

# Predicting yeast synthetic lethal genetic interactions using protein domains-Code Replication

---

## The problem :

Predicting synthetic lethality pairwise interaction upon knowledge of protein domains of the pair

gene A



gene B



Are gene A and B, SL?

# **What do we need first?**

# DATA!!!

Data

- Known SL pairs
  - Known nSL pairs
  - Protein domains + proteins
- 

- All the current knowledge on yeast genetic interaction is in **BioGrid**

## What do we need 2nd?

- to implement the features of the problem in order for the method to “learn” from

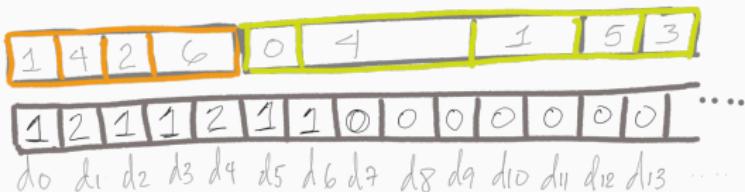
## Features of the paper

One pair example :

gene A : 

gene B : 

Pair = gene A  $\cup$  gene B

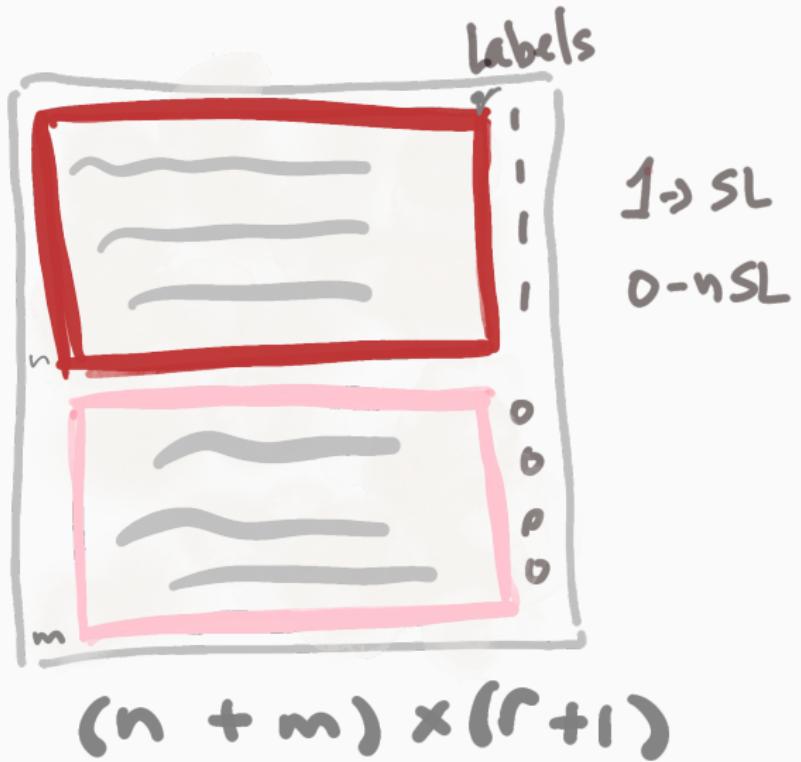
  
do d<sub>1</sub> d<sub>2</sub> d<sub>3</sub> d<sub>4</sub> d<sub>5</sub> d<sub>6</sub> d<sub>7</sub> d<sub>8</sub> d<sub>9</sub> d<sub>10</sub> d<sub>11</sub> d<sub>12</sub> d<sub>13</sub> ...

## Feature matrix

	$d_0$	$d_1$	$d_2$	$\dots$	$d_{k-1}$	$d_k$	$\dots$	$d_{r-1}$	
Pair 1	0	1	1	0	1	0	2	0	
Pair 2	2	0	0	0	2	1	1	1	
Pair 3	0	1	0	0	0	2	0	1	
	⋮			⋮		0	1	0	⋮
Pair K	0	2	0	1	0	0	0	1	
Pair n	1	0	1	⋮	2	0	2	0	

$n \times r$  feature matrix





# Everything starts! :)

- Splitting the data<sup>1</sup> for training and testing

---

<sup>1</sup>The data was a random sample of 10000 pairs taken from the given

# Everything starts! :)

- Splitting the data<sup>1</sup> for training and testing
- Train a classifier with the training data

---

<sup>1</sup>The data was a random sample of 10000 pairs taken from the given

# Everything starts! :)

- Splitting the data<sup>1</sup> for training and testing
- Train a classifier with the training data
- Make some predictions

---

<sup>1</sup>The data was a random sample of 10000 pairs taken from the given

# Everything starts! :)

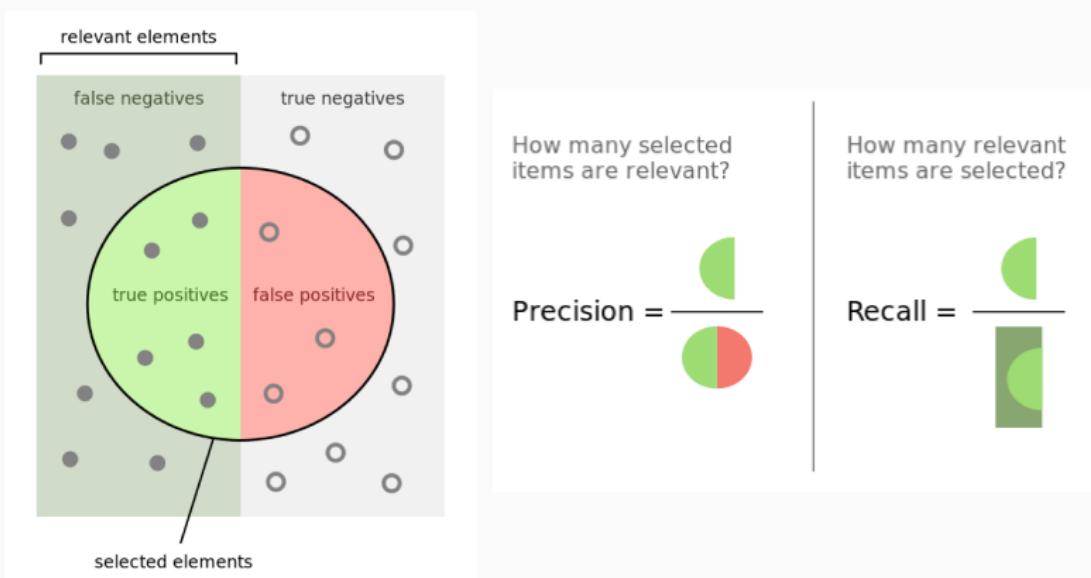
- Splitting the data<sup>1</sup> for training and testing
- Train a classifier with the training data
- Make some predictions
- Evaluate the predictions based on the testing data.

---

<sup>1</sup>The data was a random sample of 10000 pairs taken from the given

## Evaluating the method based on its prediction

- Accuracy : The percentage of in how many cases the correct class was predicted.
- Sensitivity : Percentage of correctly predicted SL interactions over the total number of SL interactions in the test dataset.
- Specificity : Percentage of correctly identified negative data over the total number of negative data.
- Precision : Percentage of correctly predicted positive data over the total number of predicted positive data.



*From Wikipedia*

## Results of the evaluation

	Paper	My replication
Accuracy	0.84	0.89
Specificity (Recall of the negative class)	0.83	0.89
Sensitivity (Recall of the positive class)	0.85	0.89
Precision	0.83	0.89
F-score	0.84	0.89
AUC	0.927	0.9

## 5-fold Cross validation study

