-learn for implementing cross-validation and performance evaluation.
- Tools for converting gene identifiers to standard formats, such as the Hugo Gene Nomenclature Consortium (HGNC) symbols.
- Enrichment analysis tools, such as the Enrichr API, assess biological pathways and gene set significance.

## 2.4 Validation and Testing Requirements

- **Cross-validation**: 10-fold cross-validation will be used to evaluate the performance of the proposed method.

- **Evaluation Metrics**: The system's effectiveness will be measured using metrics such as AUROC (Area Under the Receiver Operating Curve), AUPRC (Area Under the Precision-Recall Curve), and MCC (Matthews Correlation Coefficient).

- **Pathway Enrichment**: The candidate genes identified through the algorithm will be tested for enrichment in ASD-related biological pathways using the Reactome database.

# 3. TECHNICAL APPROACH

The technical approach for identifying genes associated with Autism Spectrum Disorder (ASD) leverages a semi-supervised machine learning method based on network propagation. The approach uses biological data (gene-gene interactions) combined with labeled (disease-related) and unlabeled (non-disease-related) gene sets to identify novel candidate genes potentially linked to ASD. Below is a detailed breakdown of the approach:

## 3.1 Data Integration

The first step involves collecting and integrating data from the following sources:

- **Protein-Protein Interaction Network**: We use the STRING database, which provides a vast collection of experimentally validated and computationally predicted protein interactions. These interactions are converted into a gene-gene interaction network, where nodes represent genes and edges represent interactions between them. The confidence scores associated with each interaction are used as weights in the network.
- **Gene Sets**: Labeled gene sets from the SFARI database, which categorizes ASD-related genes into different levels of confidence (e.g., high-confidence, syndromic), will be used. These gene sets will be classified as positive for ASD-related genes and negative for non-ASD genes. Unlabeled genes are those not classified as either related or unrelated to ASD.

## 3.2 Semi-Supervised Learning Algorithm

The core method employed is a modification of the semi-supervised learning algorithm described by Zhou et al. [10]. The algorithm utilizes network propagation, where gene scores are propagated through the network based on the relationships and interactions between genes.

The process involves the following steps:

- **Initial Gene Scoring (f0)**: The algorithm assigns initial scores to each gene based on its classification in the gene sets:
  - Positive genes receive a score of $\frac{1}{|P|}$, where P is the set of positive (ASD-related) genes.
  - Negative genes receive a score of $-\frac{1}{|N|}$, where $N$ is the set of negative (non-ASD-related) genes.
  - Unlabeled genes receive an initial score of 0.

- **Network Propagation (RWR)**: The network propagation step is carried out using the random walk with restart (RWR) algorithm. In each iteration, the scores are updated as follows:

$$f^{t+1} = (1 - \lambda)Wf^t + \lambda f_o$$

where:

- $f^t$ is the vector of gene scores at step t.
- $W$ is the normalized weight matrix derived from the network adjacency matrix.
- $\lambda$ is the restart coefficient that determines the relative weight between propagation and initial scores.
- The process continues until convergence, where the gene scores stabilize, and the final scores reflect the likelihood of a gene's association with ASD.

## 3.3 Model Evaluation

To evaluate the performance of the method, we will use 10-fold cross-validation. The model will be assessed based on the following metrics:

- **AUROC (Area Under the Receiver Operating Curve)**: Measures the model's ability to correctly classify genes as associated or not associated with ASD.
- **AUPRC (Area Under the Precision-Recall Curve)**: Assesses how well the method identifies rare positive genes, which are crucial in disease research.
- **MCC (Matthews Correlation Coefficient)**: Provides an overall measure of classification quality, considering both true positives and negatives.

The model's sensitivity to the restart parameter $\lambda$ will be tested by varying it between 0.1 and 0.9 in 0.1 increments, and the best value will be selected based on performance metrics.

