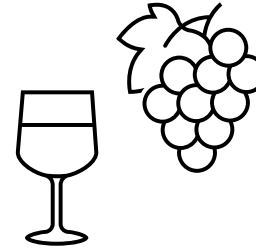


# Data Mining – Team MINEHeim

Wine Prediction 



# Agenda

Research Problem



Data & Pre-Processing



Data Mining



Conclusion



# Research Problem

## PROBLEM

How can you determine the quality of wine without actually drinking it?  
What makes a great wine ... actually great?

- Which factors/components influence the quality and taste?



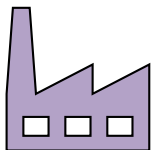
## GOAL

Objectively predict and classify the quality of wine based on its ingredients, components and features.

- e.g. sulphates, residual sugar, alcohol, pH-Value, type, ...



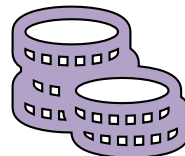
## BUSINESS VALUE



Wineries & Dealers



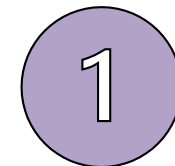
Old Wine



Reduce Costs



Save Time



Single Source of Truth

# Data Structure



Data set related to red and white variants of Portuguese wine "Vinho Verde"



6497 samples, already gathered



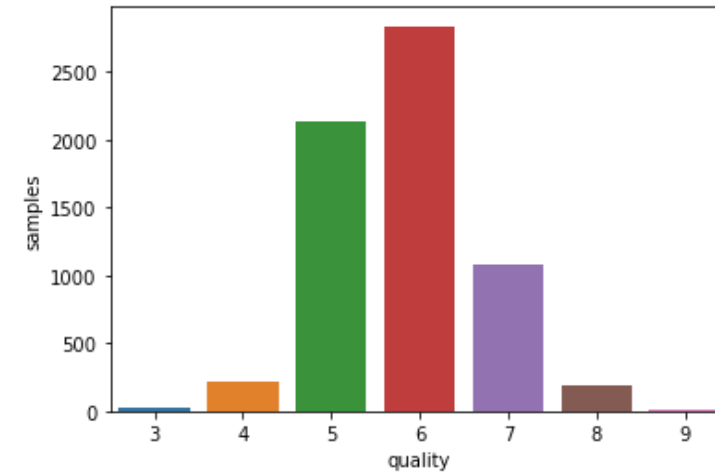
13 features, mainly physiochemical properties e.g., pH-value, alcohol percentage etc.



Quality feature ranges from 0 to 10, assessed by wine tastings



Highly unbalanced data: lots of mediocre wines, few bad and excellent wines



# Data & Pre-Processing



## Duplicate Deletion and Replacing Null-Values



- No duplicates
- 34 rows with missing values were deleted



## Column Preprocessing



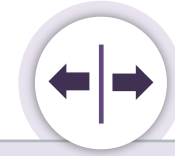
- Normalization of numerical values
- Encoding of categorical values



## Outlier Deletion



- Deleting 57 outliers on the basis of a 90 % confidence interval



## Data Separation and Splitting



- Separation of the feature and target variable
- Train and test split

# Data Mining

## ML METHODS

K-Nearest  
Neighbors  
(K-NN)

Support  
Vector  
Machine  
(SVM)

