

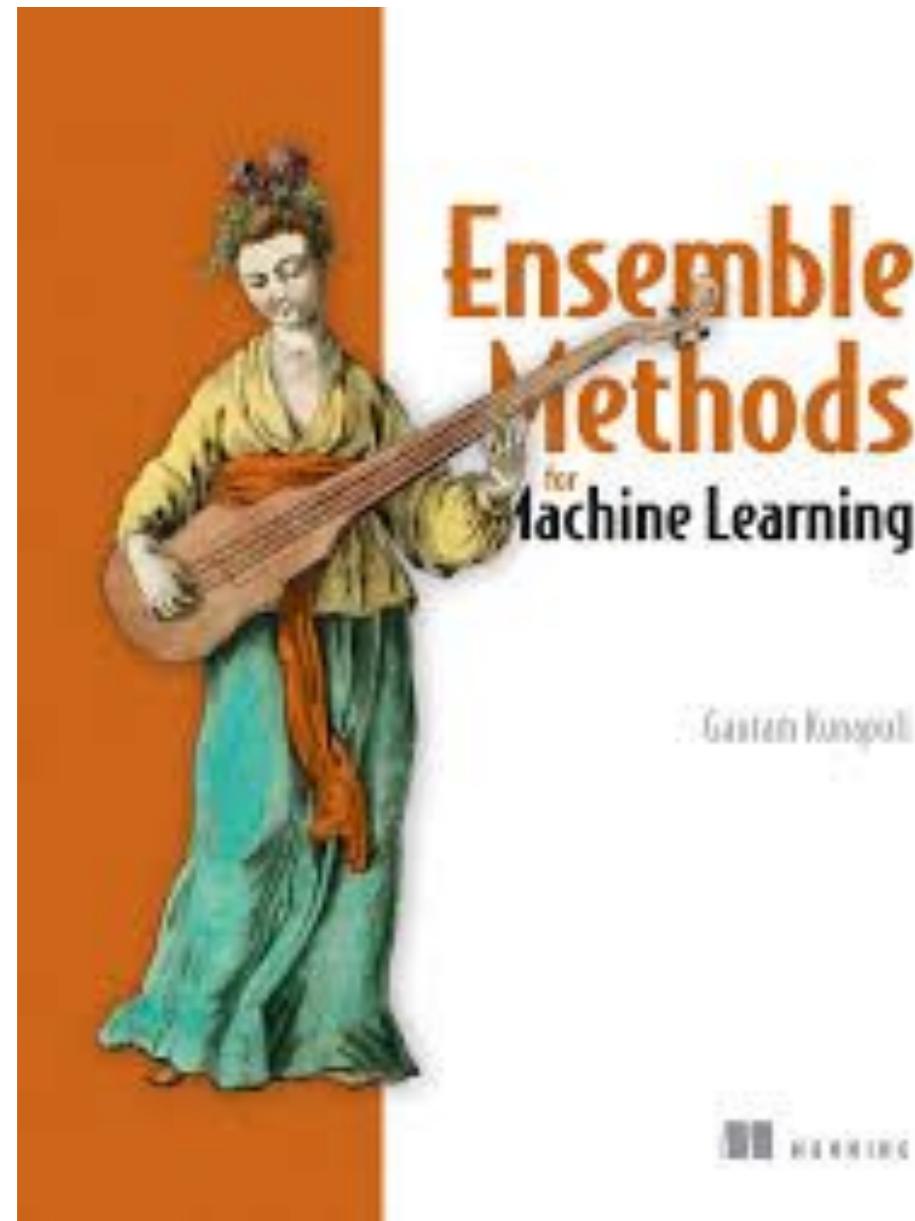


IF3270 Pembelajaran Mesin Ensemble Methods

Tim Pengajar IF3270

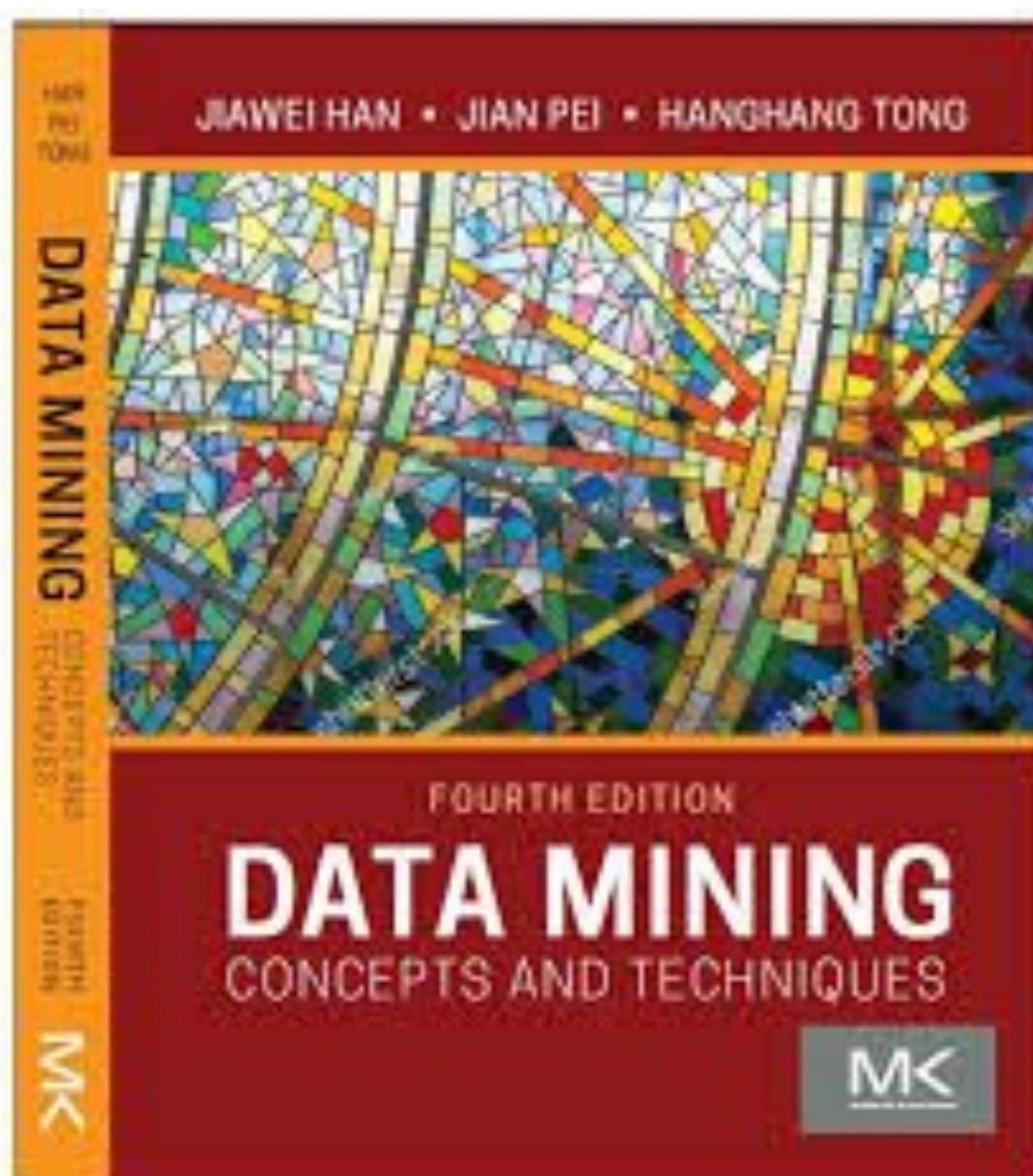


Reference



Kunapuli, G. (2023). *Ensemble methods for machine learning*.
Simon and Schuster.

<https://livebook.manning.com/book/ensemble-methods-for-machine-learning/chapter-1/19>



Han, J., Pei, J., & Tong, H. (2022). *Data mining: concepts and techniques*. Morgan kaufmann.

- 6.7 Techniques to improve classification accuracy
 - 6.7.1 Introducing ensemble methods
 - 6.7.2 Bagging
 - 6.7.3 Boosting
 - 6.7.4 Random forests

Outline

Introduction
Ensemble
Methods

Parallel
Ensemble:
Bagging &
RF

Sequential
Ensemble:
AdaBoost
& GB

+ xgboost

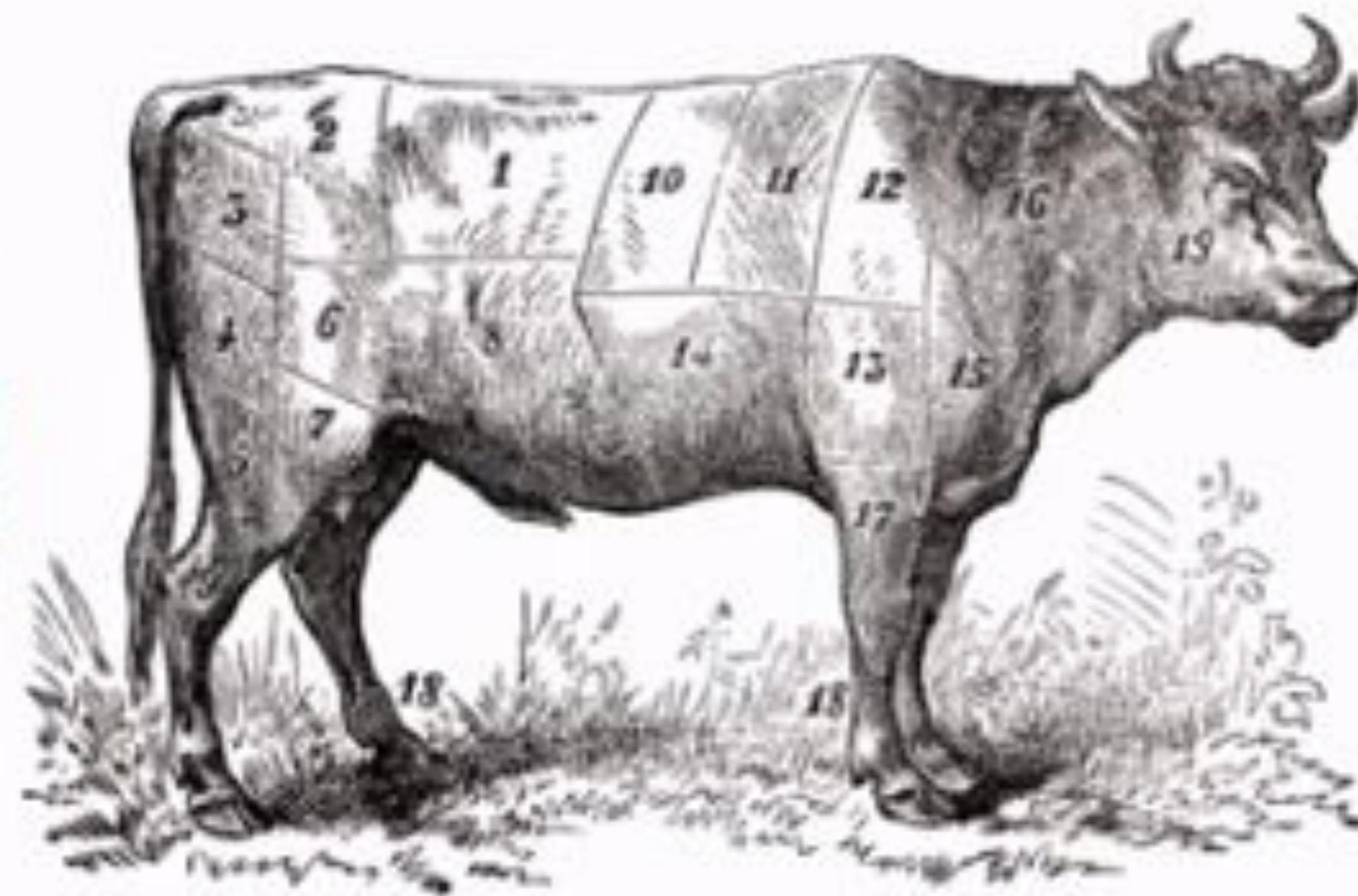
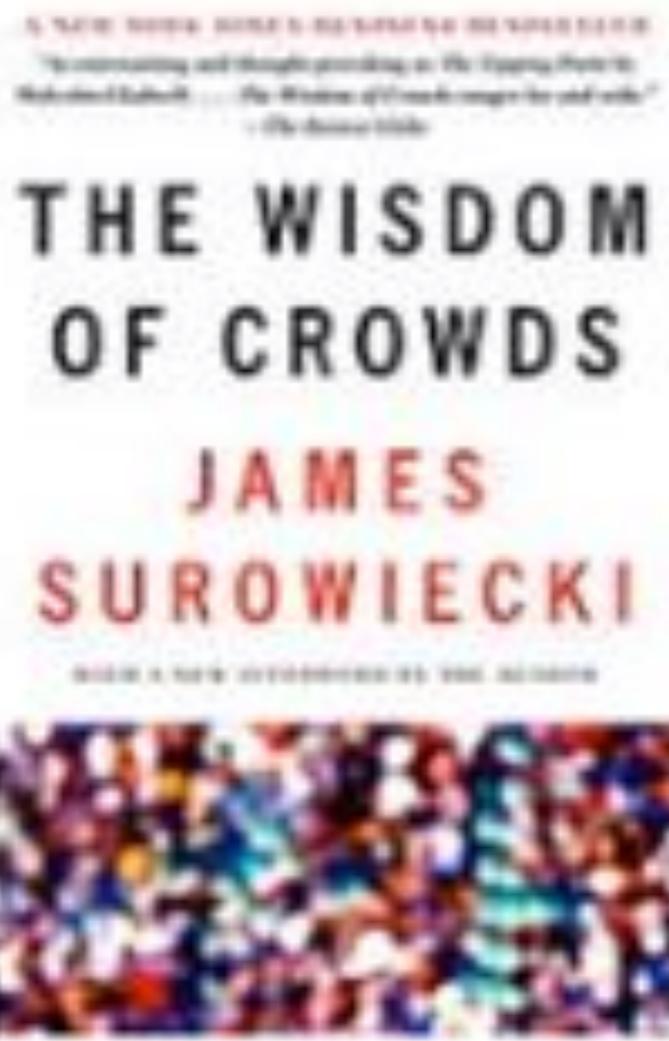


Overview Ensemble Methods

- Wisdom of the crowds: the *combined answer* of many models is often better than any one individual answer.
 - Ensemble diversity → training
 - Model aggregation into a single final decision → inference
- Parallel Ensembles:
 - Strong learners → **single model biasanya strong learner**
 - Homogeneous parallel ensembles
 - Bagging: bootstrap sampling + aggregating
 - Random Forest: bootstrap sampling + random subspace + aggregating
 - Heterogeneous parallel ensembles
 - Heterogeneous ensembles with weighting
 - Heterogeneous ensembles with meta learning
- Sequential Ensembles:
 - Weak learners (decision stump/depth 1), homogeneous learning algorithm
 - Adaptive boosting: instance weights
 - Kunapuli (2023) vs Han (2022)

→ buat model yg sgt \leq simple

The Wisdom of Crowds



*average of 800 guesses = 1,197
actual weight of the ox = 1,198*



The wisdom of the crowd is the process of taking into account the **collective opinion** of a group of individuals rather than a **single expert** to answer a question.

Ensemble Methods: The Wisdom of the Crowds

Dr. Randy Forrest works at a teaching hospital, and have a team with a diversity skills. Each resident has their own strengths.

Every time Dr. Forrest gets a new case, he solicits the opinions of his residents and then democratically **selects the final diagnosis as the most common one** from among all those proposed.

