# Machine Learning Analysis of Protein Function Prediction

Matthew Segar

Tong Wu

# Outline

Introduction

Methods

Results

Conclusion

# Outline

**Introduction**

Methods

Results

Conclusion

# A Dichotomy Exists

Sequencing has become increasingly easy

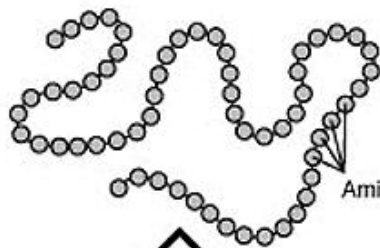Databases (SwissProt) make it easy to upload data

Data that is sequenced >>>> data that is annotated

# Protein Function

"Anything and everything that happens to or through a protein" (Syed)
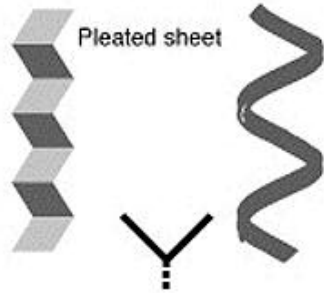
From sequence alone, the best predictors are only 90% accurate (Rost, Whisstock)
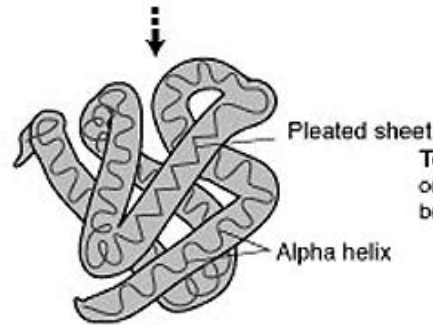
# Protein Structure



**Primary protein structure** is sequence of a chain of amino acids
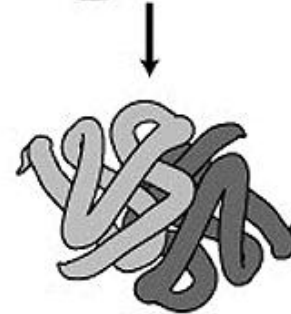
Amino Acids

Pleated sheet

Alpha helix

**Secondary protein structure** occurs when the sequence of amino acids are linked by hydrogen bonds

Pleated sheet

Alpha helix

**Tertiary protein structure** occurs when certain attractions are present between alpha helices and pleated sheets.

**Quaternary protein structure** is a protein consisting of more than one amino acid chain.

# Outline

Introduction

**Methods**

Results

Conclusion

# Data

Randomly selected 5000 protein sequences from SwissProt

Used *Reactome* to determine GO annotation

Input list of Uniprot IDs

Output broad function

# Reactome

## Analysis results, per pathway

This table provides an overview of your expression data in a pathway context. For each Reactome pathway, the total number of proteins is shown, plus the number of genes/proteins in your dataset that match. By clicking on a pathway name, you will be taken to an interactive graphical representation of the pathway, where your expression levels are represented as coloration of proteins.

Select format to download this table: microsoft xcel (tsv) | Download

| Pathway ▼▲ | Species ▼▲ | IDs in pathway (%) ▼▲ | Enrichment (pval) ▼▲ | FDR ▼▲ |
|---|---|---|---|---|
| Not assigned | Not known | 75 (0%) | | |
| Apoptosis | Homo sapiens | 73 (48%) | 2.53E-05 | 1.52E-04 |
| Binding and Uptake of Ligands by Scaveng | Homo sapiens | 5 (2%) | | |
| Cell Cycle | Homo sapiens | 168 (35%) | 2.55E-09 | 5.09E-08 |
| Cell-Cell communication | Homo sapiens | 18 (14%) | 8.69E-02 | 8.69E-02 |
| Cellular responses to stress | Homo sapiens | 36 (14%) | | |
| Chromatin organization | Homo sapiens | 4 (4%) | | |
| Circadian Clock | Homo sapiens | 4 (11%) | 1.87E-02 | 3.75E-02 |
| Developmental Biology | Homo sapiens | 62 (15%) | 1.87E-02 | 3.74E-02 |
| Disease | Homo sapiens | 280 (25%) | 3.24E-03 | 9.73E-03 |
| DNA Repair | Homo sapiens | 58 (54%) | 6.39E-02 | 6.39E-02 |
| DNA Replication | Homo sapiens | 57 (56%) | 2.28E-06 | 2.28E-05 |
| Extracellular matrix organization | Homo sapiens | 53 (21%) | 0.1 | 0.1 |
| Gene Expression | Homo sapiens | 229 (31%) | 3.98E-02 | 3.98E-02 |
| Hemostasis | Homo sapiens | 116 (27%) | 1.55E-03 | 7.46E-03 |
| Immune System | Homo sapiens | 233 (19%) | 1.86E-03 | 7.46E-03 |
| Meiosis | Homo sapiens | 18 (31%) | 4.17E-02 | 4.17E-02 |
| Membrane Trafficking | Homo sapiens | 13 (9%) | 0.5 | 0.5 |
| Metabolism | Homo sapiens | 353 (24%) | 1.51E-02 | 3.02E-02 |
| Metabolism of proteins | Homo sapiens | 65 (11%) | 0.9 | 0.9 |
| Muscle contraction | Homo sapiens | 1 (1%) | 0.2 | 0.2 |
| Neuronal System | Homo sapiens | 22 (8%) | 0.7 | 0.7 |
| Organelle biogenesis and maintenance | Homo sapiens | 8 (15%) | | |
| Reproduction | Homo sapiens | 0 (0%) | | |
| Signal Transduction | Homo sapiens | 229 (12%) | 1.82E-02 | 3.64E-02 |
| Transmembrane transport of small molecul | Homo sapiens | 43 (8%) | 0.6 | 0.6 |

