**EpiMed Open Course – Session 4**

**AI for Omics – Use case of leukemia classification – Part 2**
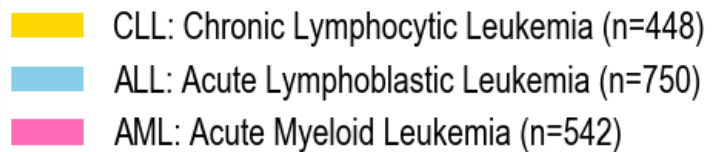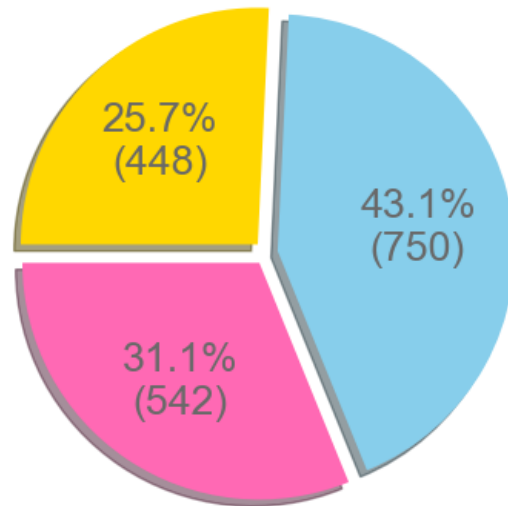
Ekaterina Flin

13/04/2020

# Use case: leukemia dataset

Goal: Predict leukemia type CLL, ALL or AML

**Leukemia Dataset - GSE13159**
**(n=1740)**

25.7%
(448)

43.1%
(750)

31.1%
(542)

CLL: Chronic Lymphocytic Leukemia (n=448)
ALL: Acute Lymphoblastic Leukemia (n=750)
AML: Acute Myeloid Leukemia (n=542)

Transcriptomic data, microarrays

Total number of samples = 1 740

Total number of genes ("features") = 21 875

3 labels to identify : CLL, ALL or AML
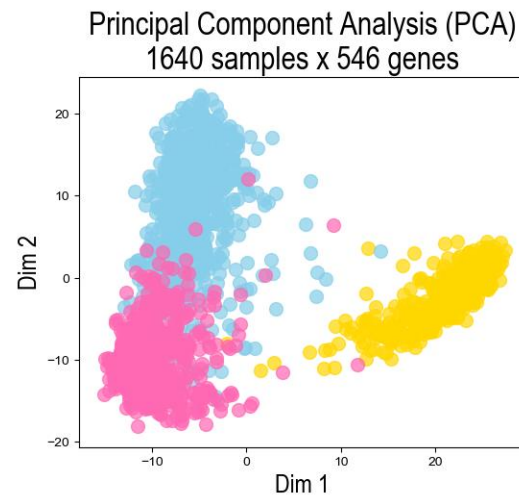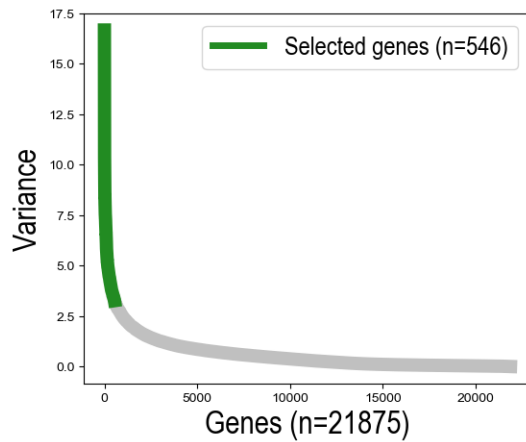
**1** Split dataset

| Training Set: 1640 | Leave-out Set: 100 |
|---|---|

**1** Split dataset

| Training Set: 1640 | Leave-out Set: 100 |
|---|---|

**2** Selection of variables



Principal Component Analysis (PCA)
1640 samples x 546 genes

# Pipeline for data preparation

**1** ## Split dataset

| Training Set: 1640 | Leave-out Set: 100 |
|---|---|

**3** ## Cross-validation

| Train | Train | Test |
|---|---|---|
| Test | Train | Train |
| Train | Test | Train |

**2** ## Selection of variables



Principal Component Analysis (PCA)
1640 samples x 546 genes

Pipeline for data preparation

**1** Split dataset

| Training Set: 1640 | Leave-out Set: 100 |

**3** Cross-validation

| Train | Train | Test |
| Test | Train | Train |
| Train | Test | Train |

**2** Selection of variables

Selected genes (n=546)

Principal Component Analysis (PCA)
1640 samples x 546 genes

Variance — Genes (n=21875) — Dim 1 — Dim 2

**4** Data scaling

| Train | Train | Test |

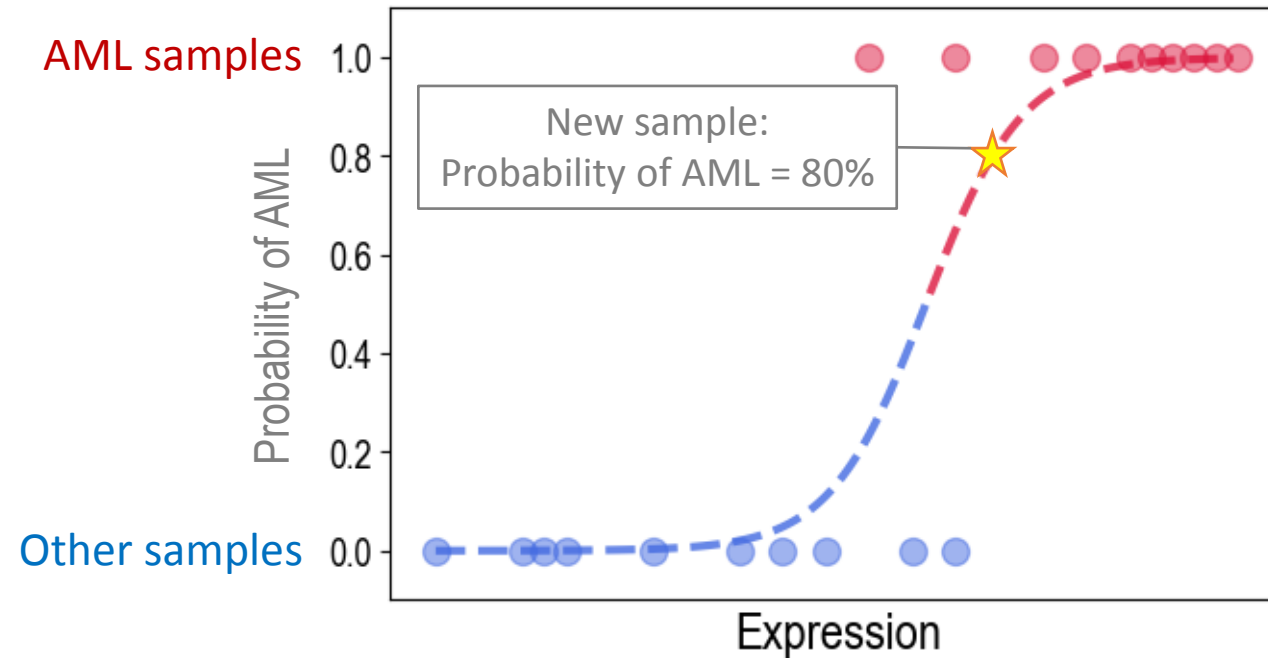$$scaled\_expression = \frac{expression - \mu}{\sigma}$$