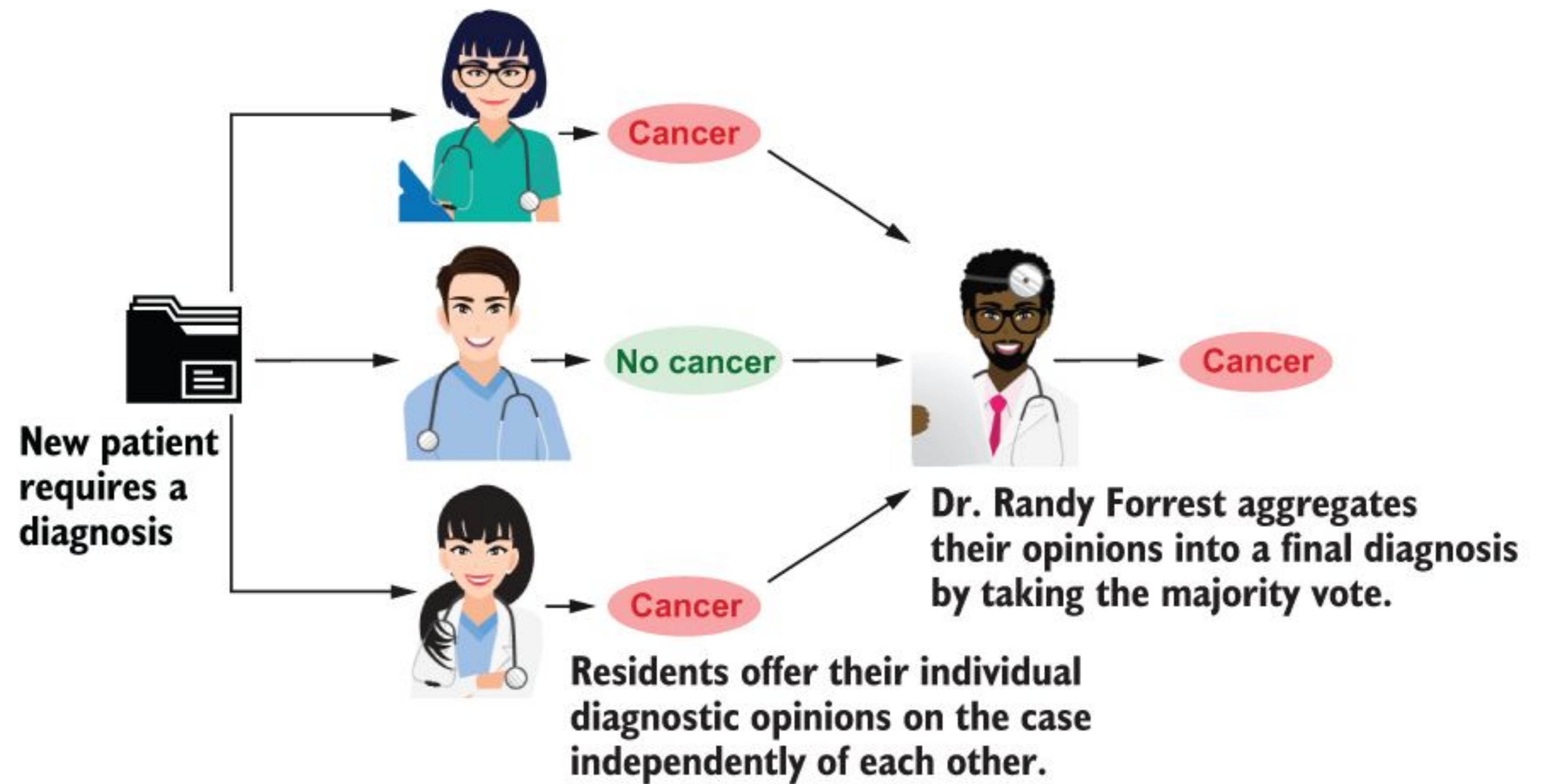


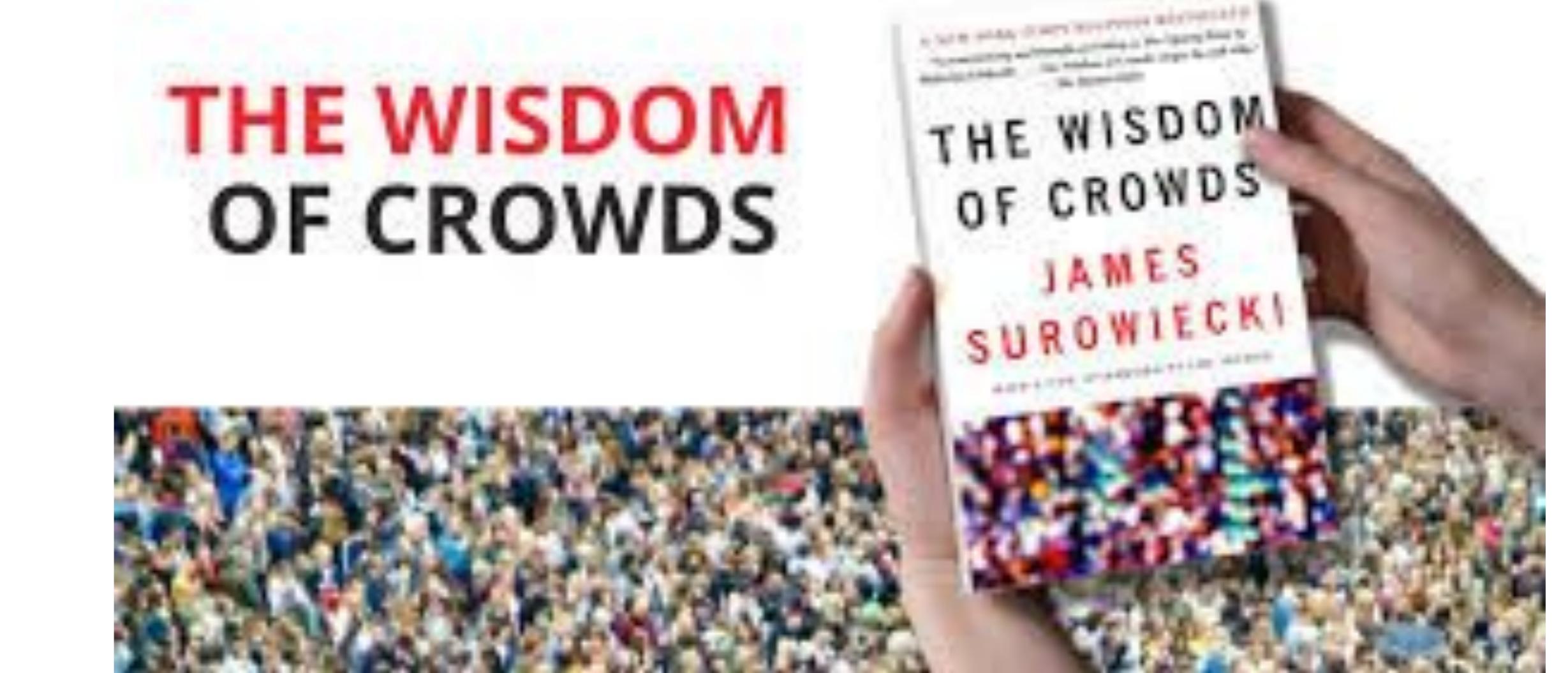
Figure 1.1 The diagnostic procedure followed by Dr. Randy Forrest every time he gets a new case is to ask all of his residents their opinions of the case. His residents offer their diagnoses: either the patient does or does not have cancer. Dr. Forrest then selects the majority answer as the final diagnosis put forth by his team.

Ensemble Methods: The Wisdom of the Crowds

An ensemble method relies on the notion of “wisdom of the crowd”: the **combined answer** of many models is often **better** than any **one** individual answer.

The secrets to the success of ensemble methods are:

- Ensemble diversity (variety of opinions to choose from) *varian orgz yg ditanya*
- Model aggregation into a single final decision *agregasi modelnya*



Diversity and independence are important because the best collective decisions are the product of disagreement and contest, not consensus or compromise.

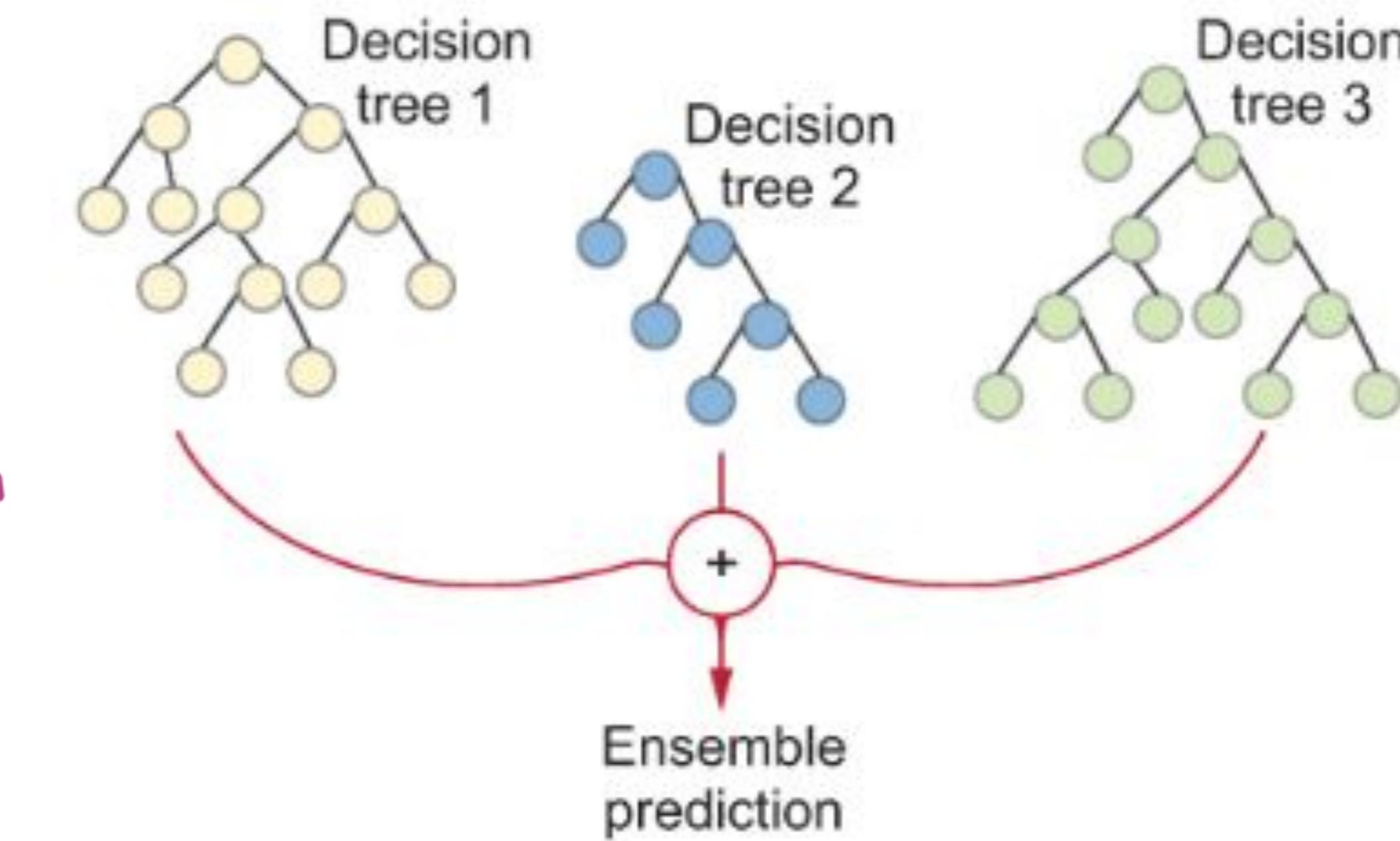
Terminology

satu model yg jd kolektif dr modelnya

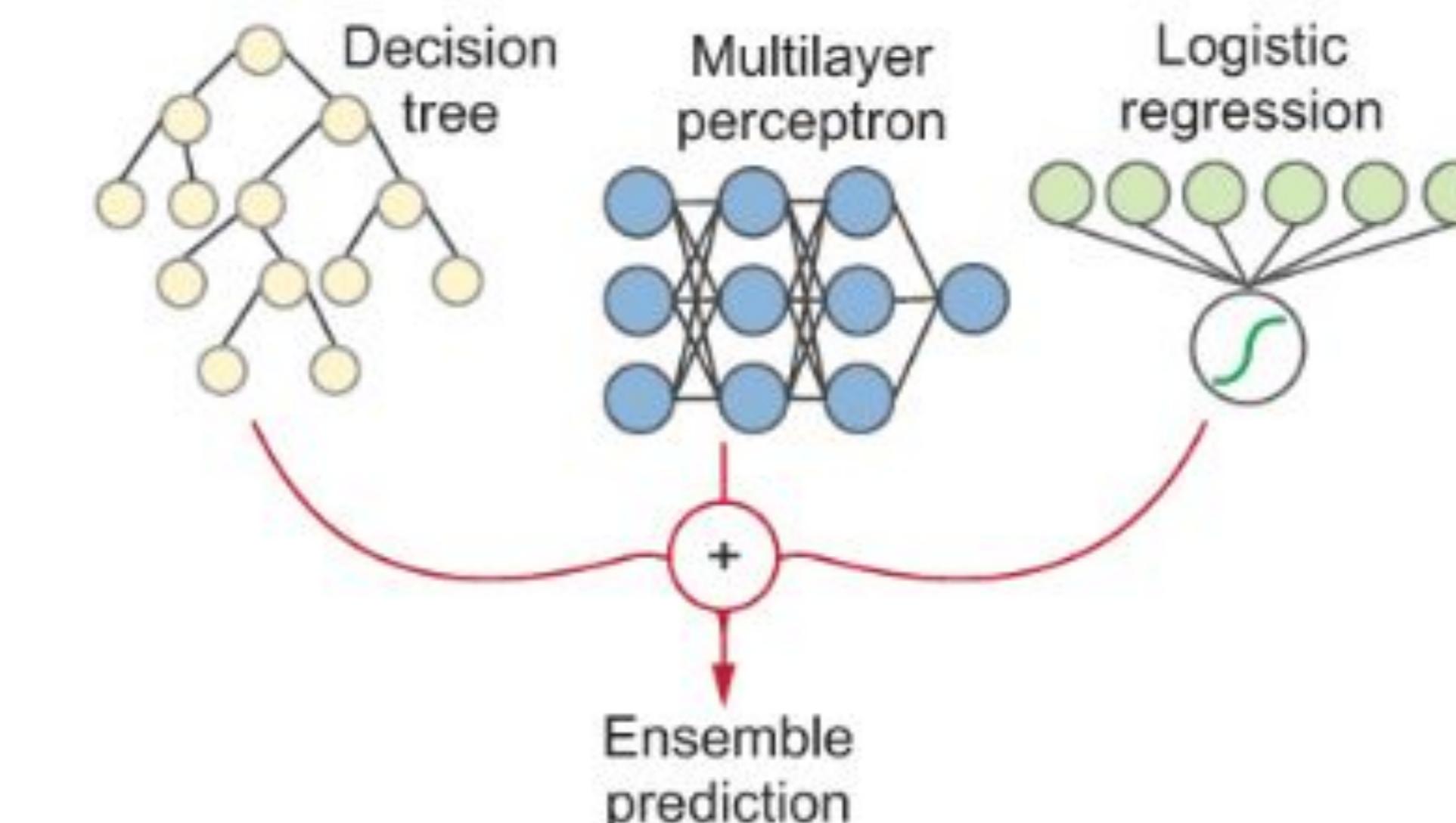
- Base models/learners/estimators: individual machine-learning models
- Parallel vs sequential ensemble method: train each base model independently (or in parallel) vs sequentially in stages.
- Homogeneous vs heterogeneous ensembles: all the base learners are trained using the same vs different base-learning algorithm.

mempengaruhi waktu training

Homogeneous parallel ensembles



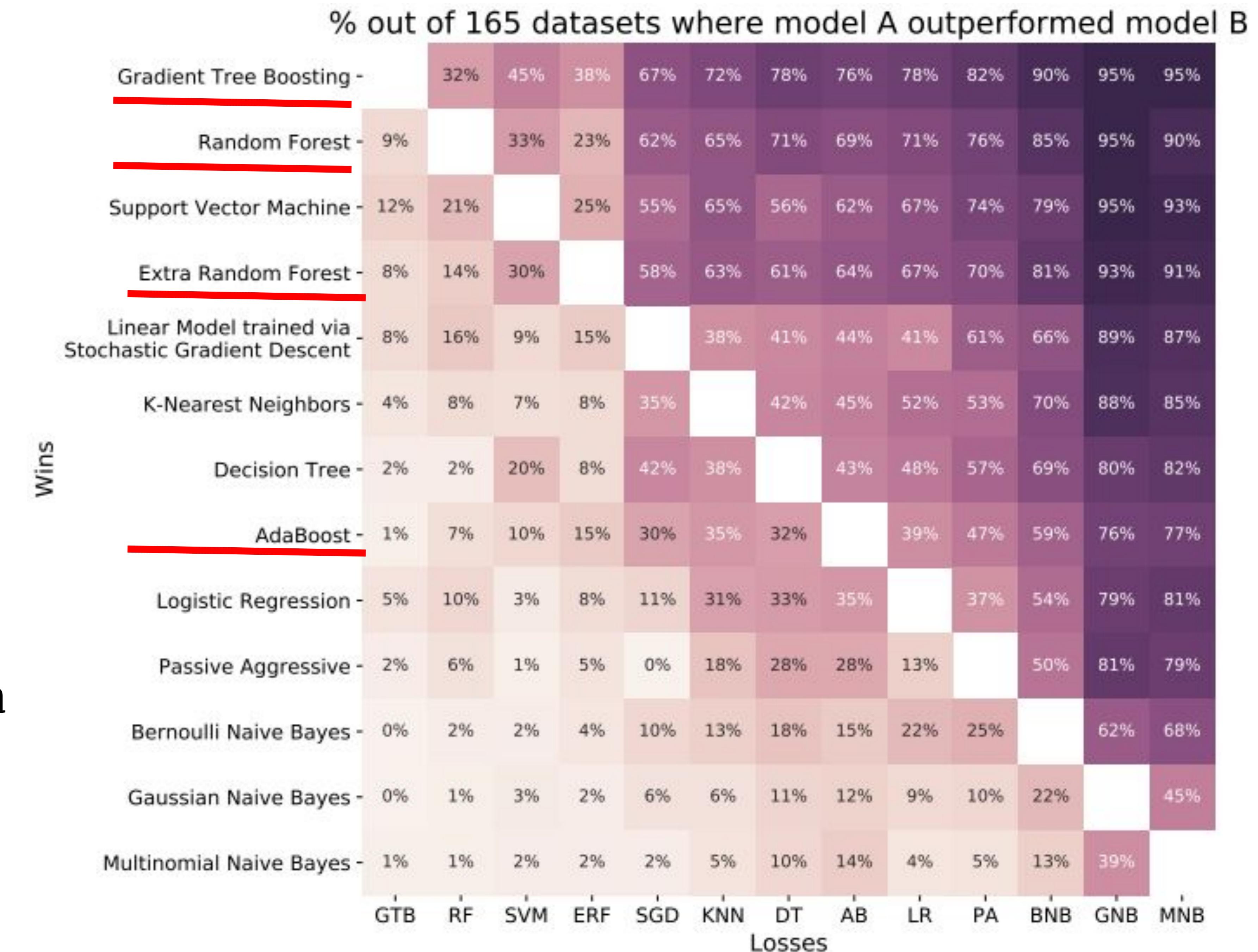
Heterogeneous parallel ensembles



Why Ensemble Methods ?

