

# Machine learning

## Decision Trees & Ensemble methods

xgboost



---

Wouter Gevaert & Marie Dewitte

# Inhoud

→ Decision trees voor classificatie

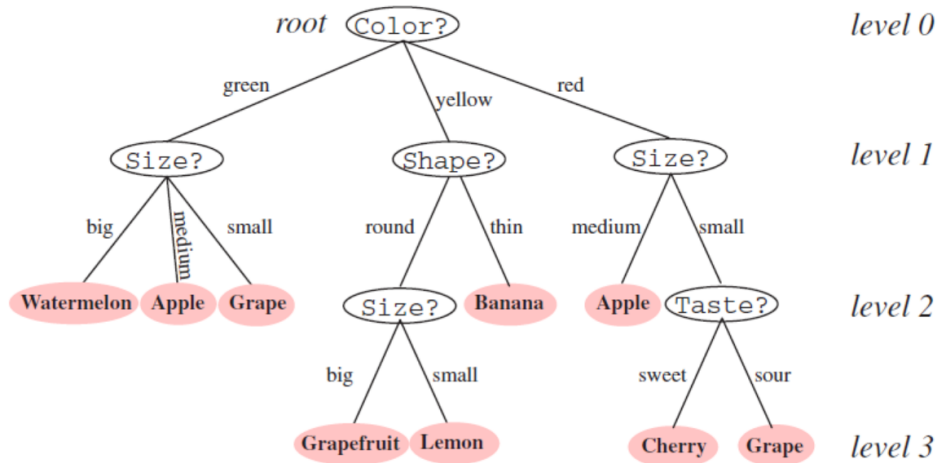
→ Decision trees voor regressie

Ensemble learning

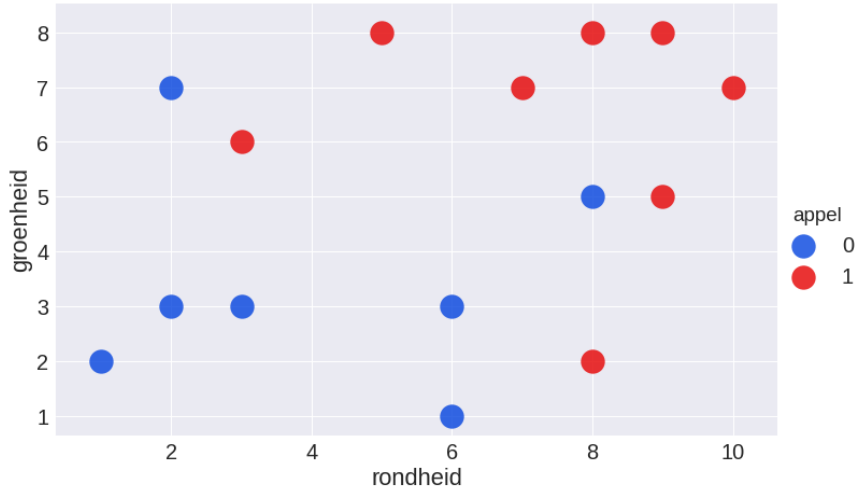
## Decision trees voor classificatie

---

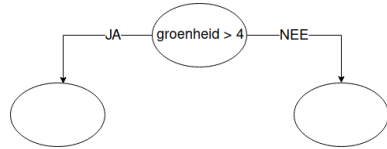
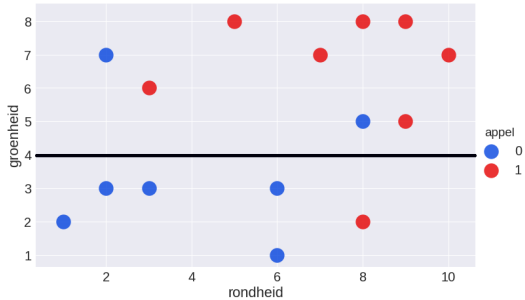
# Introductie aan de hand van een voorbeeld



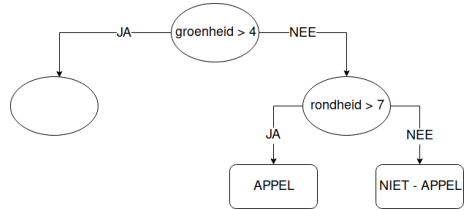
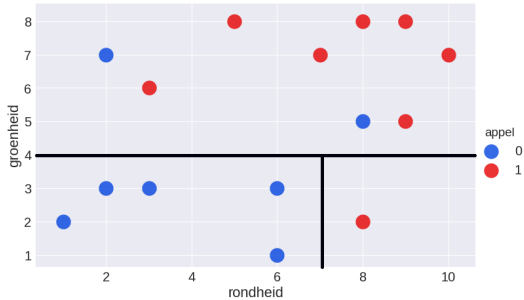
# Introductievoorbeeld - appels



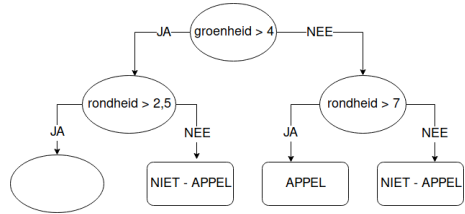
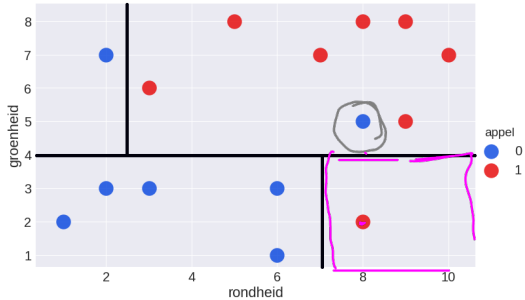
# Introductievoorbeeld - appels



