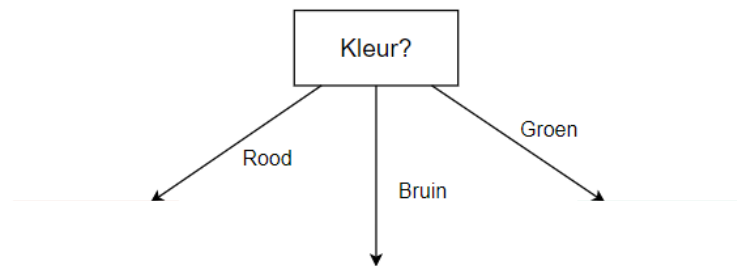


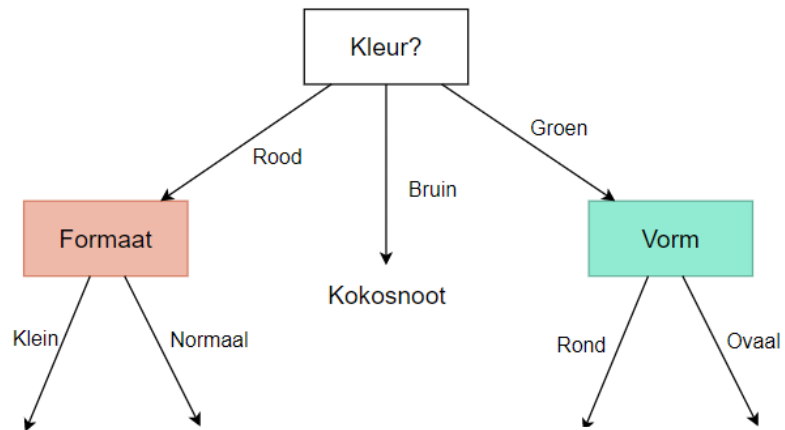
The background features a series of concentric, light gray circles centered on the slide. Overlaid on these are stylized, light blue circuit-like lines with small circles at the nodes, appearing in the corners and along the edges.

