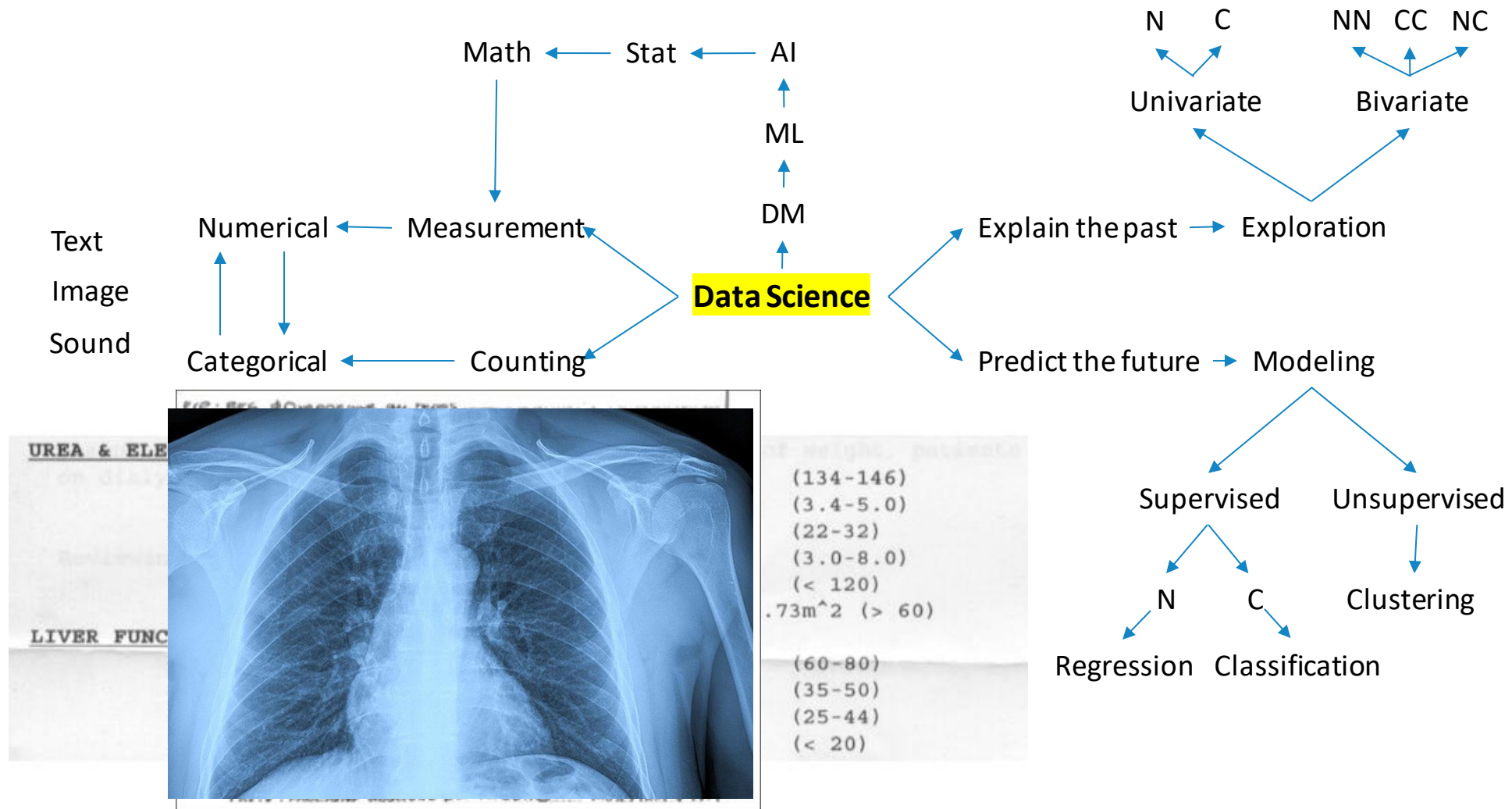


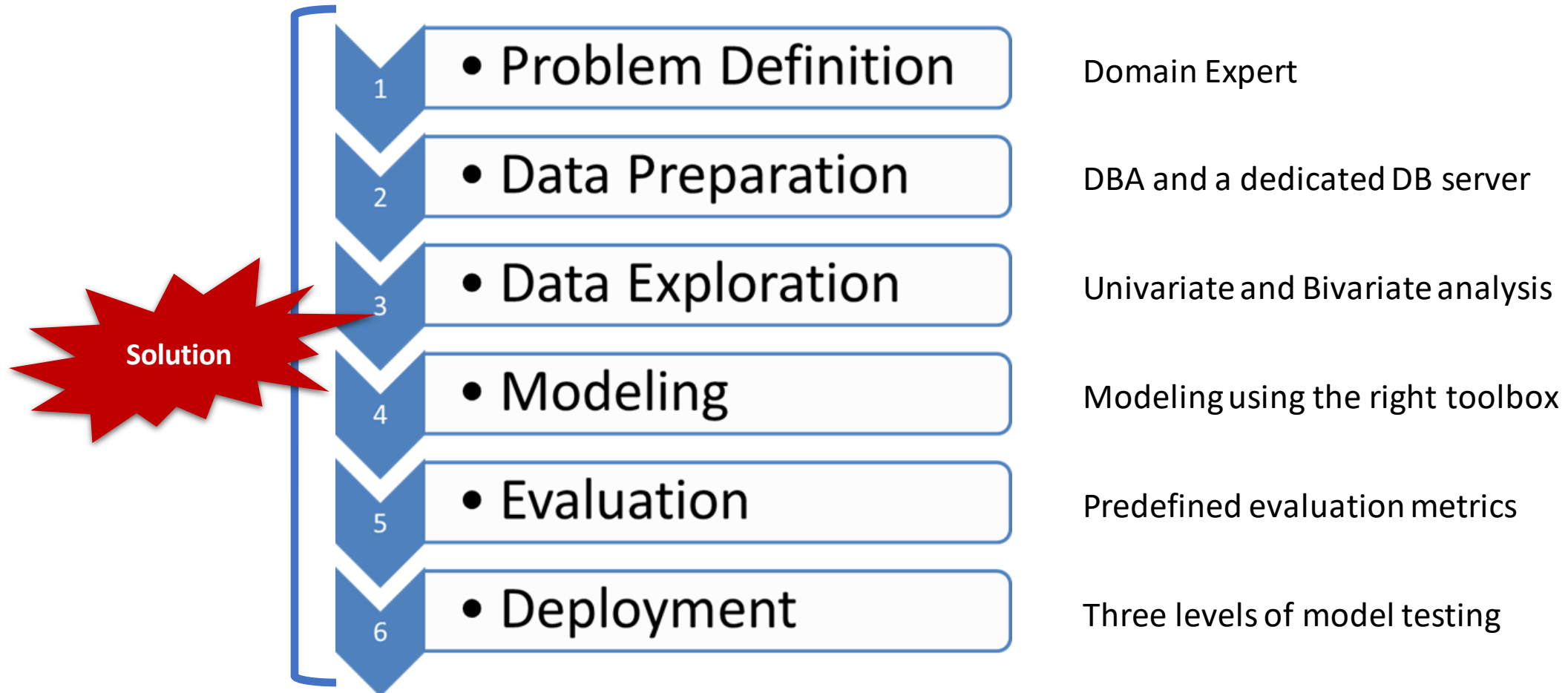


Data Science for Biomarkers Discovery

Saed Sayad M.D. Ph.D.
Bioada.com



Data Science 6-Step



1- Problem Definition

- Building a **classification model** using a novel combination of serum microRNAs for detecting breast cancer in the early stage
- Success Criteria
 - ROC chart/AUC
 - Precision/Recall
 - FDR
 - Gain/Lift chart
 - K-S chart
 - Deciles chart

