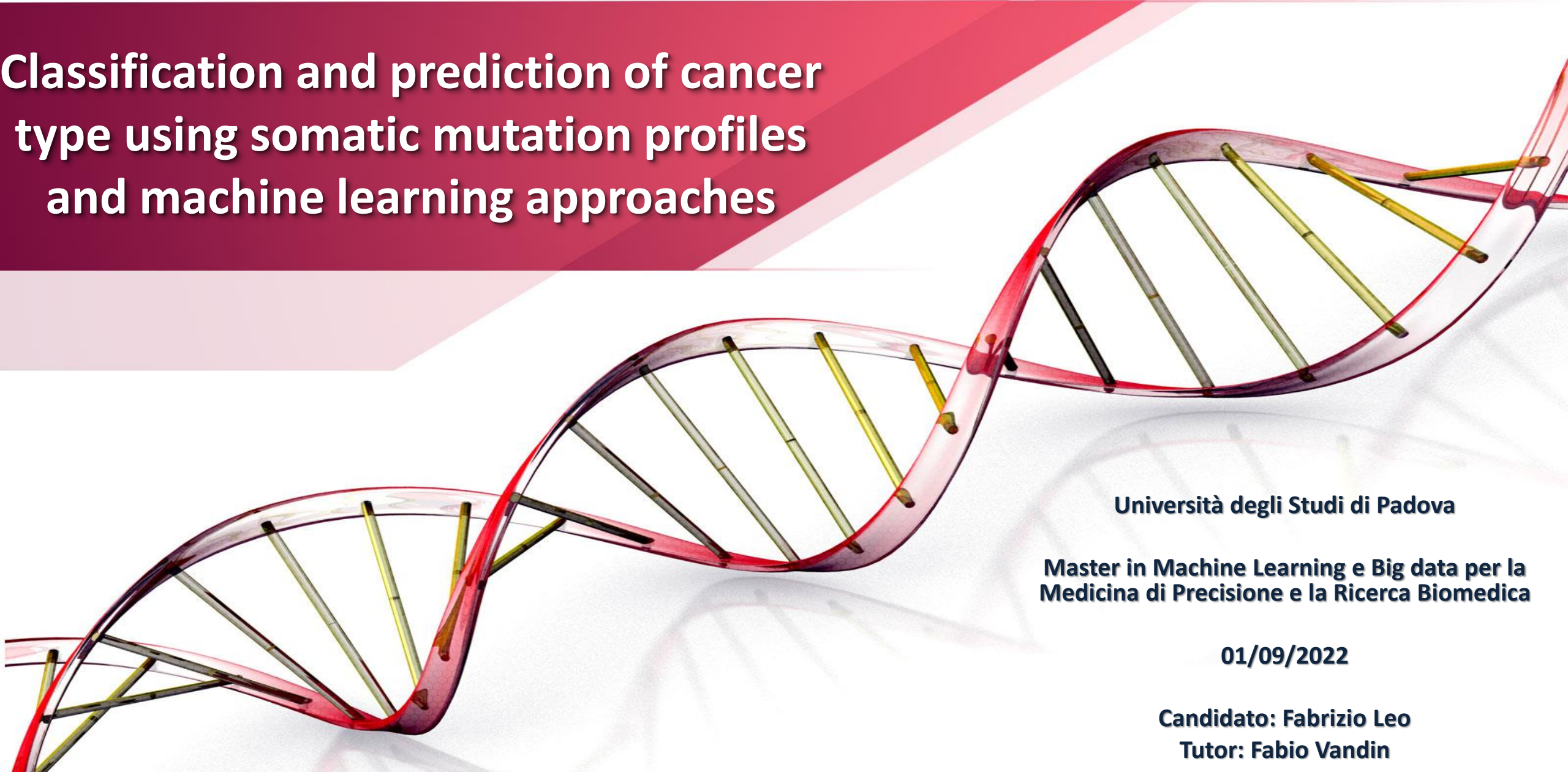# Classification and prediction of cancer type using somatic mutation profiles and machine learning approaches