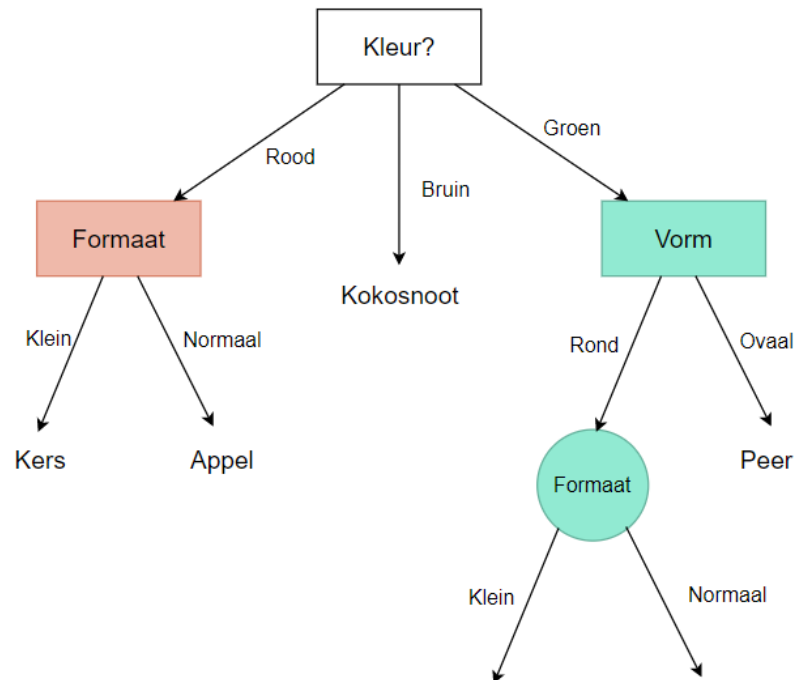
# DECISION TREES

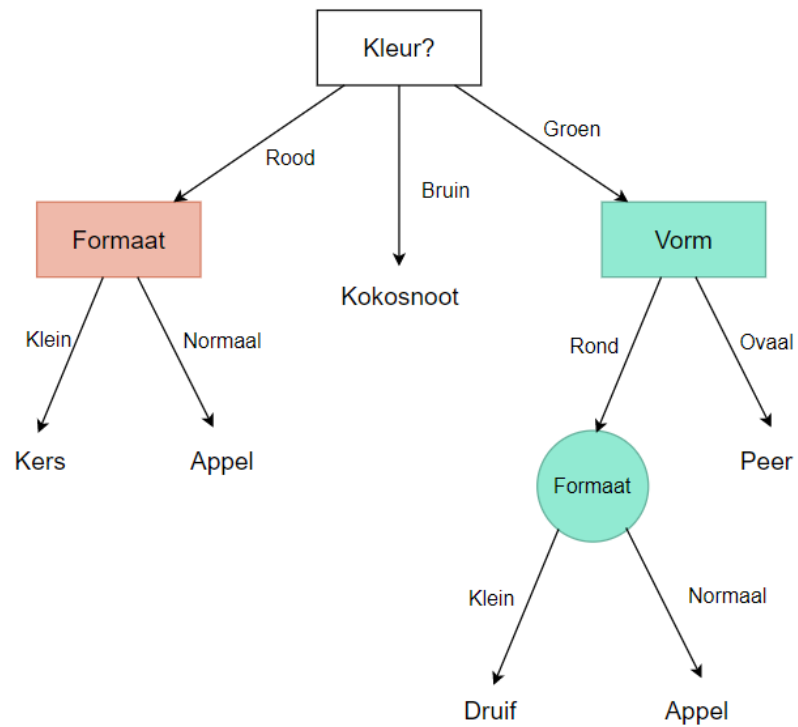
JENS BAETENS

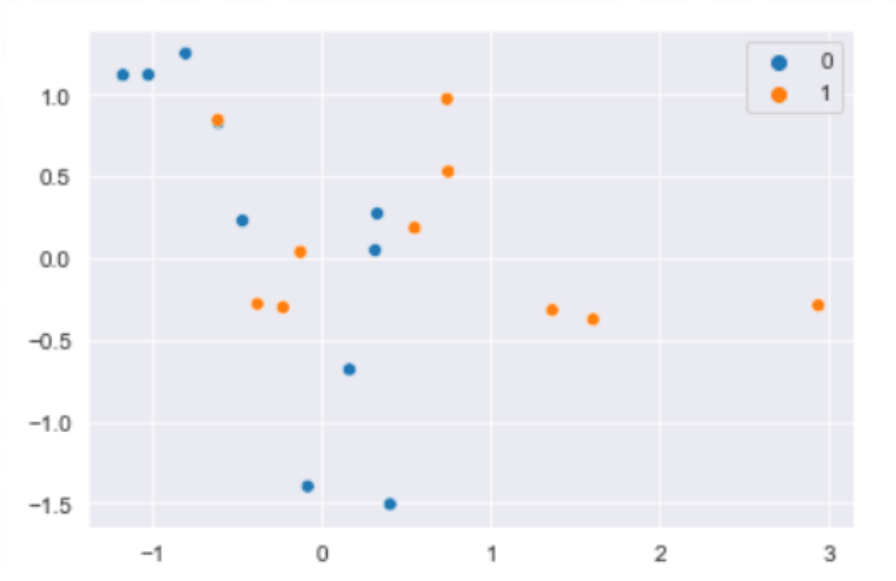


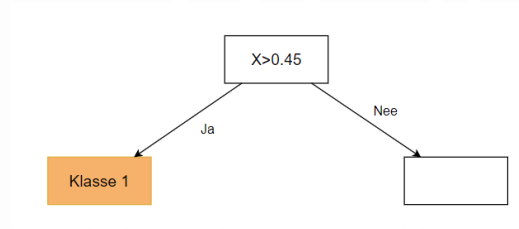
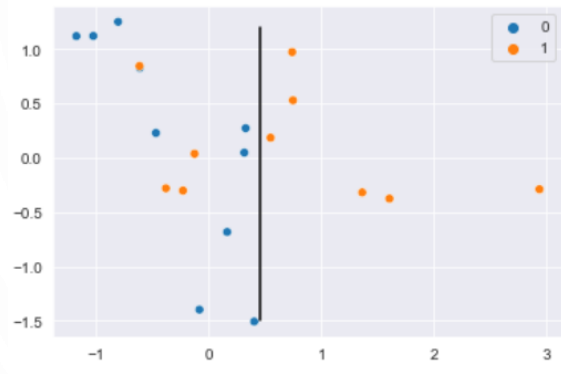