# Train a classifier: Logistic Regression



Logistic function $\quad S(x) = \dfrac{1}{1 + e^{-x}} = \dfrac{e^x}{1 + e^x}$

S-shaped curve, also called "sigmoid" function

Logistic function $\quad S(x) = \dfrac{1}{1 + e^{-x}} = \dfrac{e^x}{1 + e^x}$

S-shaped curve, also called "sigmoid" function

# Logistic Regression in Python

```python
from sklearn.linear_model import LogisticRegression
from sklearn import metrics

# Create a classifier
classifier = LogisticRegression(multi_class='multinomial',
                                penalty='none', solver='newton-cg')

# Train classifier using training dataset
classifier.fit(X_train_scaled, y_train)

# Predict labels for test dataset
y_pred_test = classifier.predict(X_test_scaled)

# Calculate accuracy comparing prediction
# with real labels in test dataset
accuracy = metrics.accuracy_score(y_test, y_pred_test)
```
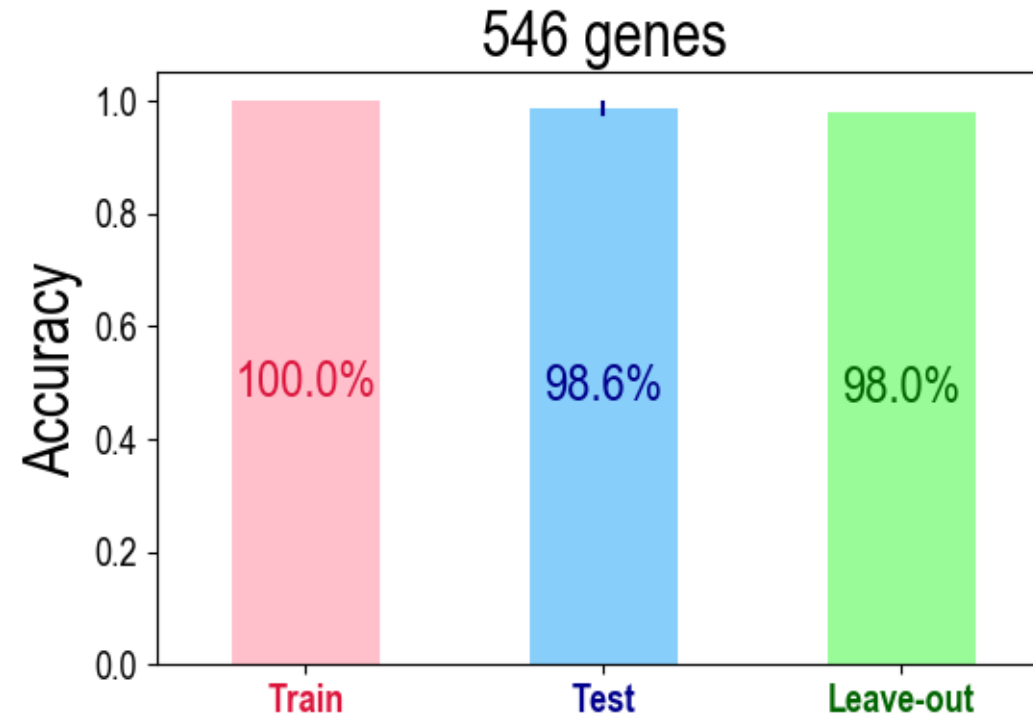
# Train a classifier: Logistic Regression

Training Set: 1640 | Leave-out Set: 100

| Train | Train | Test |
| Test | Train | Train |
| Train | Test | Train |



546 genes

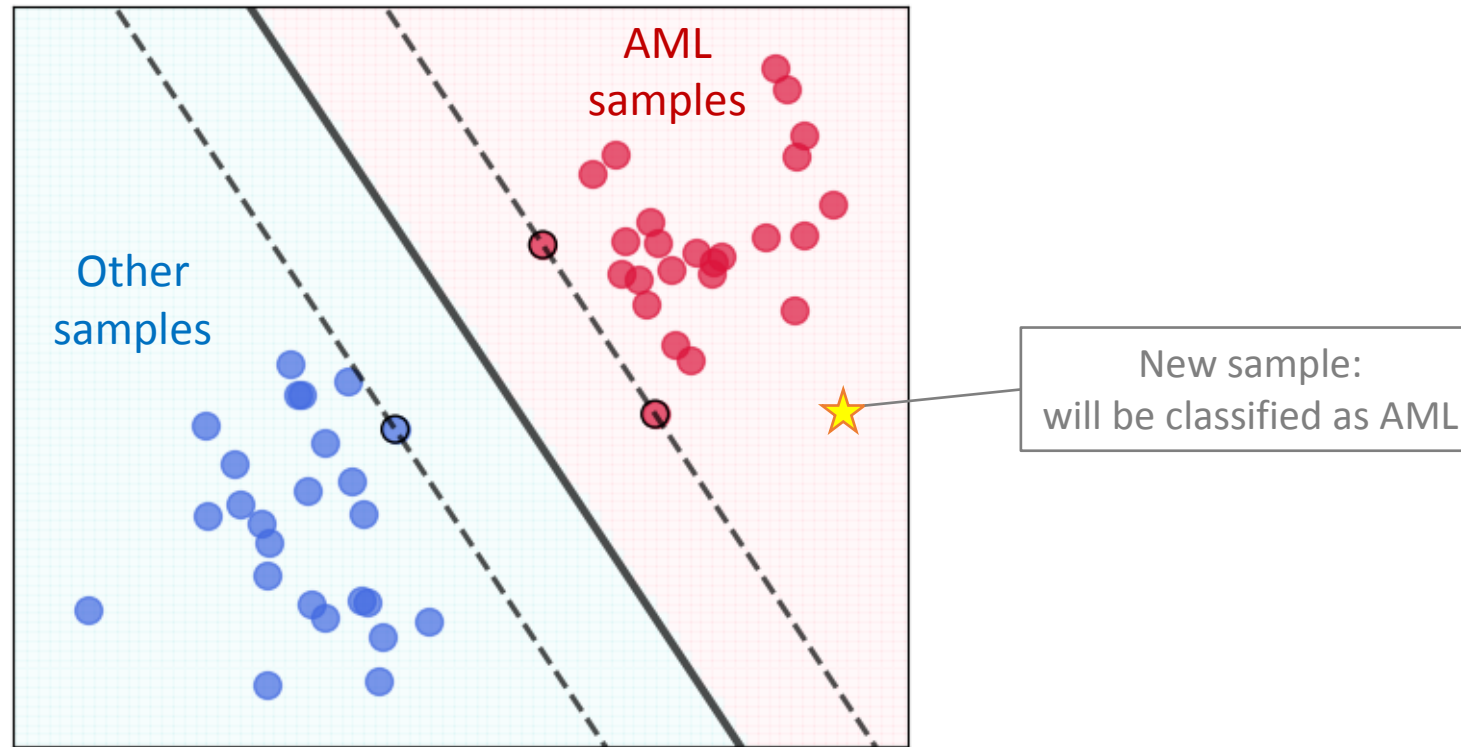Train: 100.0%  Test: 98.6%  Leave-out: 98.0%