Data

- **GSE73002**

<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE73002>

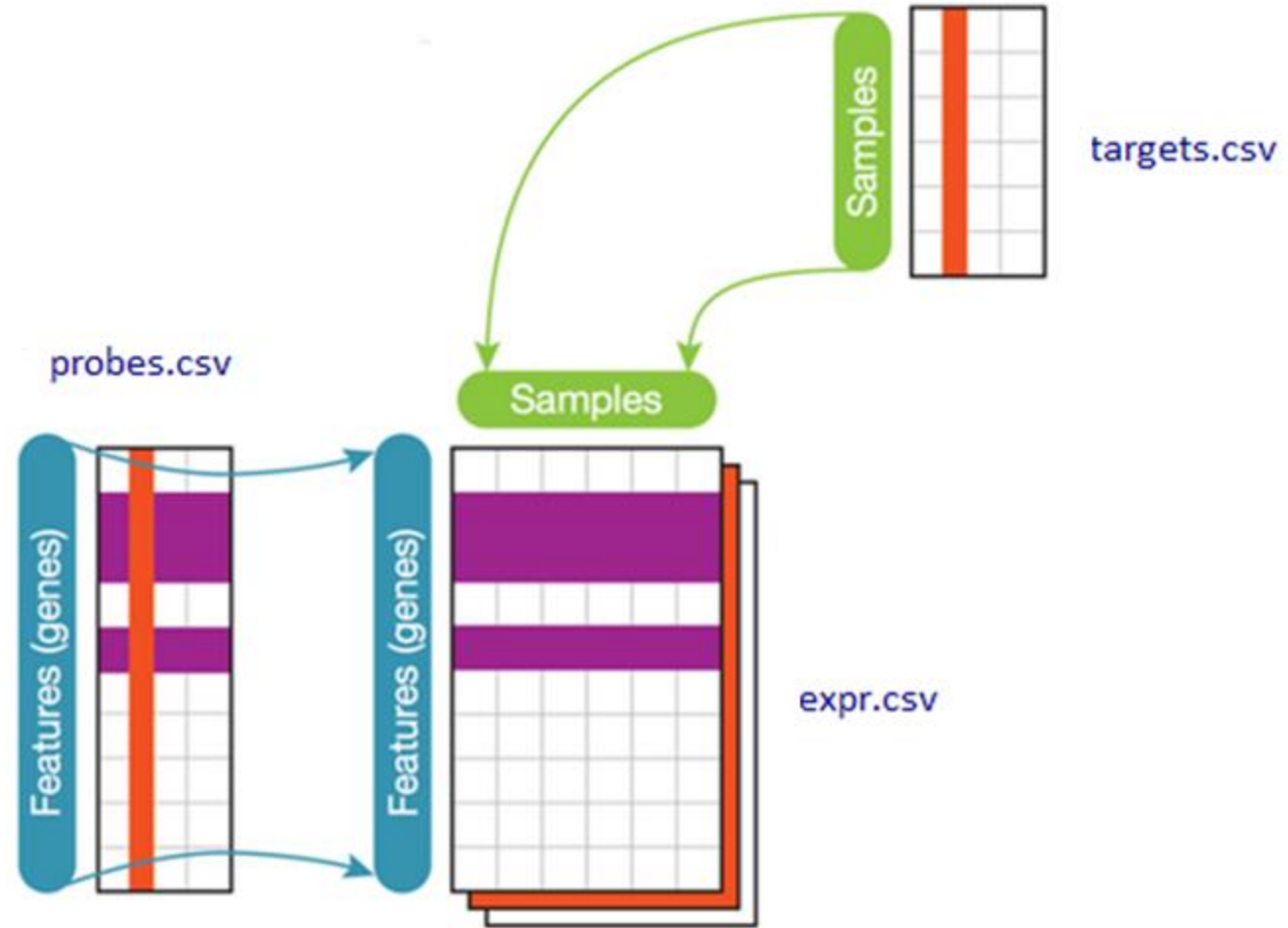
- **Title**

A novel combination of **serum microRNAs** for detecting breast cancer in the early stage

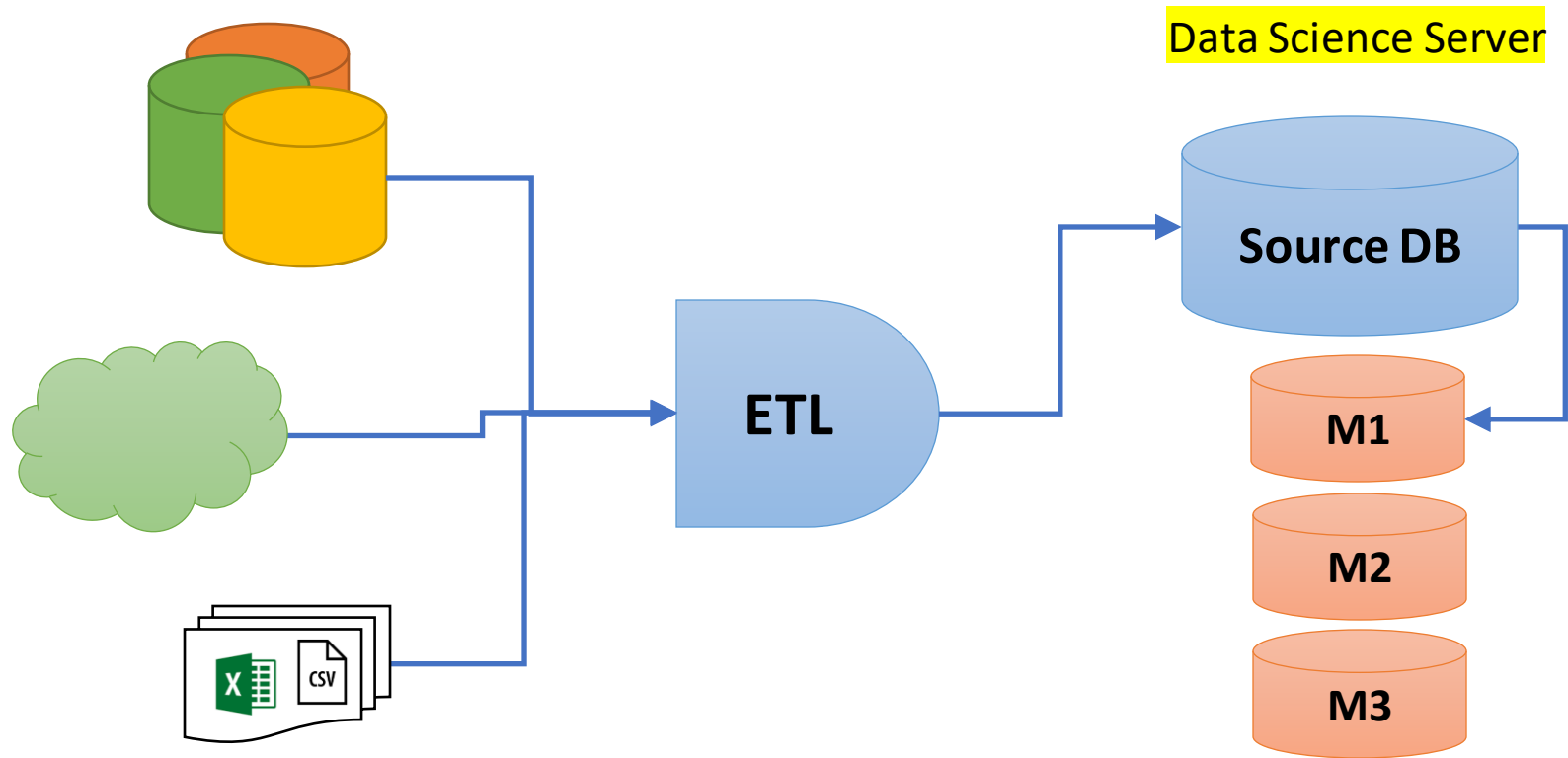
- **1280 breast cancer** and **2686 non-cancer control**