**Università degli Studi di Padova**

**Master in Machine Learning e Big data per la Medicina di Precisione e la Ricerca Biomedica**

**01/09/2022**

**Candidato: Fabrizio Leo**
**Tutor: Fabio Vandin**

# Data

- Source: The Cancer Genome Atlas (TCGA)

- Information about cancer patient somatic mutations and type of cancer for each patient

  a) «samples_labels.txt» -> list of patients with their cancer type
  b) «snvs.tsv» -> list of mutated genes for each patient
  c) «Compendium_Cancer_Genes.txt» -> list of genes considered important

# Pre-processing

- I only considered, for computational reasons, the list of 568 genes considered as important
- Built binary mutation matrix considering these genes
- Removed 42 genes not mutated in any patient

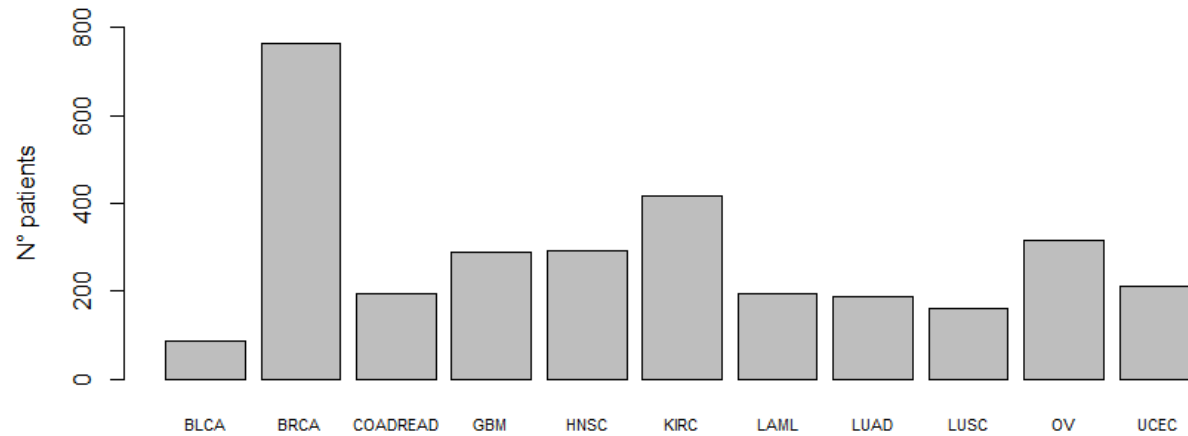- Therefore, I obtained a mutation matrix 3109 (patients) x 526 (genes)

# Outline

- Is it possible to predict cancer type based on genes with somatic mutation in a patient?

- Is there a «small» set of genes having a good predictive power, or at least as good as the entire set of genes?

- Does the patient grouping based on similarity of mutated genes reflect the grouping based on cancer type?
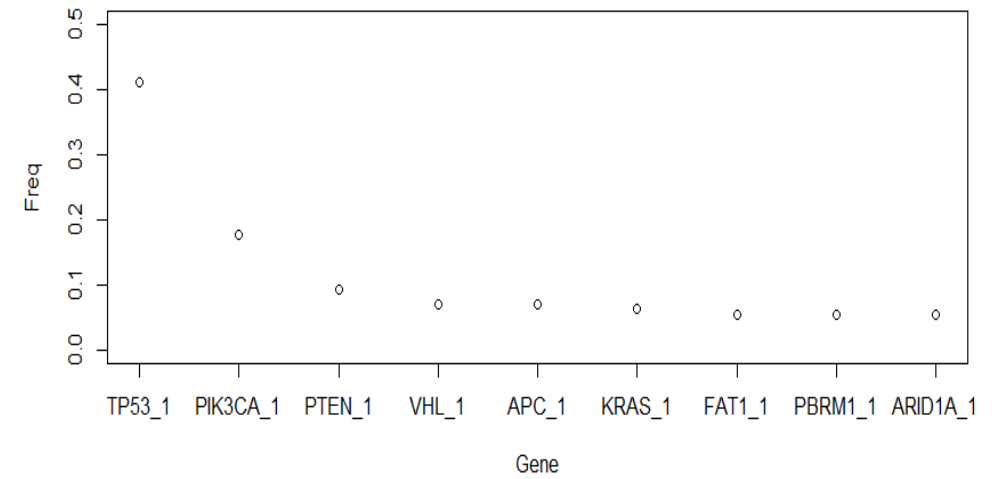
# Explorative analysis