# Introductievoorbeeld - appels

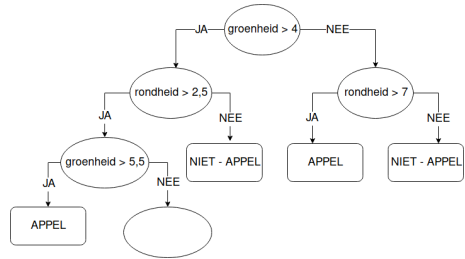
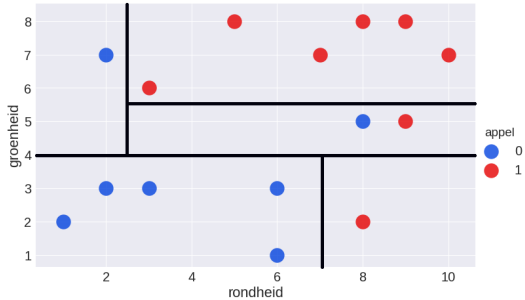


# Introductievoorbeeld - appels



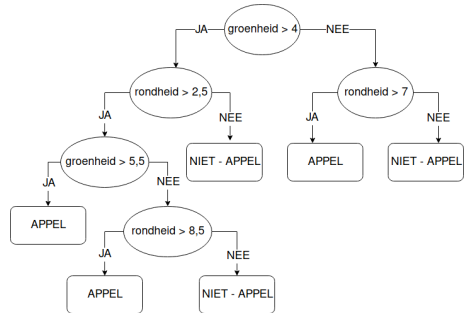
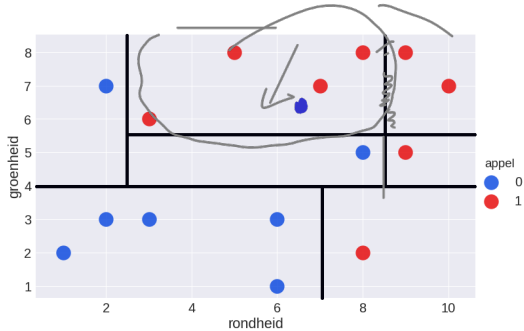


# Introductievoorbeeld - appels



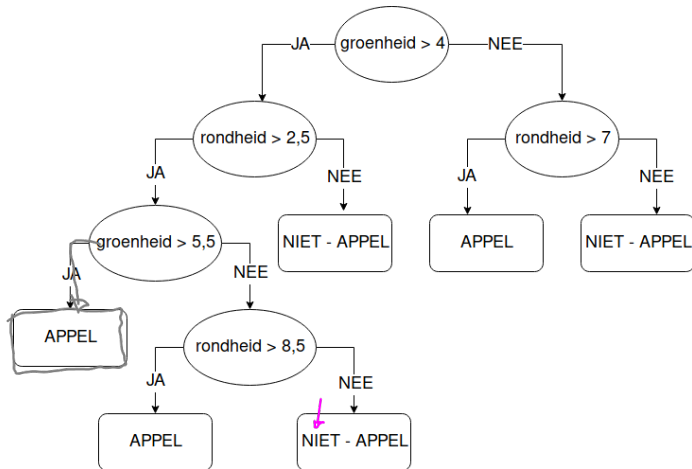
# Introductievoorbeeld - appels

$$\rightarrow \begin{cases} R = 6 \\ G = 7 \end{cases}$$



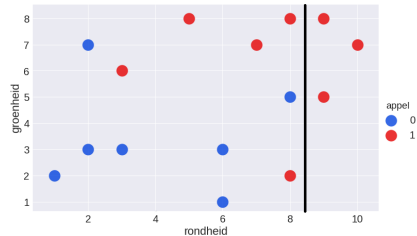
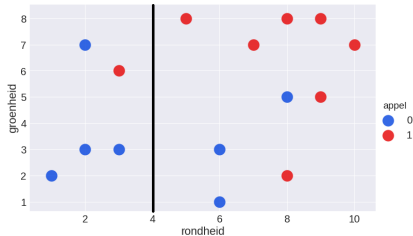
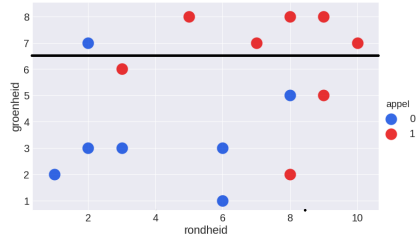
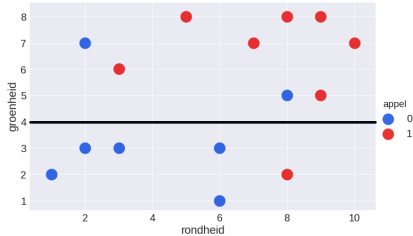
# Introductievoorbeeld - appels

Classificeer een stuk fruit met een rondheid van 6 en een groenheid van 7



# Splitsen van de boom

Hoe bepalen waar de boom te splitsen?