- **Bioconductor** R packages can be used for data preparation

Data Model



Data Management

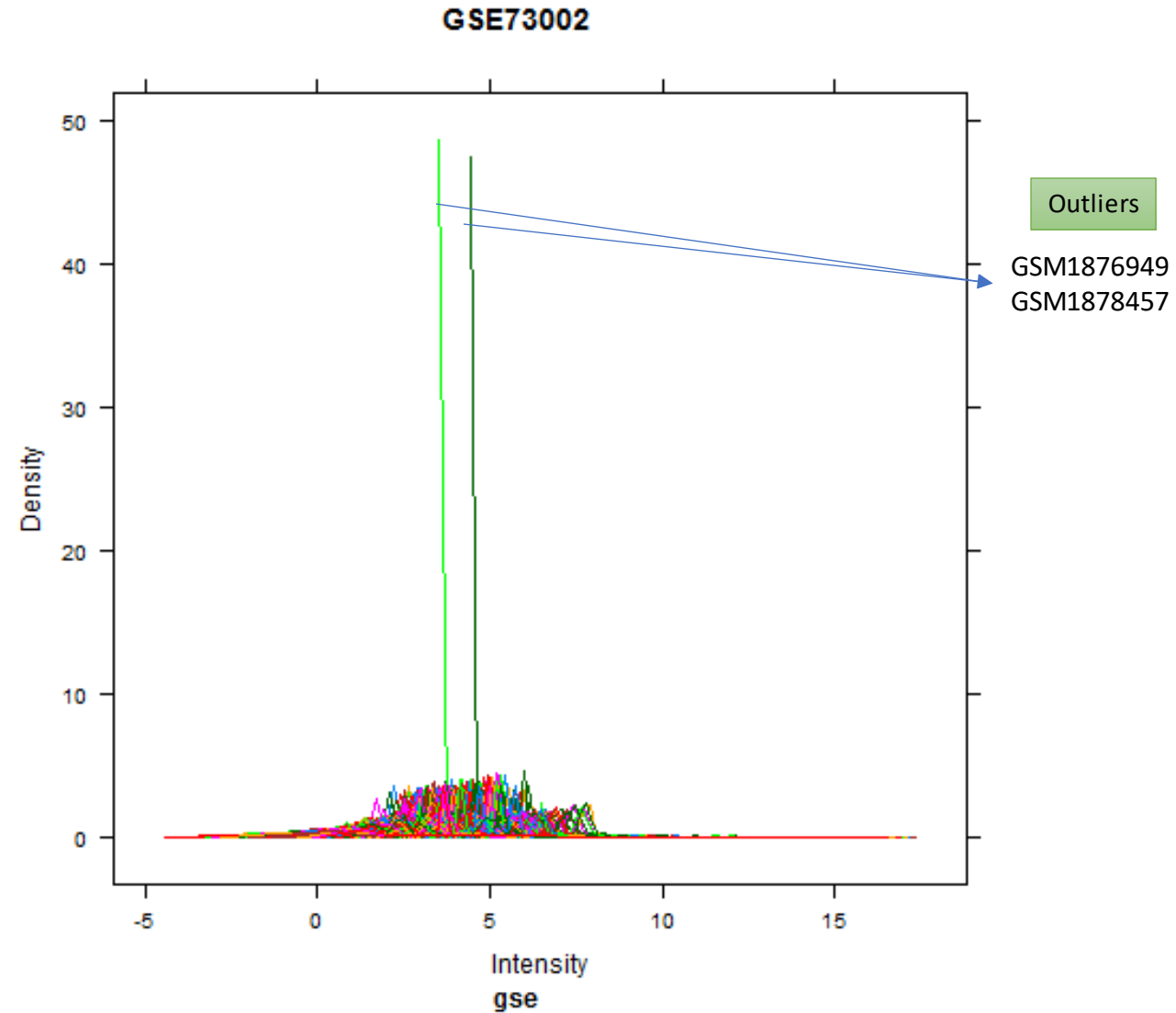


3- Data Exploration