Leave-out Set

Principal Component Analysis (PCA)
100 samples x 546 genes

- ● CLL: Chronic Lymphocytic Leukemia (nb=26)
- ● ALL: Acute Lymphoblastic Leukemia (nb=43)
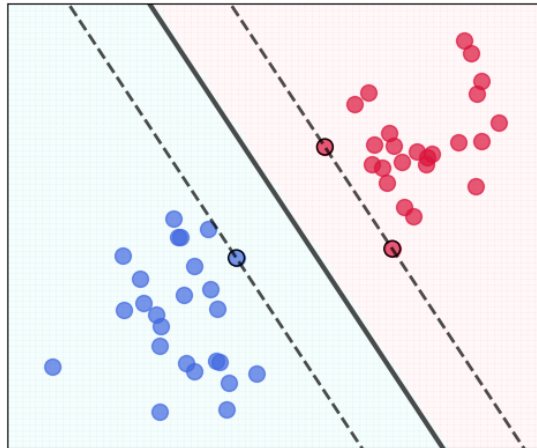- ✕ ALL predicted as AML by Logistic Regression
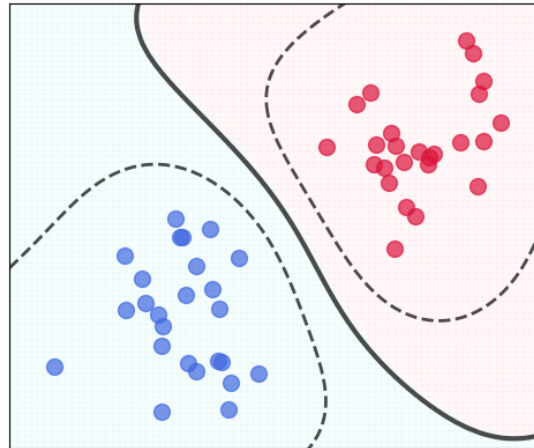- ● AML: Acute Myeloid Leukemia (nb=31)
- ✕ AML predicted as ALL by Logistic Regression

SVM method separates samples by a clear gap with the maximal possible margin.



In addition to performing linear classification, SVMs can efficiently perform a non-linear classification using a kernel, implicitly mapping their inputs into high-dimensional feature spaces.

SVM
546 genes
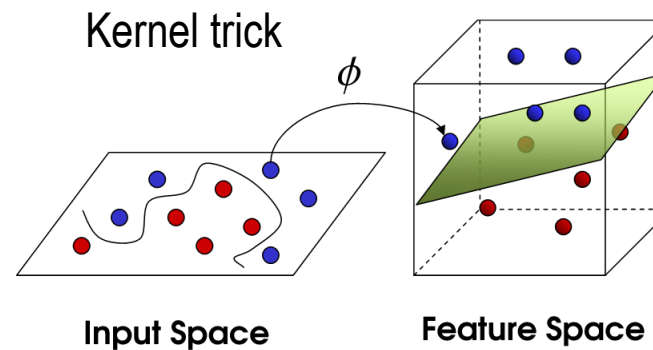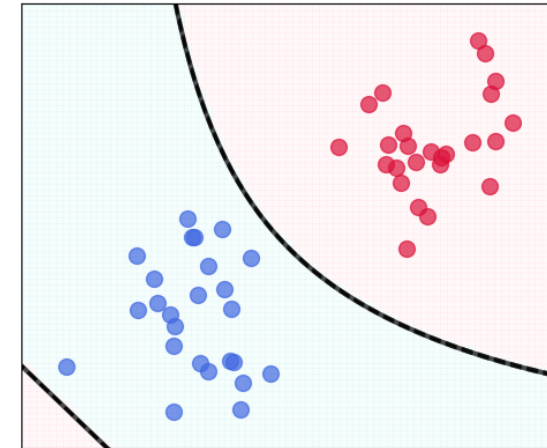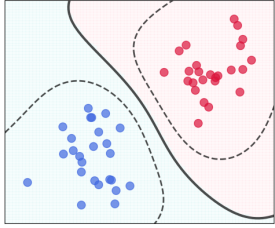
Leave-out Set

Principal Component Analysis (PCA)
100 samples x 546 genes

CLL: Chronic Lymphocytic Leukemia (nb=26)
ALL: Acute Lymphoblastic Leukemia (nb=43)
× ALL predicted as AML by Linear SVM
AML: Acute Myeloid Leukemia (nb=31)
× AML predicted as ALL by Linear SVM

In addition to performing linear classification, SVMs can efficiently perform a non-linear classification using a kernel trick, implicitly mapping their inputs into high-dimensional feature spaces.

No kernel (linear)

RBF kernel

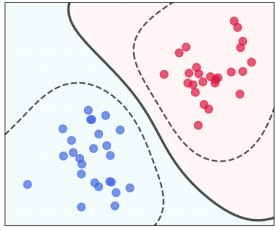Sigmoid kernel

Kernel trick

$\phi$

Input Space

Feature Space

To train a non-linear SVM classifier we need to fix two parameters:

- Parameter C: degree of ridge regularization
  *how smooth should our model be to avoid overfitting?*

- Parameter $\gamma$: degree of non-linearity
  *strongly non-linear model or almost linear?*

classifier = SVC(kernel='rbf', **C=?**, **gamma=?**)

Grid search for optimal parameters using cross-validation



classifier = SVC(kernel='rbf', C=**6**, gamma=**1e-3**)

Random Forest creates several decision trees using random subsets of data. The final result is obtained by a vote from all decision trees.
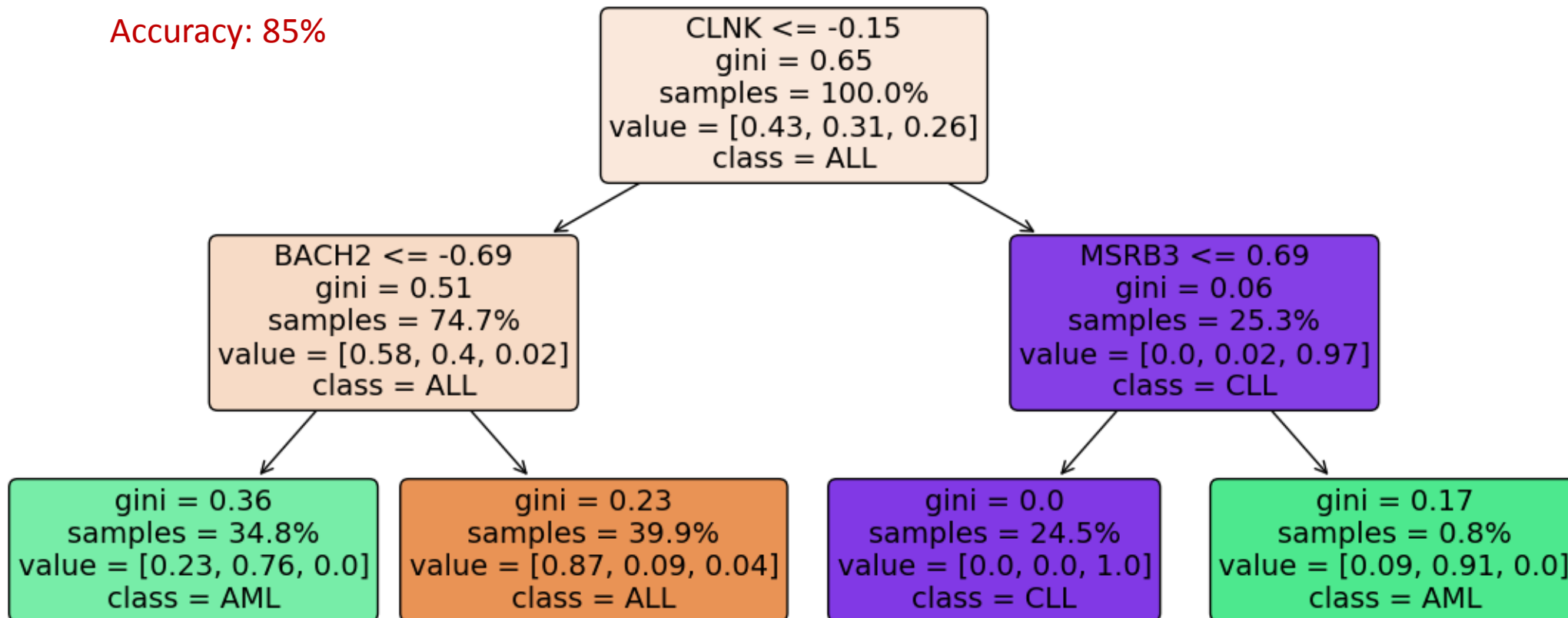