**Patients by Cancer Type**

**Genes with most frequent mutation in the entire sample**

- 763 cases BRCA (24.5%)
- 87 cases BLCA (2.8%)

# Explorative analysis

# Outline

- Is it possible to predict cancer type based on genes with somatic mutation in a patient?

- Is there a «small» set of genes having a good predictive power, or at least as good as the entire set of genes?

- Does the patient grouping based on similarity of mutated genes reflect the grouping based on cancer type?

# Is it possible to predict cancer type?

- Compared four models: CIDT, KNN, SVM and RandomForest
- 5-folds CV (25% data test)
- Tuning parameters: mincriterion, k, cost and mtry



**RandomForest Training**



**Model Comparison**

*** , ACC > NIR
p < .001

# Is it possible to predict cancer type?



Classification Sensitivity- Test Data - Best Model (RF)

$$\frac{TP}{TP + FN}$$

Balanced Accuracy - Test Data - Best Model (RF)

$$\frac{\left(\frac{TP}{TP + FN}\right) + \left(\frac{TN}{TN + FP}\right)}{2}$$
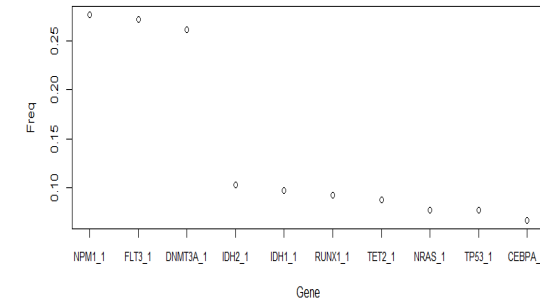
Confusion Matrix - Test Data - Best Model (RF)

Genes with most frequent mutation in the LUSC Cancer Type

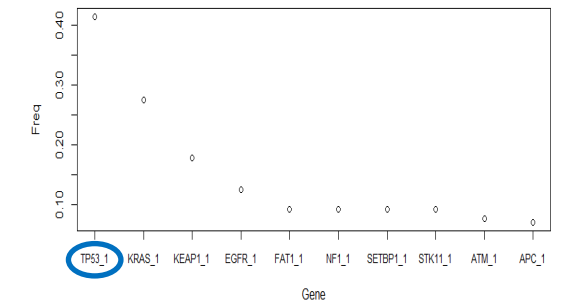Genes with most frequent mutation in the HNSC Cancer Type