ketika data terstruktur & independen
biasanya bagus pakai ensemble

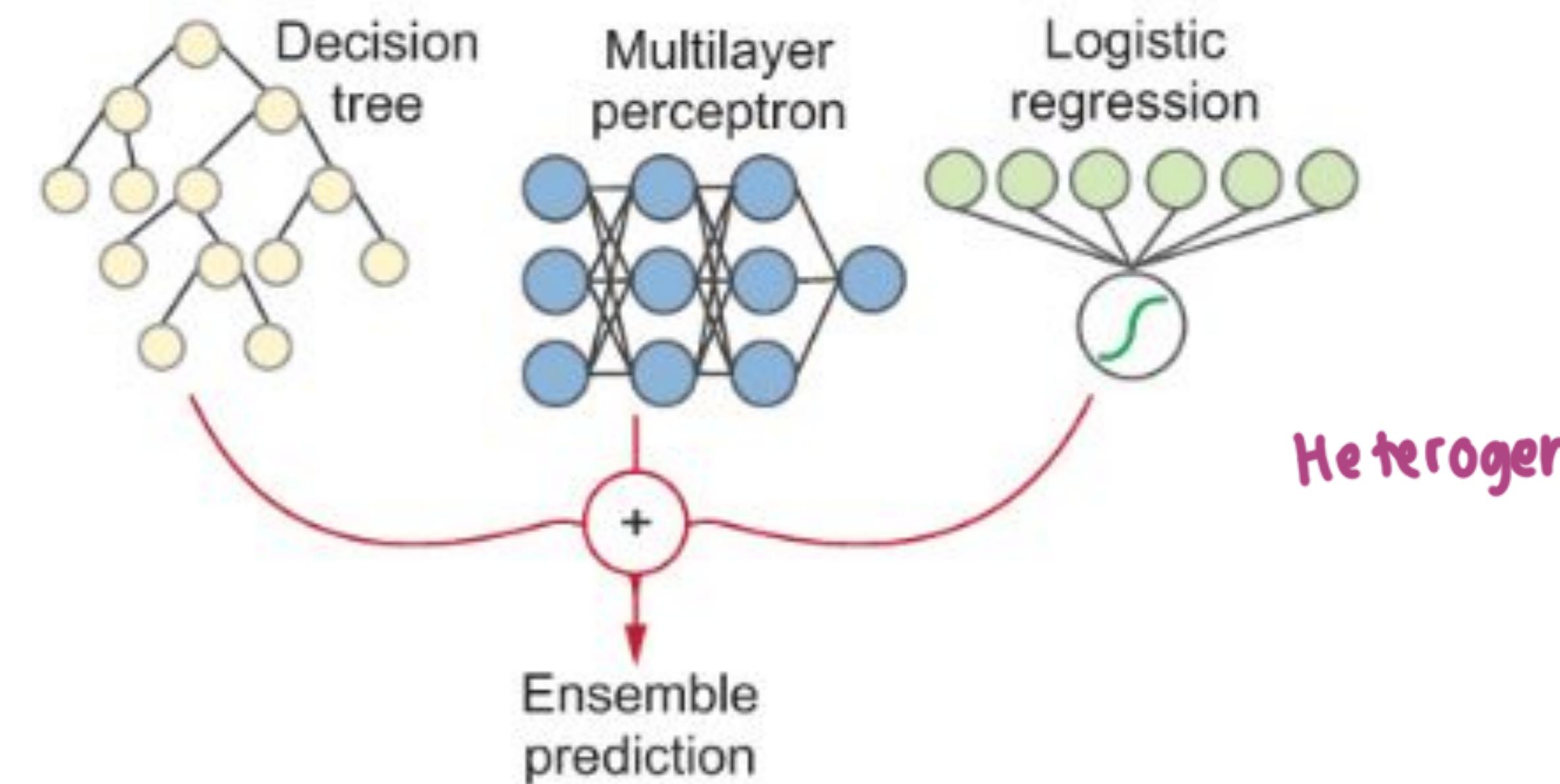
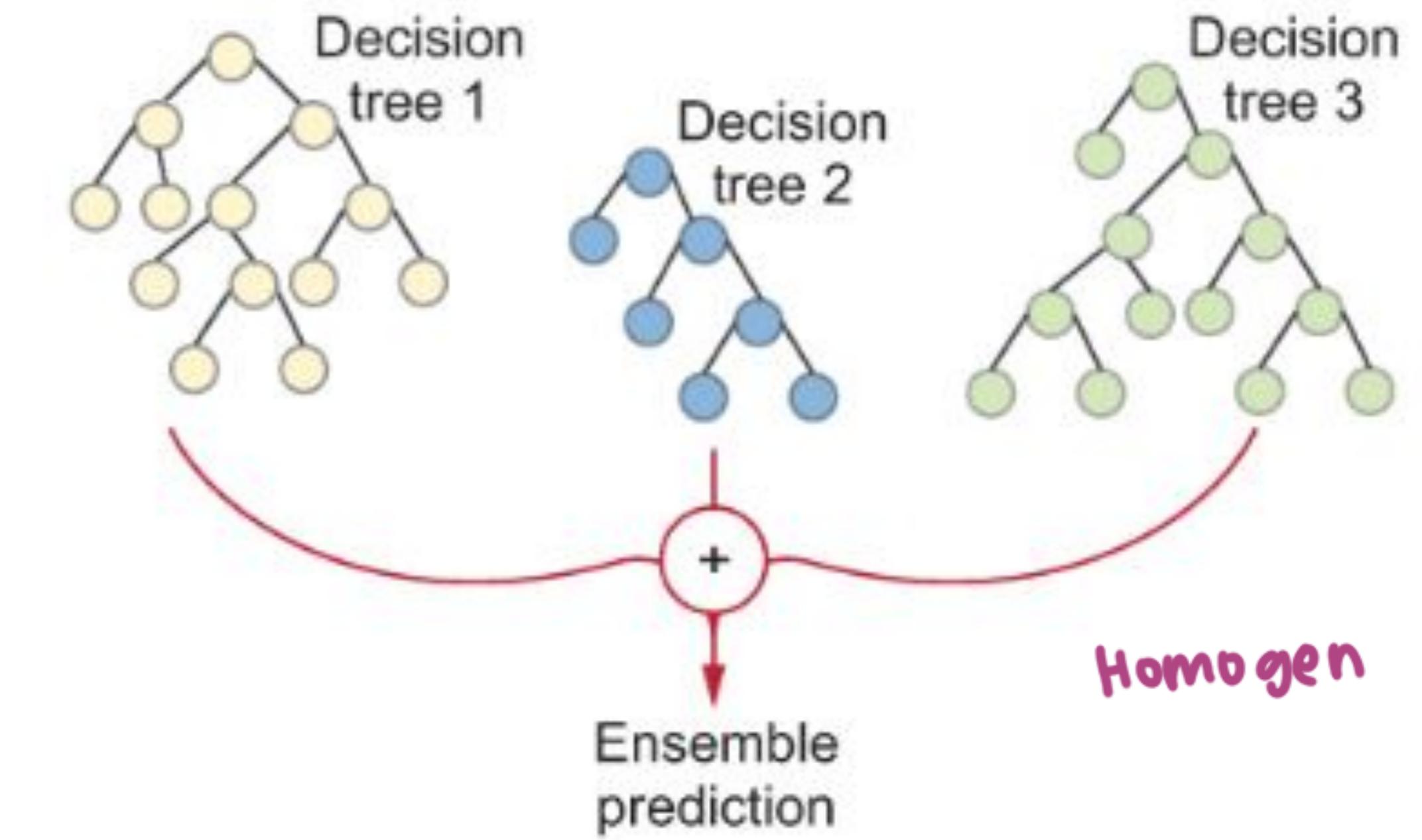
- Ensemble methods have proven to be very successful on **structured data**.
 - Anthony Goldbloom, CEO of Kaggle, revealed in 2015 that the three **most successful** algorithms for **structured problems** were XGBoost, random forest, and gradient boosting, **all ensemble methods**.
 - Olson et al. (2018) ranked machine-learning algorithm based on their performance on all benchmark data sets. Ensemble methods **outperformed** other methods, that's why ensemble methods (specifically, **tree-based ensembles**) are considered **sota** for structured dataset.



sota

Parallel Ensembles

exploit the *independence* of each base estimator

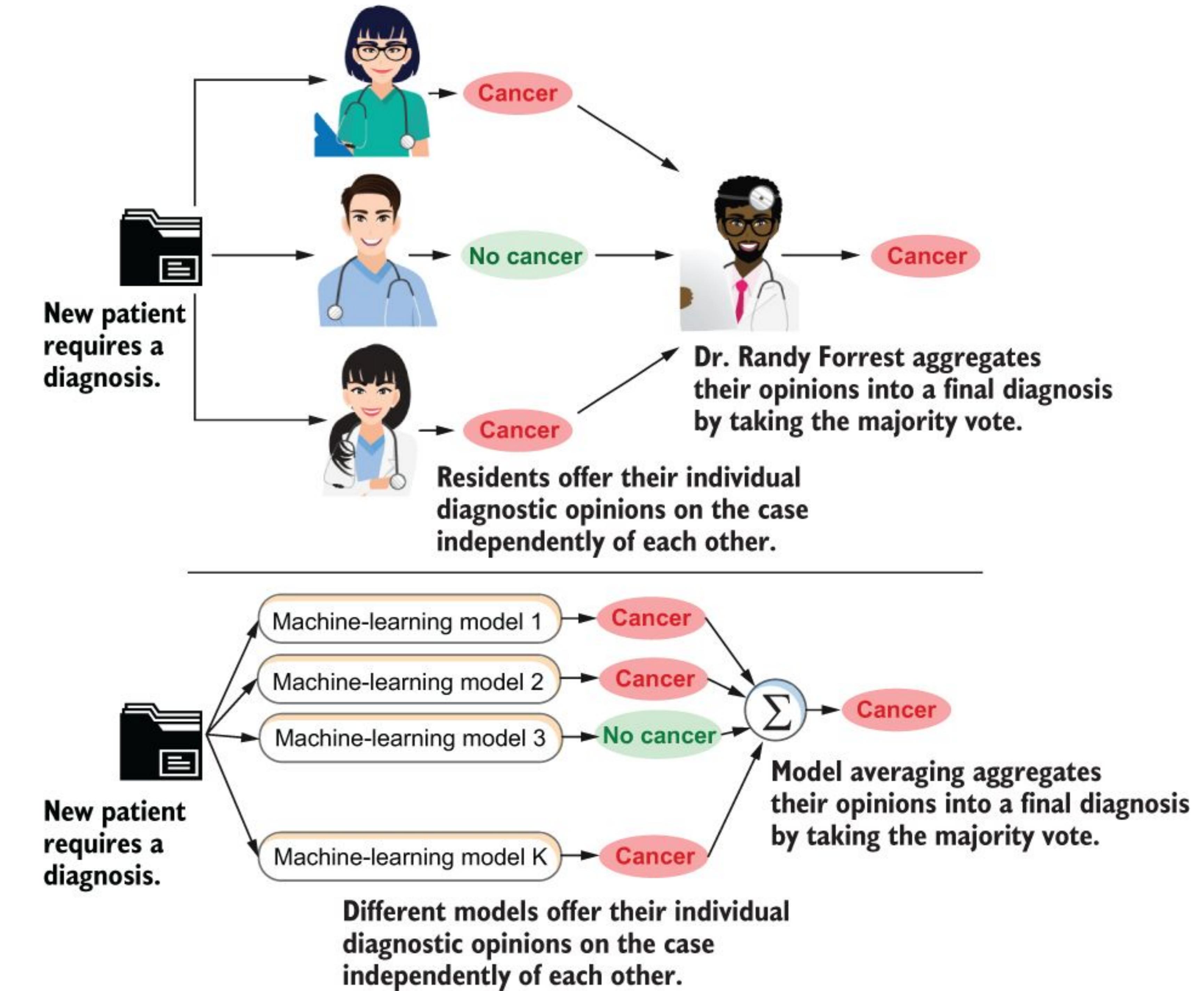


Homogeneous Parallel Ensembles

Homogeneous: component models are all trained using the **same** machine-learning algorithm

Parallel: train each component base estimator **independently** of the others (in parallel)

Dr. Randy Forrest's diagnostic process is an analogy of a parallel ensemble method



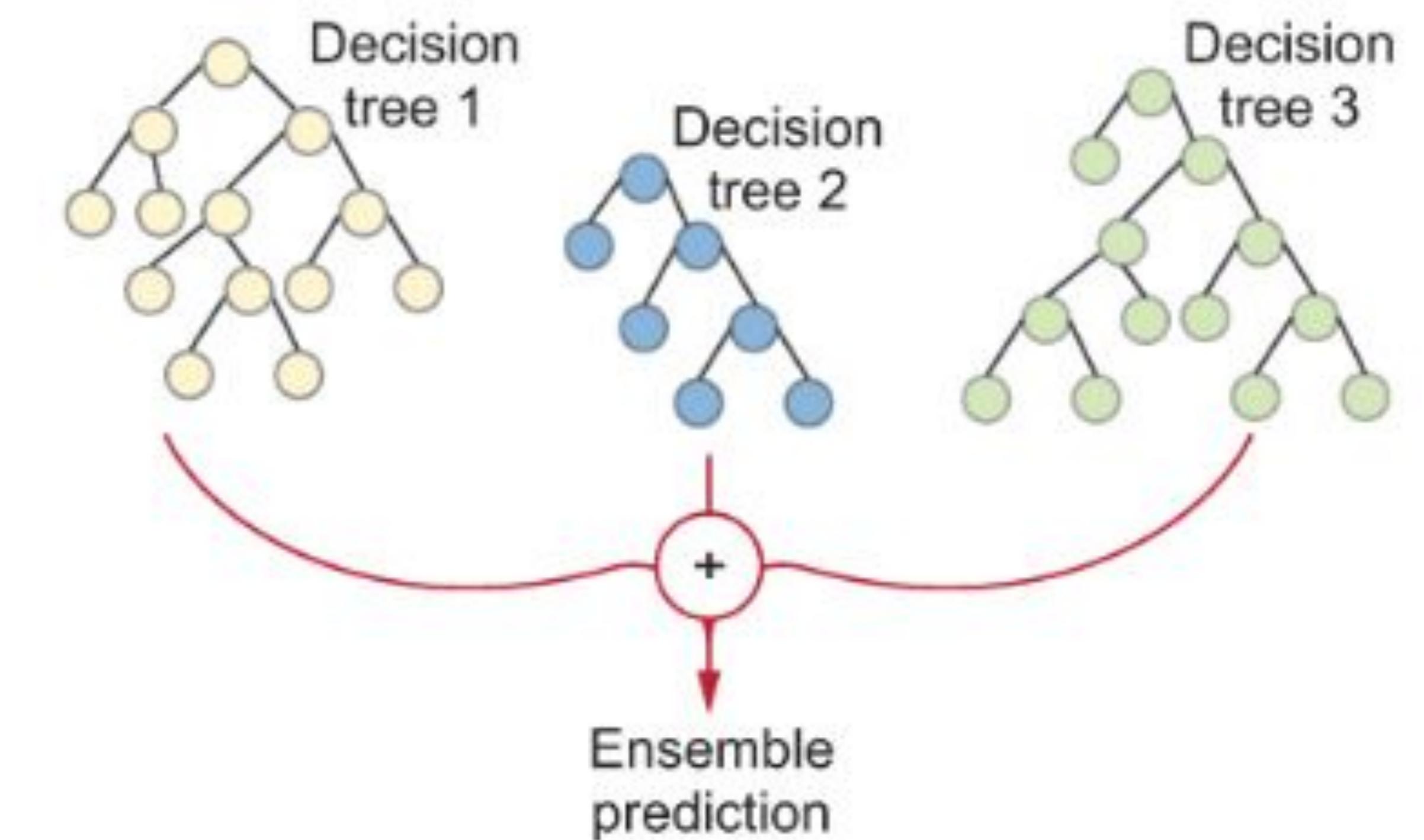
Homogeneous Parallel Ensembles: Ensemble Diversity



We have a single data set, so how do we obtain different base models with the same learning algorithm in homogeneous ensemble ?



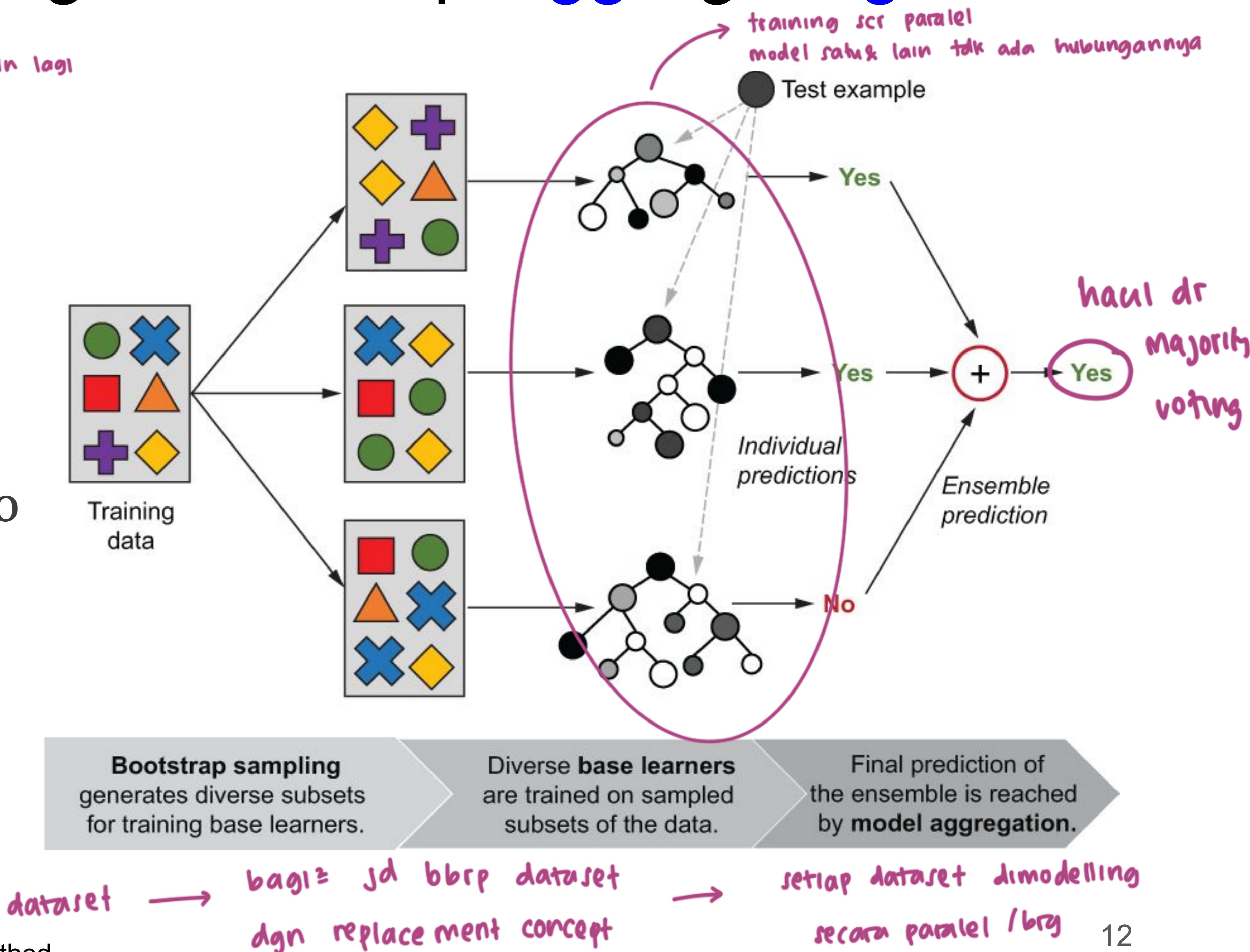
Tentuin best attribute yg bs jd node pakai information gain (AI)