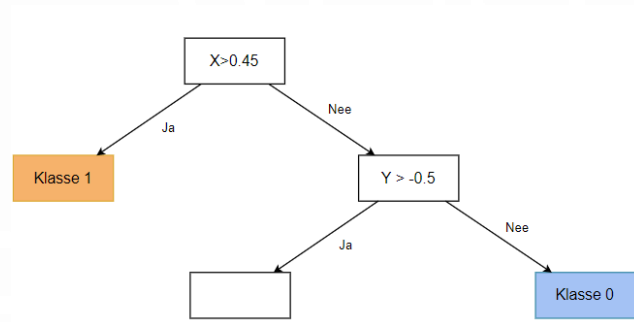
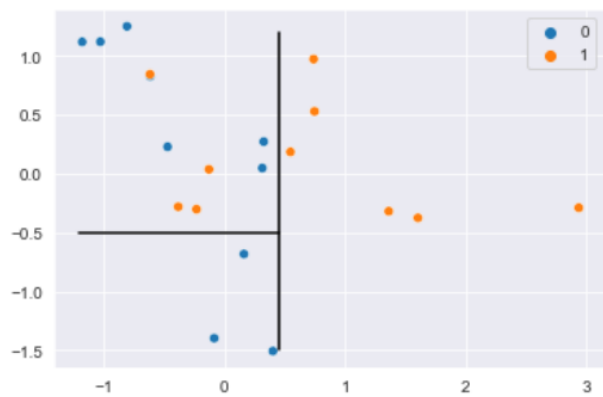


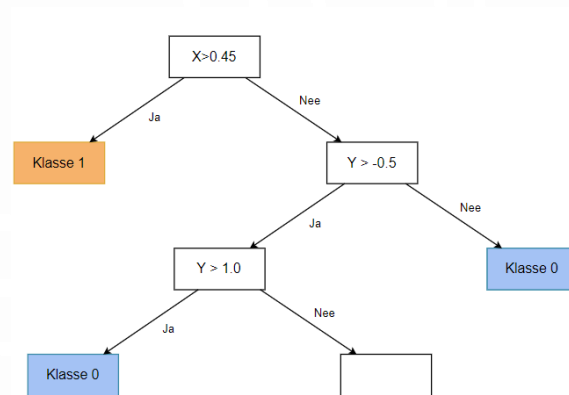


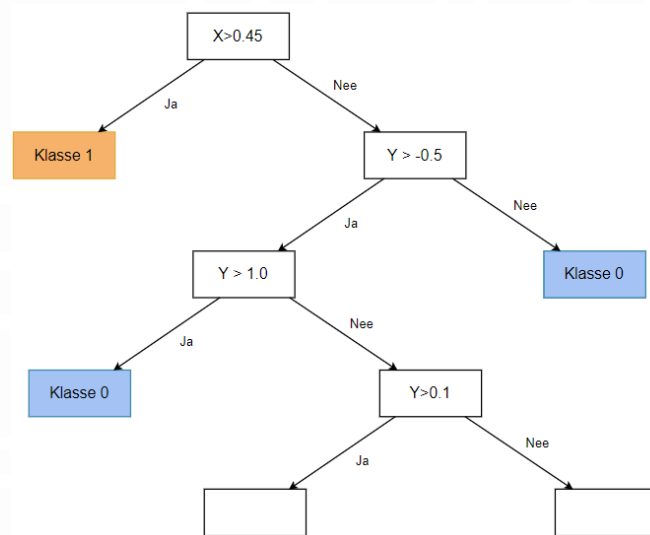
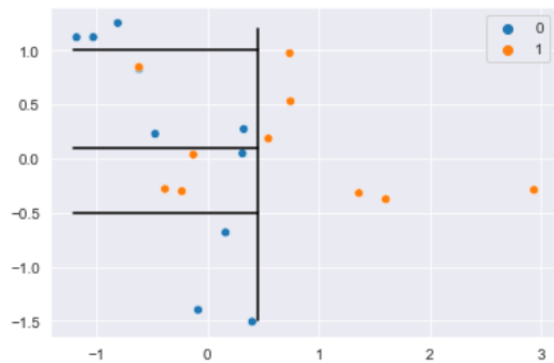