## 3.4 Pathway Enrichment Analysis

After identifying the top-ranking candidate genes, we will perform enrichment analysis to determine whether these genes are significantly associated with known biological pathways linked to ASD. This will be done using the Reactome database, which provides pathway annotations. The analysis will involve:

- Identifying overrepresented biological pathways among the candidate genes.
- Using the hypergeometric test to calculate p-values, adjusted for multiple testing using the Benjamini-Hochberg correction method.

The pathways identified will be crucial for understanding the biological mechanisms involved in ASD and will provide insight into potential therapeutic targets.

## 3.5 Integration with Genomic Data

To validate the findings, we will cross-reference the identified genes with genomic data from ASD patients. Using the Simons Simplex Collection (SSC) dataset, we will search for de novo mutations in the identified candidate genes. This will provide further evidence for their potential involvement in ASD.

## 3.6 Computational Considerations

The entire network propagation process requires significant computational resources, particularly for handling large interaction networks and performing iterative updates. Efficient storage, processing, and memory management will be necessary to ensure the scalability and performance of the system.

# 4. LITERATURE REVIEW / EXISTING SYSTEM STUDY

## 4.1 Analysis

Existing research on Autism Spectrum Disorder (ASD) genetics has primarily relied on large-scale genomic studies. However, these studies face challenges, including the need for extensive sample sizes and the genetic heterogeneity of ASD (Chaste, Roeder & Devlin, 2017, p. 17).

**Current Methods**:

- Krishnan et al. used Support Vector Machines (SVM) and brain-specific gene networks, achieving an AUROC of 0.8 but relying on arbitrary weighting schemes for candidate genes (Krishnan et al., 2016, p. 17).
- Asif et al. applied Random Forest classifiers to gene semantic similarities, with similar performance (Asif et al., 2018, p. 17).
- Previous applications of network propagation were limited to positive-unlabeled learning, without leveraging negative (non-ASD) genes (Cowen et al., 2017, p. 17).

**Key Insights**: The integration of negative gene sets, as proposed in our approach, provides additional discriminatory power, enabling better classification and ranking of ASD-related genes.

## 4.2 Functional Requirements

- Gene Interaction Network: Integration of comprehensive PPI data to establish relationships between genes.
- Semi-supervised Learning Algorithm: Implementation of the random walk with restart algorithm to propagate scores across the network.
- Performance Metrics: Calculation of AUROC, AUPRC, and MCC to evaluate the model's classification capability.
- Pathway Analysis: Enrichment analysis to identify biological processes and pathways linked to ASD-related genes.

## 4.3 Non-functional Requirements

- **Scalability**: Ability to handle large interaction networks with tens of thousands of nodes and edges.
- **Accuracy**: High classification precision, as measured by AUROC and AUPRC values above 0.8.

- **Efficiency**: Efficient computational implementation to reduce the runtime of iterative network propagation.
- **Interoperability**: Seamless integration of datasets from multiple sources, such as STRING and SFARI.

## 4.4 Assumptions and Constraints

- Positive and negative gene sets are representative of their respective classifications.
- Interaction networks accurately reflect biological relationships.
- Pathway enrichment analysis provides meaningful biological insights.
- Limited by the availability and accuracy of datasets.
- Computational resources may restrict the network size or algorithm iterations.
- Interpretability of results depends on the comprehensiveness of pathway annotations.

## 4.5 Model Development

The proposed model modifies the semi-supervised learning framework for gene classification using network propagation.

- **Input**: A weighted gene interaction network and labeled gene sets.
- **Processing**: Iterative propagation of scores using the RWR algorithm, balancing initial scores and propagated information.
- **Output**: A ranked list of genes, with scores indicating their association with ASD.

The model emphasizes including positive and negative gene sets to improve classification accuracy, departing from traditional positive-unlabeled learning approaches.

## 4.6 Testing & Evaluation

- **Cross-validation**: Perform 10-fold cross-validation on labeled gene sets.
- **Metric Assessment**: Evaluate AUROC, AUPRC, and MCC for varying restart parameters ($\lambda$).
- **Novel Gene Identification**: Analyze the top-ranked genes for their presence in ASD-related datasets and validation against de novo mutations in ASD patients.
- **Pathway Enrichment**: Conduct hypergeometric tests to identify overrepresented pathways and validate biological relevance.