*A simple decision tree*

# Train a classifier: Decision Tree

A simple forest composed of a unique tree:

Accuracy: 85%

```
CLNK <= -0.15
gini = 0.65
samples = 100.0%
value = [0.43, 0.31, 0.26]
class = ALL
```

```
BACH2 <= -0.69
gini = 0.51
samples = 74.7%
value = [0.58, 0.4, 0.02]
class = ALL
```

```
MSRB3 <= 0.69
gini = 0.06
samples = 25.3%
value = [0.0, 0.02, 0.97]
class = CLL
```

```
gini = 0.36
samples = 34.8%
value = [0.23, 0.76, 0.0]
class = AML
```

```
gini = 0.23
samples = 39.9%
value = [0.87, 0.09, 0.04]
class = ALL
```

```
gini = 0.0
samples = 24.5%
value = [0.0, 0.0, 1.0]
class = CLL
```

```
gini = 0.17
samples = 0.8%
value = [0.09, 0.91, 0.0]
class = AML
```
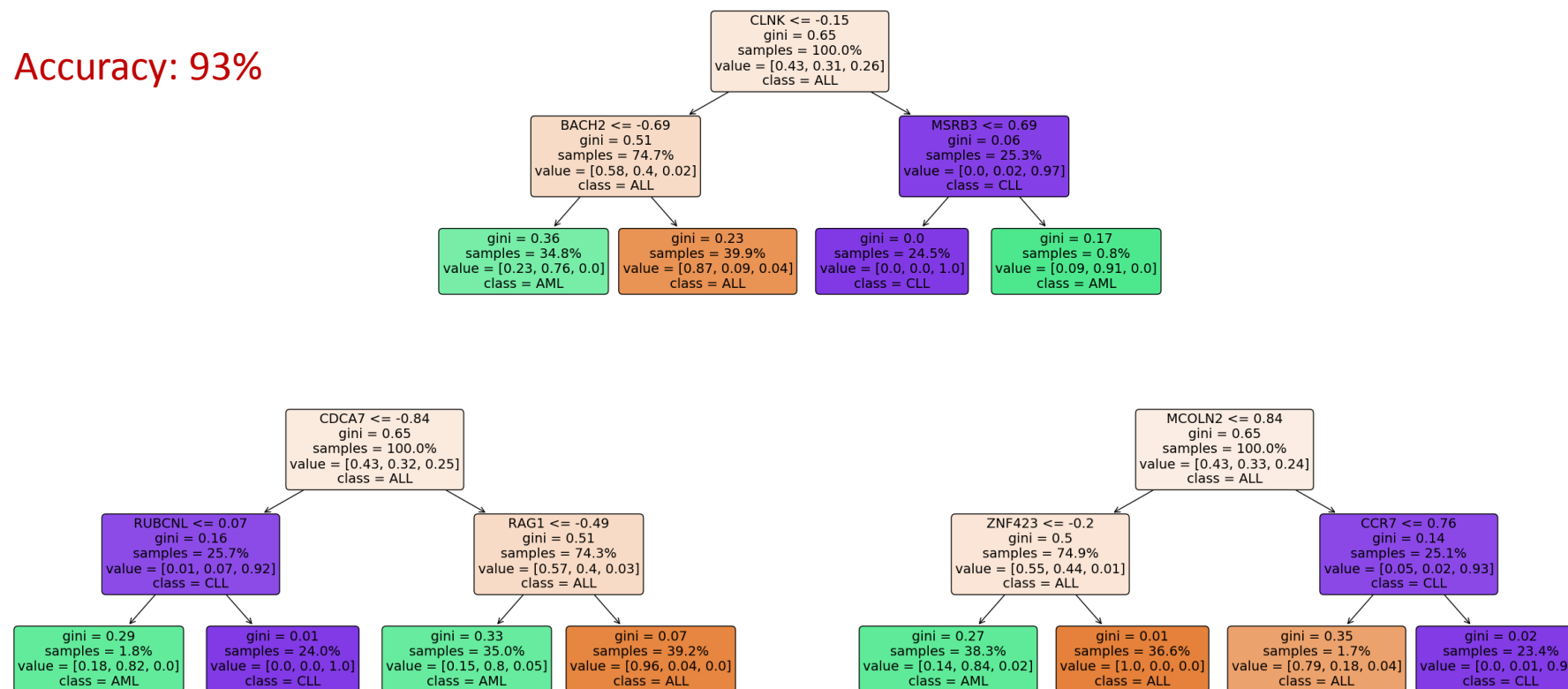
```
classifier = RandomForestClassifier(n_estimators=1, max_depth=2)
```