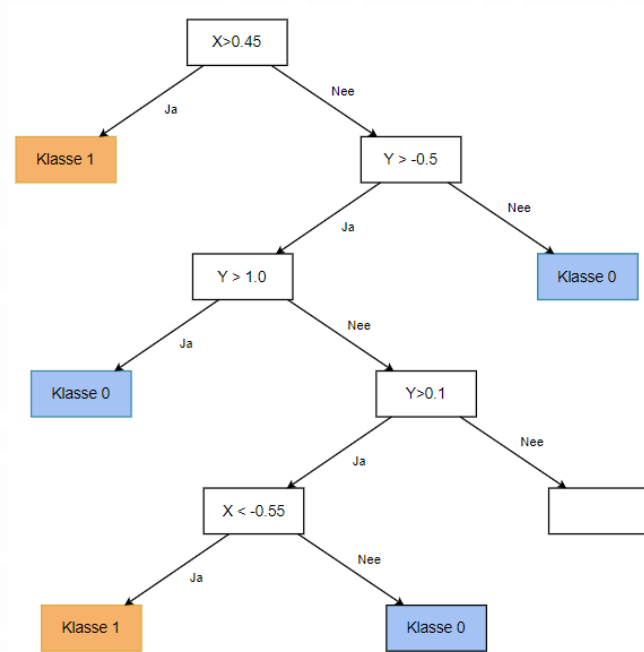
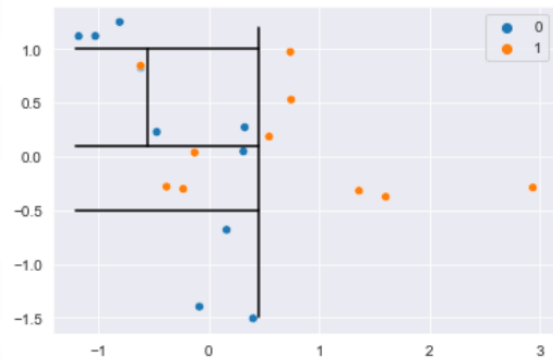


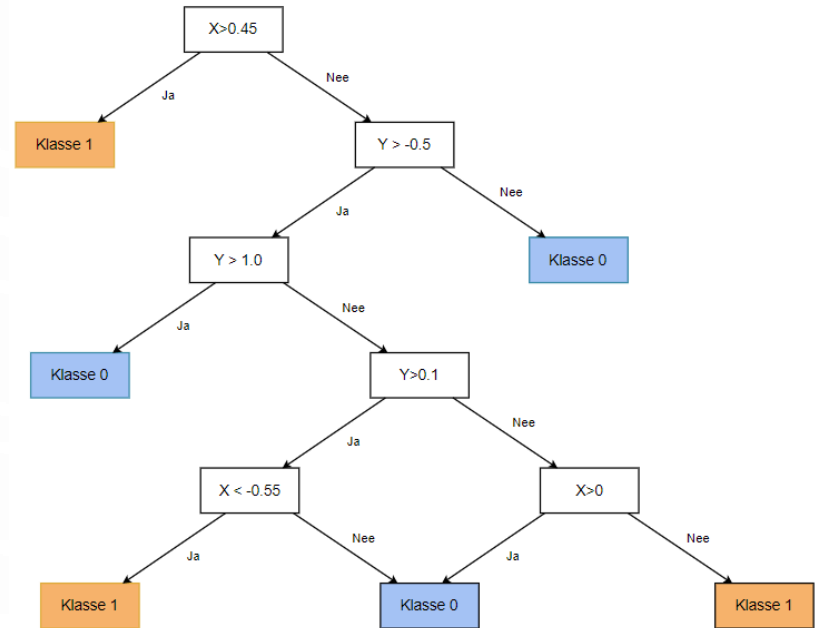
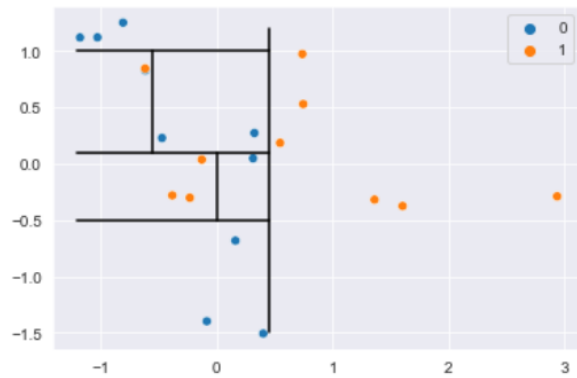