# Splitsen van de boom

Entropy = mate van  
wanorde

## Entropy en information gain

Een goede split is deze waarbij de wanorde afneemt en beide kanten 'zuiverder (pure) worden'

**Entropy H:** maat voor de wanorde (gemiddelde informatieinhoud)

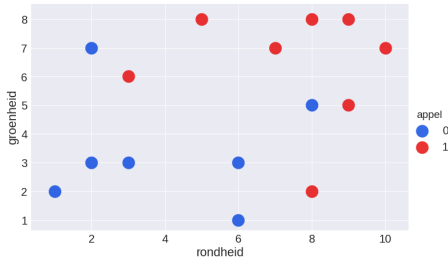
$$H = \sum_{i=1}^N p_i \log_2 \left( \frac{1}{p_i} \right)$$

Entropy  $\swarrow$

$$H = - \sum_{i=1}^N p_i \log_2 (p_i)$$

# Splitsen van de boom

## Entropy en information gain



Kans op een appel:  $\frac{8}{15} = 0,53$

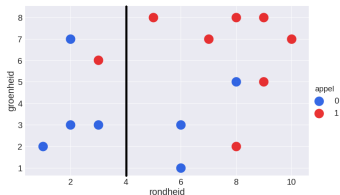
Kans op een niet-appel:  $\frac{7}{15} = 0,47$

$$H = -(0,53 \times \log_2(0,53) + 0,47 \times \log_2(0,47)) = -(-0,99740) = \boxed{0,9974}$$

# Splitsen van de boom

## Entropy en information gain

Information gain = entropy voor de split - entropy na de split



### Entropy links

$$p(\text{appel}) = \frac{1}{5} = 0,20$$

$$p(\text{niet} - \text{appel}) = \frac{4}{5} = 0,80$$

$$H = -(0,20 \log_2(0,20) + 0,80 \log_2(0,80)) = \boxed{0,72}$$

### Entropy rechts

$$p(\text{appel}) = \frac{7}{10} = 0,70$$

$$p(\text{niet} - \text{appel}) = \frac{3}{10} = 0,30$$

$$H = -(0,70 \log_2(0,70) + 0,30 \log_2(0,30)) = \boxed{0,88}$$

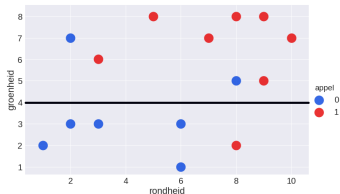
$$\text{Information gain} = 0,9974 - \left( \frac{5}{15} \times 0,72 + \frac{10}{15} \times 0,88 \right) = \boxed{0,17073}$$

*Handwritten notes: 0,182 and inf. gain*

# Splitsen van de boom

## Entropy en information gain

Information gain = entropy voor de split - entropy na de split



Entropy boven

$$p(appel) = \frac{7}{9} = 0,78$$

$$p(niet - appel) = \frac{2}{9} = 0,22$$

$$H = -(0,78 \log_2(0,78) + 0,22 \log_2(0,22)) = 0,76$$

Entropy onder

$$p(appel) = \frac{1}{6} = 0,17$$

$$p(niet - appel) = \frac{5}{6} = 0,83$$

$$H = -(0,17 \log_2(0,17) + 0,83 \log_2(0,83)) = 0,66$$

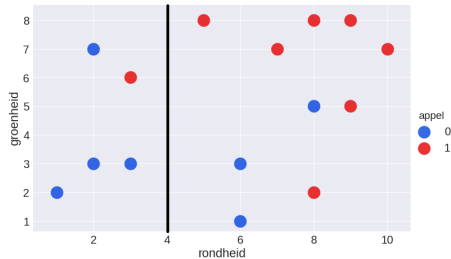
$$\text{Information gain} = 0,9974 - \left( \frac{9}{15} \times 0,76 + \frac{6}{15} \times 0,66 \right) = \boxed{0,27740} \approx 0,28$$



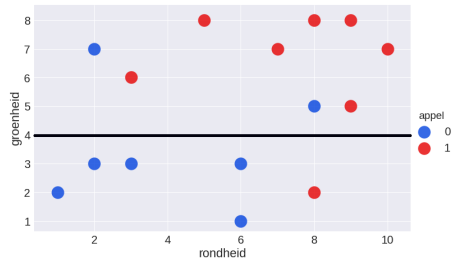
# Splitsen van de boom

## Entropy en information gain

Splits telkens bij de split die je de **hoogste information gain** oplevert.



Information gain = 0,17073



Information gain = 0,27740

# Splitsen van de boom

## → Gini impurity

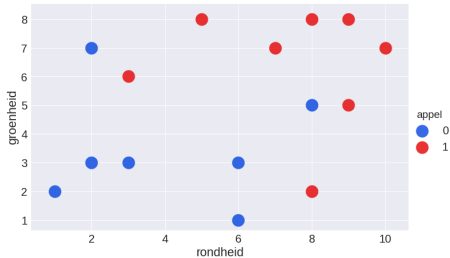
Alternatief voor het gebruik van de entropy.

De Gini impurity is een maat voor hoe vaak een willekeurig gekozen element van de set verkeerd gelabeld zou worden als het willekeurig werd gelabeld volgens de distributie van de labels in de subset.

$$G = 1 - \sum_{i=1}^N p_i^2$$

# Splitsen van de boom

## Gini impurity



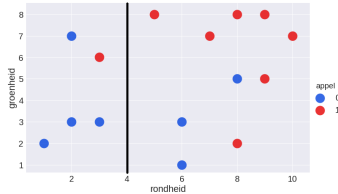
Kans op een appel:  $\frac{8}{15} = 0,53$

Kans op een niet-appel:  $\frac{7}{15} = 0,47$

$$G = 1 - (0,53^2 + 0,47^2) = 0,4982$$

# Splitsen van de boom

## Gini impurity



### Gini impurity links

$$p(\text{appel}) = \frac{1}{5} = 0,20$$

$$p(\text{niet} - \text{appel}) = \frac{4}{5} = 0,80$$

$$G = 1 - (0,20^2 + 0,80^2) = 0,32$$

### Gini impurity rechts

$$p(\text{appel}) = \frac{3}{10} = 0,30$$

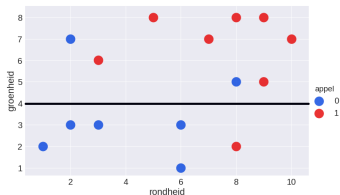
$$p(\text{niet} - \text{appel}) = \frac{7}{10} = 0,70$$

$$G = 1 - (0,30^2 + 0,70^2) = 0,42$$

$$\text{Gewogen Gini impurity} = \frac{5}{15} \times 0,32 + \frac{10}{15} \times 0,42 = 0,39$$

# Splitsen van de boom

→ Gini impurity → *defout*.