# Train a classifier: Random Forest

A simple Random Forest:
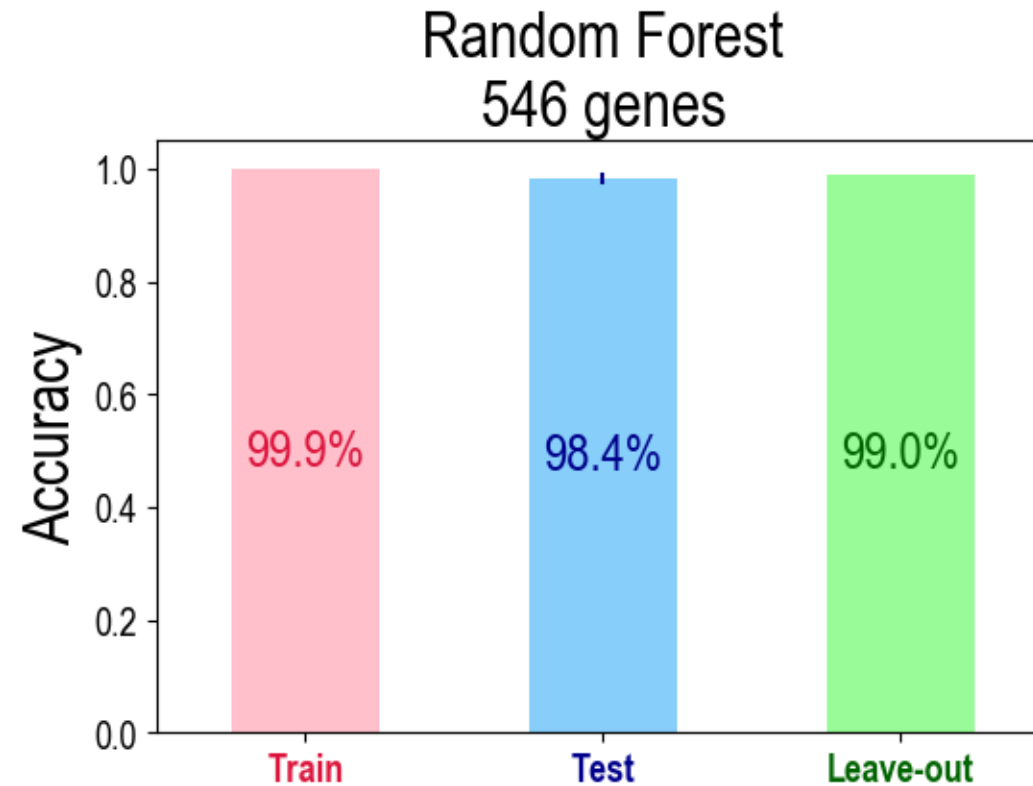number of trees = 3, depth of each tree = 2

Accuracy: 93%



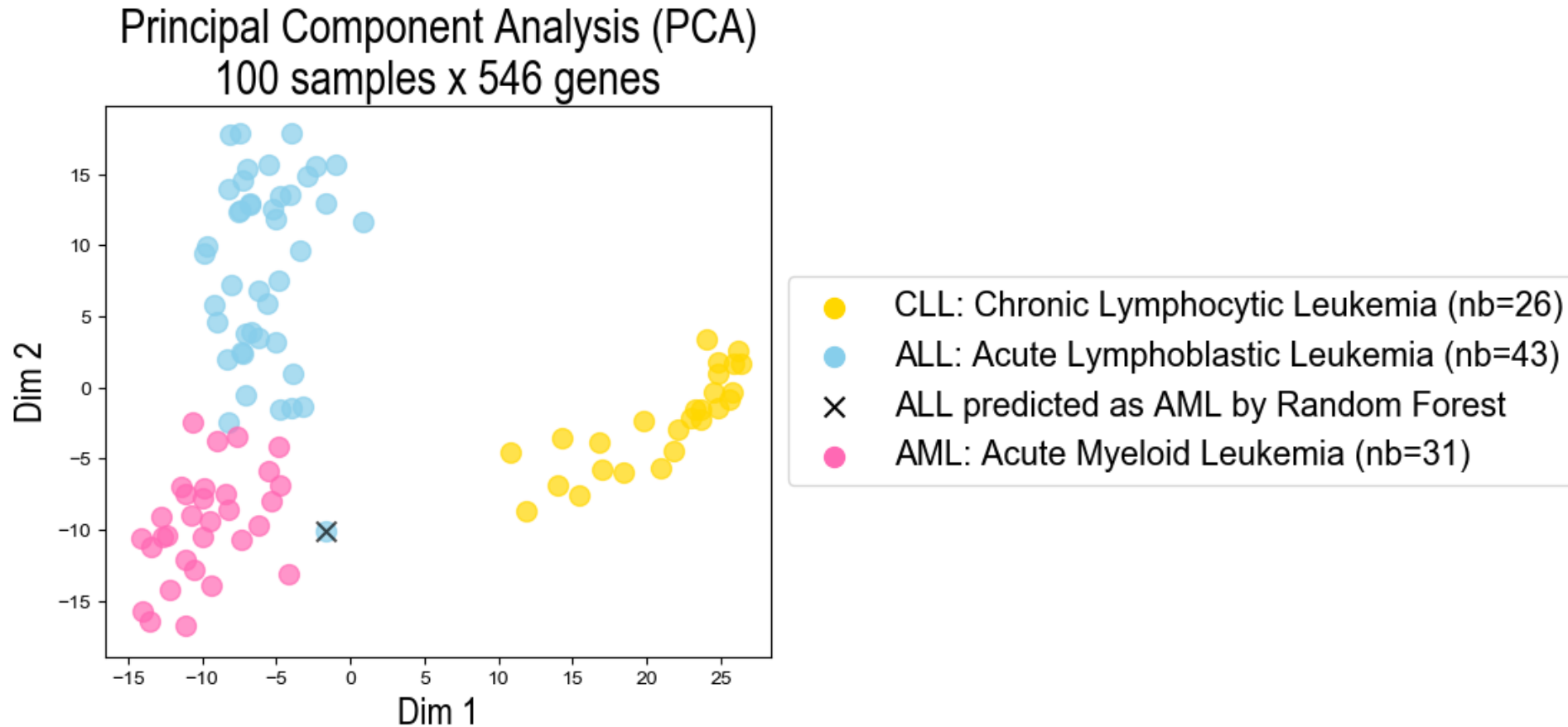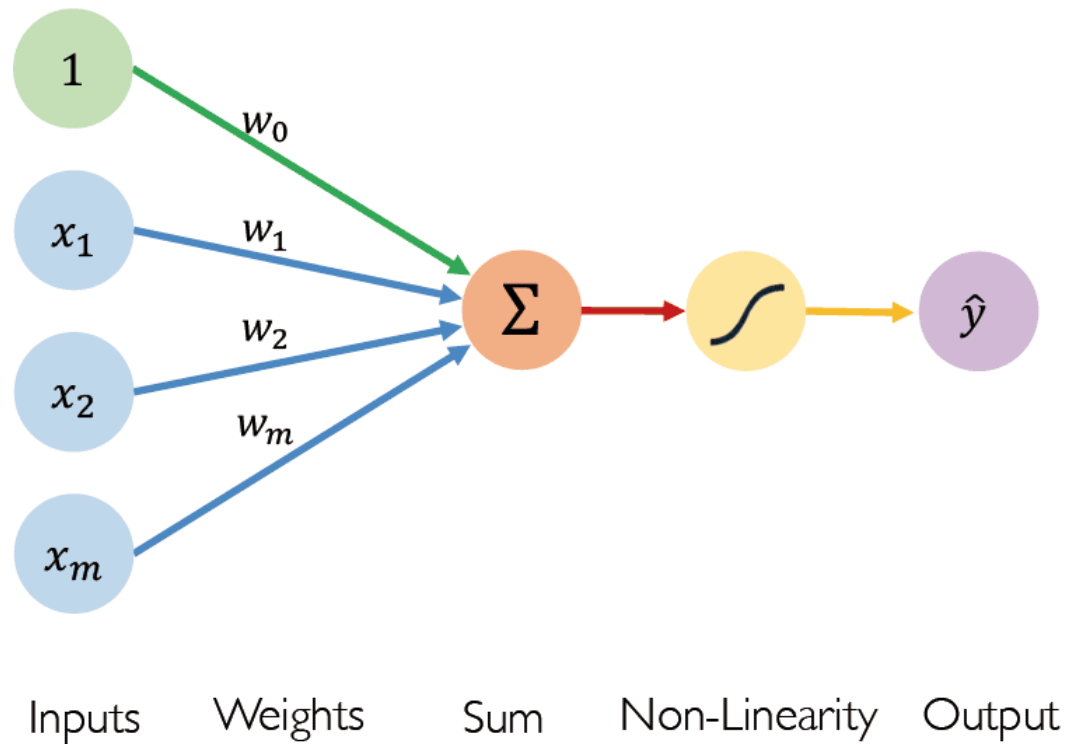classifier = RandomForestClassifier(n_estimators=3, max_depth=2)