# Outline

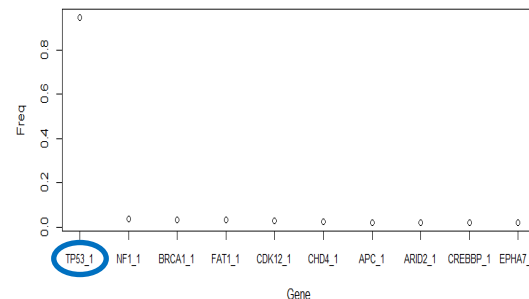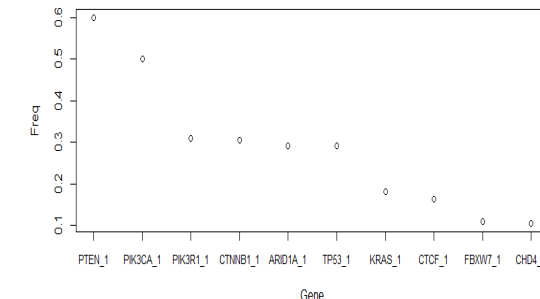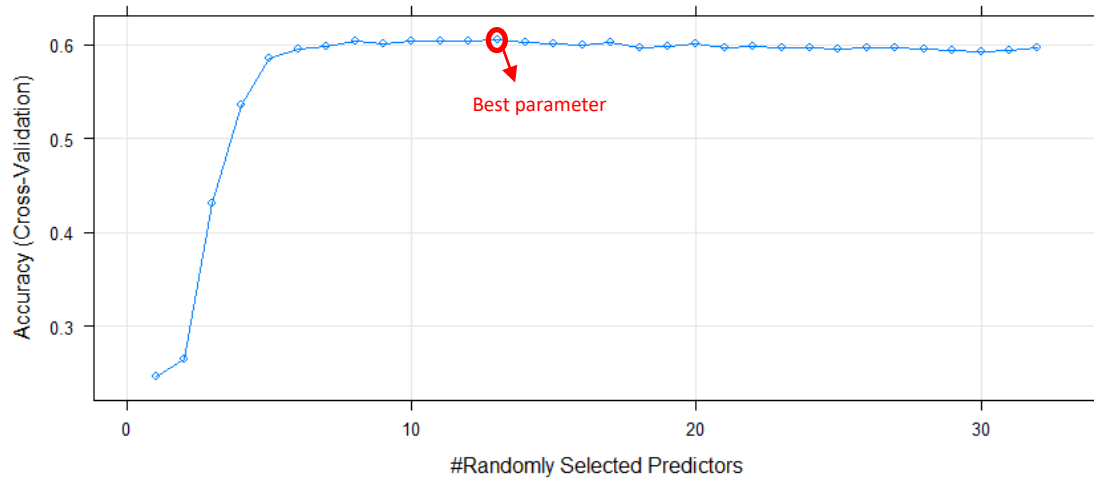- Is it possible to predict cancer type based on genes with somatic mutation in a patient?

- Is there a «small» set of genes having a good predictive power, or at least as good as the entire set of genes?