26 rows

# Filter

Only 3,085 protein sequences had functions

Removed functions with < 40 samples

Each function needed to constitute 1% of the total dataset

Resulted in 19 functions and 3,039 sequences

# Data

| Row | Function | Count | Percentage |
| --- | --- | --- | --- |
| 1 | Apoptosis | 52 | 1.66% |
| 2 | Binding and uptake of ligands by scavenger | 110 | 3.51% |
| 3 | Cell Cycle | 167 | 5.33% |
| 4 | Cell-Cell communication | 167 | 5.33% |
| 5 | Cellular response to stress | 71 | 2.26% |
| 6 | Chromatin organization | 49 | 1.56% |
| 7 | Developmental Biology | 124 | 3.95% |
| 8 | Disease | 308 | 9.83% |
| 9 | DNA repair | 32 | 1.02% |
| 10 | Extracellular matrix organization | 118 | 3.76% |
| 11 | Gene expression | 139 | 4.43% |
| 12 | Hemostatis | 198 | 6.32% |
| 13 | Immune system | 366 | 11.68% |
| 14 | Membrane trafficking | 42 | 1.34% |
| 15 | Metabolism | 535 | 16.76% |
| 16 | Metabolism of proteins | 143 | 4.56% |
| 17 | Neuronal system | 138 | 4.41% |
| 18 | Signal transduction | 277 | 8.84% |
| 19 | Transmembrane transport of small molecule | 106 | 3.38% |

# Protein Sequence Features

Extracted 32 features for each protein

Wrote custom C++ script to calculate most features

R was used for the rest

Output into ARFF format

# Features

| Dimensions | Sequence Feature | Description |
|---|---|---|
| 1 | Number of amino acids | Total number of amino acids. |
| 2 | Molecular weight | Total molecular weight of the protein. |
| 3 | Theoretical pI | The isoelectric point. |
| 4-23 | Amino Acid Composition | Percentage of each amino acid. |
| 24 | Positively charged residue 1 | Percentage of lysine and arginine. |
| 25 | Positively charged residue 2 | Percentage of histidine, lysine, and arginine. |
| 26 | Number of atoms | Total number of atoms. |
| 27 | Carbon | Total number of carbon atoms. |
| 28 | Hydrogen | Total number of hydrogen atoms. |
| 29 | Nitrogen | Total number of nitrogen atoms |
| 30 | Oxygen | Total number of oxygen atoms. |
| 31 | Sulphur | Total number of sulphur atoms. |
| 32 | Non standard residue | Does it contain non-standard residue? |

# ML Analysis

1. Decision Tree

2. Random Forest

3. Support Vector Machine

4. Neural Network

# *WEKA*

# *WEKA*

# Outline

Introduction

Methods

**Results**

Conclusion

# Results

| | Decision Tree | Random Forest | Support Vector Machine | Neural Network |
|---|---|---|---|---|
| Correct classifications | 857 (27%) | **1091 (36%)** | 990 (33%) | 1056 (35%) |
| Incorrect classifications | 2182 (72%) | **1948 (64%)** | 2049 (67%) | 1983 (65%) |
| TP Rate | 0.282 | **0.359** | 0.326 | 0.347 |
| FP Rate | **0.065** | 0.070 | 0.105 | 0.067 |
| Precsion | 0.287 | **0.359** | 0.263 | 0.338 |
| Recall | 0.282 | **0.359** | 0.326 | 0.347 |
| F-measure | 0.284 | **0.350** | 0.255 | 0.338 |
| AUC | 0.620 | **0.737** | 0.729 | 0.729 |

# ROC – Binding and Uptake

# ROC – Signal Transduction

# Outline

Introduction

Methods

Results

**Conclusion**

# Conclusion

Random Forest performed best

Best ML application only had 36% correctness

AUC was over 70%

# Improvements

Need more data (upwards of 10,000)

More even distribution of functions

Need more protein characteristics

# GitHub

Available online at:


www.github.com/msegar/IT529-Protein-Prediction