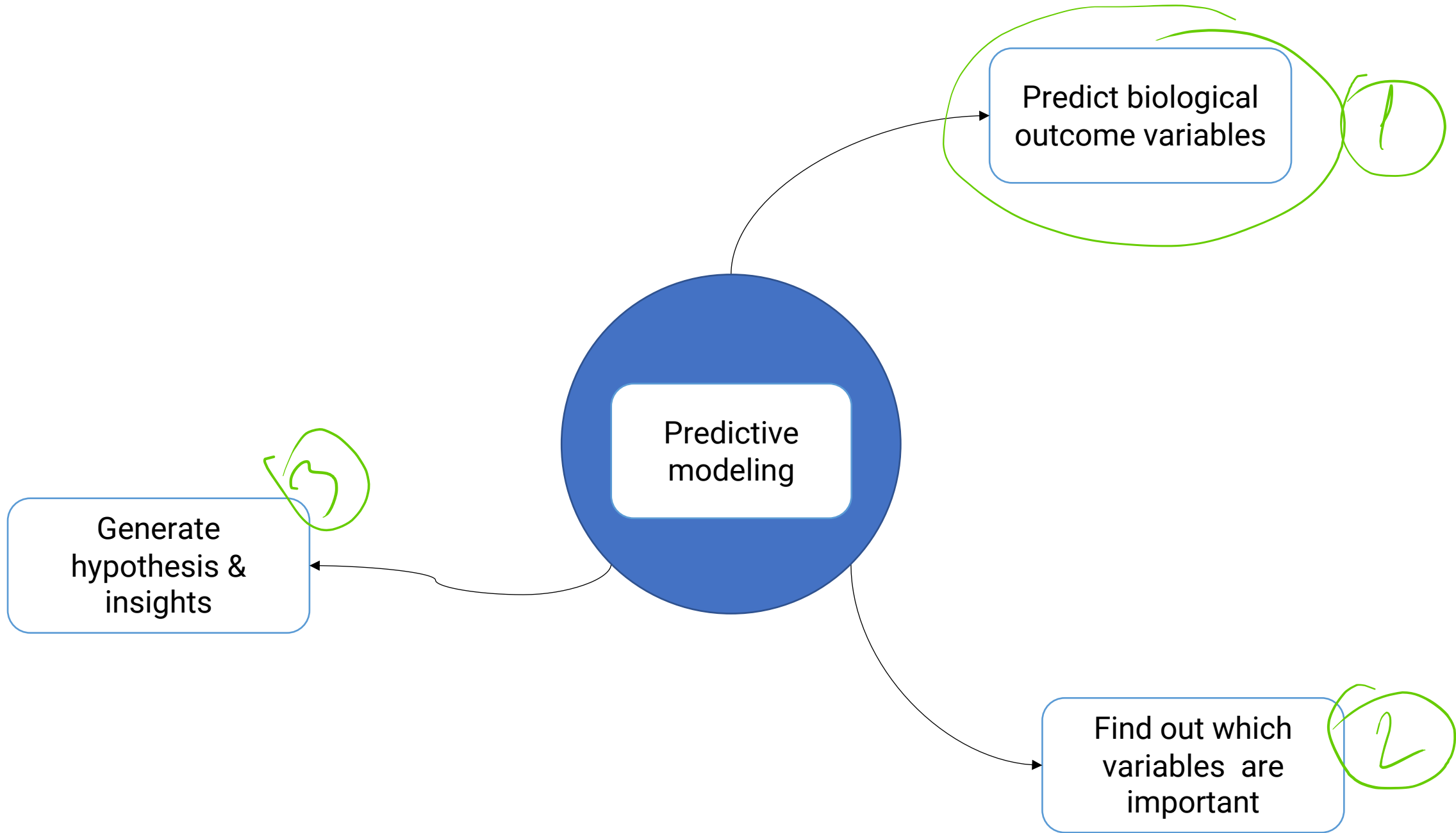# Computational genomics: hands on course

## Predictive modeling with supervised learning

# How are machine learning models fit ?

- Define a prediction function or method f(X) → LR SVM RF

- Devise a loss/cost function: such as $\sum(Y-f(X))^2$ org pred.

- Apply optimization & find best parameters for $\sum(Y-f(X))^2$ f(t)

X-WH

k-means
cost nax TSS cluster

# Steps for supervised learning

1) Pre-processing data → transform, filter, normalize

2) Data split → Training & test

3) Training → specific model / algo.