Ensemble Diversity: Bagging - Bootstrap Aggregating

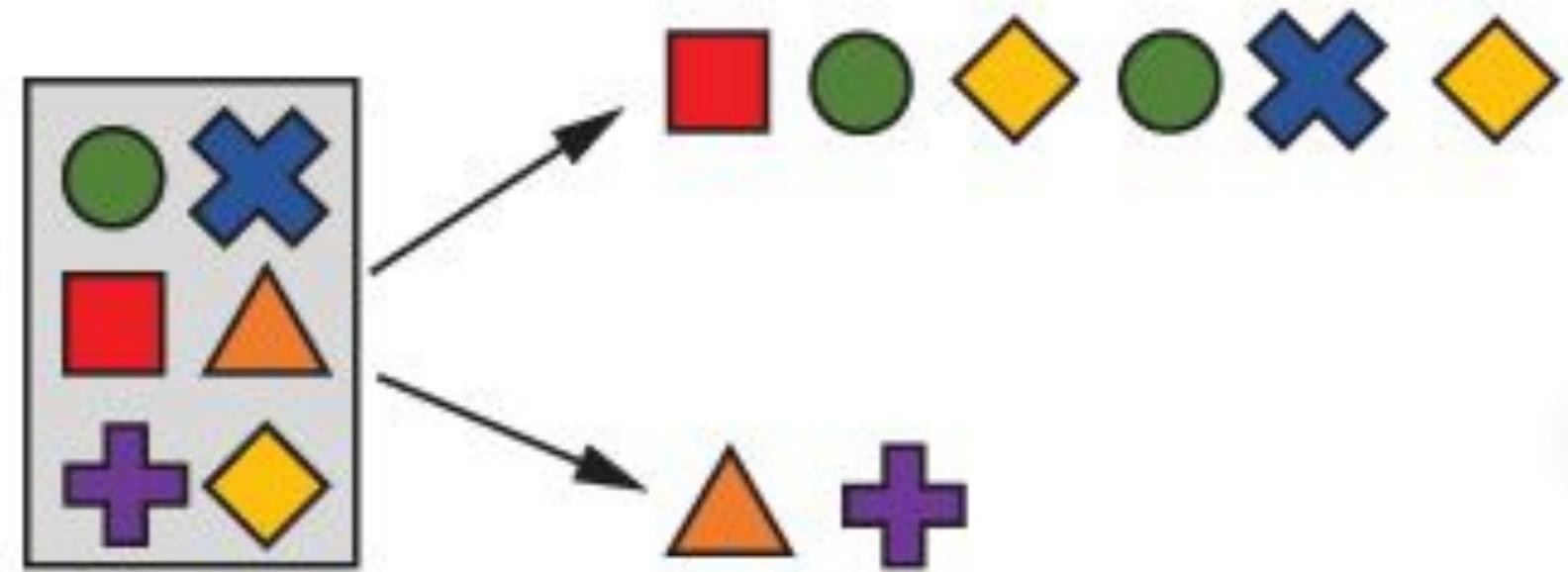
Training: bootstrap sampling (or sampling with replacement) is used to generate replicates of the training data set that are different from each other but drawn from the original data set. This ensures ensemble diversity.

Prediction: model aggregation is used to combine the predictions of the individual base learners into one ensemble prediction (classification: majority voting; regressing: simple averaging).



Bootstrap Sampling

Original set of objects
(e.g., training data)



Bootstrap sample: Sampling with replacement allows some objects to be selected more than once.

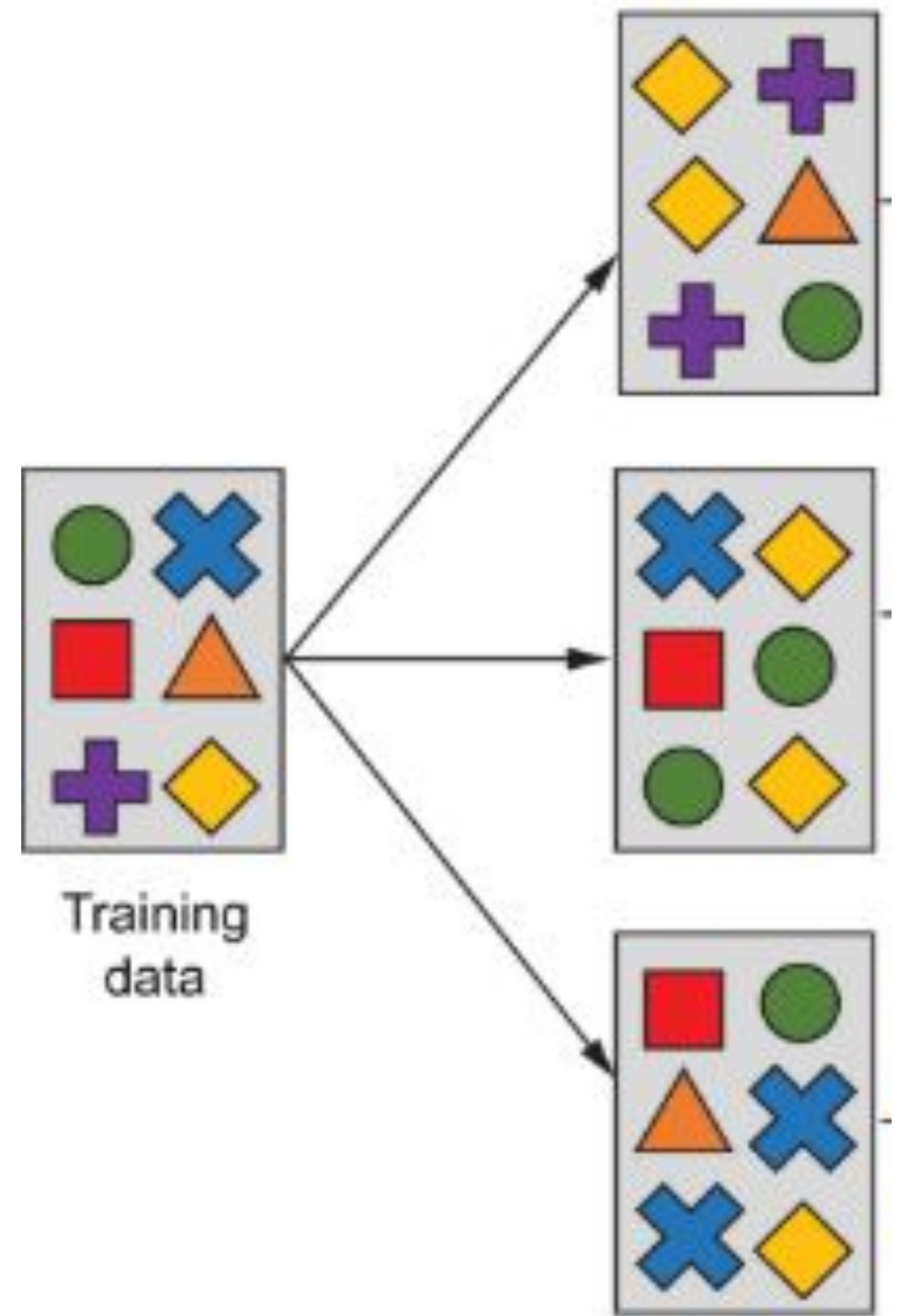
Out-of-bag sample: Sampling with replacement means some objects will not be selected even once.

```
import numpy as np
bag = np.random.choice(range(0, 50), size=50, replace=True)
np.sort(bag)
```

```
array([ 1,  3,  4,  6,  7,  8,  9, 11, 12, 12, 12, 14, 14, 14, 15, 15,
       21, 21, 21, 21, 24, 24, 25, 25, 26, 26, 26, 29, 29, 29, 31, 32, 32, 33, 33, 34,
       34, 34, 35, 35, 35, 37, 37, 39, 39, 39, 40, 43, 43, 43, 44, 44, 46, 46, 46, 48, 48, 48, 48, 49, 49, 49])
```

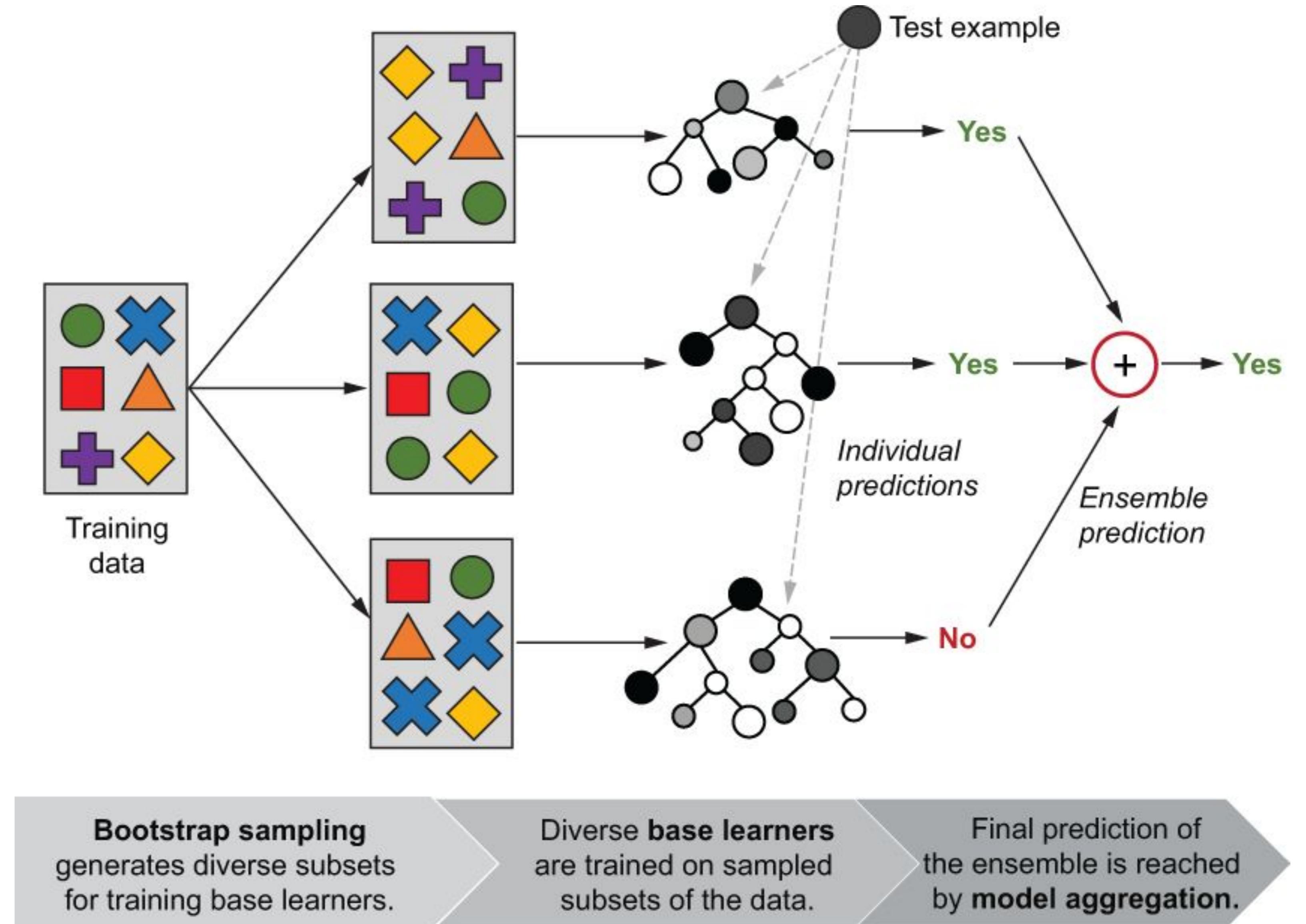
```
oob = np.setdiff1d(range(0, 50), bag)
oob
```

```
array([ 0,  2,  5, 10, 13, 16, 17, 18, 19, 20, 22, 23, 27, 28, 30, 36, 38,
       41, 42, 45, 47])
```



Training with Bootstrap Sampling

Because different bootstrap samples will contain different examples repeating a different number of times, each base estimator will turn out to be somewhat different from the others.



Bagging: Bootstrap Aggregation (Han et al., 2022)

Algorithm: Bagging. The bagging algorithm—create an ensemble of classification models for a learning scheme where each model gives an equally weighted prediction.

Input:

- D , a set of d training tuples;
- k , the number of models in the ensemble;
- a classification learning scheme (e.g., decision tree algorithm, naïve Bayesian, etc.).

Output: The ensemble—a composite model, M^* .

Method:

- (1) **for** $i = 1$ to k **do** // create k models:
- (2) create bootstrap sample, D_i , by sampling D with replacement;
- (3) use D_i and the learning scheme to derive a model, M_i .
- (4) **endfor**

To use the ensemble to classify a tuple, X :

let each of the k models classify X and return the majority vote;

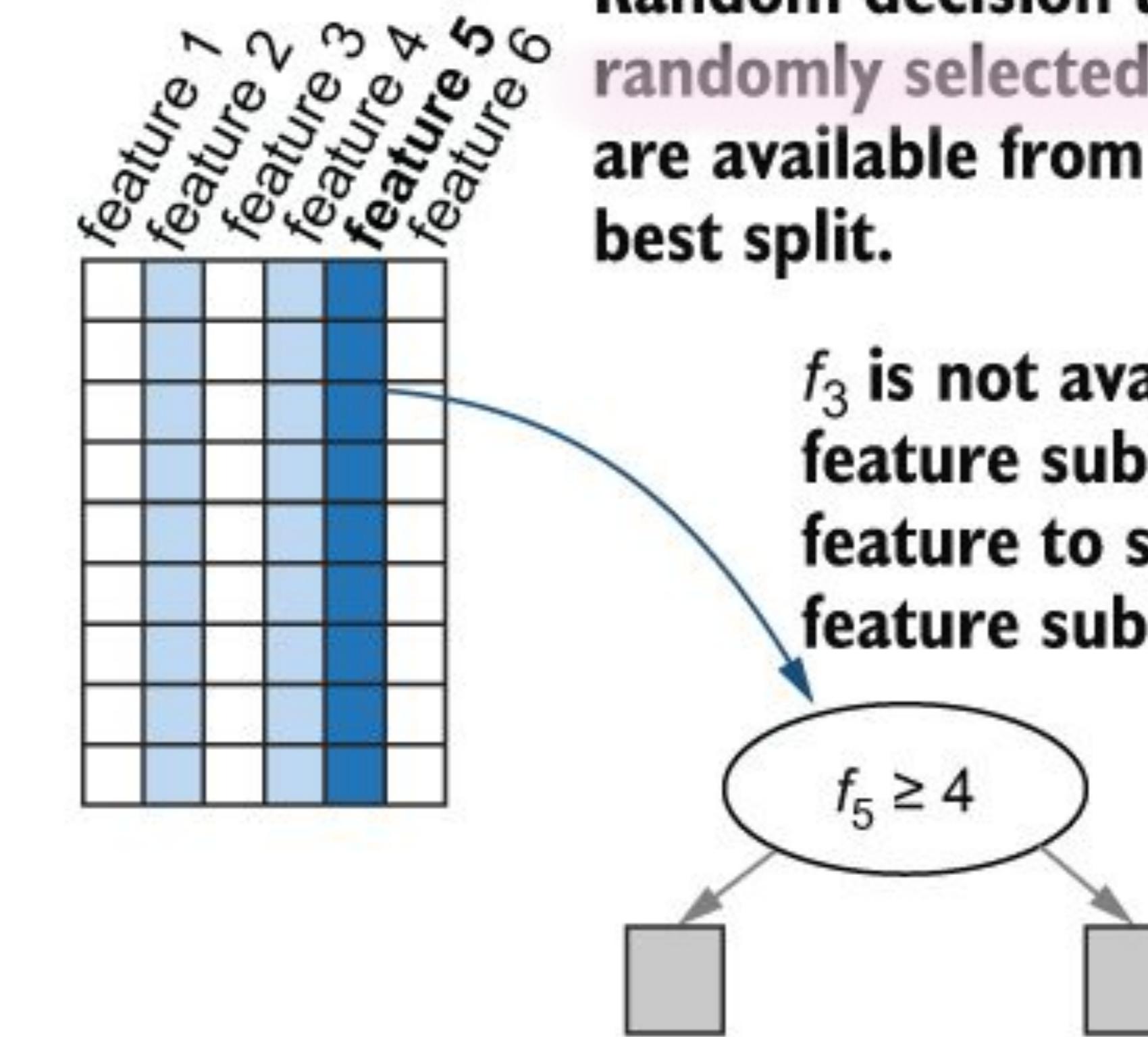
The bagged classifier, M^* , counts the votes and assigns the class with the most votes to X .