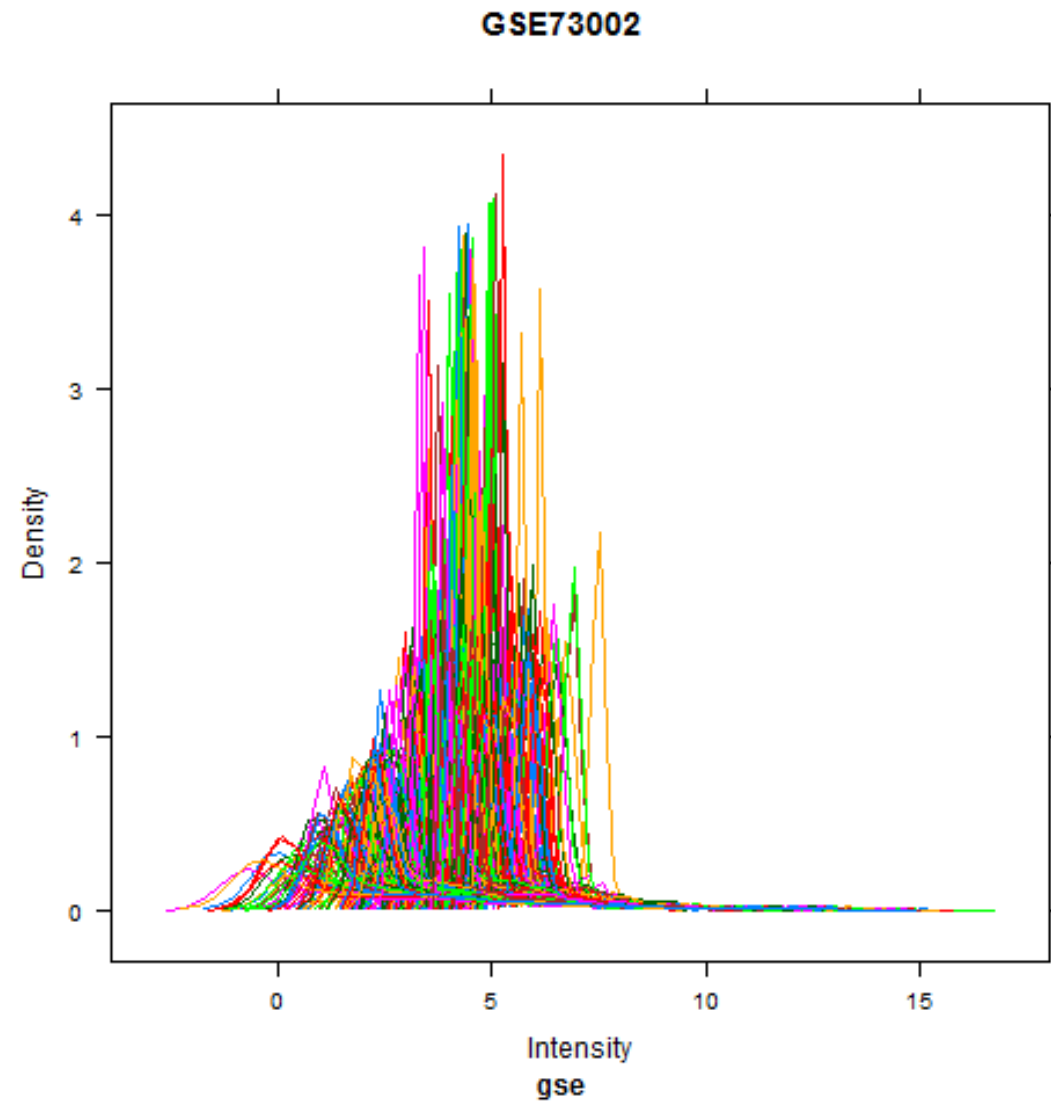
- **Univariate Analysis** on all variables/probes
- **Bivariate Analysis** at least against the target