4) Performance estimation → error

5) Model tuning → to further optimize parameters

R packages: mlr and caret
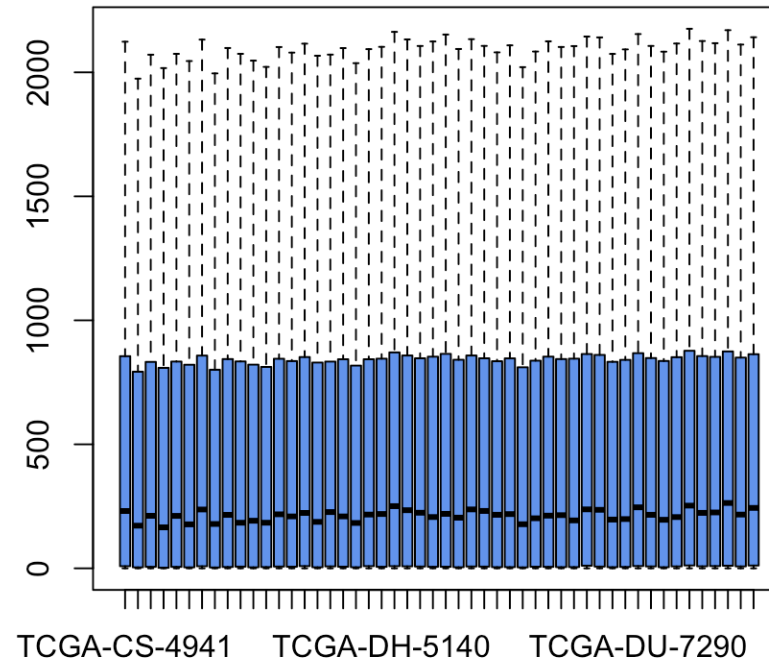
# Use case: Disease subtype from genomics data

# Data pre-processing
## Data transformation

Systematic differences between samples and outliers are a problem for fitting ML models