Gini impurity boven

$$p(\text{appel}) = \frac{7}{9} = 0,78$$

$$p(\text{niet} - \text{appel}) = \frac{2}{9} = 0,22$$

$$G = 1 - (0,78^2 + 0,22^2) = 0,34$$

Gini impurity onder

$$p(\text{appel}) = \frac{1}{6} = 0,17$$

$$p(\text{niet} - \text{appel}) = \frac{5}{6} = 0,83$$

$$G = 1 - (0,17^2 + 0,83^2) = 0,28$$

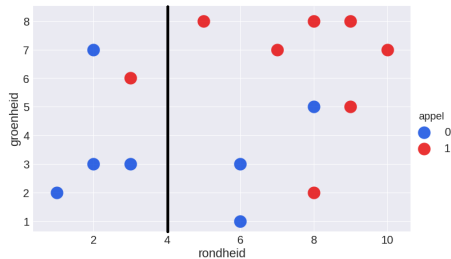
| 0,32

$$\text{Gewogen Gini impurity} = \frac{9}{15} \times 0,34 + \frac{6}{15} \times 0,28 = \boxed{0,32}$$

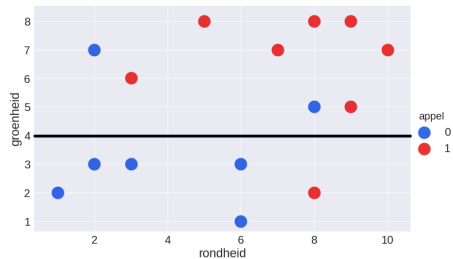
# Splitsen van de boom

## Gini impurity

Splits telkens bij de split die je de **laagste Gini impurity** oplevert.



Gewogen Gini impurity = 0,39



Gewogen Gini impurity = 0,32

→ dataset → niet met. lineaire  
→ DT

→ lineair dataset → Lineaire  
regressie

## Decision trees voor regressie

---

- niet-gescaleerde data
- performantie
- Transparant

$$\frac{25 + 30 + 35 + 38 + 48}{5}$$

= ....

Day	Outlook	Temp.	Humidity	Wind	Golf Players
1	Sunny	Hot	High	Weak	25
2	Sunny	Hot	High	Strong	30
3	Overcast	Hot	High	Weak	46
4	Rain	Mild	High	Weak	45
5	Rain	Cool	Normal	Weak	52
6	Rain	Cool	Normal	Strong	23
7	Overcast	Cool	Normal	Strong	43
8	Sunny	Mild	High	Weak	35
9	Sunny	Cool	Normal	Weak	38
10	Rain	Mild	Normal	Weak	46
11	Sunny	Mild	Normal	Strong	48
12	Overcast	Mild	High	Strong	52
13	Overcast	Hot	Normal	Weak	44
14	Rain	Mild	High	Strong	30

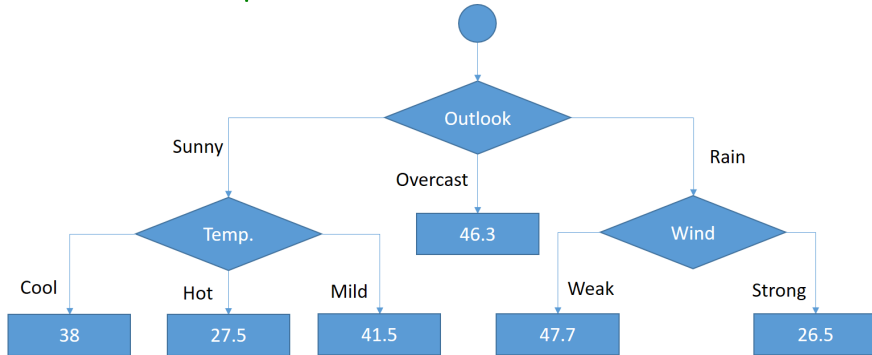
$$VAR = \frac{\sum (x_i - \bar{x})^2}{5}$$

→ VAR

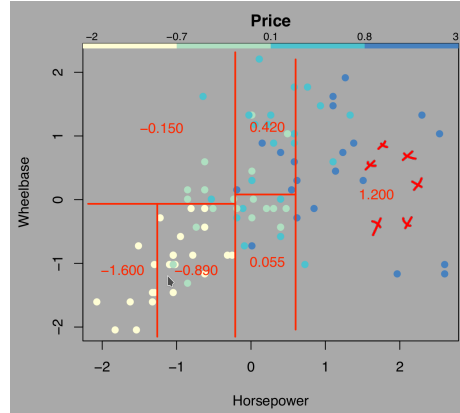
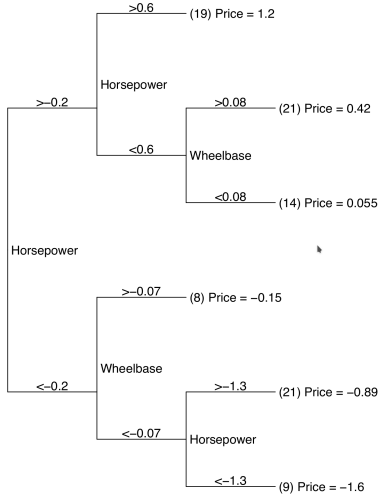
→ VAR