FIGURE 6.25

Random Forest

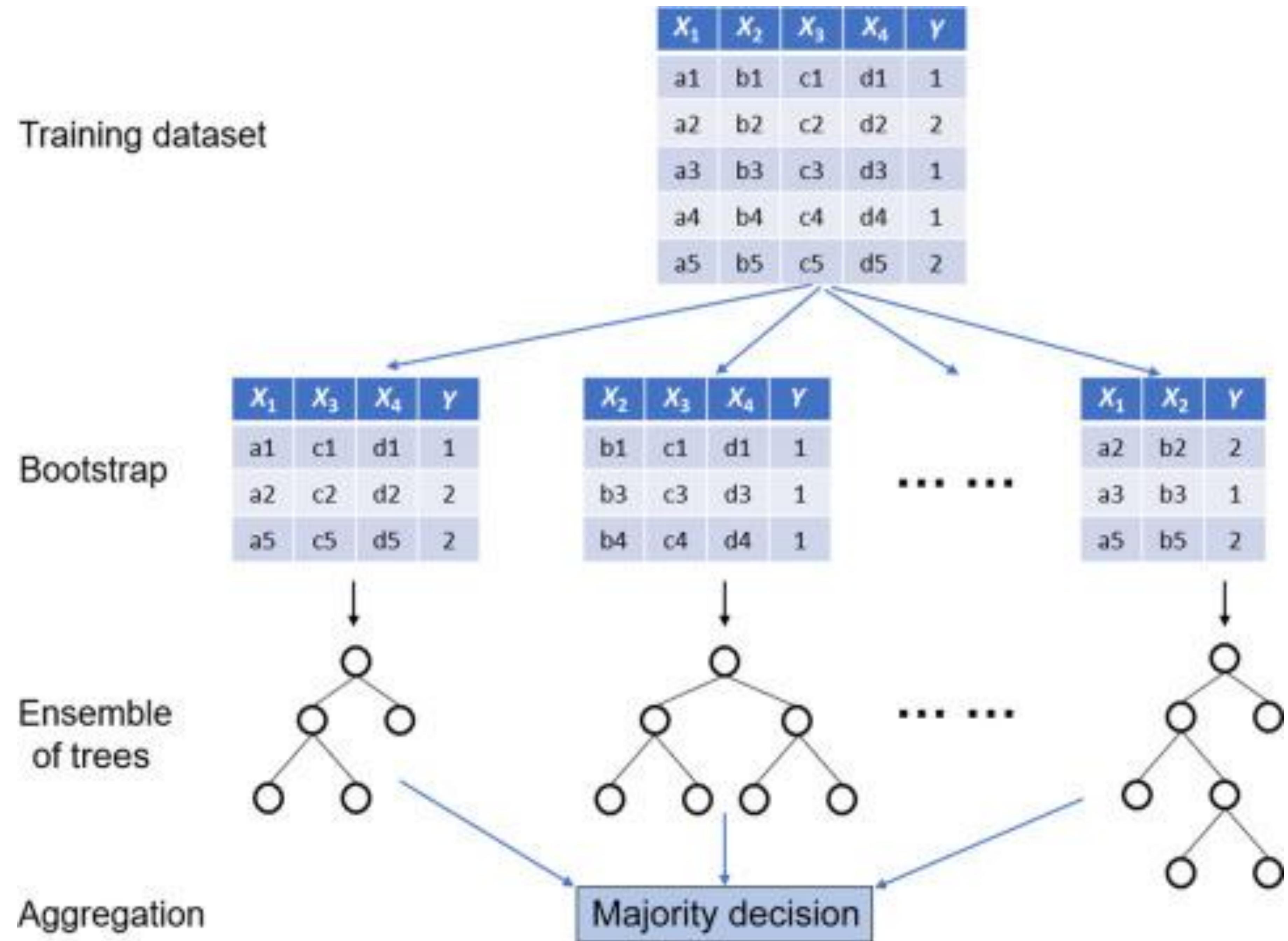
bagging ensemble tp tdk semua
feature nya digunakan

- an ensemble of **randomized decision trees** constructed using a special extension of **bagging**.
- Randomized decision trees are trained using a **modified decision-tree learning algorithm**
 - Problem: the same small number of **dominant features** are **repeatedly used** in different trees. This makes the ensemble **less diverse**.
 - Solution: **randomly samples features** before creating a decision node
- This additional source of randomness increases ensemble diversity and generally leads to better predictive performance.



Random decision tree: only some randomly selected subsets of features are available from which to select the best split.

Random Forest



Random forest consists of using randomly selected inputs (**bootstrap sample**) or combinations of inputs at each node (**random subspace**) to grow each tree.

sample dr featurenya

variasi banyak karena ambil random features x random data

nentuinya random

Algorithm 1: Pseudo code for the random forest algorithm

To generate c classifiers:

for $i = 1$ to c **do**

 Randomly sample the training data D with replacement to produce D_i

 Create a root node, N_i , containing D_i

 Call BuildTree(N_i)

end for

BuildTree(N):

if N contains instances of only one class **then**

return

else

 ga semua atribut dihitung information gain nya,
 random cm brp % atribut yg dihitung

 Randomly select $x\%$ of the possible splitting features in N

 Select the feature F with the highest information gain to split on

 Create f child nodes of N , N_1, \dots, N_f , where F has f possible values (F_1, \dots, F_f)

for $i = 1$ to f **do**

 Set the contents of N_i to D_i , where D_i is all instances in N that match

F_i

 Call BuildTree(N_i)

end for

end if

Misal ada 10 atr \rightarrow disini ambil 50%

\rightarrow random 5 atr x hitung gain nya \rightarrow baru pilih best atr

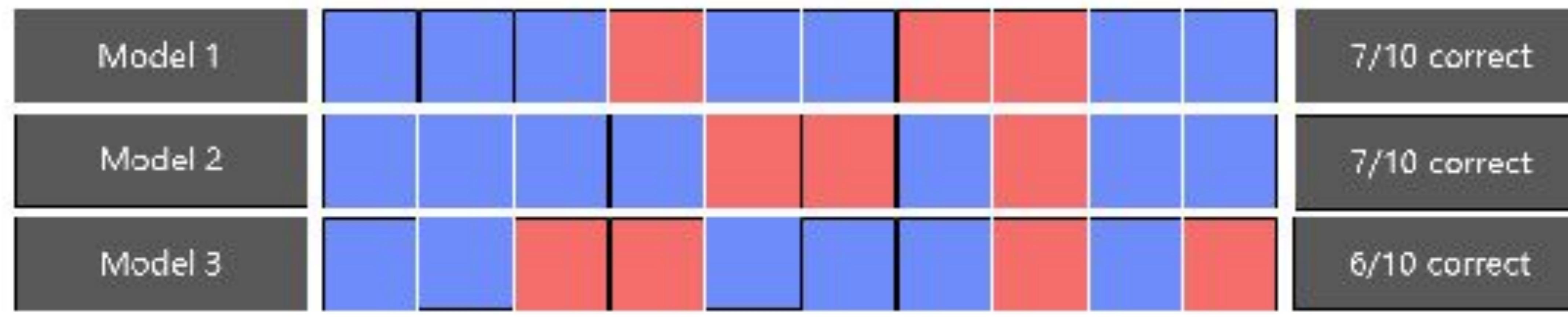
<https://scikit-learn.org/stable/modules/ensemble.html#forest>

- Each tree in the ensemble is built from a sample drawn with replacement (i.e., a bootstrap sample) from the training set.
- When splitting each node during the construction of a tree, the best split is found either from all input features or a **random subset** of size `max_features`.
- the scikit-learn implementation combines classifiers by **averaging their probabilistic prediction**, bukan majority instead of letting each classifier vote for a single class.

Kalo ga ada info di ujian
ambil yg kanan

Majority Voting

task regress → average



Ensemble Model
(Majority Voting)

! or more tutorials: algobeans.com

The predicted class probabilities of an input sample are computed as the mean predicted class probabilities of the trees in the forest.

Why Random Forest ?

- Ensemble methods have been a very effective tool to **improve the performance** of multiple existing models (Breiman, 2021)
- Random forest is also well known as one of **the most accurate** and as a **fast learning** method **independently** from the nature of the datasets, as shown by various recent comparative studies (Fernández-Delgado et al. 2014; Caruana and Niculescu-Mizil 2006; Rokach 2016).
- Ensemble learning algorithms like random forests are well suited for medium to large datasets.

data jelek pake ensemble → hasilnya bagus

Random Forest: $d \gg n$

- When the number of independent variables (d) is **larger than** the number of observations (n), linear regression and logistic regression algorithms will not run, because the number of parameters to be estimated exceeds the number of observations. Random forest works because not all predictor variables are used at once.

Schonlau, M., & Zou, R. Y. (2020). The random forest algorithm for statistical learning. *The Stata Journal*, 20(1), 3-29.

Random Forest: Summary

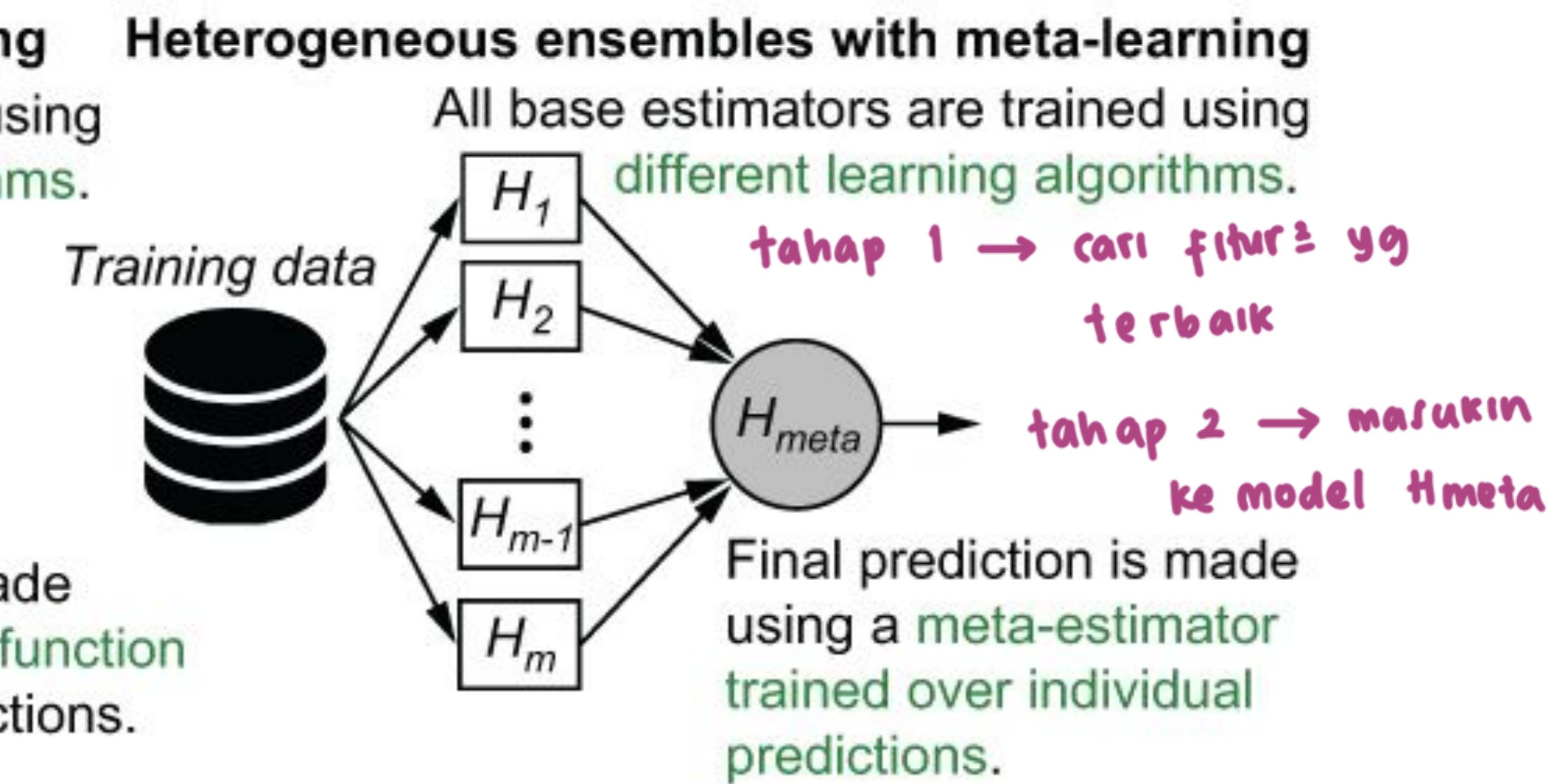
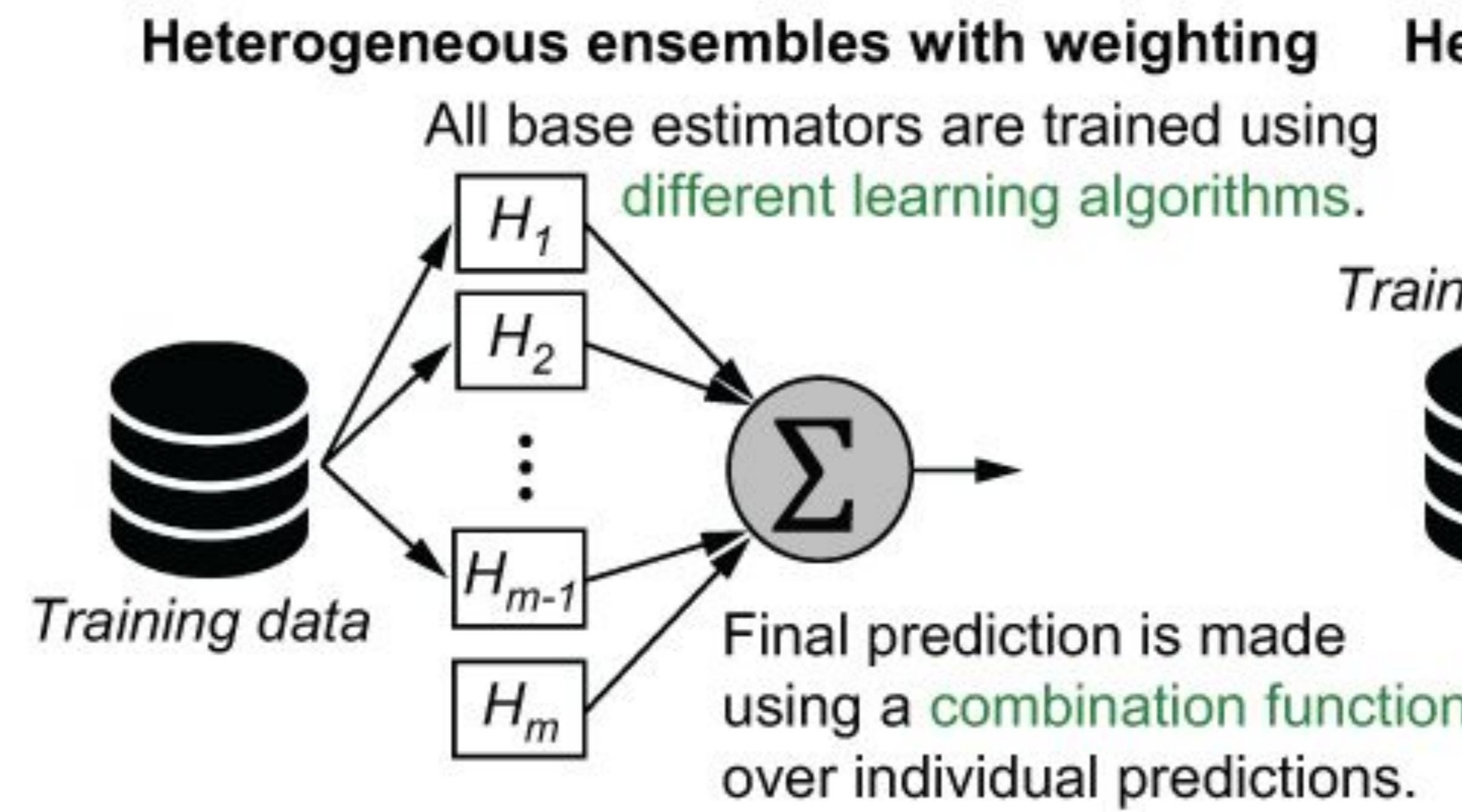
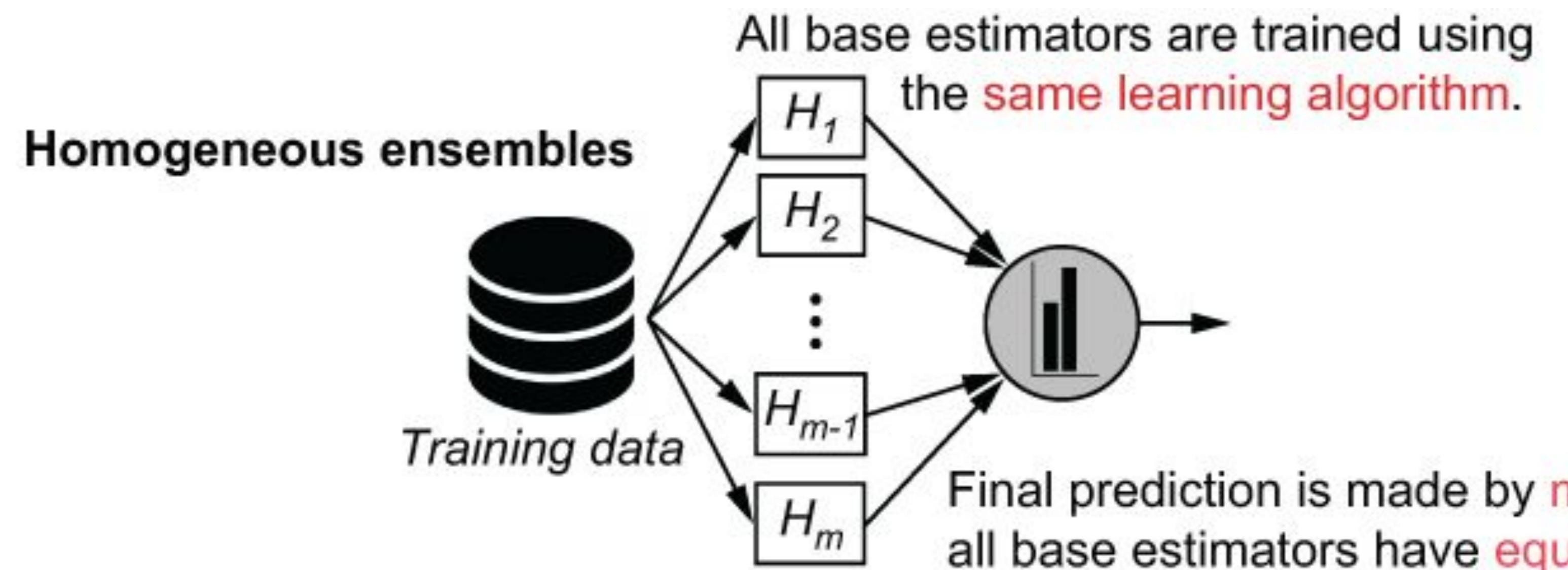
Advantages

- one of the most accurate learning algorithms available
- It runs efficiently on large databases. (parallel)
- It has methods for balancing error in class population unbalanced data sets.