Decision  
Tree

Random  
Forest

Neural  
Network

## EVALUATION METRICS



Precision



Recall



**F<sub>1</sub>-Score**

## APPROACH

Default Values

1.  
Hyperparameter  
Tuning

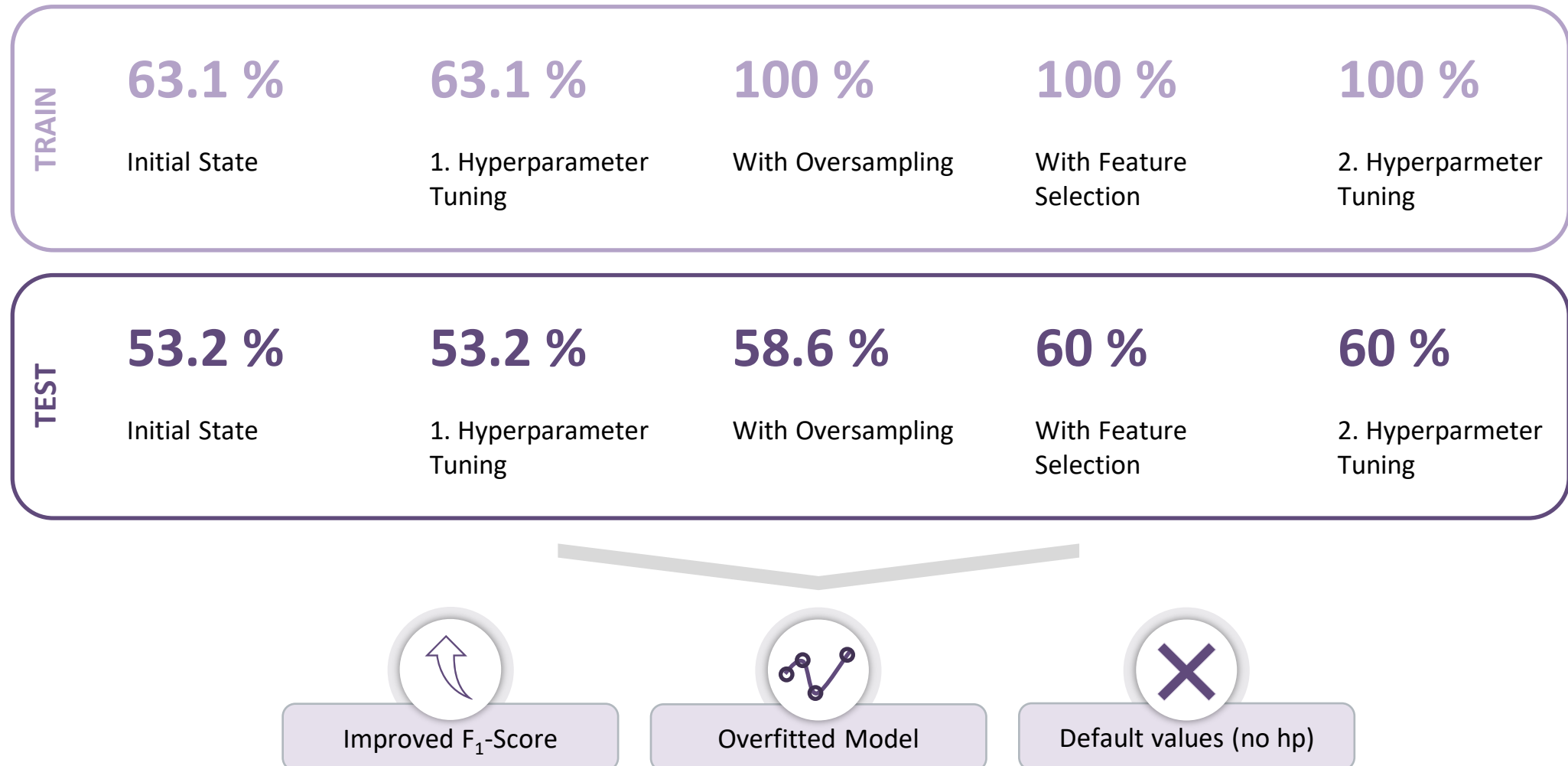
Balancing  
(SMOTE)