Definitive Random Forest model:
number of trees = 50, depth of each tree = 8



```
classifier = RandomForestClassifier(n_estimators=50, max_depth=8)
```
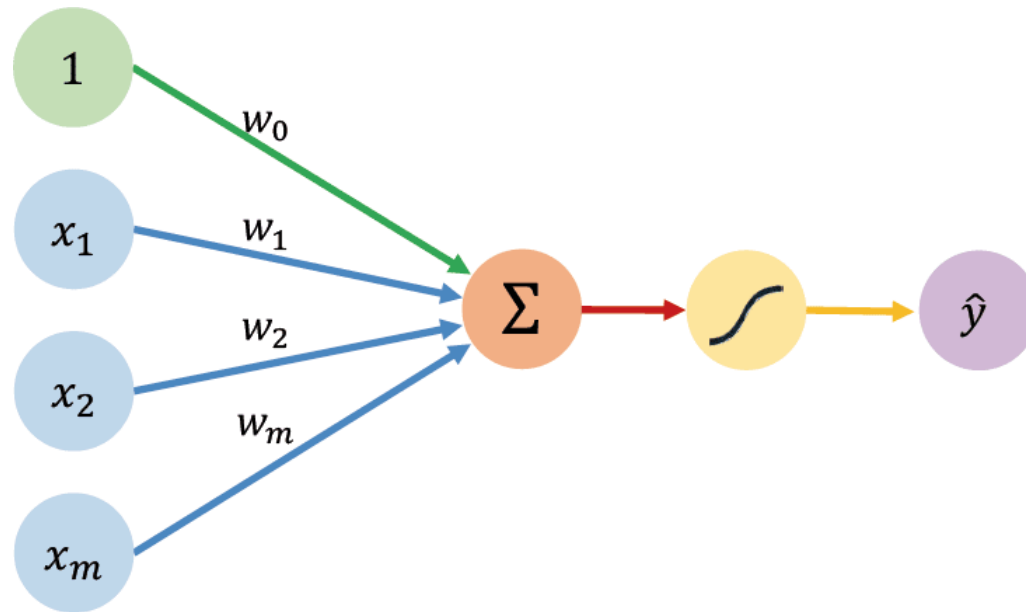
$$\hat{y} = g\left(w_0 + \sum_{i=1}^{m} x_i \, w_i\right)$$

Output

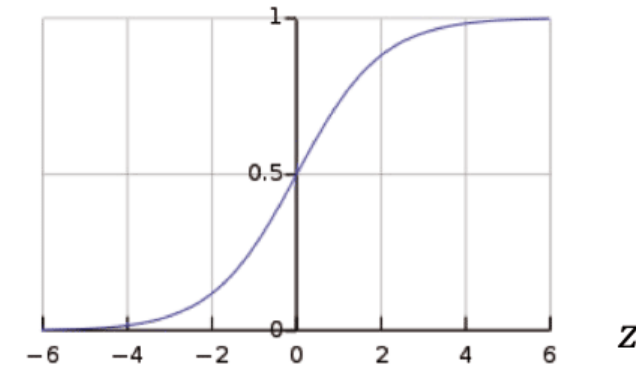Linear combination of inputs

Non-linear activation function

Bias

Inputs   Weights   Sum   Non-Linearity   Output

## Activation Functions

$$\hat{y} = g\left( w_0 + X^T W \right)$$

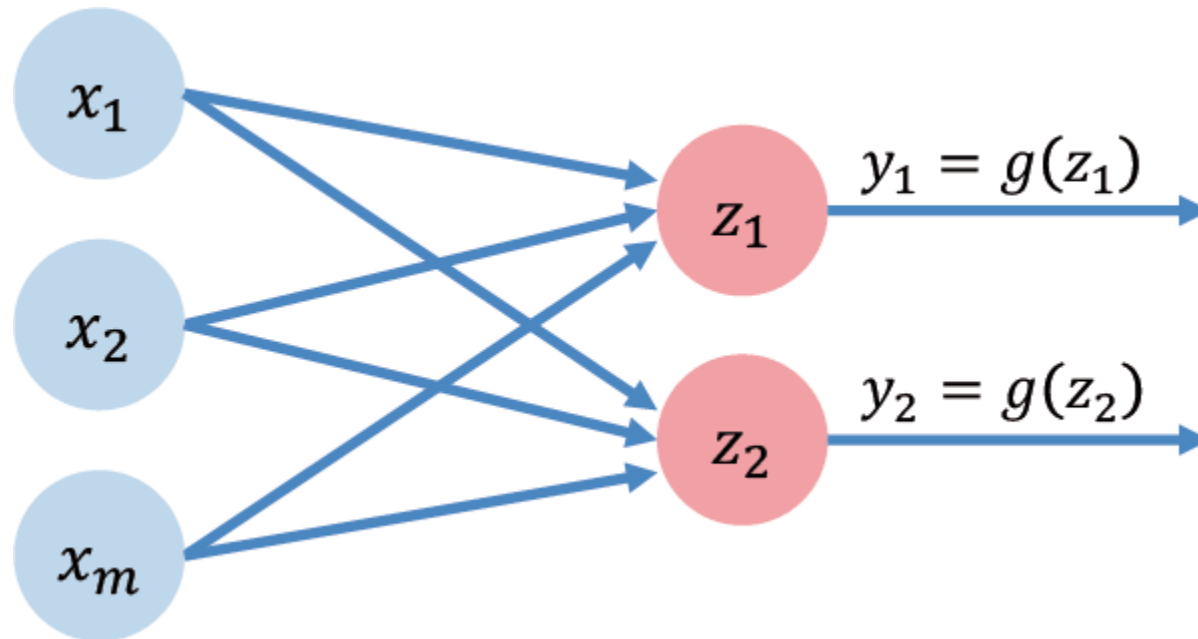- Example: sigmoid function

$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Train a perceptron with sigmoid activation function is mathematically equivalent to train a Logistic Regression.

Source: MIT Course 6.S191 Introduction to Deep Learning

# Multi-layer perceptron (Neural Network)

Simplified representation of a perceptron

Perceptron with two outputs

# Multi-layer perceptron (Neural Network)

Single layer neural network



Inputs             Hidden layer           Output

# Multi-layer perceptron (Neural Network)

Deep neural network



Inputs

... Many hidden layers ...