Gene expression values from glioblastoma samples

# Data pre-processing
## Data transformation

Some of the data transformation operations
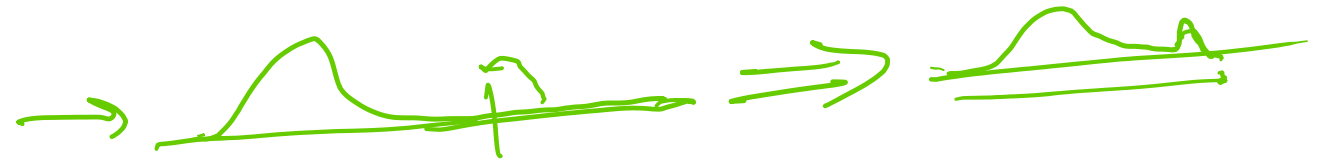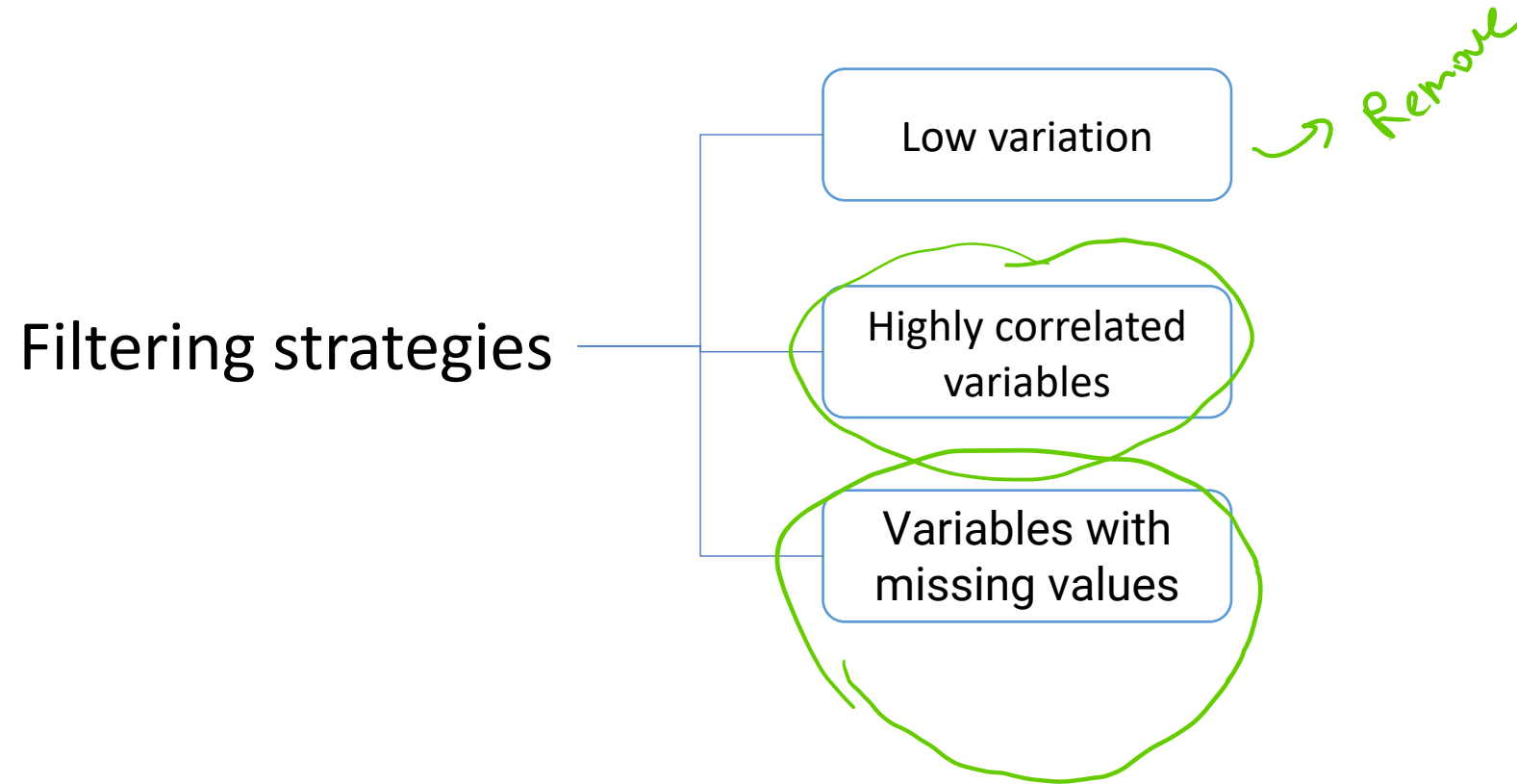
- Normalization /scaling
- Log transform
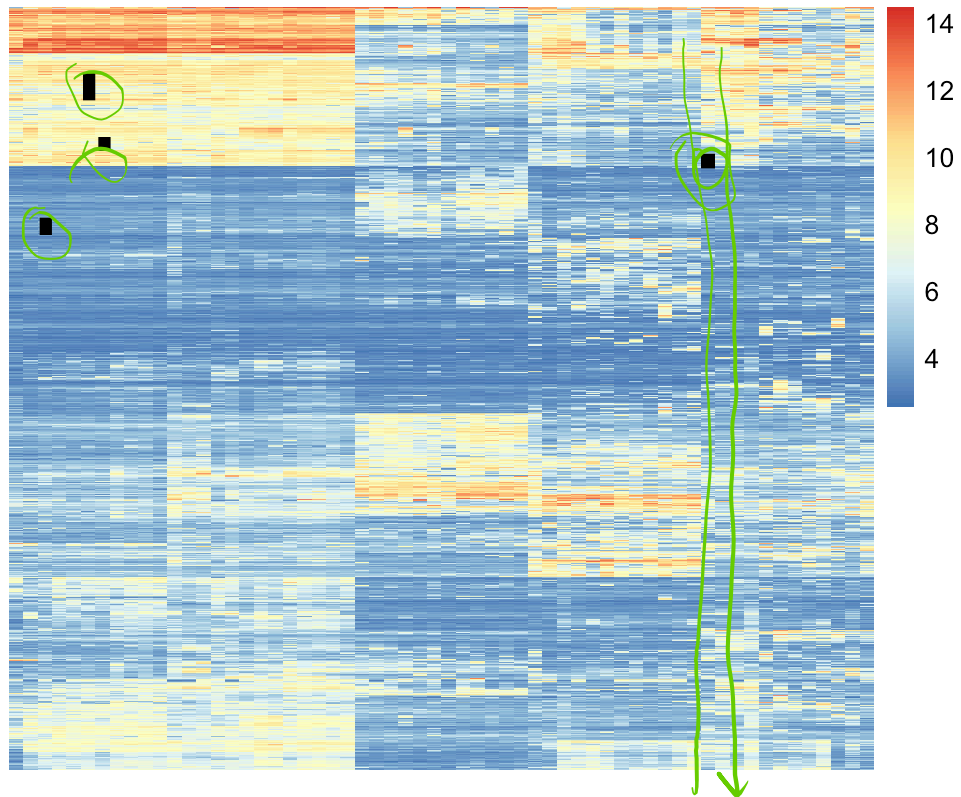- Box Cox transform
- winsorizing