# Voorbeeld



# Voorbeeld



- Bagging
- Boosting
- Stacking

## Ensemble learning

---

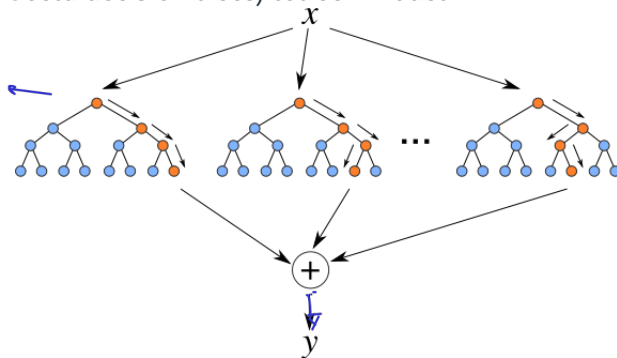
# Ensemble learning

## Problematiek van decision trees

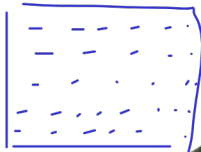
Decision trees hebben de neiging tot overfitting.

Oplossing: **combineer** de voorspellingen van verschillende gerandomiseerde modellen (bijvoorbeeld decision trees) tot één model

Weak  
Combination



# Ensemble learning



$$\rightarrow \boxed{JA = 60\%}$$
$$NFE = 40\%$$

Given a jury of voters and assuming independent errors. If the probability of each single person in the jury of being correct is above 50% then the probability of the jury being correct tends to 100% as the number of persons increase.



Nicolas de Condorcet (1743 - 1794)

# Ensemble learning

## Overzicht Ensemble learning methodes

- • Bootstrap aggregating (bagging)
- • Boosting
- ~~Bayes optimal classifier~~
  - ~~Bayesian model averaging/combination~~
  - ~~Bucket of models~~
- • Stacking

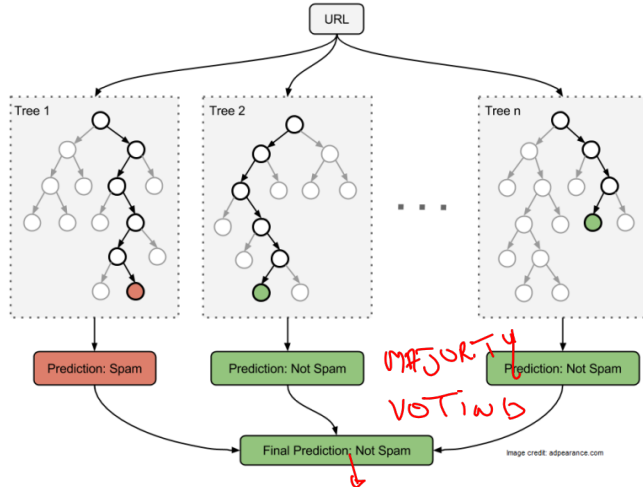
$$p(y|x) =$$

- - -

# 2) Bagging

Concept van bagging - Majority voting

→ Random Forest Trees  
(Bagging)



# Bagging

## Concept van bagging - Majority voting

	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Sample 7	Sample 8	Sample 9	Sample 10	
→ Tree 1	Correct	Correct	Correct	Wrong	Correct	Correct	Wrong	Wrong	Correct	Correct	7/10 correct
→ Tree 2	Correct	Correct	Correct	Correct	Wrong	Wrong	Correct	Wrong	Correct	Correct	7/10 correct
→ Tree 3	Correct	Correct	Wrong	Wrong	Correct	Correct	Correct	Wrong	Correct	Wrong	6/10 correct
Majority voting	Correct	Correct	Correct	Wrong	Correct	Correct	Correct	Wrong	Correct	Correct	8/10 correct

Correct  
Wrong

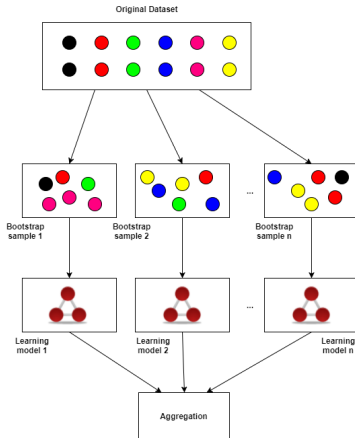
80%



# Bagging

## De techniek van bagging

Methode om de variantie te reduceren en overfitting te vermijden.



Kies verschillende subsets (bags) uit de trainingset.

Willekeurige selectie met teruglegging

Train op elke subset een classifier

Max. features = 4  
to features

RAIN

1000

TEST

1996

63

200.