# HOE BEPALEN VAN DE SCHEIDINGSLIJN?

Bereken entropie (maatstaf voor de wanorde) voor elke mogelijke verdeling

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

$P_i$  is het percentage van elke klasse in een gebied

Bereken het verschil met de entropie voor en na de verdeling

=> Information gained => Hoe groter hoe beter

# ALTERNATIEF VOOR ENTROPIE

Logaritmes voor berekenen entropie zijn rekenintensief

Alternatieve manier: Gini impurity

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

Beste scheidingslijn heeft de kleinste Gini Impurity

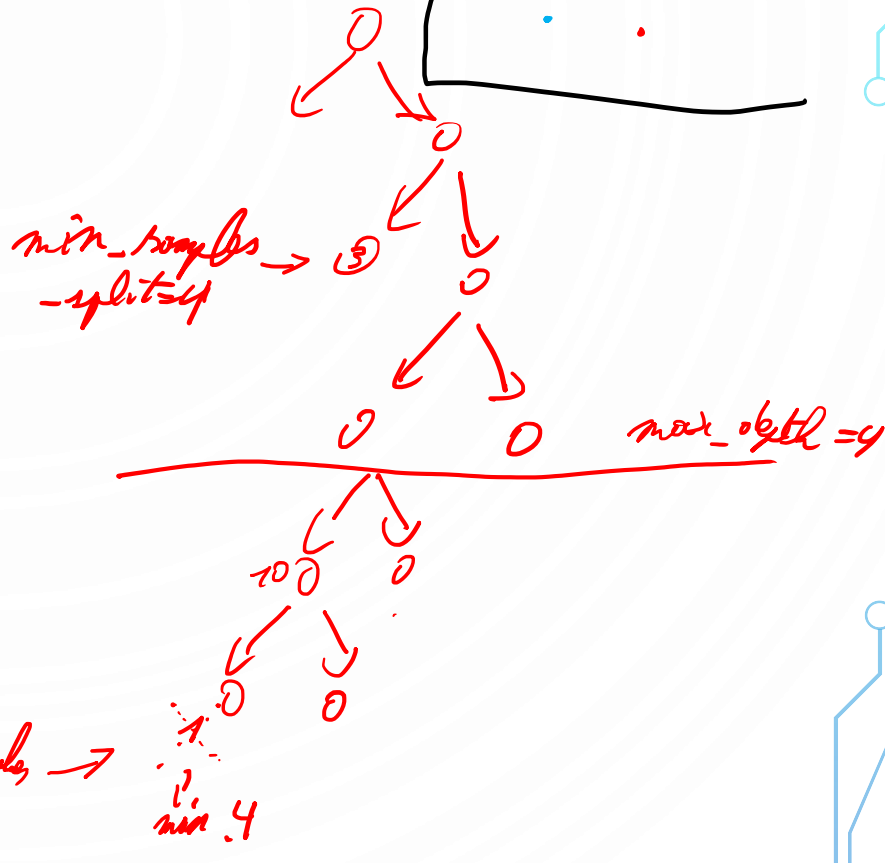
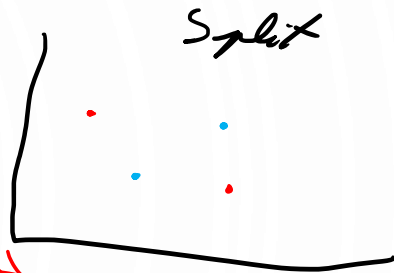
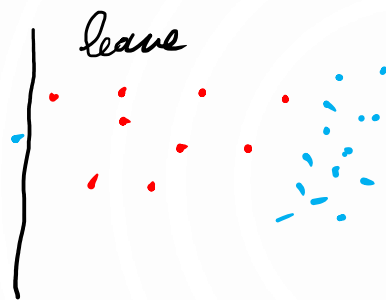
# NADELEN