Feature Selection

2.  
Hyperparameter  
Tuning

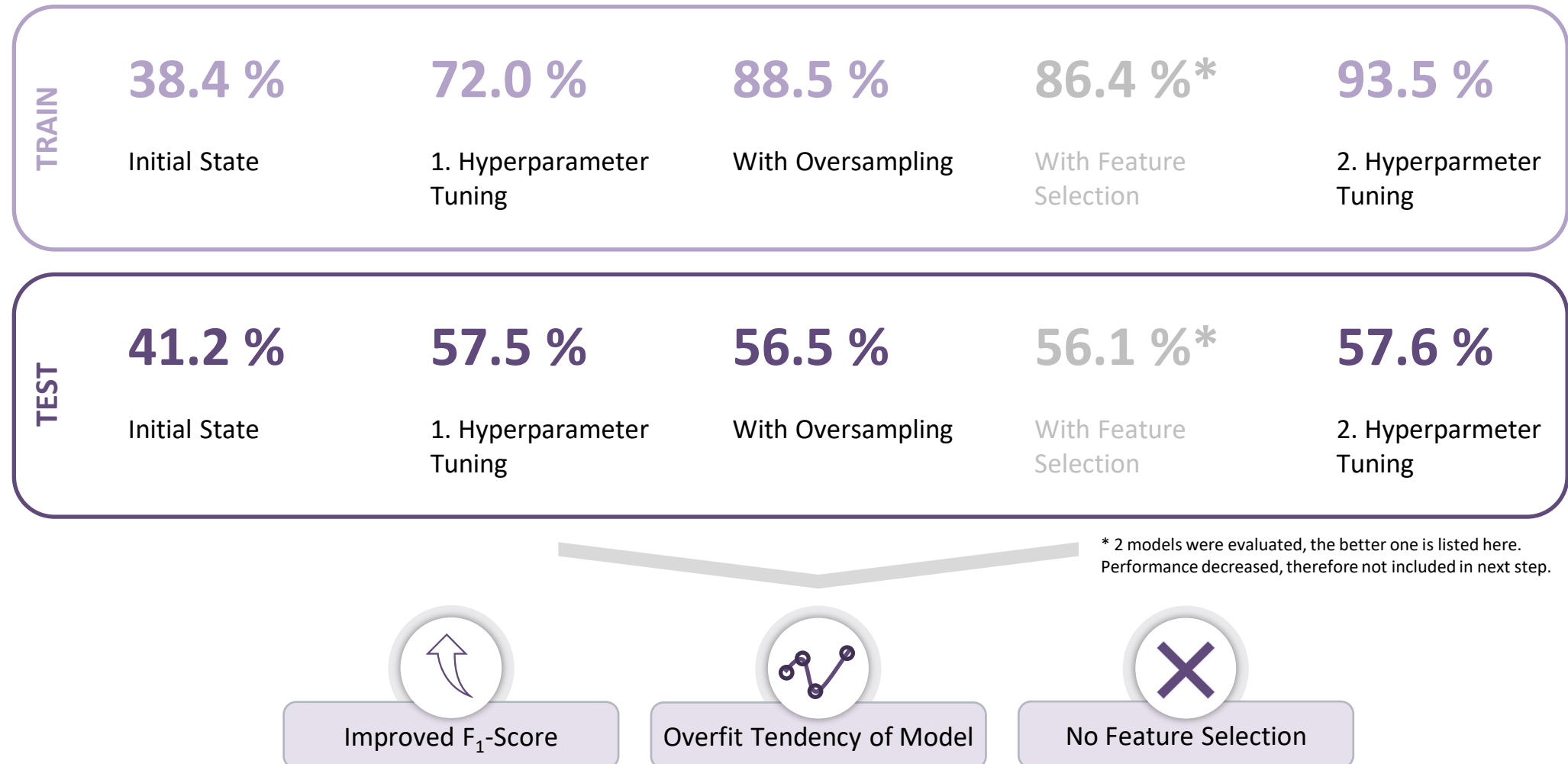
# K-Nearest Neighbors (K-NN)