# Data pre-processing
## Data filtering

Filtering strategies

Low variation → *Remove*

Highly correlated variables

Variables with missing values

# Data pre-processing
## Dealing with missing values
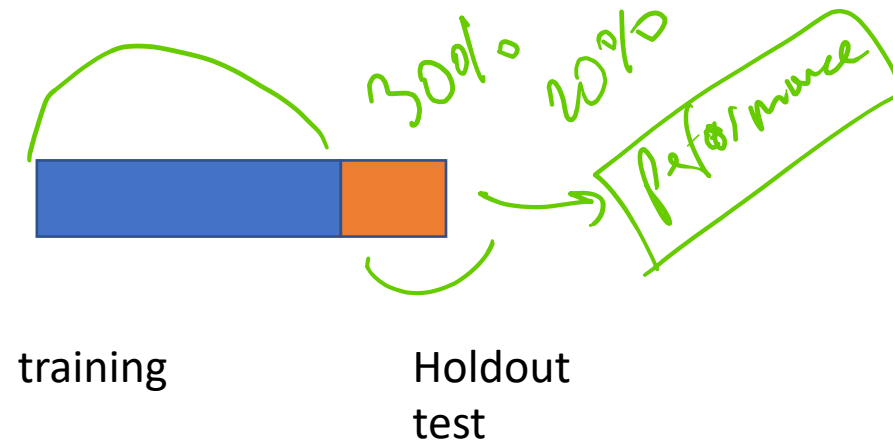


Choices:
1) Remove samples/variables with missing values

2) Assign the mean/median value

3) Try to predict missing values

```
knnImpute=prePprocess(missing_tgexp,method="knnImpute")
```
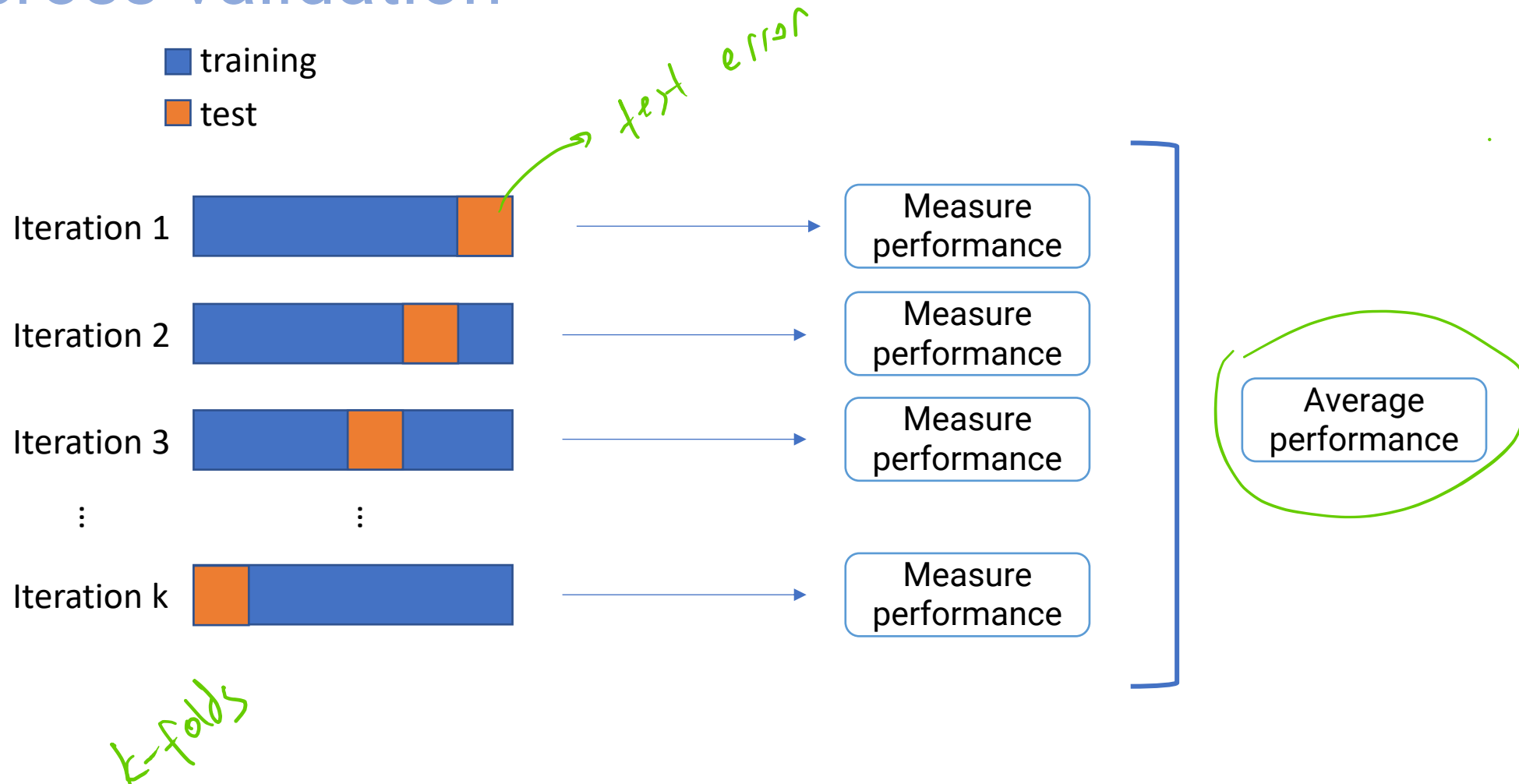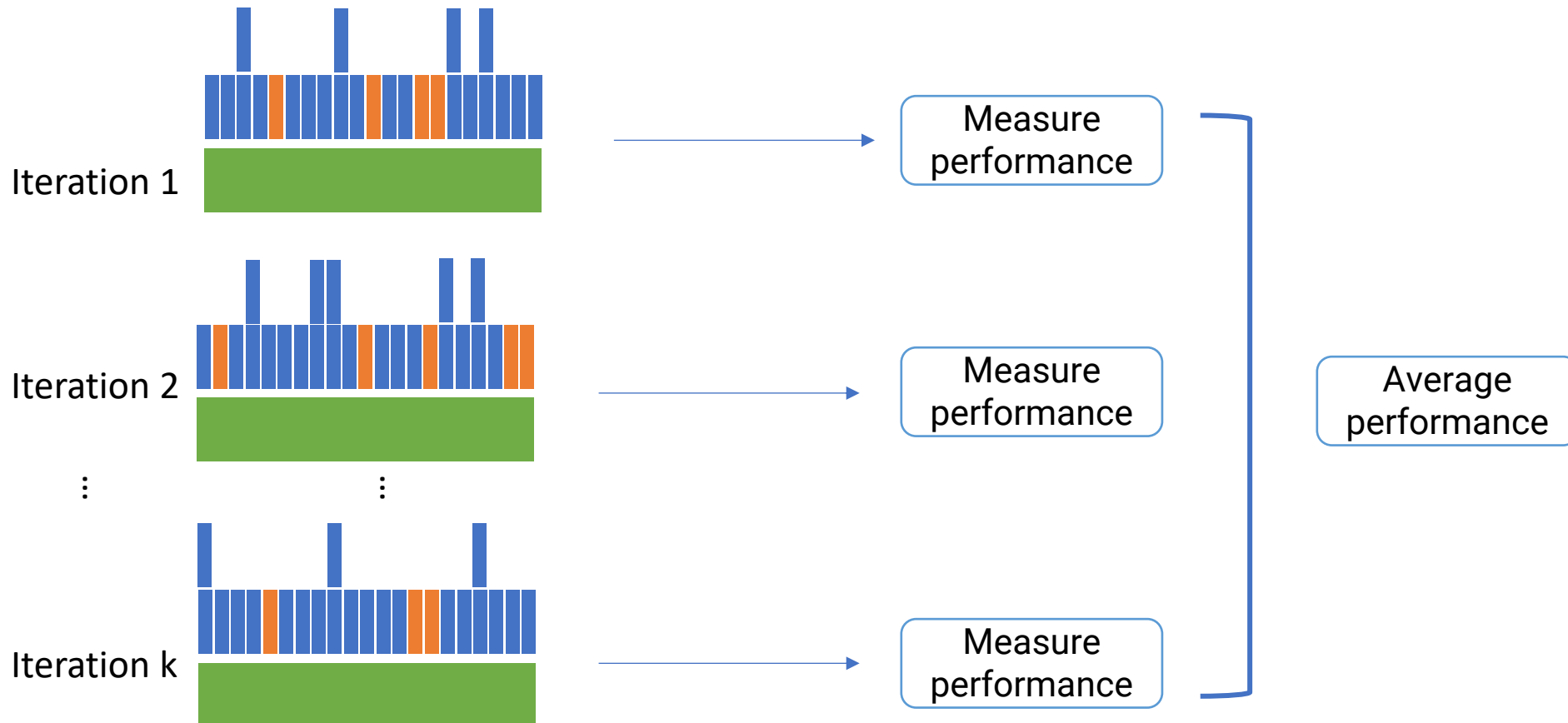
# Data split
## Holdout test dataset