Gevoelig aan ruis

Sterke neiging tot overfitting

Hyperparameters voor regularisatie:

- max\_depth
- min\_samples\_split
- min\_samples\_leave





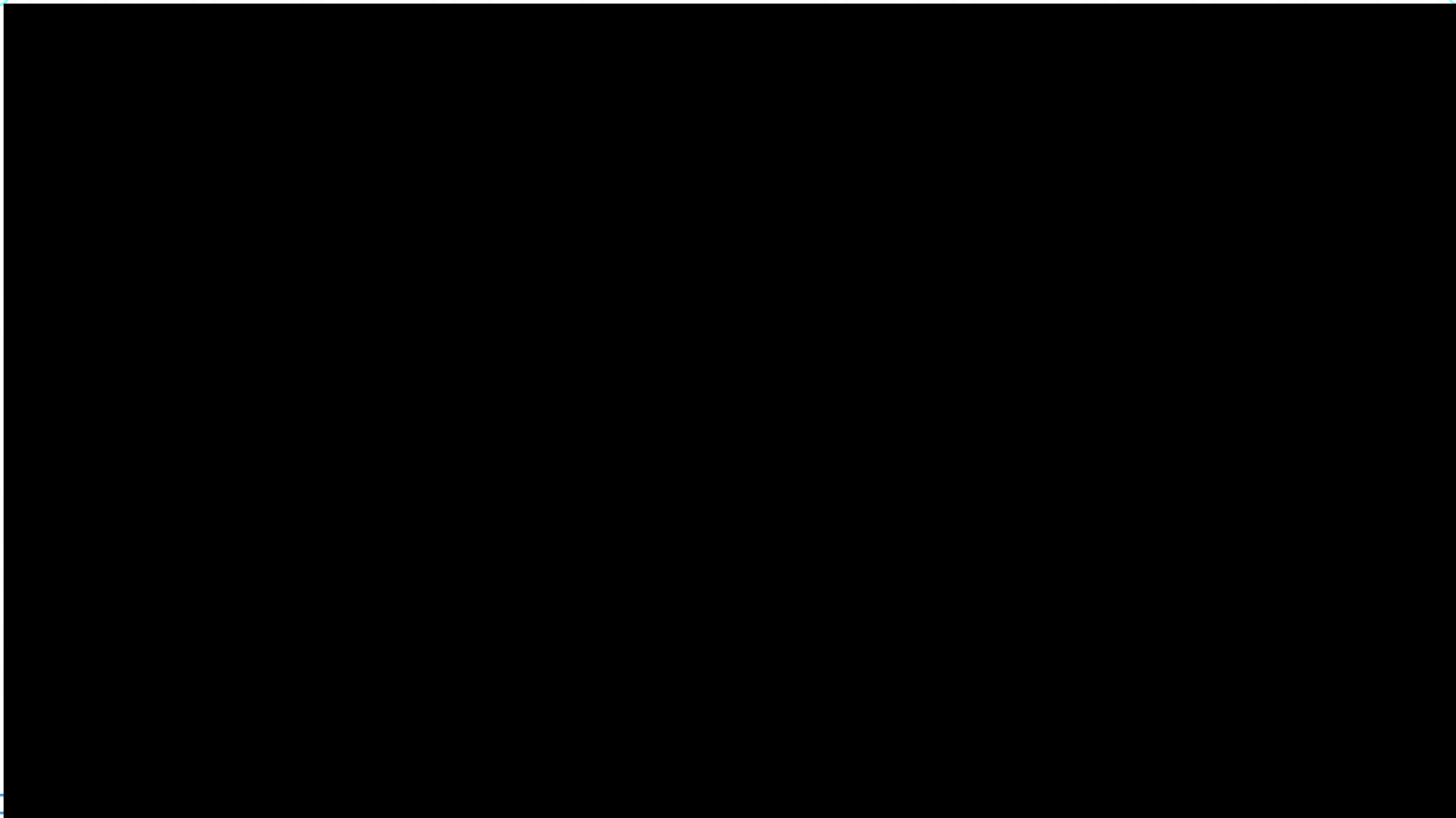
The slide features a light blue background with a subtle pattern of concentric circles. In each of the four corners, there are decorative circuit-like lines in a slightly darker blue, consisting of straight segments and small circles, resembling a stylized electronic board.