- Does the patient grouping based on similarity of mutated genes reflect the grouping based on cancer type?
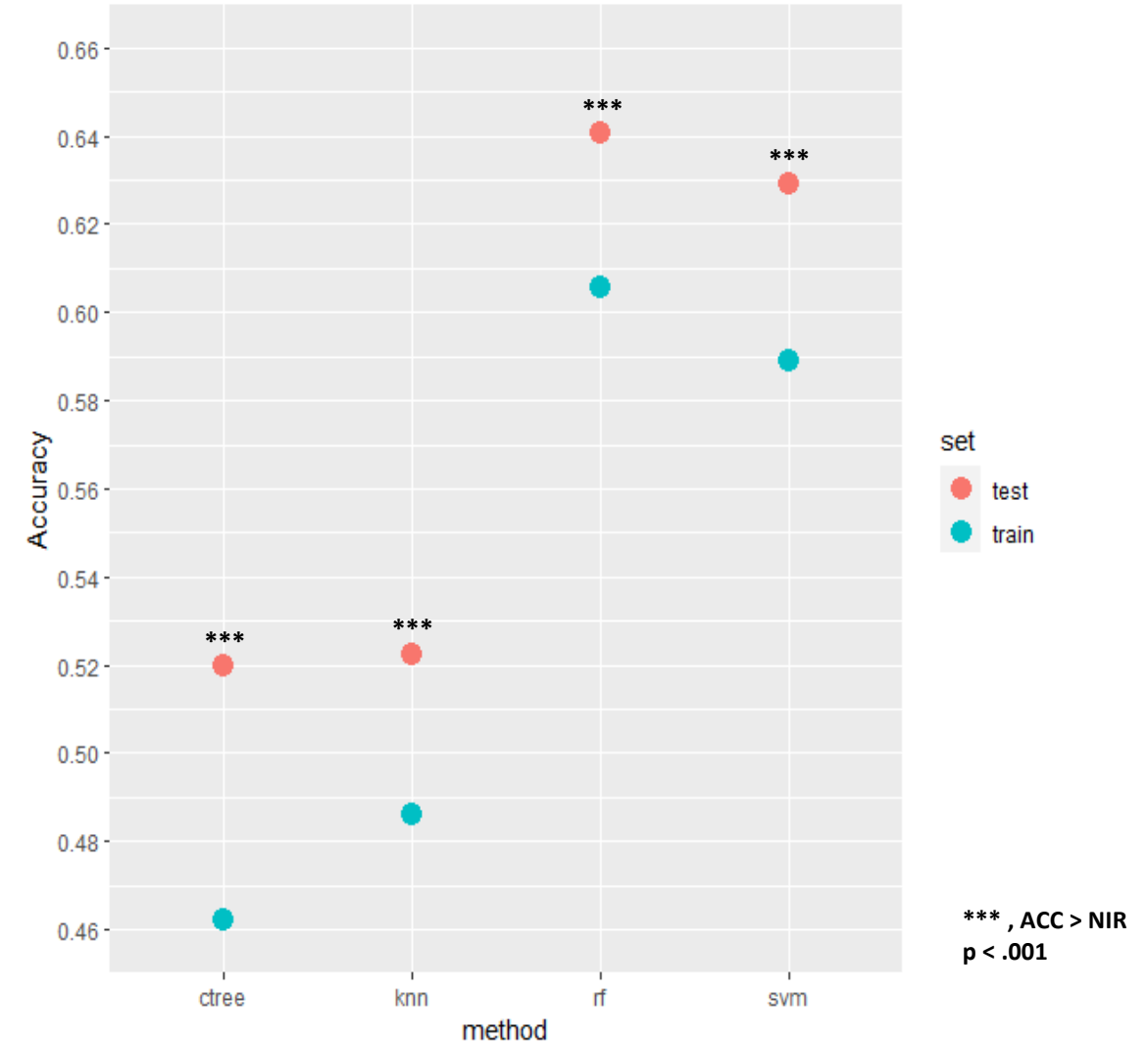
# Is there a small set of genes having a good predictive power?

- Used two approaches for gene selection:

  1. **RandomForest**, using different threshold values for Gini score (>2, >3, >4, >5, >6, >10)
  2. **Boruta**: a) selection of relevant genes in all CV runs;

     b) selection of genes defined as relevant in at least one CV run

- Training RandomForest with selected genes
- Model evaluation using the test set

# Is there a small set of genes having a good predictive power?

**Feature Selection Performance Comparison in Test set**



- Good accuracy for set of 30-40 genes and above

- Performance decreases significantly under 20 genes

# Outline

- Is it possible to predict cancer type based on genes with somatic mutation in a patient?

- Is there a «small» set of genes having a good predictive power, or at least as good as the entire set of genes?

- Does the patient grouping based on similarity of mutated genes reflect the grouping based on cancer type?

# Does the clustering based on mutated gene similarity reflect cancer type?

- **NMF approach** (Lee algorithm, 10 runs)



- Optimal K > 24, therefore higher than the number of cancer types

# Does the clustering based on mutated gene similarity reflect cancer type?

- Even when setting k = 11 (Lee algorithm, 30 runs), the clustering is different than real classes

- E.g., BRCA cancer is clusterized as follows:

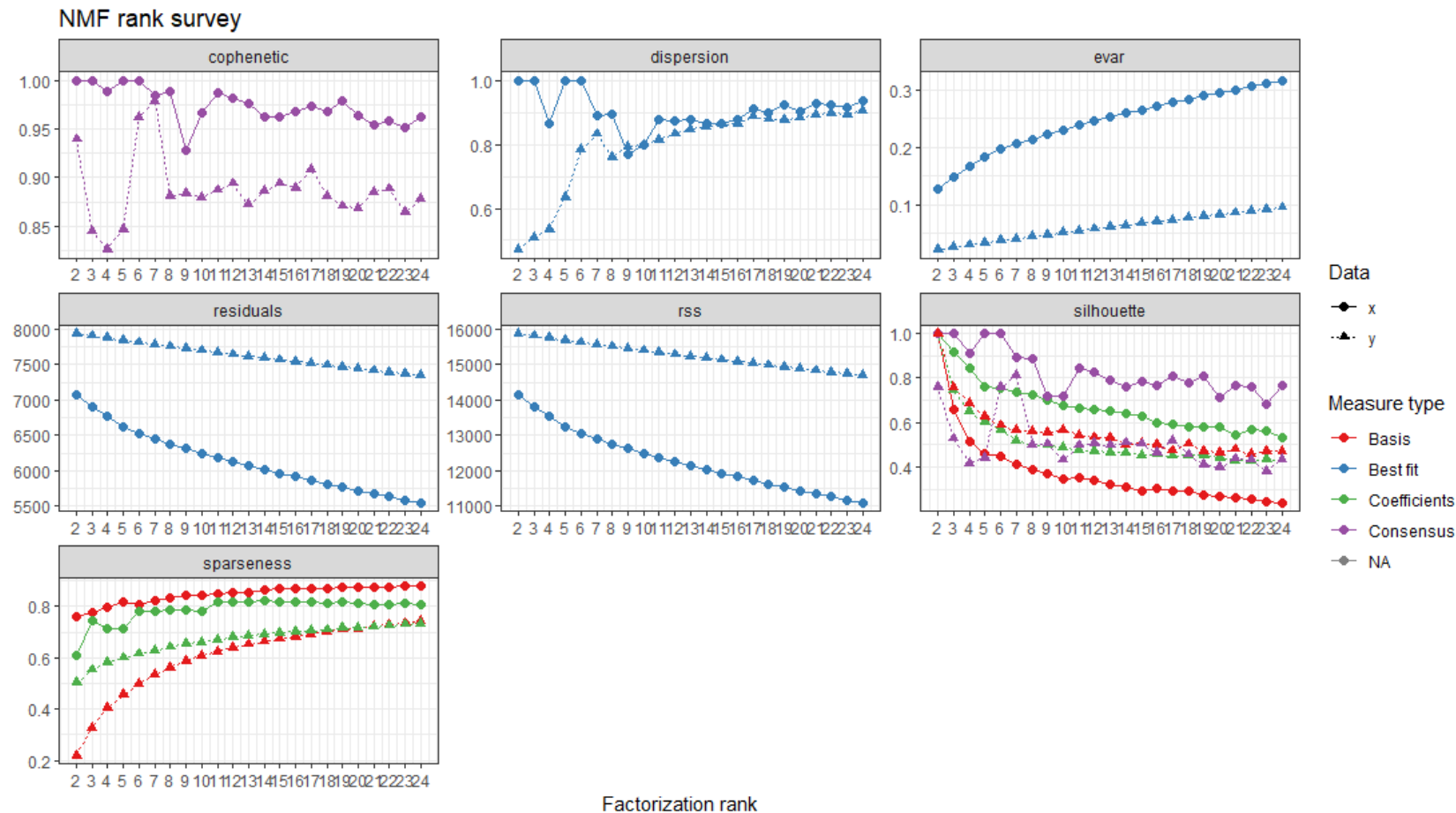| k | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|---|---|---|---|---|---|---|---|---|---|----|----|
| cases | **180** | 27 | 31 | **181** | 5 | 4 | 13 | **81** | 8 | **227** | 6 |

# Conclusions

- Is it possible to predict cancer type based on genes with somatic mutation in a patient?

  Is is possible to obtain a good predictive accuracy for specific cancer types whereas model performance is low for others

- Is there a «small» set of genes having a good predictive power, or at least as good as the entire set of genes?

  Yes, is is possible to select a set of 30-40 genes having a predictive power similar to the one obtained with all the genes

- Does the patient grouping based on similarity of mutated genes reflect the grouping based on cancer type?

  No, several cancer types seem to be composed by subsets of patients with different somatic mutation profiles

**Thanks for your attention!**