training                 Holdout
test

30%    20%

performance

# Data split
## cross-validation

training
test

*test error*

| Iteration 1 | → | Measure performance |
| Iteration 2 | → | Measure performance |
| Iteration 3 | → | Measure performance |
| ⋮ | | |
| Iteration k | → | Measure performance |

Average performance

*k-folds*

# Data split
## bootstrap resampling



- ■ Training (bootstrap sample)
- ■ Test: out-of-the-bag (OOB)
- ■ All the data

Iteration 1 → Measure performance

Iteration 2 → Measure performance

⋮ ⋮

Iteration k → Measure performance

Average performance

# Predicting the subtype with k-nearest neighbors

k-NN in a nutshell: find similar patients and use their labels

Class B

unkown

Class A

class A

New point is assigned to class A by k-NN where k=3

# Assessing the performance of our model

| | Actual CIMP | Actual noCIMP |
|---|---|---|
| Predicted as CIMP | True Positives (TP) | False Positive (FP) |
| Predicted as noCIMP | False Positives (FN) | True negatives (TN) |

Precision, $TP/(TP + FP)$

Sensitivity, $TP/(TP + FN)$

Specificity, $TN/(TN + FP)$

# Assessing the performance of our model
## Receiver Operating Characteristic (ROC) Curves

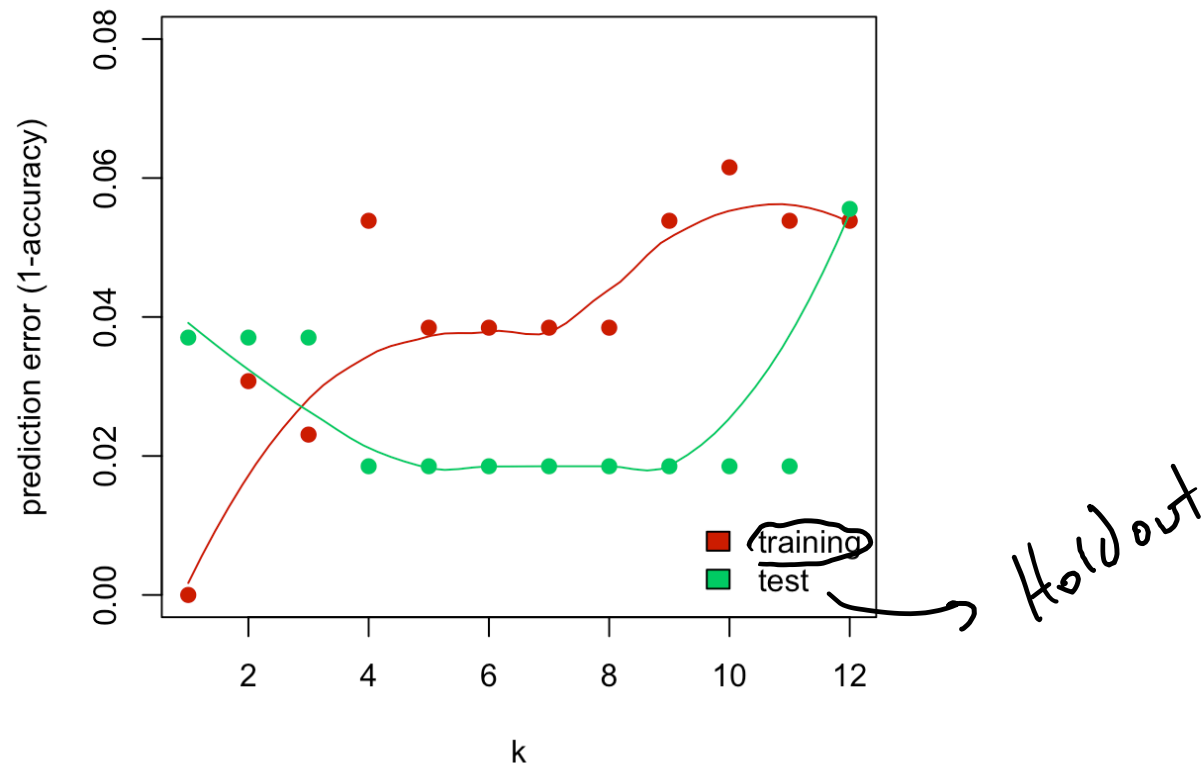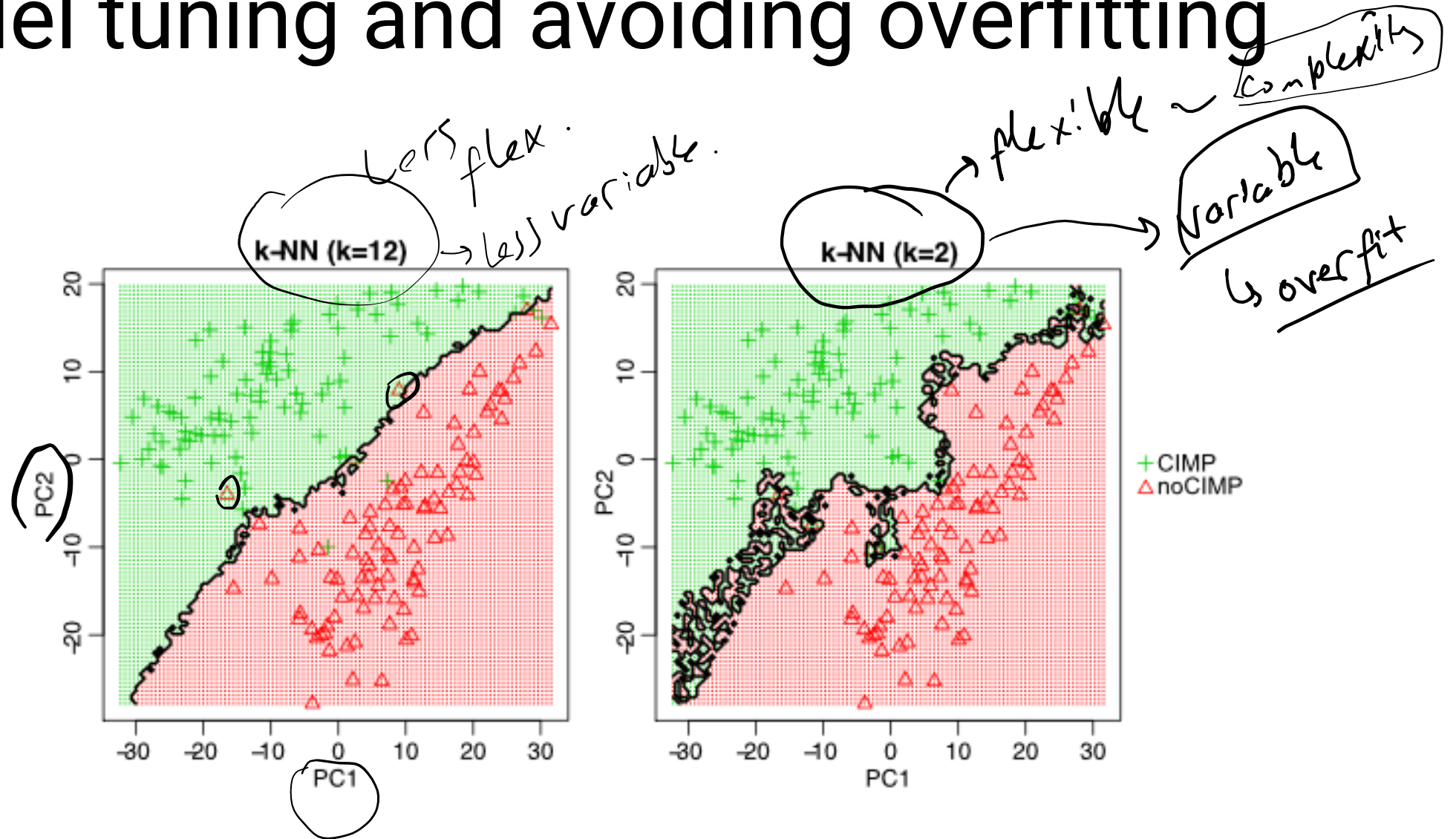# Model tuning and avoiding overfitting