003

## Variants

## BAGGING

is

60E D

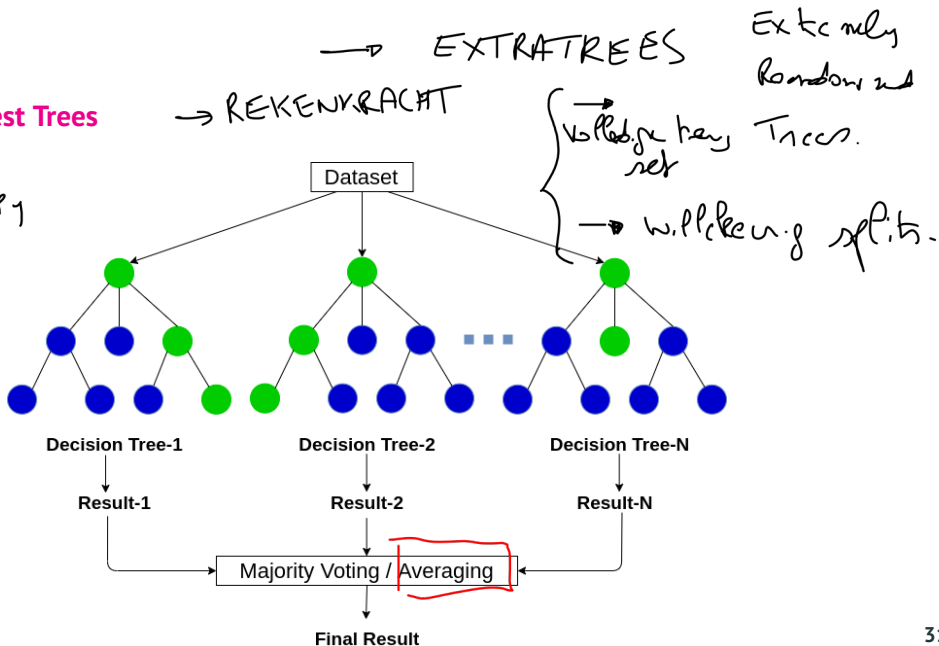
Bij

# WITSCHMIETER

# Bagging

## Random Forest Trees

- Gini
- ENTROPY



# Bagging

## Random Forest Trees - hyperparameters

**n\_estimators:** number of trees in the forest. Meestal hoe meer hoe beter.

**Criterion:** Gini of Entropy (default Gini)

**Maximum number of features:** het maximum aantal features per boom.

- int : aantal features
- float: percentage
- 'auto': max\_features = vierkantswortel van totaal aantal features
- 'sqrt': max\_features = vierkantswortel van totaal aantal features
- 'log2': log van het totaal aantal features
- Default worden alle features gebruikt.

# Random forest trees

## Random Forest Trees - hyperparameters

**max\_depth:** de maximale diepte van de boom. Als je te maken hebt met noisy data is het aan te raden de maximale diepte beperkt te houden.

**min\_samples\_split:** het minimum aantal samples nodig om binnen een boom te blijven splitsen.

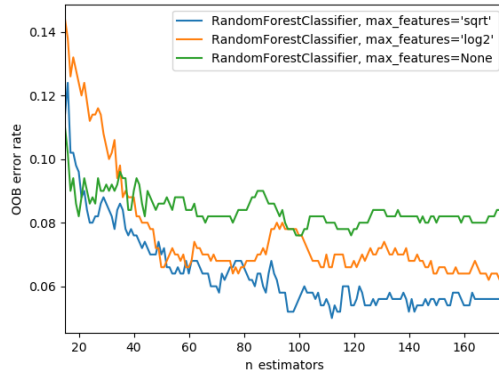
**min\_samples\_leaf:** het minimum aantal samples dat zich aan een blad van de boom moet bevinden. Hoe groter deze waarde, hoe minder vatbaar voor overfitting.

**Bootstrap aggregating:** Bagging (zie verder). Staat default op True.

# Bagging

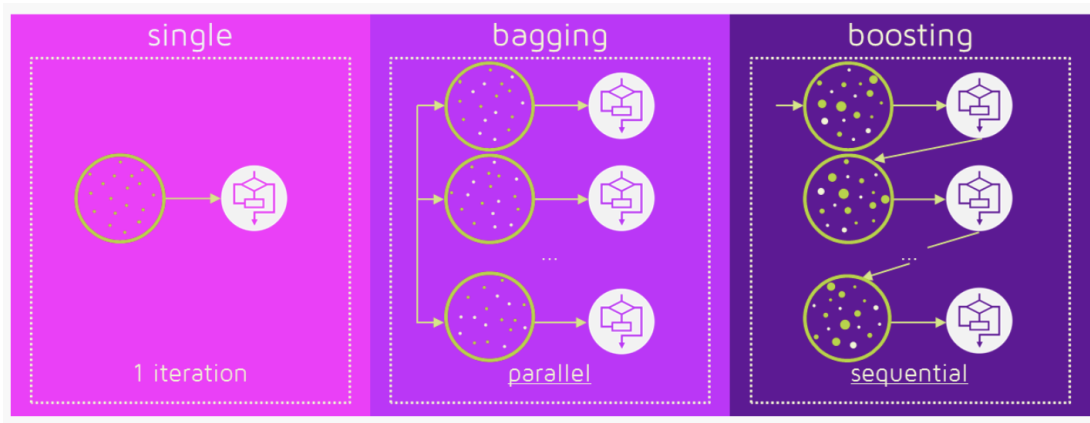
## Random Forest Trees - hyperparameters

**oob\_score**: de gemiddelde error bij het testen van een sample op bomen die niet op deze sample getraind zijn geweest.