Output

# Multi-layer perceptron in Python

```python
from sklearn.neural_network import MLPClassifier
from sklearn import metrics

# Create a classifier
classifier = MLPClassifier(early_stopping=True, alpha=1.0,
                           hidden_layer_sizes=(100, 100, 100, 100))

# Train classifier using training dataset
classifier.fit(X_train_scaled, y_train)

# Predict labels for test dataset
y_pred_test = classifier.predict(X_test_scaled)

# Calculate accuracy comparing prediction
# with real labels in test dataset
accuracy = metrics.accuracy_score(y_test, y_pred_test)
```
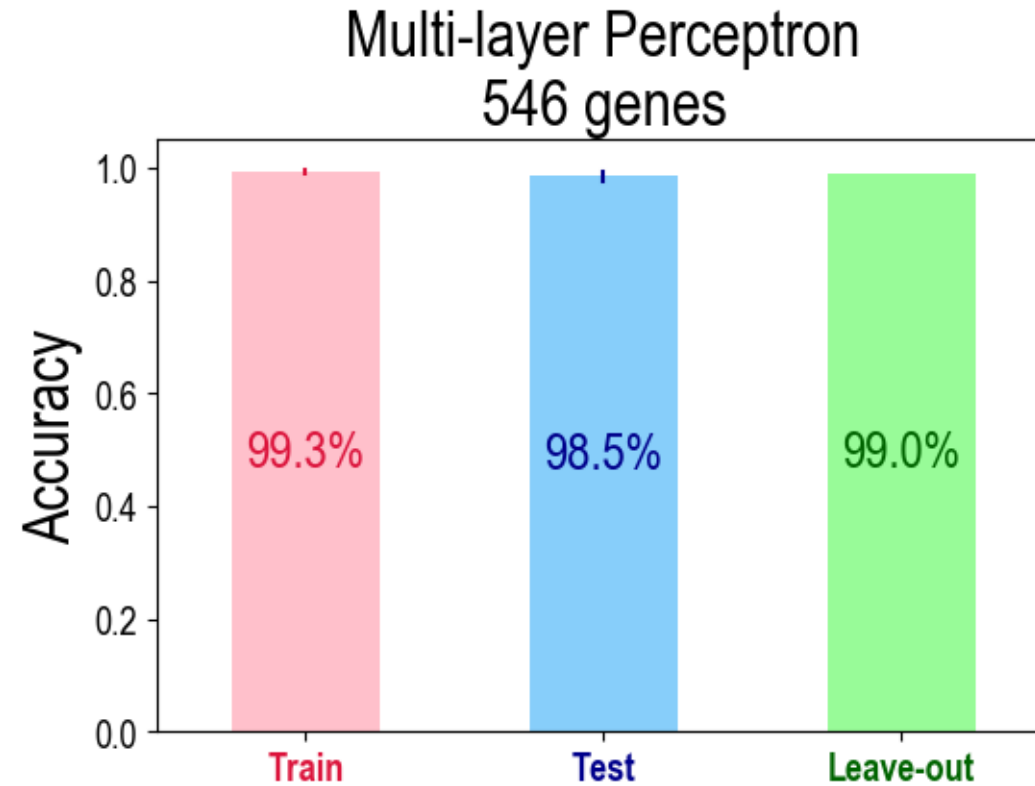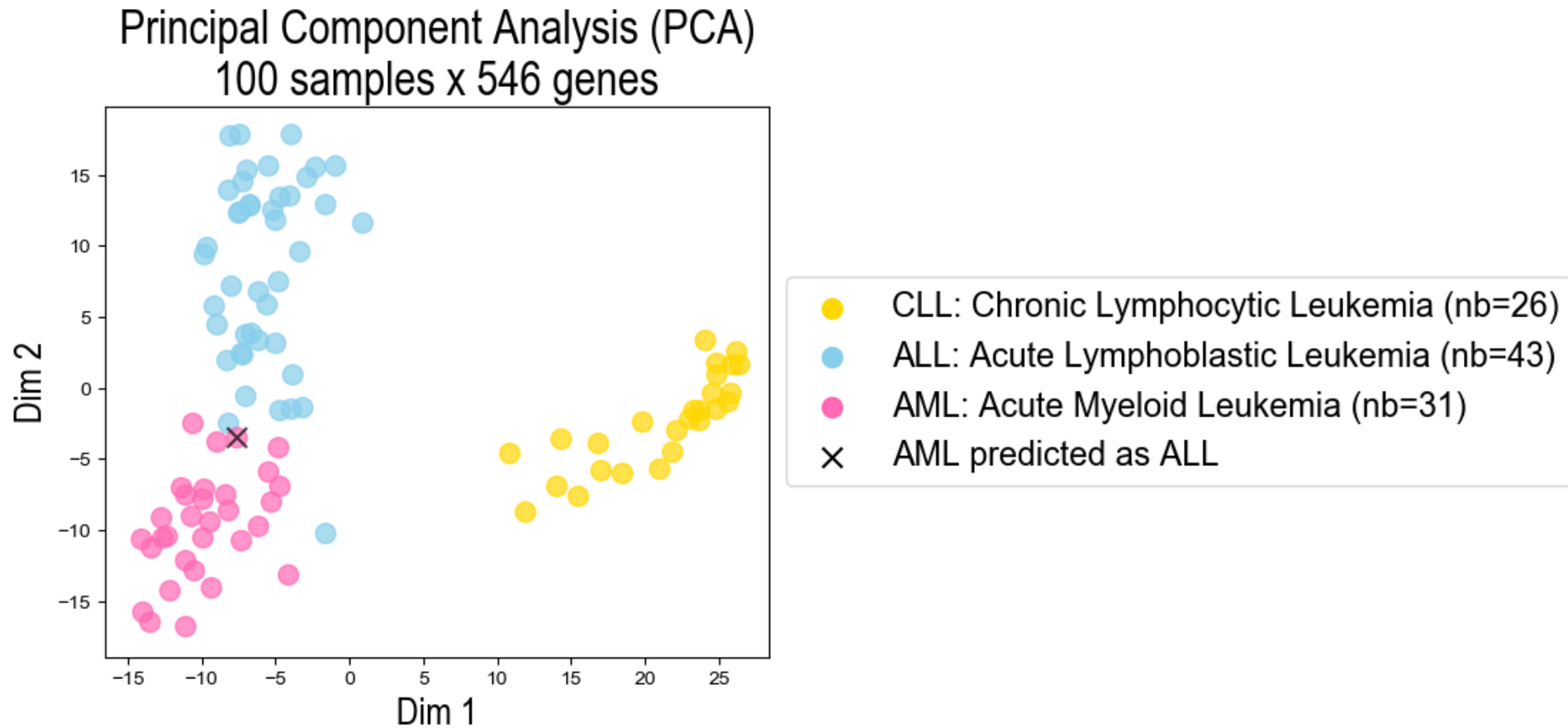
Principal Component Analysis (PCA)
100 samples x 546 genes