# Model tuning and avoiding overfitting

# Model tuning and avoiding overfitting

$$E[Error] = Bias^2 + Variance + \sigma_e^2$$

$\downarrow$ decrease complexity

Var(Prediction)

# variables parameters

Optimal complexity

# Trees and forests
## Decision trees

# Trees and forests
## Decision trees



Is Male

yes — no

Age > 60

Is PIGX exp > 100

yes    no

yes    no

75% class A
25% class B

5% class A
95% class B

Gini Impurity → 0.5

0.235

→ 0.095

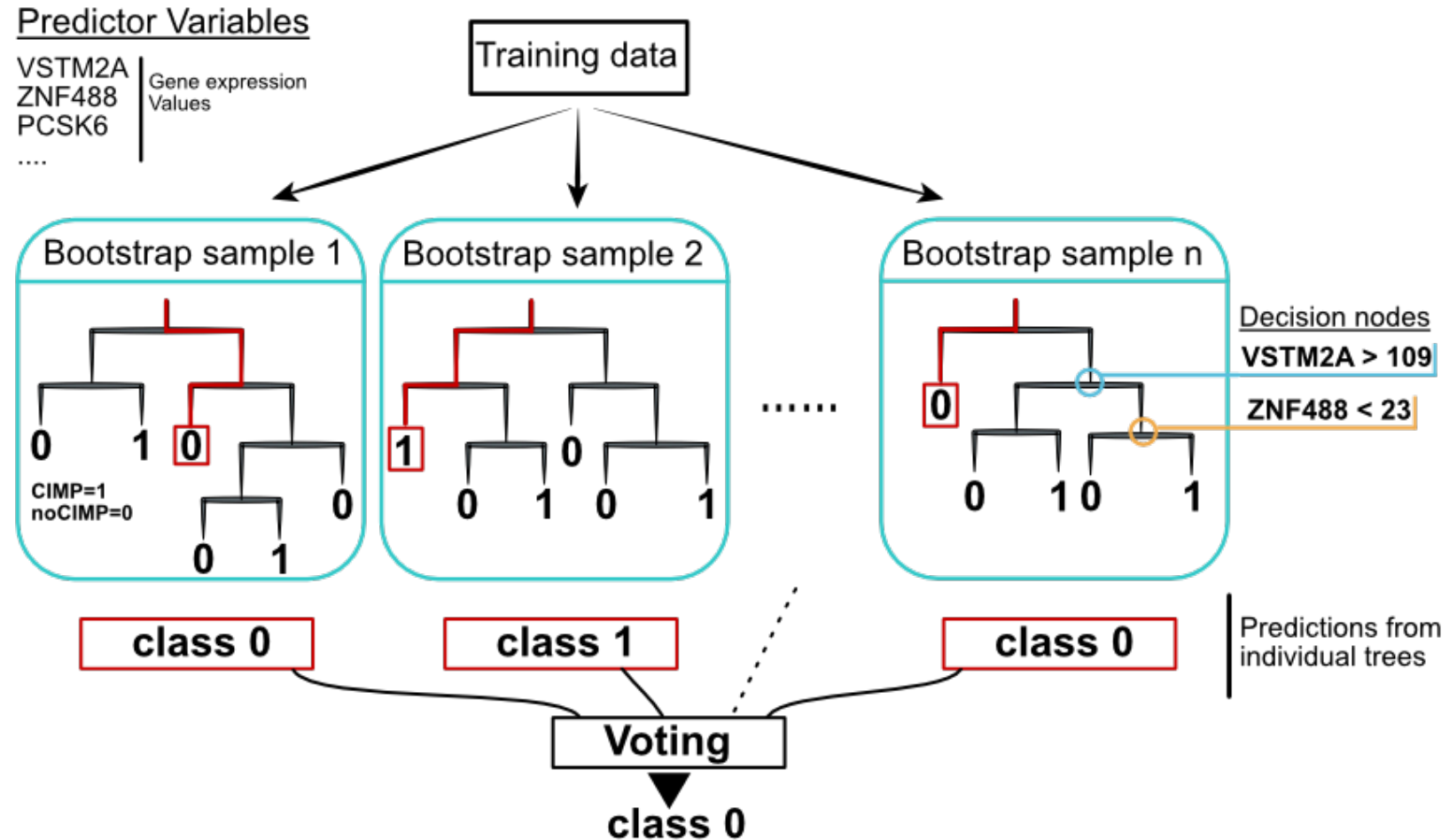$$I_G(p) = \sum_{i=1}^{K} p_i(1-p_i) = \sum_{i=1}^{K} p_i - \sum_{i=1}^{K} p_i^2 = 1 - \sum_{i=1}^{K} p_i^2$$

$$1 - (0.75^2 + 0.25)^2 = 0.375$$

# Trees and forests:
## Random Forests

# Variable importance for RF

Given a variable:

| Permute variable | → | OOB Test with permuted variable | → | Measure decrease in accuracy |
|---|---|---|---|---|

| Calculate decrease in Gini impurity | → | sum up the decrease across trees | → | Divide by total number of trees |
|---|---|---|---|---|

$$I = G_{parent} - G_{split1} - G_{split2}$$

# Variable importance
## Method agnostic

Given a predictor variable:

Drop-out loss

Permute predictor variable → Measure loss in performance

Permute response variable → Measure loss in performance

baseline loss

→ (Drop-out loss ) / (baseline loss)

DALEX package in R implements this strategy

# Regression using random forests
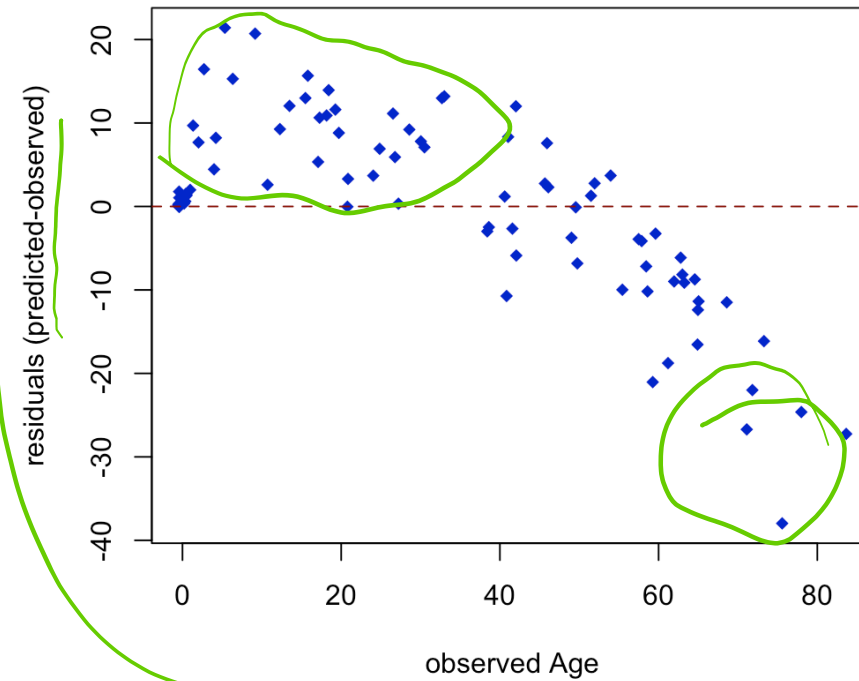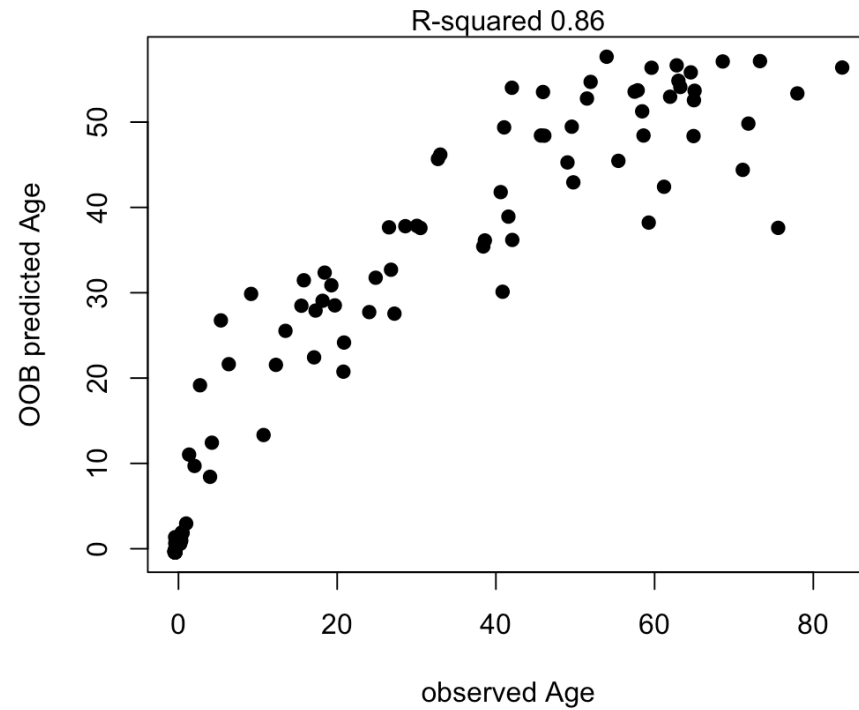
$$SSE = \sum (Y - f(x))^2$$

DNA methylation ▶ Age

↳ [0|1] CpG

108 subjects

~27000 CpGs

# Regression using random forests

# Supervised learning for Genomics:
## Recap

**Key concepts:**

1) Data prep → model performance

2) Overfitting → model complexity → resampling → test sets

3) Variable importance

4) Practical applications: classification & regression → ML