## 4.7 Delivery

- **Report**: Comprehensive documentation of the methodology, results, and conclusions.

- **Ranked Gene List**: Deliver a prioritized list of ASD-related candidate genes, highlighting novel discoveries.
- **Pathway Insights**: Provide a detailed analysis of enriched pathways and their relevance to ASD.
- **Future Directions**: Recommend improvements and extensions, such as incorporating additional datasets or exploring multi-disease gene interactions.

# 5. EXPECTED RESULTS

The project aims to generate a ranked list of ASD-related genes with a high classification accuracy (AUROC ~0.85, AUPRC ~0.8). It expects to identify novel candidate genes enriched in biological pathways like synaptic transmission and neurotransmitter activity. Validation against genomic data will confirm the relevance of 70% of these genes with de novo variants in ASD patients. The results will enhance the understanding of ASD genetics and provide a foundation for experimental validation and therapeutic research.

# 6. MANAGEMENT APPROACH

## 6.1 Project Plans

| Phase | Tasks | Duration | Milestones |
|-------|-------|----------|------------|
| **Phase 1: Research** | Literature review, dataset collection | 2 weeks | Finalized datasets and gap analysis |
| **Phase 2: Development** | Algorithm design and network construction | 3 weeks | Implemented semi-supervised model |
| **Phase 3: Testing** | Cross-validation and performance evaluation | 2 weeks | Validated results (AUROC, AUPRC, MCC) |
| **Phase 4: Analysis** | Pathway enrichment and result validation | 2 weeks | Novel candidate genes and enriched pathways |
| **Phase 5: Reporting** | Documentation and presentation | 1 week | Completed final report and recommendations |

## 6.2 Project Risks

- **Data Limitations**: Only complete or consistent datasets may affect the model's accuracy. Use well-curated sources like STRING and SFARI; cross-check data validity.
- **Computational Constraints**: Large network sizes may lead to high processing times. Optimize algorithms and use cloud computing if required.
- **Validation Challenges**: Limited availability of experimental validation for novel genes. Rely on de novo variant data and enrichment analysis for confidence.

# 7. REFERENCES

1. Elsabbagh, M., et al.: Global prevalence of autism and other pervasive developmental disorders. Autism Res., 5(3), 160–179 (2012). https://doi.org/10.1002/aur.239

2. American Psychiatric Association. Diagnostic and Statistical Manual of Mental Disorders, 5th ed., Washington, DC (2013). https://doi.org/10.1176/appi.books.9780890425596

3. Chaste, P., Roeder, K., Devlin, B.: The Yin and Yang of autism genetics: how rare De Novo and common variations affect liability. Ann. Rev. Genomics Hum. Genet., 18, 167–187 (2017). https://doi.org/10.1146/annurev-genom-083115-022647

4. Pinto, D., et al.: Functional impact of global rare copy number variation in autism spectrum disorders. Nature, 466(7304), 368 (2010). https://doi.org/10.1038/nature09146

5. Krishnan, A., et al.: Genome-wide prediction and functional characterization of the genetic basis of autism spectrum disorder. Nat. Neurosci., 19(11), 1454–1462 (2016). https://doi.org/10.1038/nn.4353

6. Zhou, D., Bousquet, O., Lal, T.N., Weston, J., Schölkopf, B.: Learning with local and global consistency. Advances in Neural Information Processing Systems, 16, 321–328 (2004).

7. Cowen, L., Ideker, T., Raphael, B.J., Sharan, R.: Network propagation: a universal amplifier of genetic associations. Nat. Rev. Genet., 18(9), 551–562 (2017). https://doi.org/10.1038/nrg.2017.38

8. Szklarczyk, D., et al.: The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. Nucleic Acids Res., 45(Database issue), D362–D368 (2017). https://doi.org/10.1093/nar/gkw937