## F<sub>1</sub>-Score Overview



# Support Vector Machine (SVM)

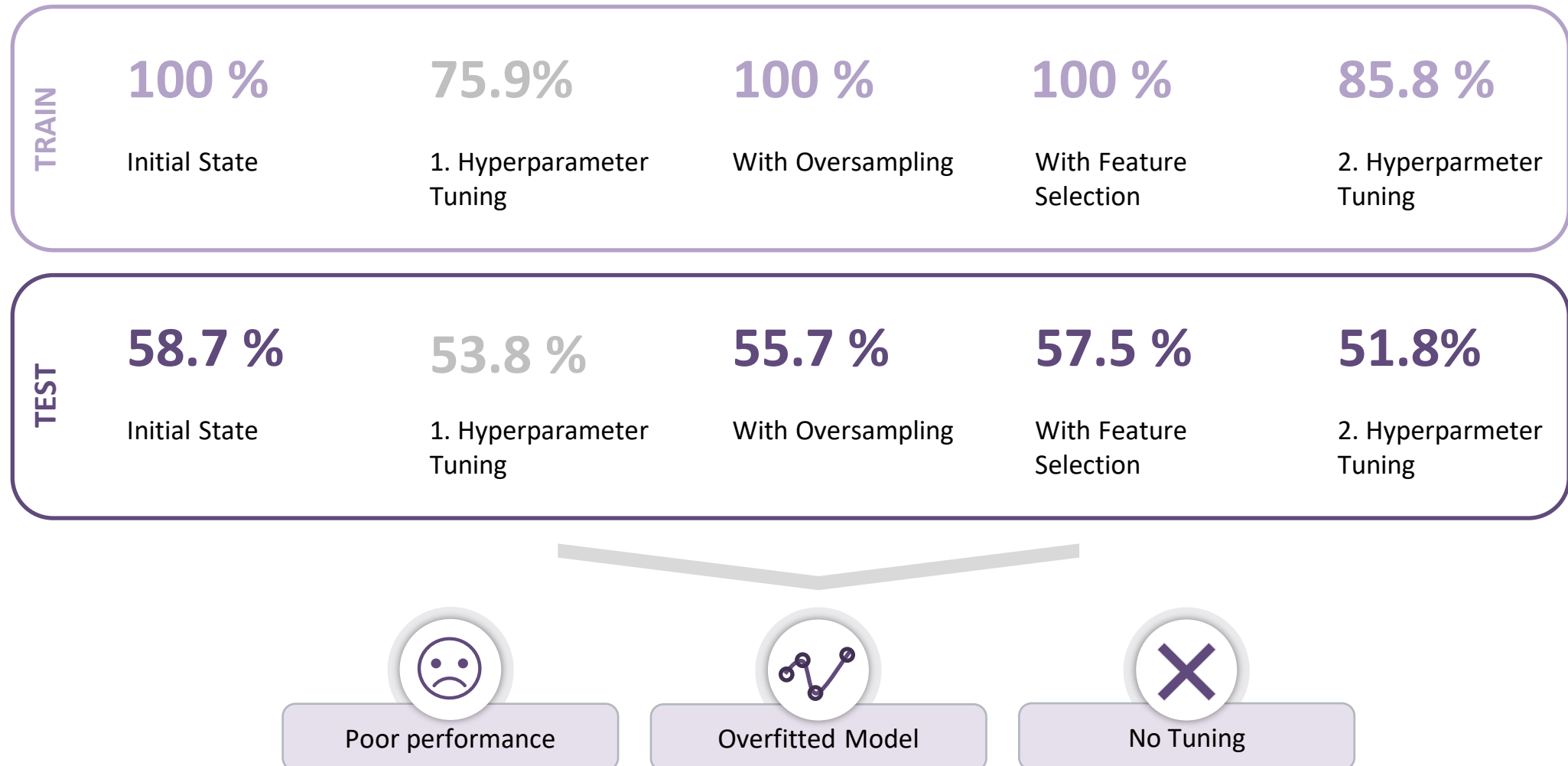
## F<sub>1</sub>-Score Overview





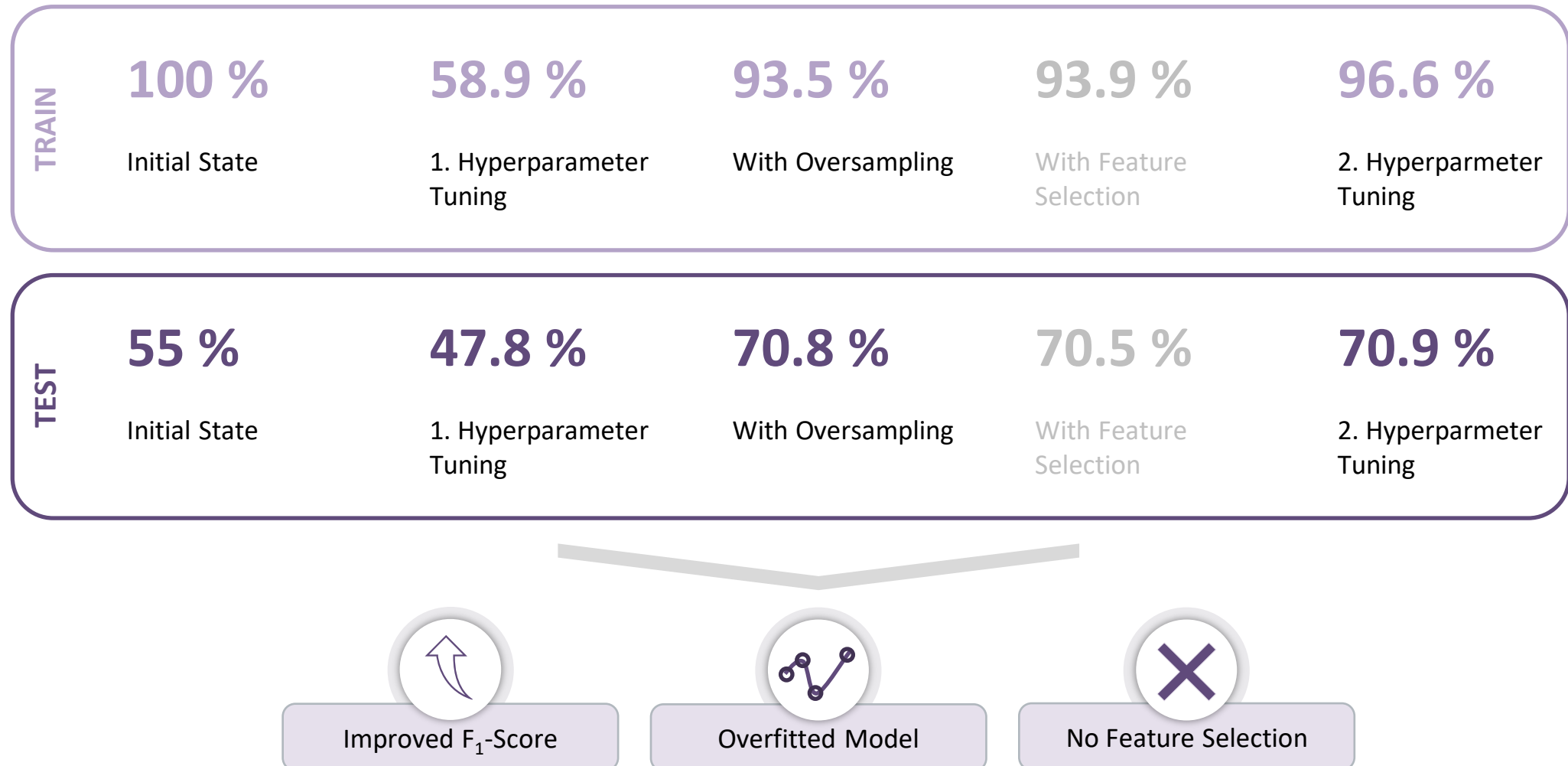
# Decision Tree

## F<sub>1</sub>-Score Overview



# Random Forest

## F<sub>1</sub>-Score Overview



# Neural Network

## F<sub>1</sub>-Score Overview

### Legend

Batch

Layers

Epochs

Neurons



UNIVERSITY  
OF MANNHEIM  
School of Business Informatics  
and Mathematics

VALIDATION

38.3 %

No Balancing

64

3

30

100

72.8 %

64

4

30

300

76.8 %

64

5

30

400

80.9 %

16

7

40

1000

84.6 %

16

8

80

2048

TEST

41.4 %

No Balancing

64

3

30

100

51.9 %

64

4

30

300

49.9 %

64

5

30

400

55.7 %

16

7

40

1000

56.6 %

16

8

80

2048



Improved F<sub>1</sub>-Score



Overfitted Model



High effort, low reward

# Evaluation & Conclusion

## Business Use Case



Predict **wine quality** with the help of contents.

## Results

### KNN

*as baseline did a reasonable job*

### Decision Tree

*mediocre performance on its own*

### Random Forest

*performed well as an aggregation of decision trees*

### SVM

*also performed well as a classification algorithm*

### Neural Network

*Efforts to create and were not rewarded with good predictions → more resources, better experience needed*

Overfitted Model

Random Forest 

Other Features besides chemical characteristics might be important for the wine quality