CLL: Chronic Lymphocytic Leukemia (nb=26)
ALL: Acute Lymphoblastic Leukemia (nb=43)
AML: Acute Myeloid Leukemia (nb=31)
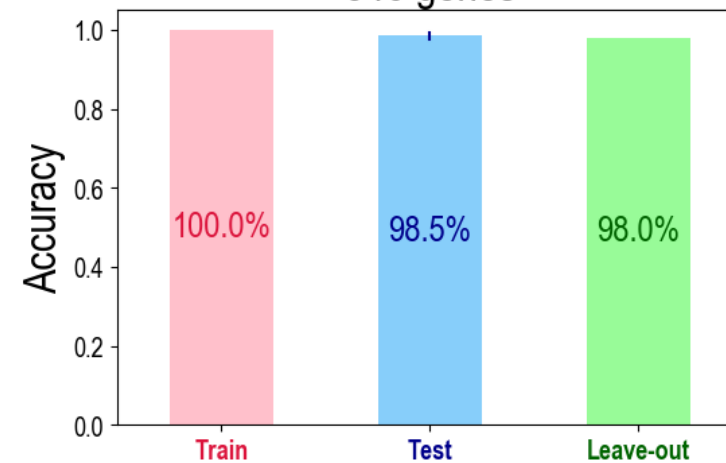× AML predicted as ALL
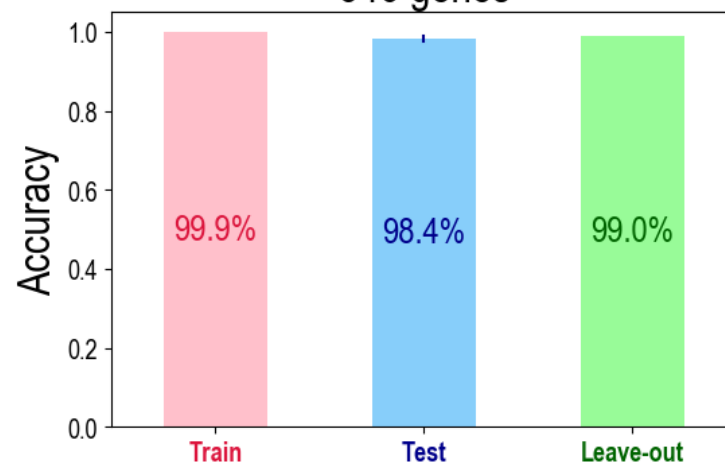
# Benchmark

Logistic Regression
546 genes

SVM
546 genes
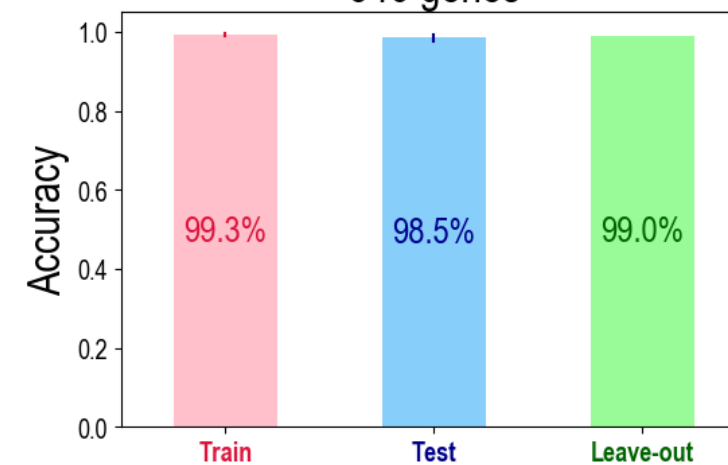
SVM with RBF kernel
546 genes

Random Forest
546 genes

Multi-layer Perceptron
546 genes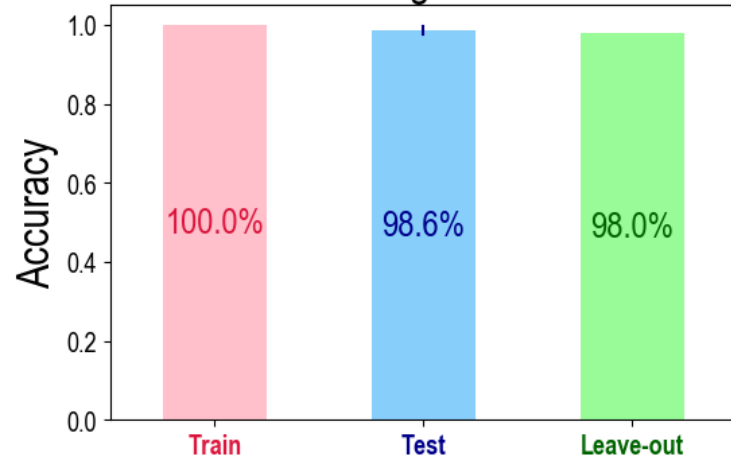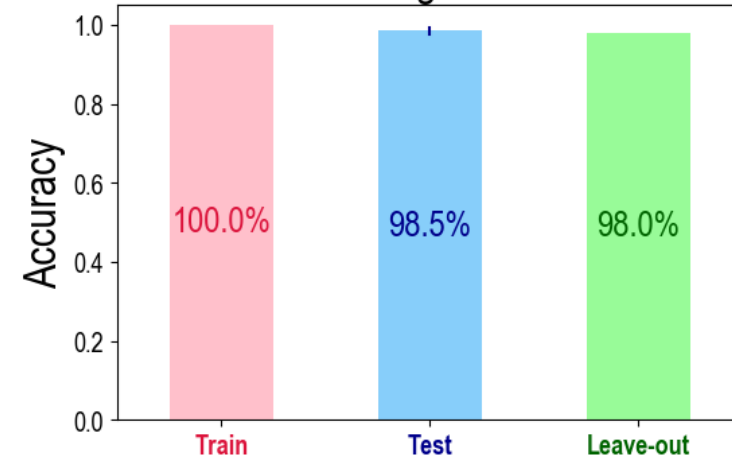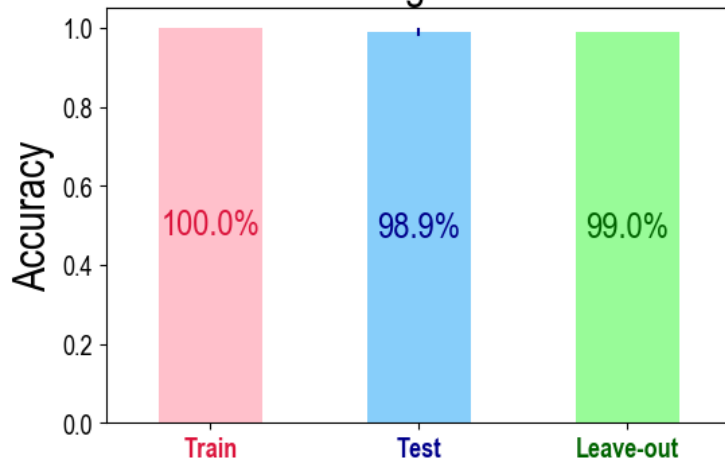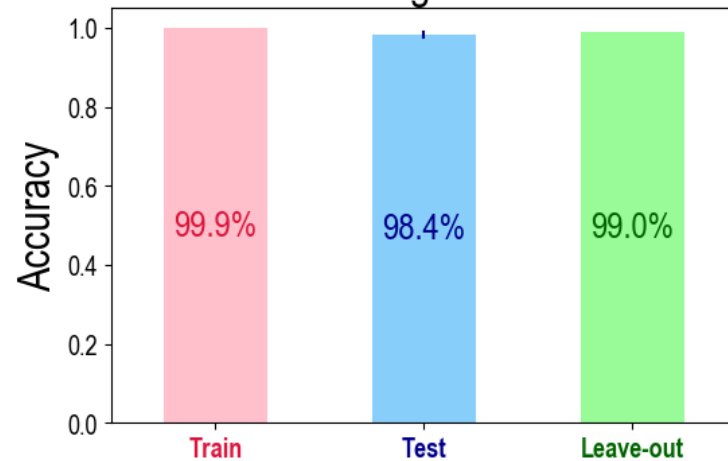