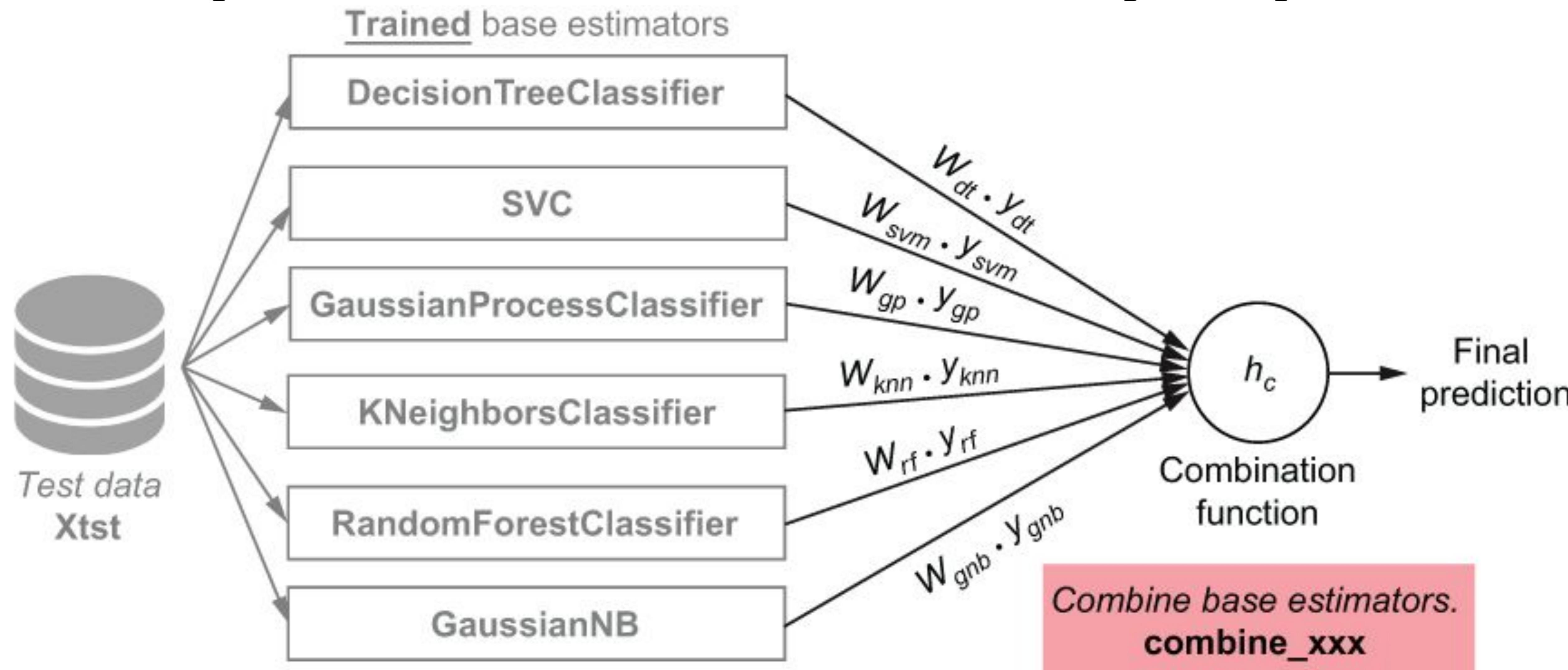
Disadvantages

- Random forests have been observed to overfit for some datasets with noisy classification
- Prediction time may be slow

Heterogeneous Parallel Ensembles

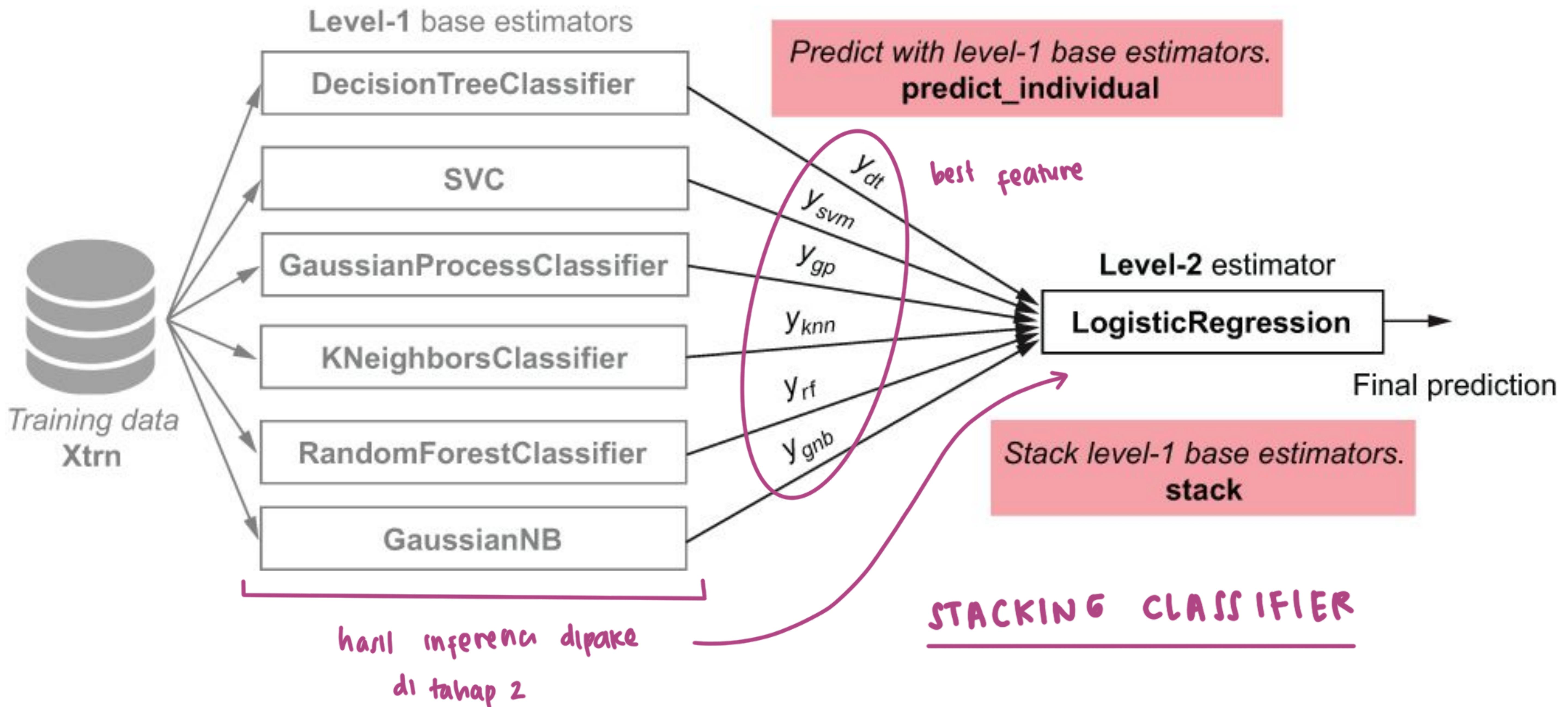


Heterogeneous Ensembles with Weighting

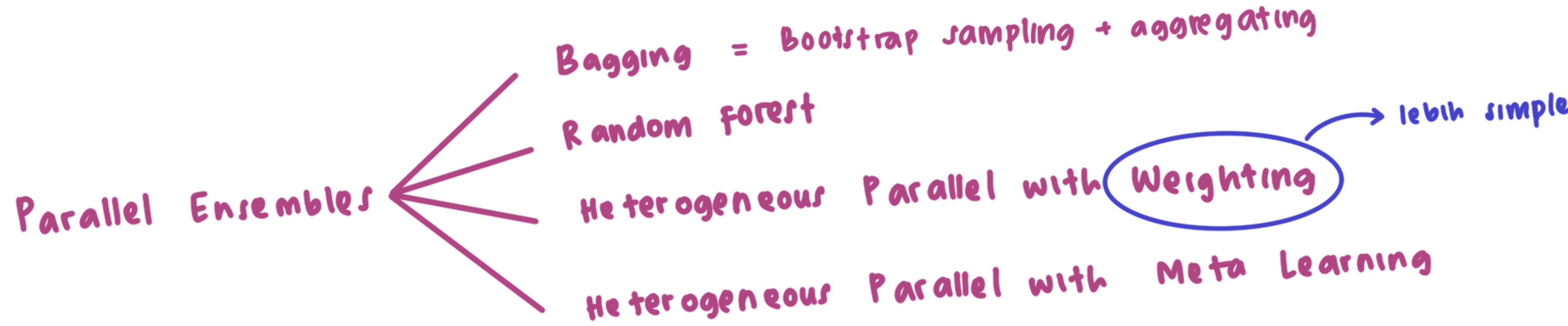


$$y_{final} = w_1 \cdot y_1 + w_2 \cdot y_2 + \cdots + w_m \cdot y_m = \sum_{t=1}^m w_t \cdot y_t$$

Heterogeneous Ensembles with Meta Learning



CONCLUSION

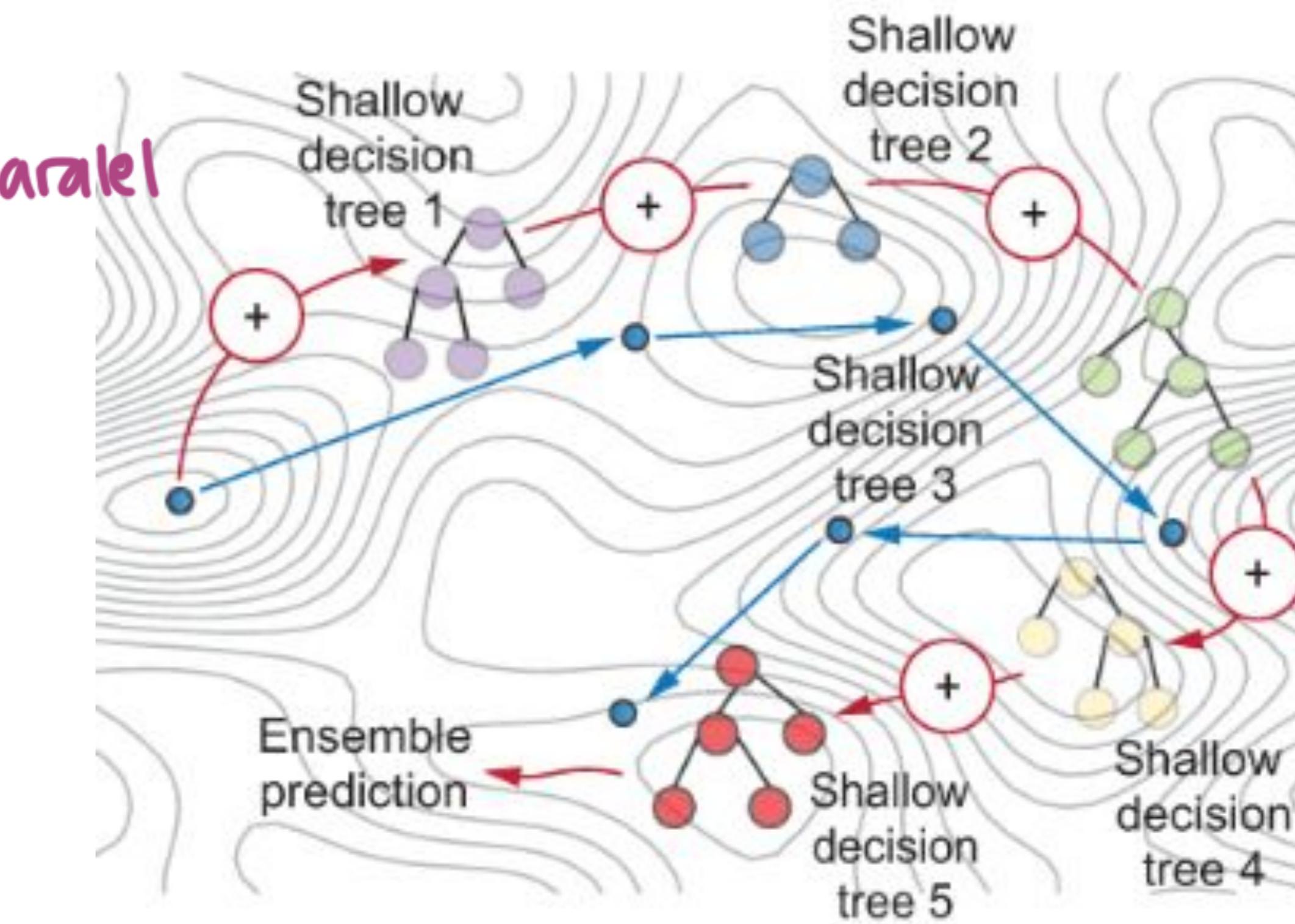
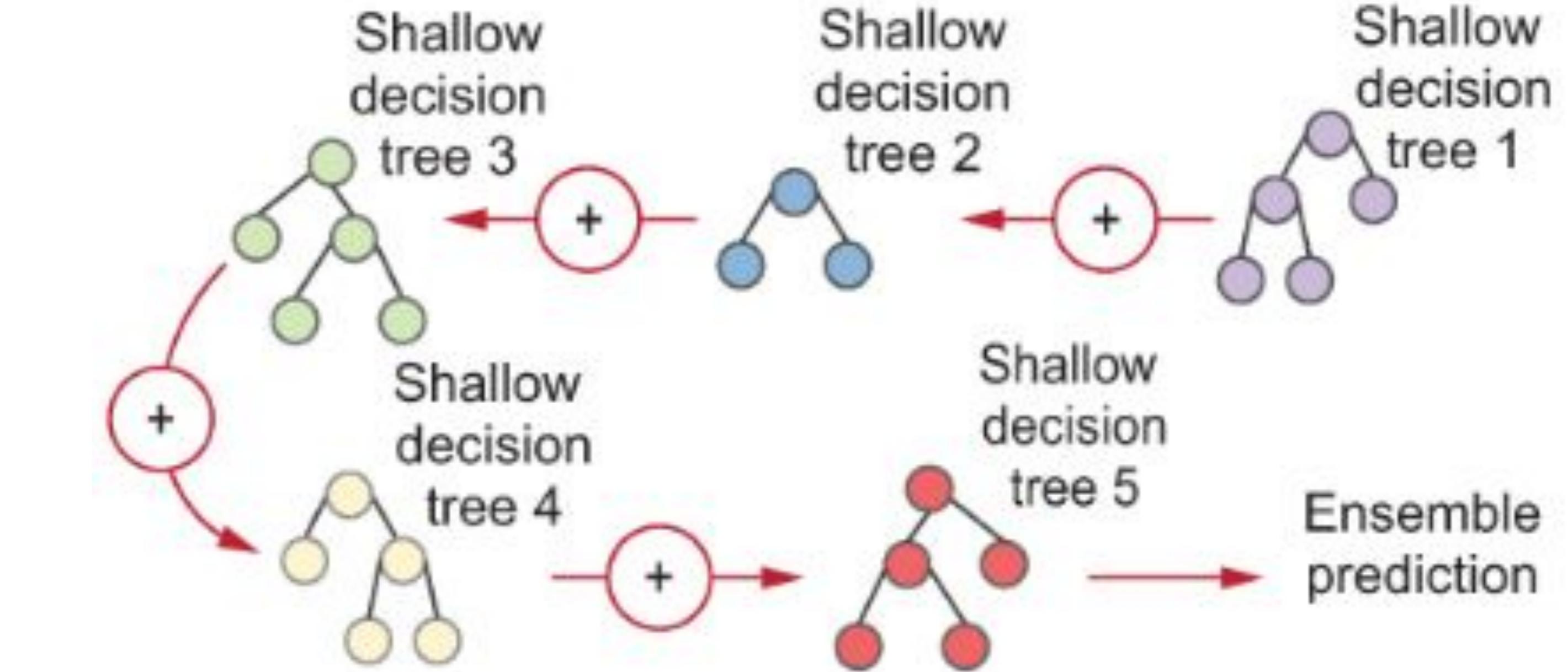