# VOORDELEN

Eenvoudig algoritme en daardoor ook snel

Men kan redeneren op het resultaat

- Niet mogelijk bij SVM, LR, Naïve Bayes, ...



<https://www.youtube.com/watch?v=Jcl5E2Ng6r4>

# **RANDOM FORESTS**

# RANDOM FORESTS

Is Rusland groter dan Afrika?

Ja: ...

Nee: ...







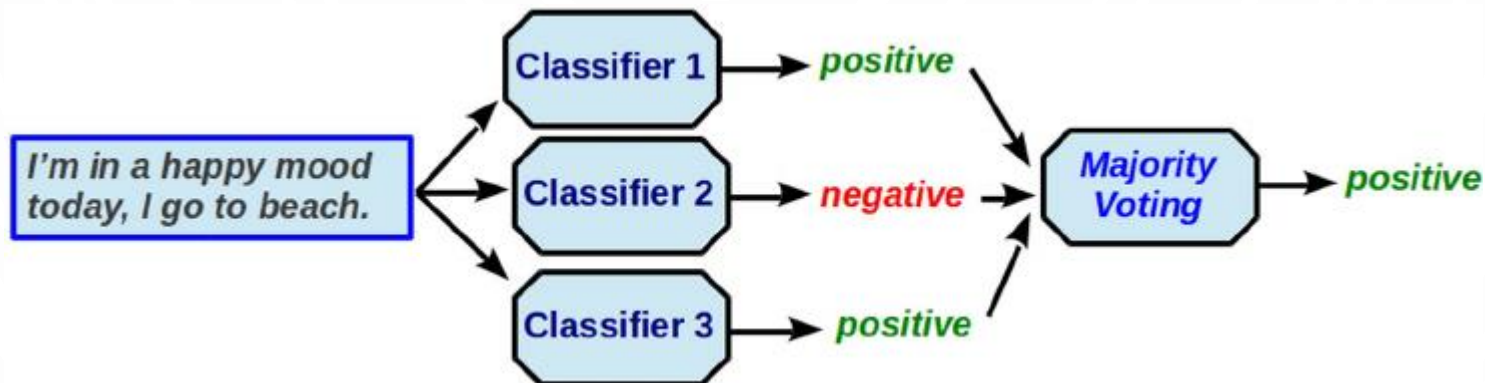
# RANDOM FORESTS

## Condorcet's jury theorem

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)

# RANDOM FORESTS



# RANDOM FORESTS - PARAMETERS

Aantal bomen

Aantal features per boom: int, float, sqrt / auto, log2, default

→ splitst op kolom

Bagging: Gebruik slechts een deel van de data om elke boom te trainen

→ splitst op rijen  
(bootstraping)

Oob\_score: Aangezien niet alle data gebruikt wordt om elke boom te trainen. Gebruik deze data om de bomen te valideren. Vooral handig als er niet genoeg data is om een aparte validatieset te behouden.