Expression data - density plot

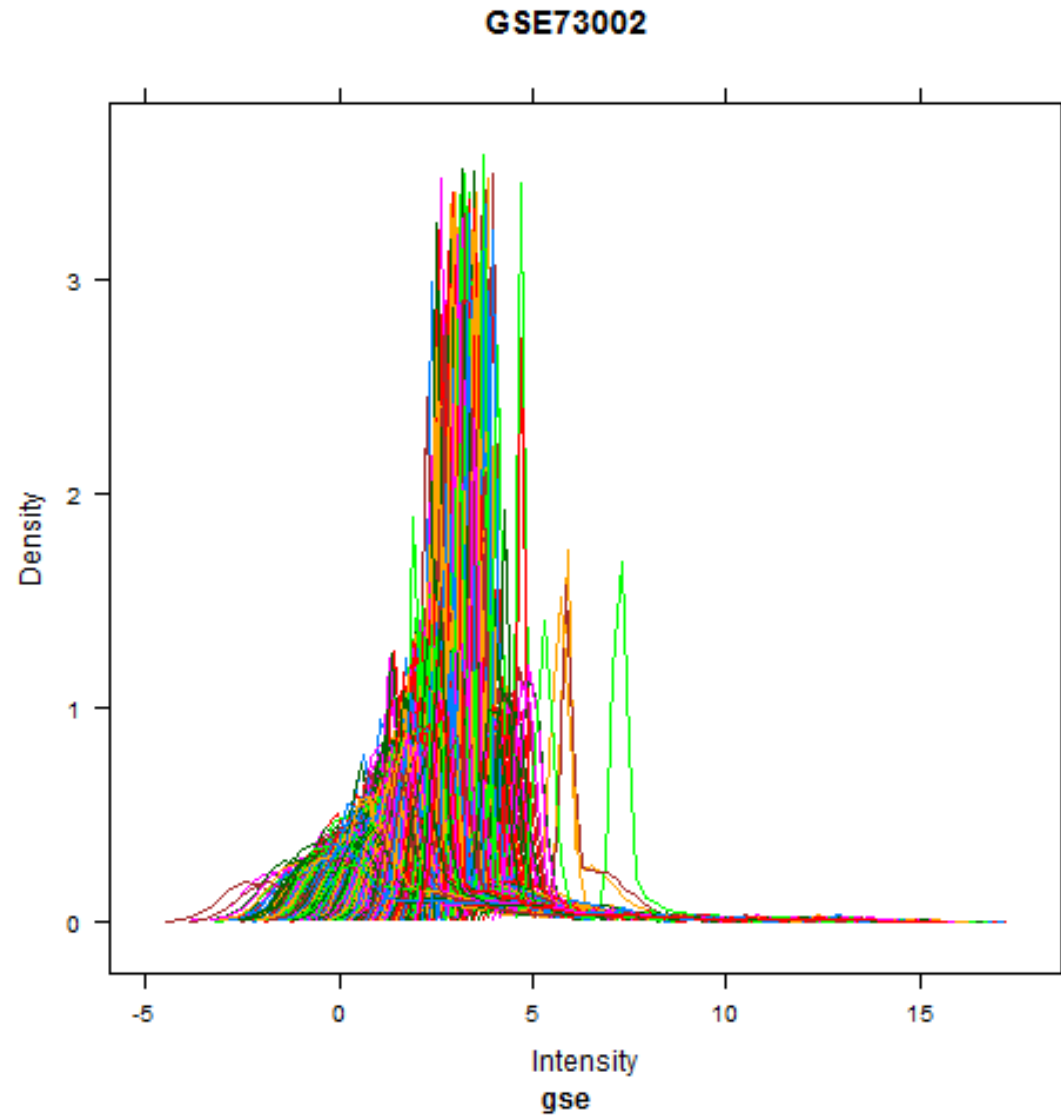
breast cancer	1280
non-cancer	2686



Exploration data - Breast Cancer



Exploration data – Normal Breast

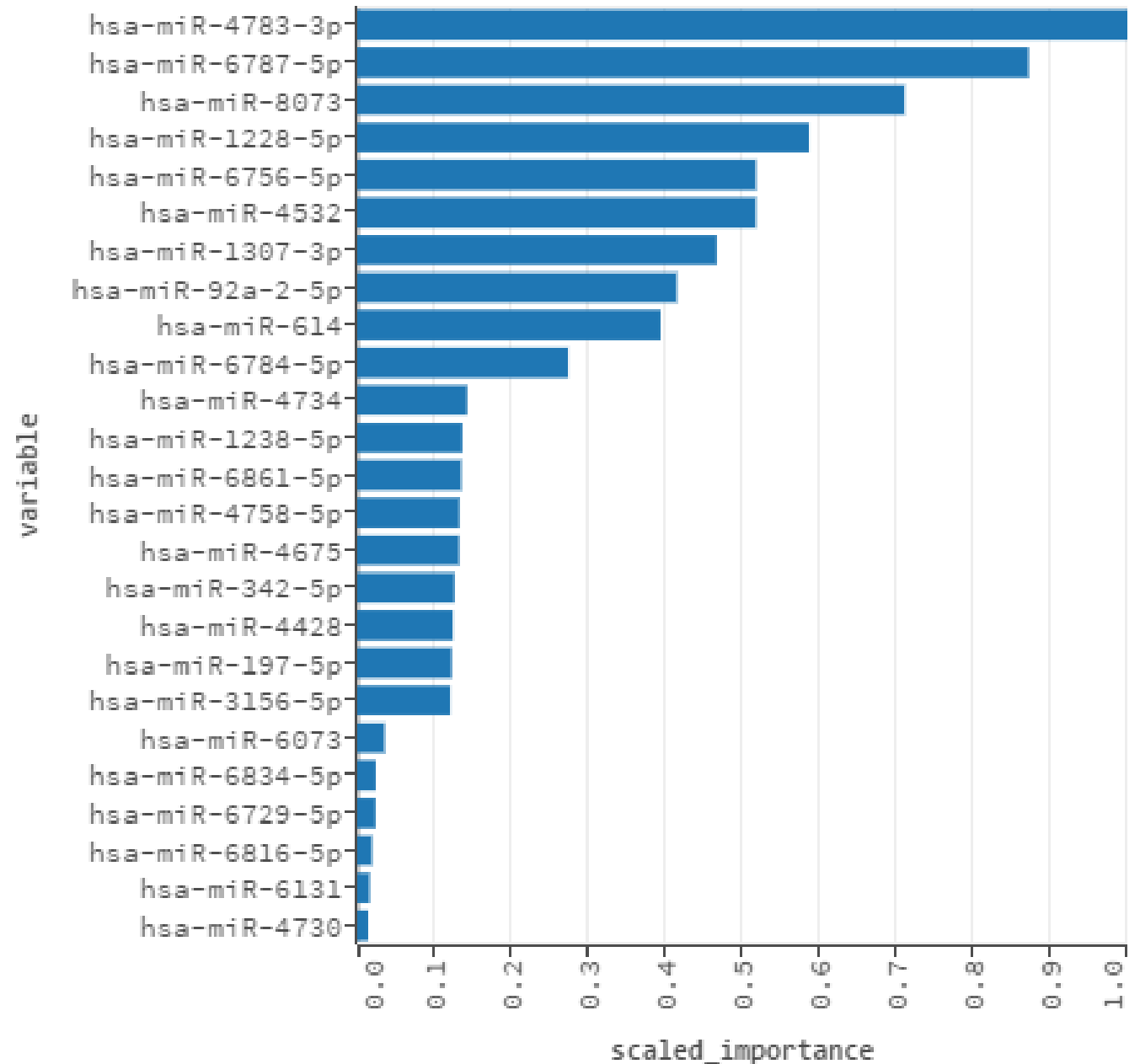


4- Modeling

1. Dividing data to **trainset** and **testset**
2. Selecting **top N probes** using **trainset only** and a proper statistical test. N is usually between 10-100
3. Building predictive models to choose the final set of probes (**biomarkers**)
 - Building a **classification model** to find the best subset of probes (biomarkers) to separate the breast cancers from the non-cancers cases