# Boosting

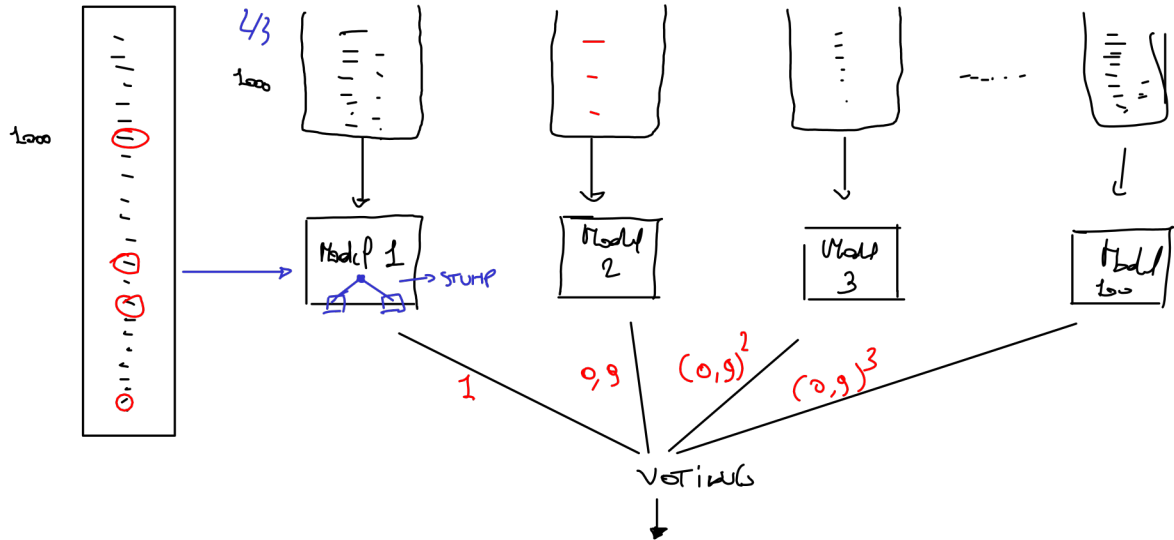
## Principle



# Adaboosting

Learn. rate. = 0,9

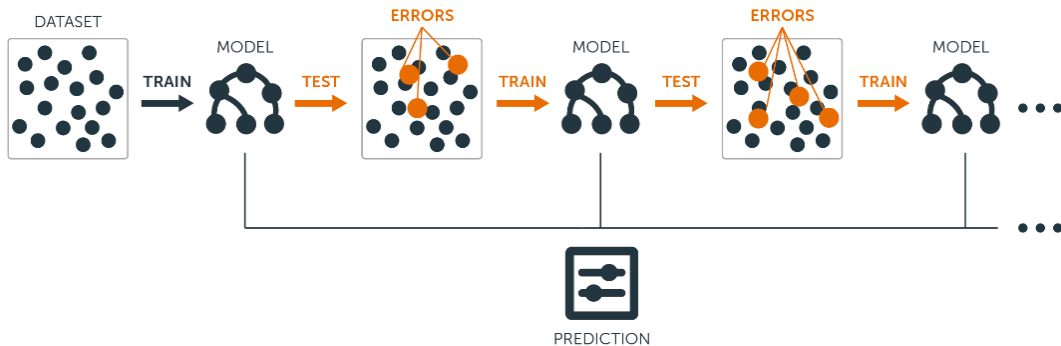
Accuracy





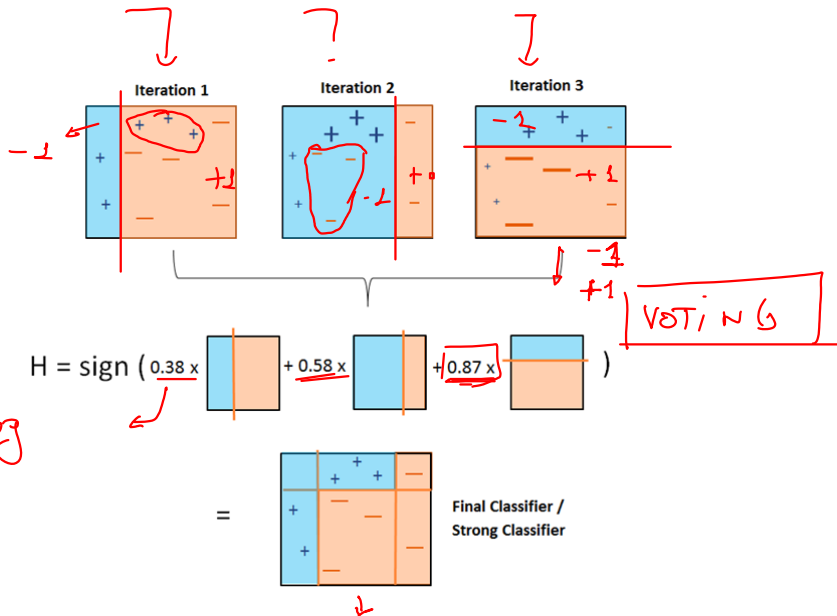
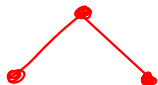
# Boosting