Sequential Ensembles

exploit the *dependence* of base estimators

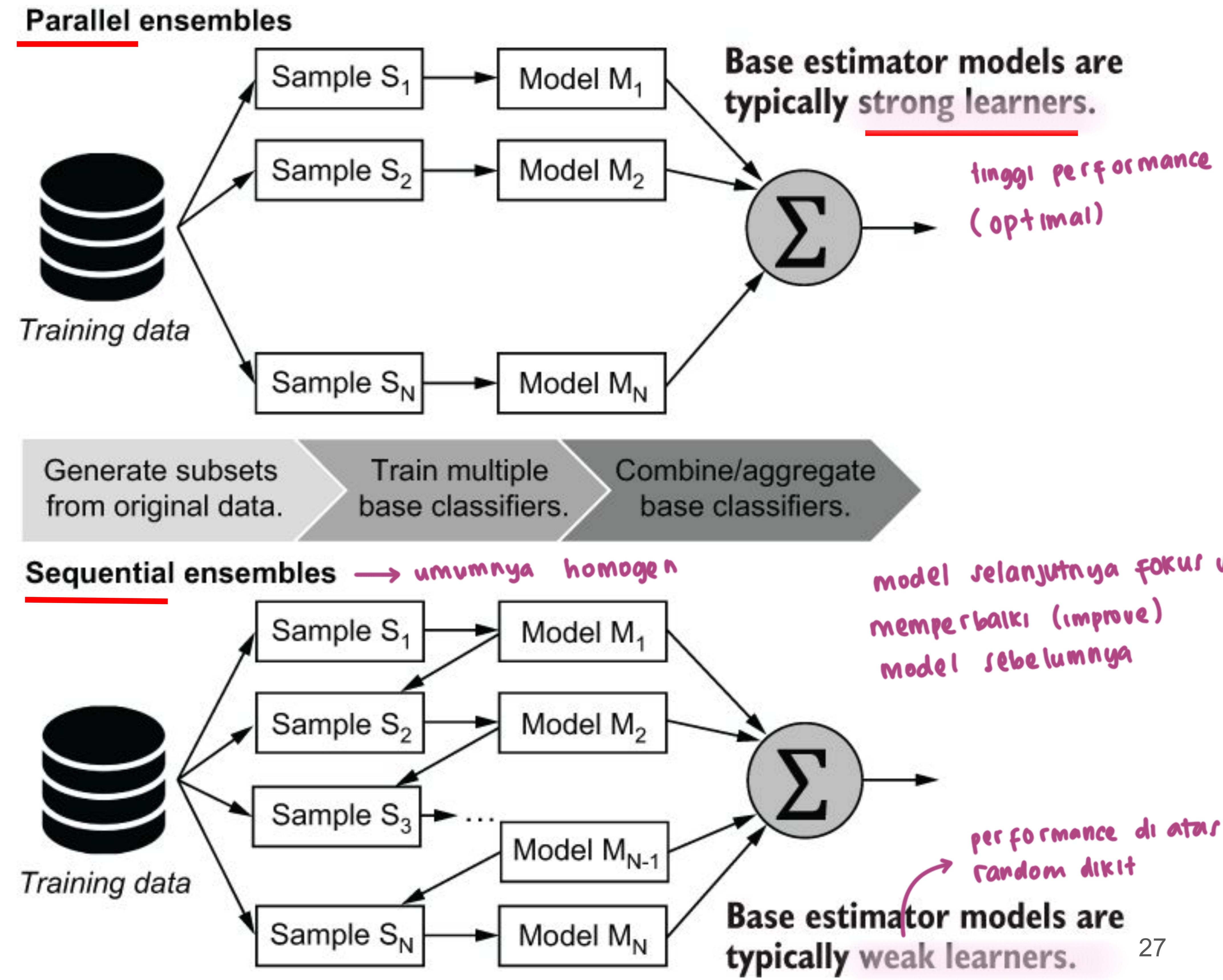
hanus weak learner
(yg performanya kurung bagus)

lebih aman dr parallel



Parallel vs Sequential Ensembles

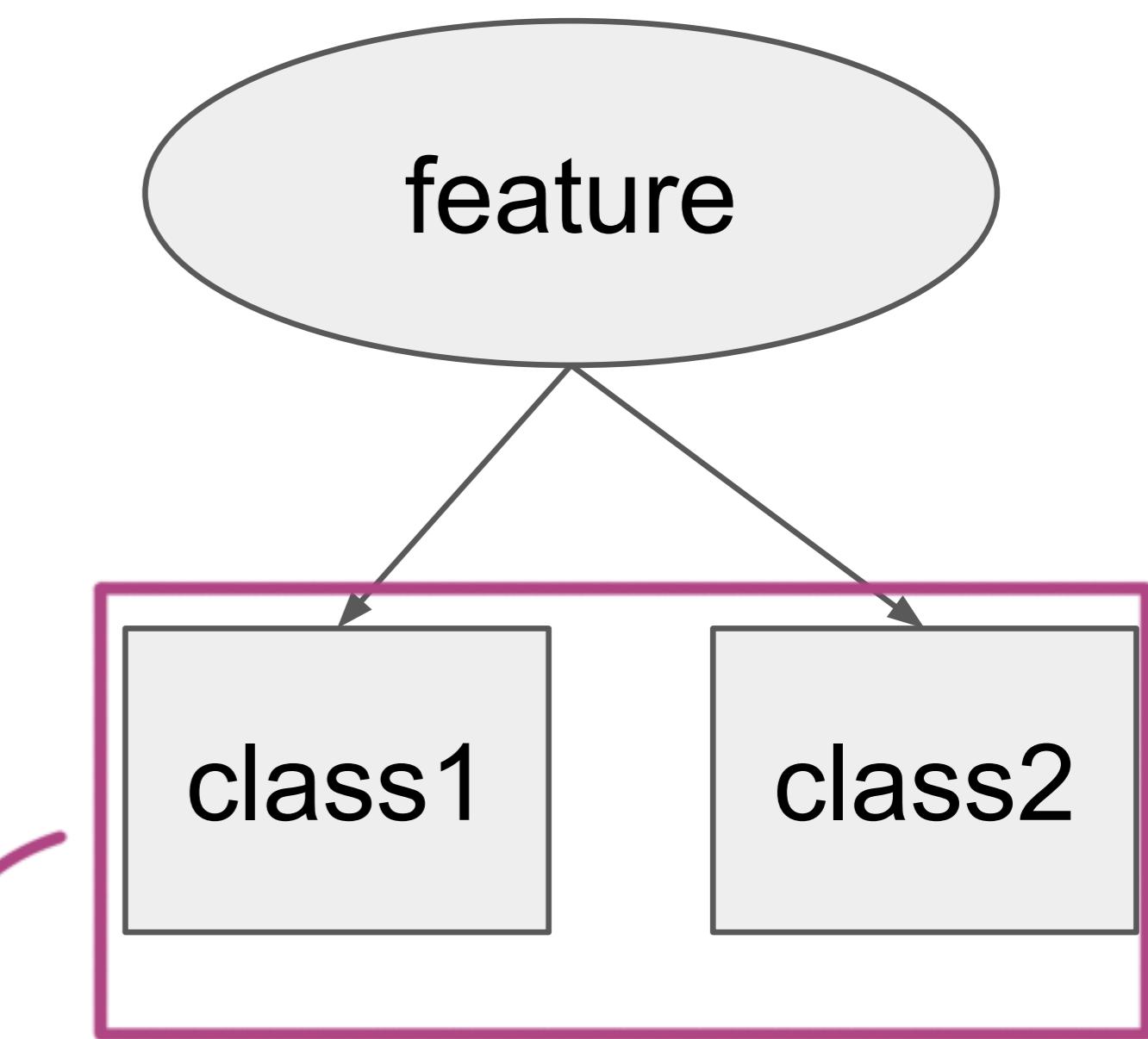
Figure 4.1 Differences between parallel and sequential ensembles: (1) base estimators in parallel ensembles are trained independently of each other, while in sequential ensembles, they are trained to improve on the predictions of the previous base estimator; (2) sequential ensembles typically use weak learners as base estimators.



Boosting

- Sequential ensemble methods such as boosting aim to combine several weak learners into a single strong learner.
 - Boosting literally aims to **boost the performance** of a collection of **weak learners** (very simple model that has performance slightly better than random chance)
- Decision trees are often used as base estimators for sequential ensembles. Boosting algorithms typically use **decision stumps**, or decision trees of depth 1
- Methods: AdaBoost, Gradient Boosting

Decision stumps



harusnya ada lg milik fitur
tp lgsg pake majority voting
aja buat pilih kelas

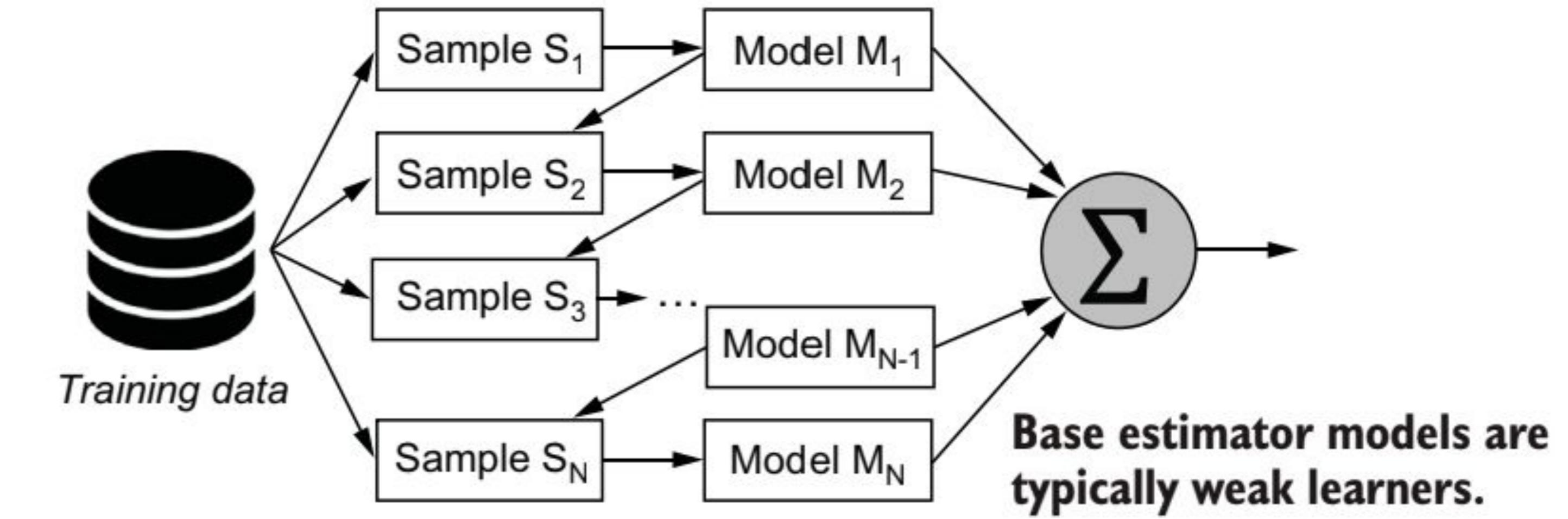
Adaptive Boosting (AdaBoost)

AdaBoost is an adaptive algorithm: at every iteration, it trains a new base estimator h_i that fixes the mistakes made by the previous base estimator h_{i-1} → prioritizes misclassified training examples.

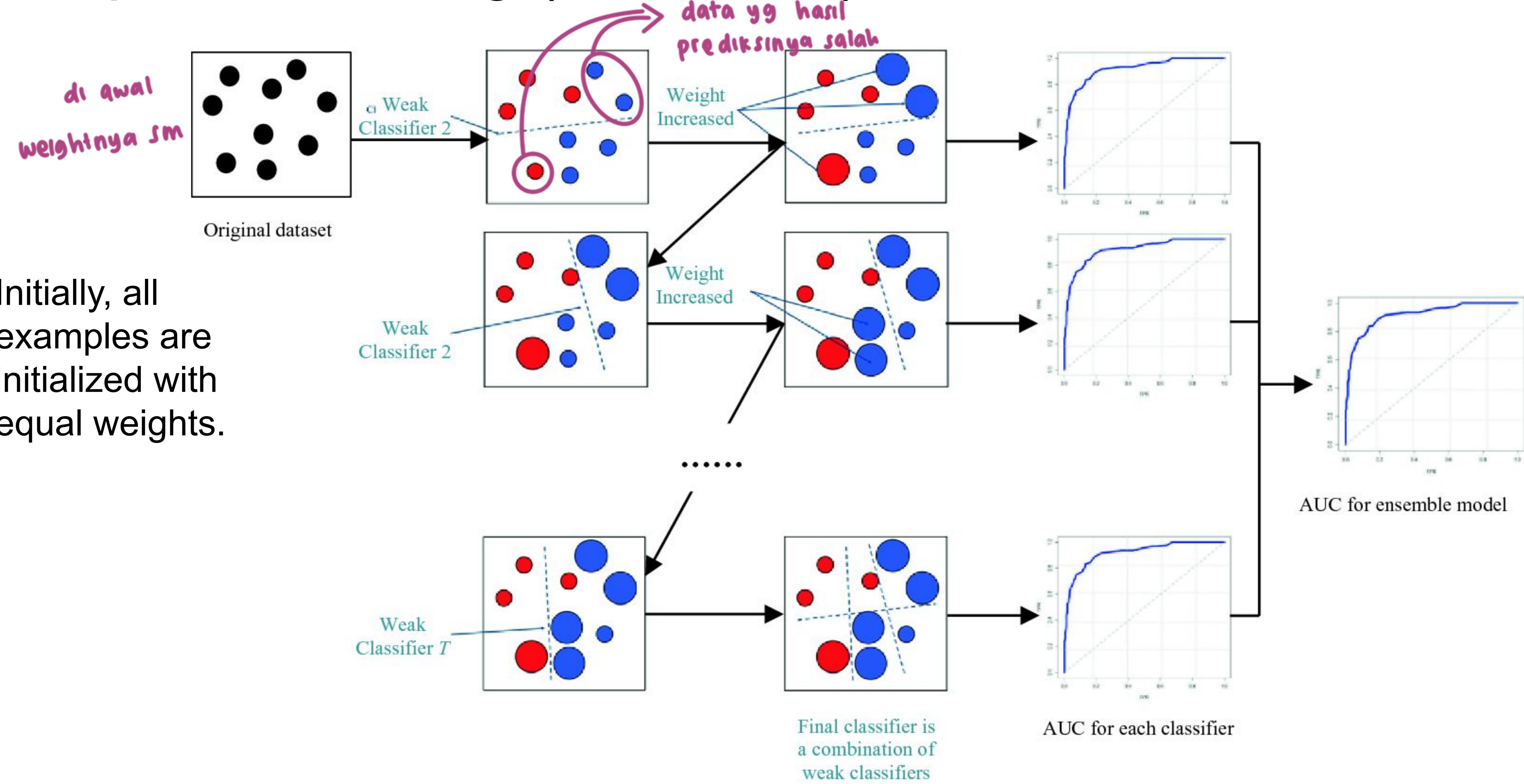
setiap data ada bobotnya

AdaBoost does this by maintaining weights over individual training examples (weights reflect the relative importance of training examples).

The final boosted classifier, H , combines the votes of each individual classifier, where the weight of each classifier's vote is a function of its accuracy. This forms a weighted linear combination of base estimators.



Adaptive Boosting (AdaBoost): Illustration



Instance Weights in Adaptive Boosting: Illustration

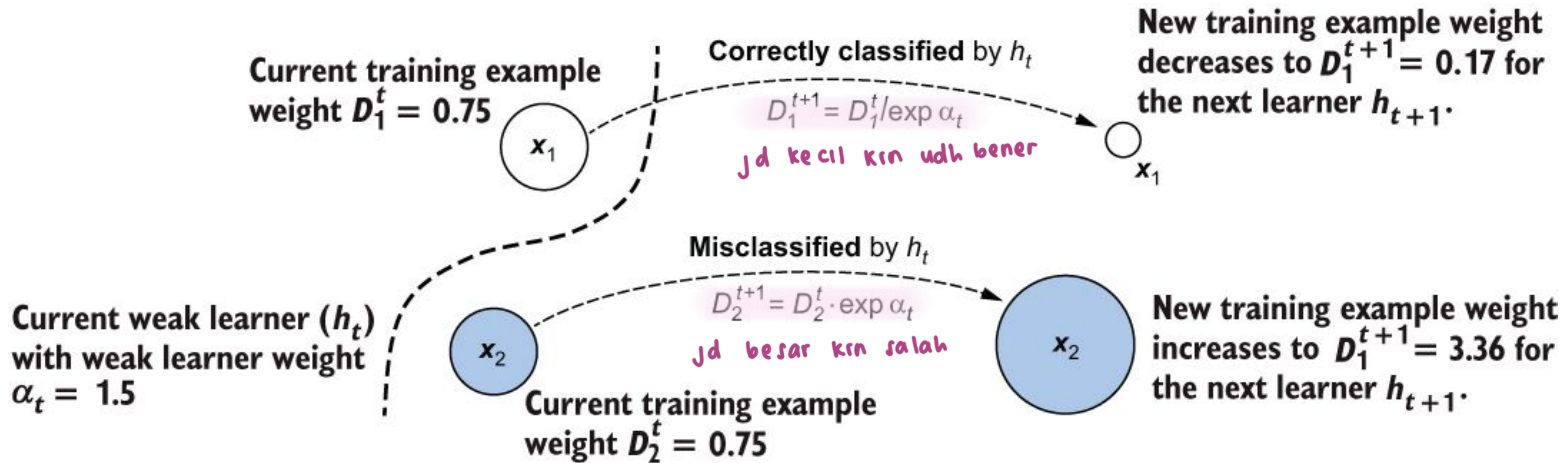


Figure 4.8 In iteration t , two training examples, x_1 and x_2 , have the same weights. x_1 is correctly classified, while x_2 is misclassified by the current base estimator h_t . As the goal in the next iteration is to learn a classifier h_{t+1} that fixes the mistakes of h_t , AdaBoost increases the weight of the misclassified example x_2 , while decreasing the weight of the correctly classified example x_1 . This allows the base-learning algorithm to prioritize x_2 during training in iteration $t+1$.

Instance Weights in Adaptive Boosting: Illustration

- Iteration t: $\alpha_t = 1.5$

Weight $D_1 = D_2 = 0.75$

$h_t(x_1) = c(x_1)$ #correctly classified, $c(x_1)$: label of x_1

$h_t(x_2) \neq c(x_2)$ #misclassified, $c(x_2)$: label of x_2

$$D_i^{t+1} = D_i^t \cdot \begin{cases} e^{\alpha_t}, & \text{if misclassified,} \\ e^{-\alpha_t}, & \text{if correctly classified} \end{cases}$$

- Update weights:

$$D_1 = D_1 / e^{\alpha_t} = 0.75 / e^{1.5} = 0.17$$

$$D_2 = D_2 \cdot e^{\alpha_t} = 0.75 \cdot e^{1.5} = 3.36$$

AdaBoost: Kunapuli (2023)

Kunapuli, G. (2023). *Ensemble methods for machine learning*. Simon and Schuster.