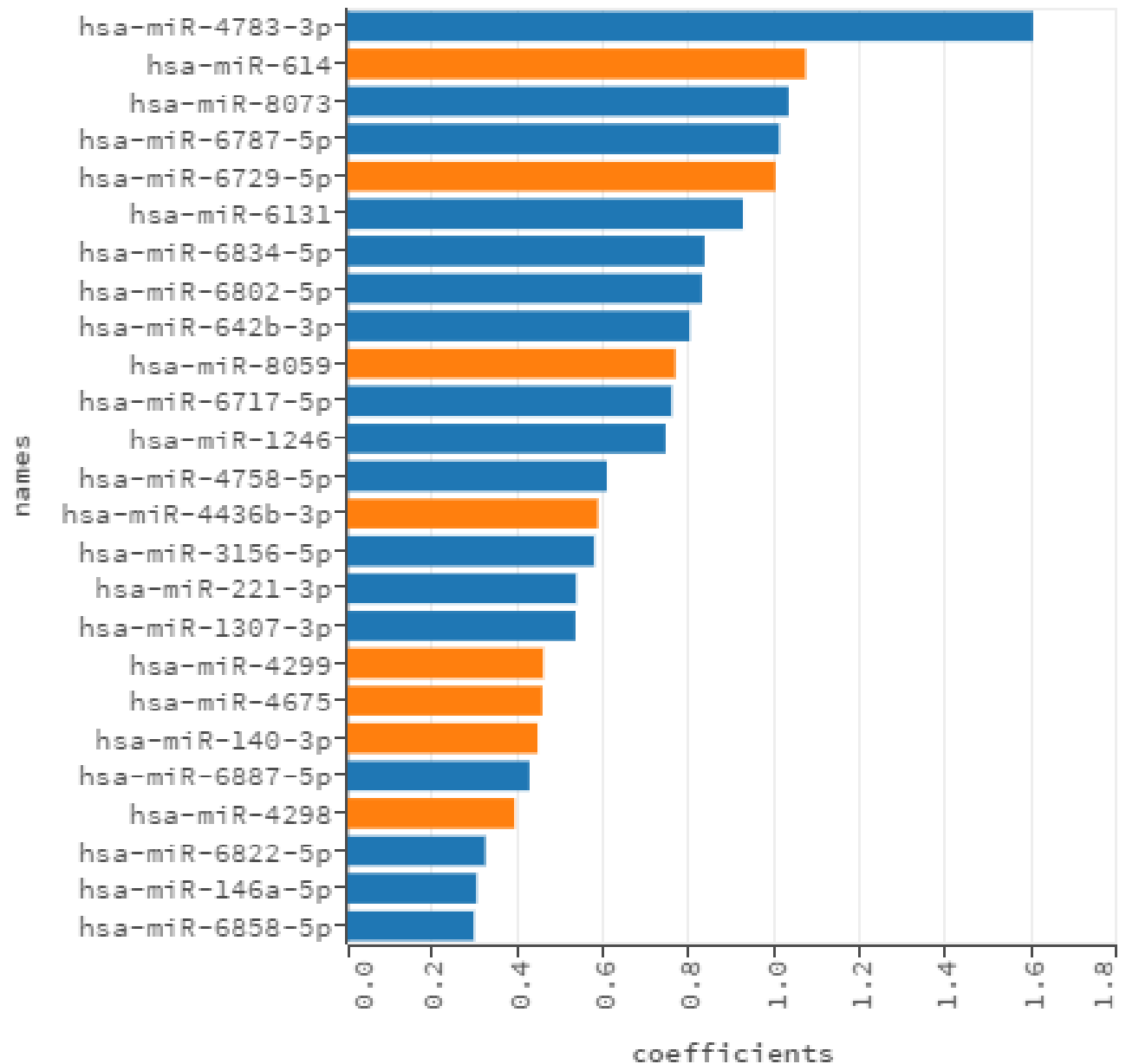
Modeling - RandomForest

Accuracy 98%



Modeling - GBM

Accuracy 99%



Modeling – Deep Learning

Accuracy 99%



DeepLearning does not have a straightforward variable selection capability.

Modeling – Linear Model (LDA)

Accuracy 99%

Xarang [1.ilm]

Project Data Learner Explorer **Modeler** Predictor Deploy

Classification Target: ☐ OnTheFly Iterations: 5 D: 0.5 Reducer Model Script

target

☒ Input: [1]

- ☐ hsa-mir-4538
- ☐ hsa-mir-4646-5p
- ☐ hsa-mir-4652-5p
- ☐ hsa-mir-4675
- ☐ hsa-mir-4690-5p
- ☐ hsa-mir-4700-5p
- ☐ hsa-mir-4727-3p
- ☐ hsa-mir-4730
- ☐ hsa-mir-4732-5p
- ☐ hsa-mir-4734
- ☐ hsa-mir-4758-5p
- ☐ hsa-mir-4771
- ☐ hsa-mir-4783-3p
- ☐ hsa-mir-483-5p
- ☐ hsa-mir-6073
- ☐ hsa-mir-608
- ☐ hsa-mir-6131
- ☐ hsa-mir-614
- ☐ hsa-mir-642b-3p
- ☐ hsa-mir-654-5p
- ☐ hsa-mir-663a
- ☐ hsa-mir-6717-5p
- ☐ hsa-mir-6729-5p

Variables	Coefficient	Contribution
hsa-mir-1307-3p	6.65893	5.02700598

Overall Contribution	Probability 1	Probability 2	Base Score
5.027006	0.6773	0.3227	46.770766

Classification Regression MultiClass Segmentation Time Series Retargeting

0.080 sec. 1.ilm

LDA selected only one probe

5- Model Evaluation

	RandomForest	GLM	DeepLearning	LDA
Accuracy	✓	✓	✓	✓
Probe Importance	✓	✓	✗	✓
Federated Learning	✗	✗	✗	✓

6- Model Deployment

➤ 3 levels of model testing

1. Testing models on the **test platform** by the modeling team
2. Testing models by the **QA team**
3. **A/B test** on the **production platform**

Summary

- Hire a **Solution Architect**
- Implement **Data Science 6-Step** and be religious about it
- Have a **clear understanding of the question/problem** and how to **measure success**
- **Own the data!** Data is your reserve currency. If you do not want to become bankrupt always make sure you have full control of your data
- Always do **univariate** and **bivariate analysis**
- **A/B Test!**
- Choose the **right ML toolbox**

Thank You!

saed@bioada.com