## Adaboost



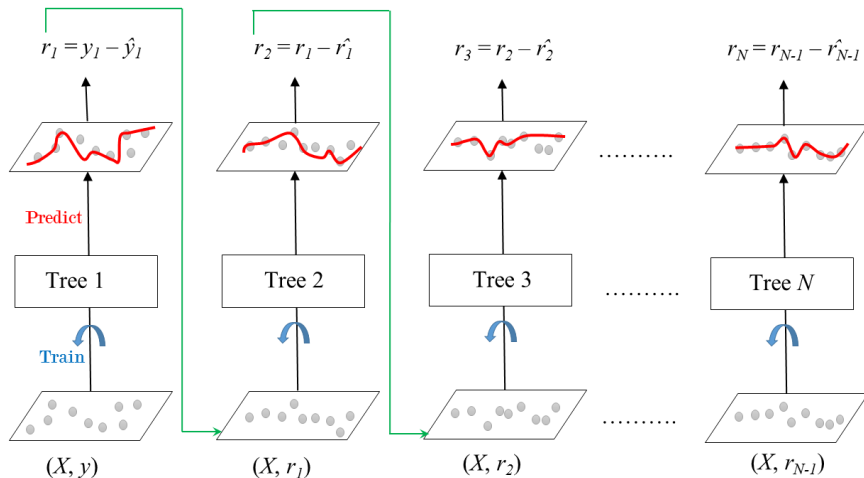
# Boosting

## Adaboost

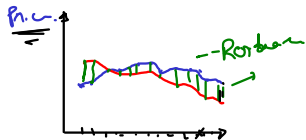


# Boosting

## → Gradient Boosting

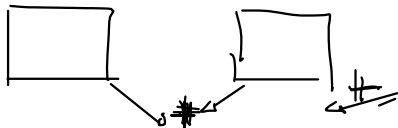


Regression



Model  
1

Model  
2



Model  
3

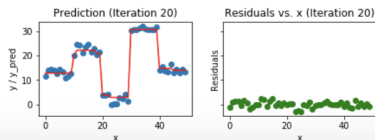
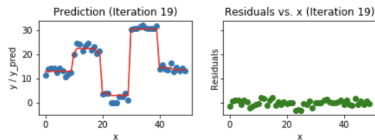
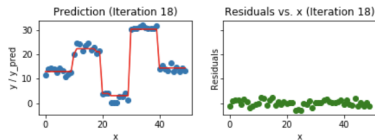
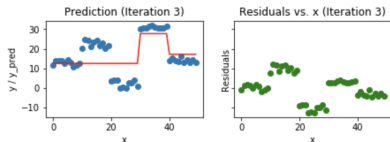
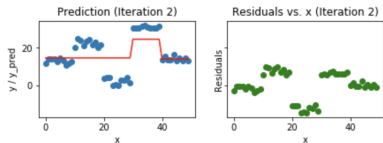
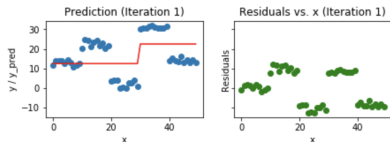
Model  
4

→ CLASSIFICATION

→ probability  
→ probability.

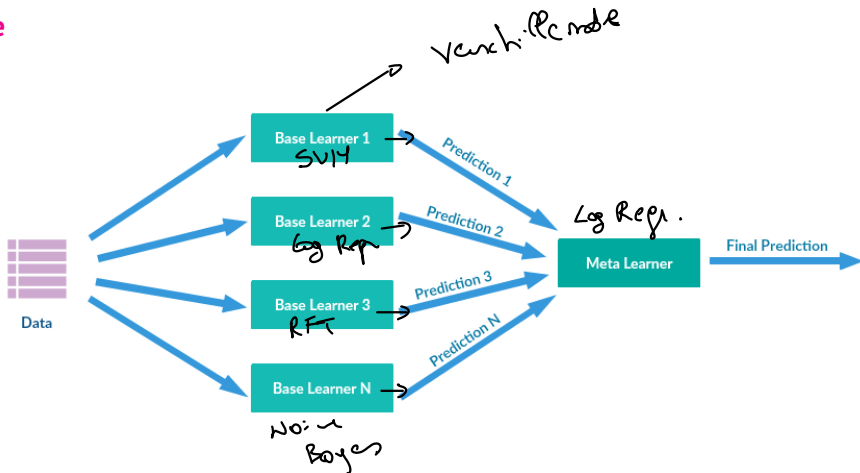
# Boosting

## Gradient Boosting



# Stacking

## Principle

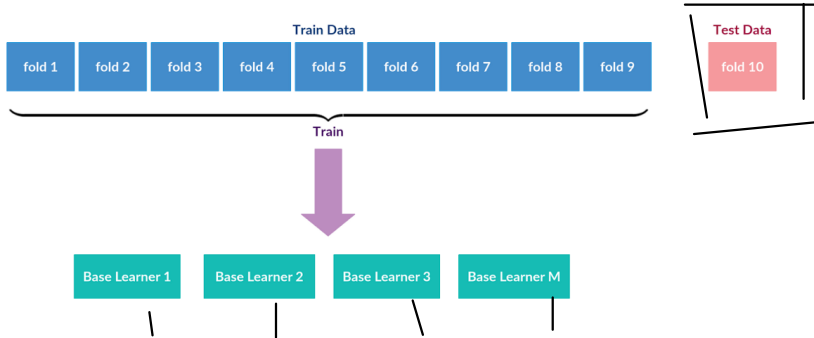


# Stacking

## Principe

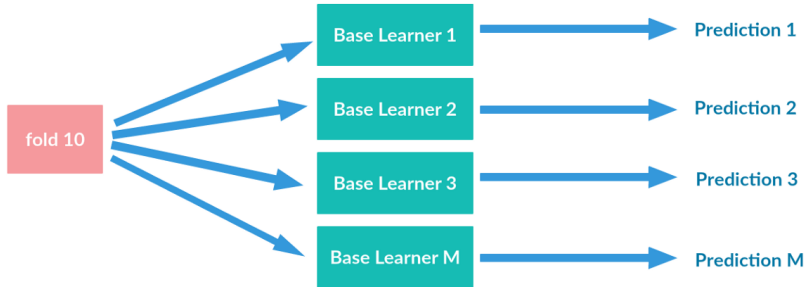
$$K = 10$$

fold size =  $N/10$



# Stacking

## Principe





# Stacking

## Principe

Data Point #	prediction from base learner 1	prediction from base learner 2	prediction from base learner 3	prediction from base learner M	actual
1	$y_{11}^{\hat{}}$	$y_{12}^{\hat{}}$	$y_{13}^{\hat{}}$	$y_{1M}^{\hat{}}$	$y_1$
2	$y_{21}^{\hat{}}$	$y_{22}^{\hat{}}$	$y_{23}^{\hat{}}$	$y_{2M}^{\hat{}}$	$y_2$
...	...	...	...	...	...
N	$y_{N1}^{\hat{}}$	$y_{N2}^{\hat{}}$	$y_{N3}^{\hat{}}$	$y_{NM}^{\hat{}}$	$y_N$

# Ensemble learning

## Samenvattend overzicht

### Bagging

- Reduceert de variantie
- Robuust tegen uitschieters en noisy data
- Dikwijls via Random Forest Trees

### Boosting

- Verhoogt de accuraatheid
- Niet robuust tegen uitschieters of noisy data

### Stacking

- Ensemble van strong learners
- Metalearner die de optimale combinatie leert van de base learners