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
import numpy as np

def fit_boosting(X, y, n_estimators=10):
    n_samples, n_features = X.shape
    D = np.ones((n_samples, ))
    estimators = []
    for t in range(n_estimators):
        D = D / np.sum(D)           ← Nonnegative weights,
                                    ← initializes an empty ensemble
                                    ← so they sum to 1
        h = DecisionTreeClassifier(max_depth=1)
        h.fit(X, y, sample_weight=D) ← Normalizes the weights
                                    ← Trains a weak learner ( $h_t$ )
                                    ← with weighted examples
        ypred = h.predict(X)
        e = 1 - accuracy_score(y, ypred, sample_weight=D) ← Computes the training error ( $\epsilon_t$ )
                                                    ← and the weight ( $\alpha_t$ ) of the weak learner
        a = 0.5 * np.log((1 - e) / e) ← Updates the example weights:
                                    ← increase for misclassified
                                    ← examples, decrease for
                                    ← correctly classified examples
        m = (y == ypred) * 1 + (y != ypred) * -1
        D *= np.exp(-a * m)
        estimators.append((a, h))   ← Saves the weak
                                    ← learner and its weight
    return estimators
```

gapake sampling with replacement

gada normalisasi

1) Train a weak learner $h_t(\mathbf{x})$ using weighted dataset $\langle \mathbf{x}_i, y_i, D_i \rangle^*$

2) Compute training error e_t dari $h_t(\mathbf{x})$

3) Computer model weight α_t

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

4) Update weight

$$D_i^{t+1} = D_i^t \cdot \begin{cases} e^{\alpha_t}, & \text{if misclassified,} \\ e^{-\alpha_t}, & \text{if correctly classified} \end{cases}$$

Adaptive Boosting (AdaBoost) - Training: Han et al. (2022)

Algorithm: AdaBoost. A boosting algorithm—create an ensemble of classifiers. Each one gives a weighted vote.

Input:

- D , a set of d class-labeled training tuples;
- k , the number of rounds (one classifier is generated per round);
- a classification learning scheme.

Output: A composite model.

Method:

- (1) initialize the weight of each tuple in D to $1/d$;
- (2) **for** $i = 1$ to k **do** // for each round:
- (3) sample D with replacement according to the tuple weights to obtain D_i ;
- (4) use training set D_i to derive a model, M_i ;
- (5) compute $\text{error}(M_i)$, the error rate of M_i (Eq. (6.34))
- (6) **if** $\text{error}(M_i) > 0.5$ **then**
- (7) abort the loop;
- (8) **endif**
- (9) **for** each tuple in D that was correctly classified **do**
- (10) multiply the weight of the tuple by $\text{error}(M_i)/(1 - \text{error}(M_i))$; // update weights
- (11) normalize the weight of each tuple.
- (12) **endfor**

Initial weights

ada normalisasi

pakai model biasa, ga perlu yg bs handle weighted model

$$\text{error}(M_i) = \sum_{j=1}^d w_j \times \text{err}(X_j)$$

Update weights if correctly classified:
 $D_i^{t+1} = D_i^t \cdot \text{error}(M_i) / (1 - \text{error}(M_i))$

To normalize a weight, we multiply it by the sum of the old weights, divided by the sum of the new weights.

Adaptive Boosting (AdaBoost) - Inference: Han et al. (2022)

To use the ensemble to classify tuple, X :

- (1) initialize weight of each class to 0;
- (2) **for** $i = 1$ to k **do** // for each classifier:
- (3) $w_i = \log \frac{1 - \text{error}(M_i)}{\text{error}(M_i)}$; // weight of the classifier's vote
- (4) $c = M_i(X)$; // get class prediction for X from M_i
- (5) add w_i to the weight for class c
- (6) **endfor**
- (7) return the class with the largest weight.

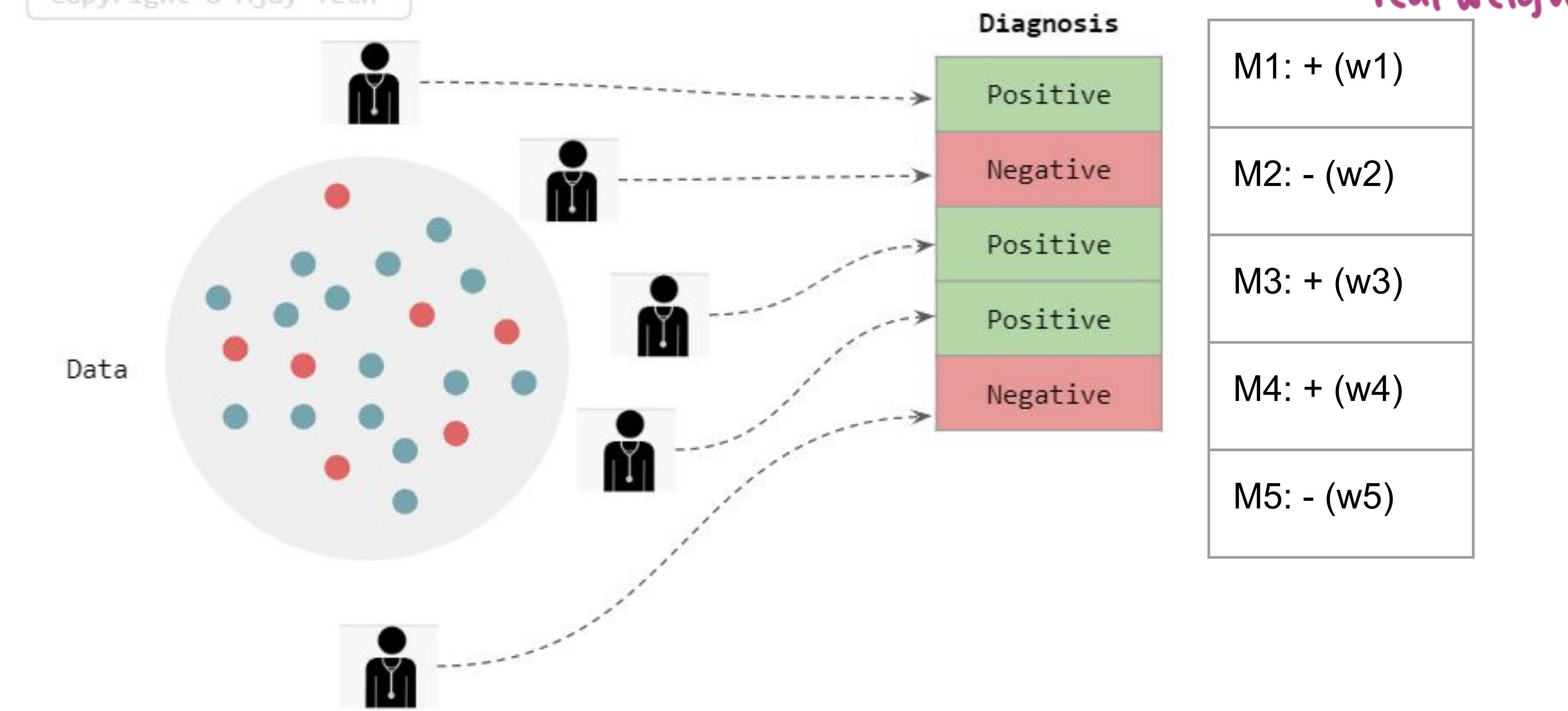
Model weights: $\alpha_t = \log[(1 - \text{error}(M_i)) / \text{error}(M_i)]$

ga efisien
hrinya nilai inferensi bisa dihitung saat training
kl ini pas train udh selesai br hitung nilai inferensi

AdaBoost Inference: Han et al. (2022)

Base Model	Error(Mi)	wi
M1	e1	$\log[(1-e1)/e1]$
M2	e2	$\log[(1-e2)/e2]$
...
Mn	en	$\log[(1-en)/en]$

Copyright © Ajay Tech



To use the ensemble to classify tuple, X :

- (1) initialize weight of each class to 0;
- (2) **for** $i = 1$ to k **do** // for each classifier:
- (3) $w_i = \log \frac{1 - \text{error}(M_i)}{\text{error}(M_i)}$; // weight of the classifier's vote
- (4) $c = M_i(X)$; // get class prediction for X from M_i
- (5) add w_i to the weight for class c
- (6) **endfor**
- (7) return the class with the largest weight.

$$w+ : w_1 + w_3 + w_4$$

$$w- : w_2 + w_5$$

Class: $\text{argmax} \{w+, w-\}$

AdaBoost Inference: Kunapuli (2023) → default pake ini

Base Model	Error(Mi)	ai
M1	e1	$a_1 = \log[(1-e_1)/e_1]$
M2	e2	$a_2 = \log[(1-e_2)/e_2]$
...
Mn	en	$a_n = \log[(1-e_n)/e_n]$

M1: + 1
M2: - 1
M3: + 1
M4: + 1
M5: - 1

$\text{pred} = a_1 * 1 + a_2 * -1 + a_3 * 1 + a_4 * 1 + a_5 * -1$

Class: $\text{sign}(\text{pred})$

Assumption: $\text{predict}(X)$ return +1 or -1

weight nya

direpresentasi sm +1 & -1

```
def predict_boosting(X, estimators):
    pred = np.zeros((X.shape[0], ))
    ← Initialize all the predictions to 0
    for a, h in estimators:
        pred += a * h.predict(X)
        ← Makes weighted prediction for each example
    y = np.sign(pred)
    ← Converts weighted predictions to -1/1 labels
    return y
```

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

AdaBoost vs Gradient Boosting

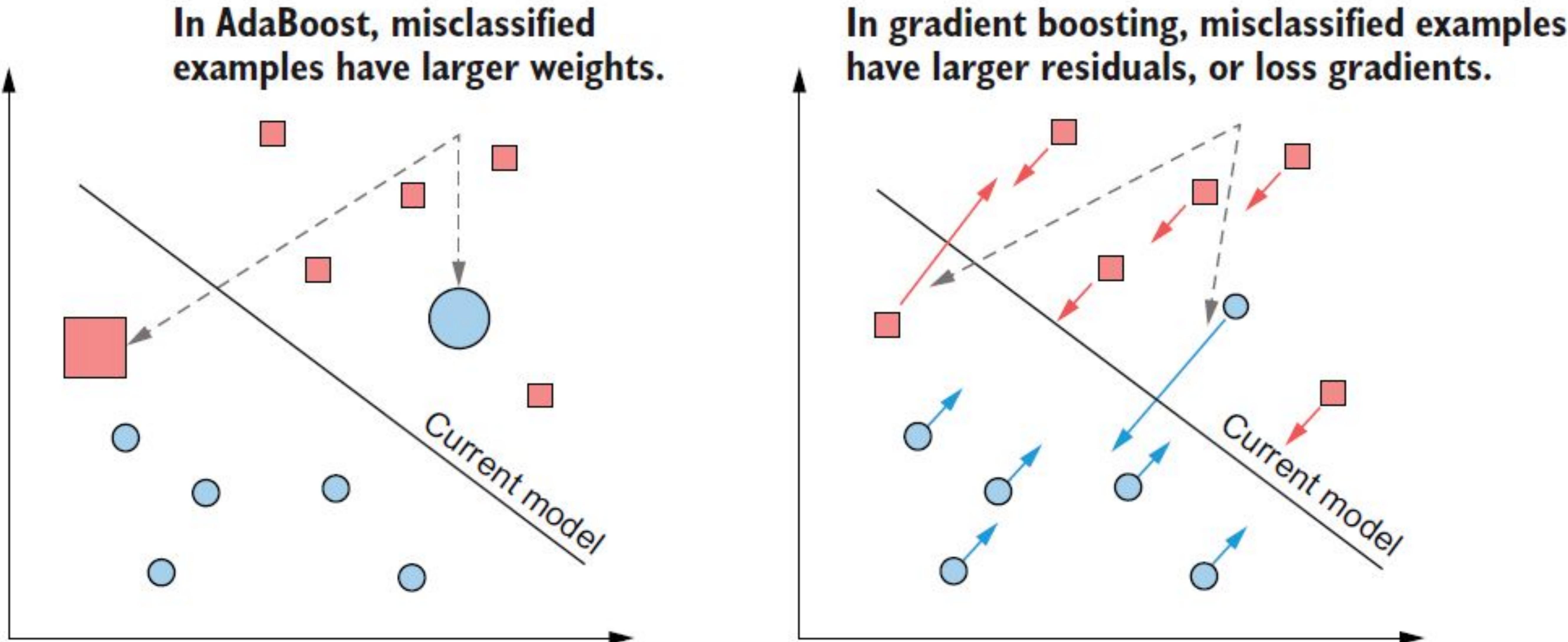
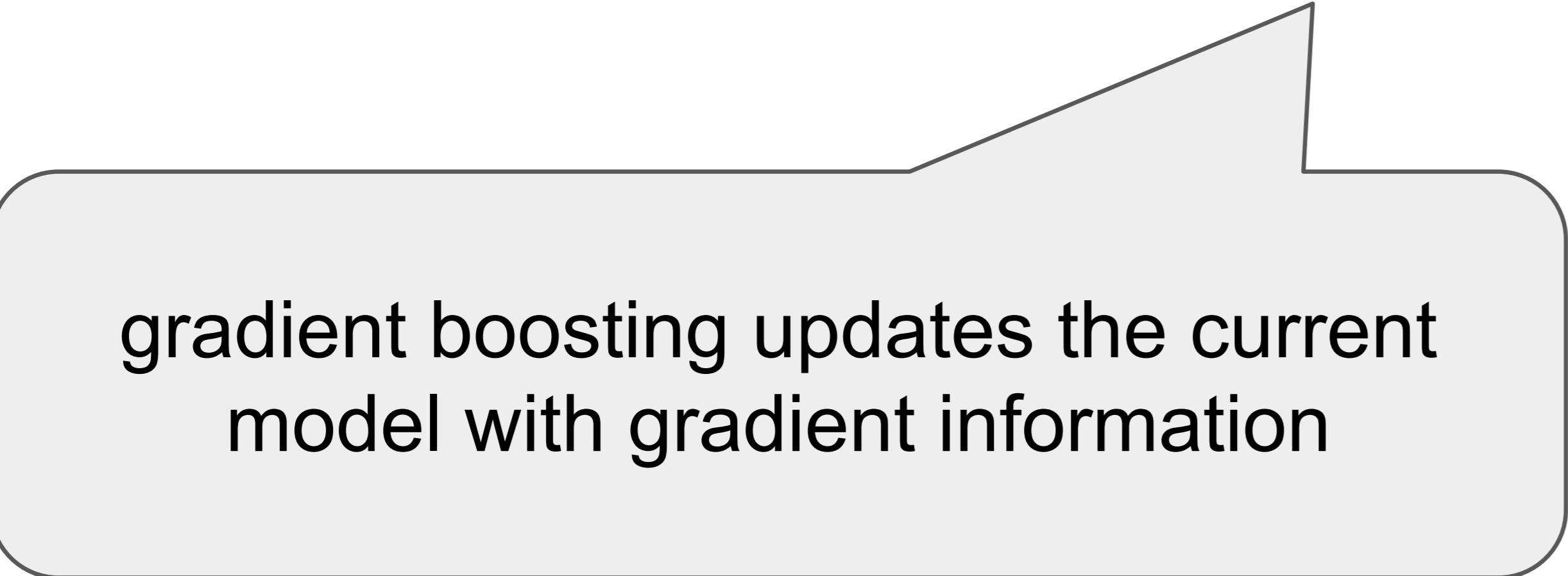


Figure 5.8 Comparing AdaBoost (left) to gradient boosting (right). Both approaches train weak estimators that improve classification performance on misclassified examples. AdaBoost uses weights, with misclassified examples being assigned higher weights. Gradient boosting uses residuals, with misclassified examples having higher residuals. The residuals are nothing but negative loss gradients.

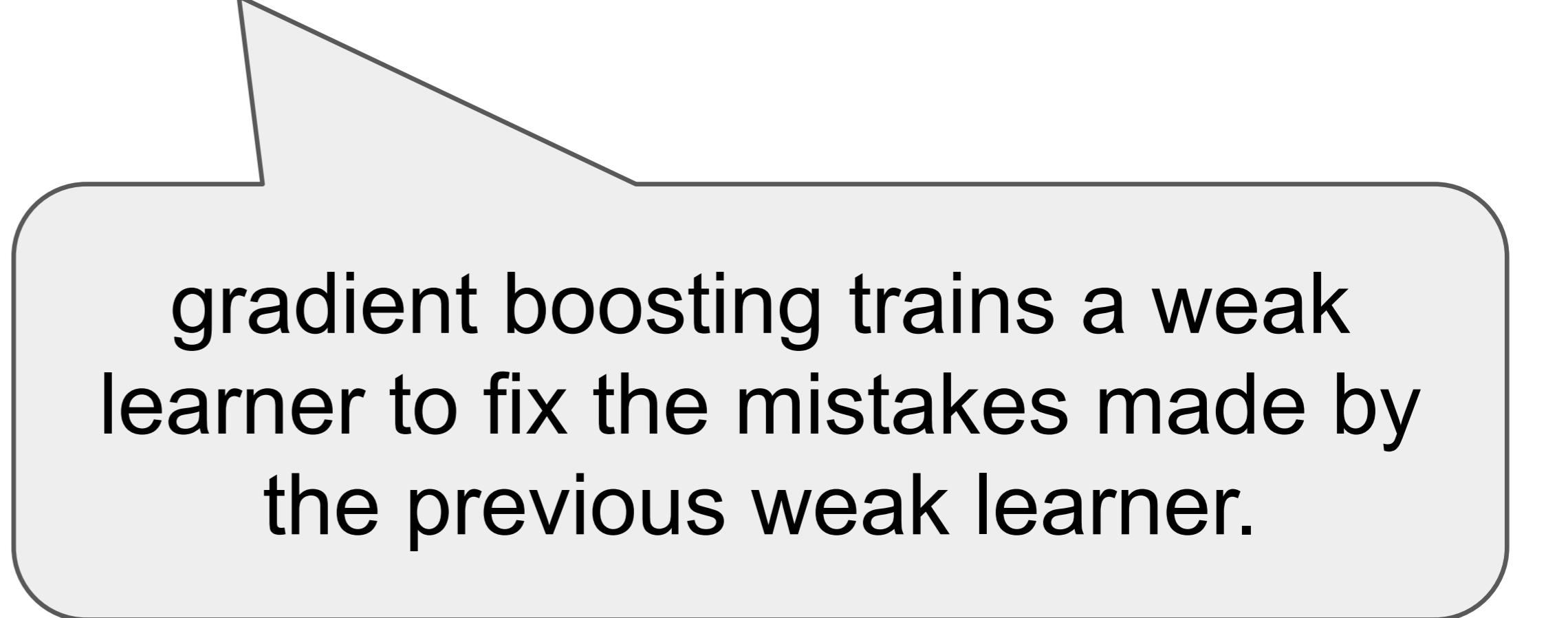
Gradient Boosting

=

Gradient Descent + Boosting



gradient boosting updates the current model with gradient information



gradient boosting trains a weak learner to fix the mistakes made by the previous weak learner.

Hill Climbing Steepest Ascent
naik

local search
value (state) \rightarrow diskrit
better neighbor

Gradient Descent
turun

state continuous
linear regression
back propagation

bisa secara iteratif pake Loss (\vec{b}) = $\sum (y_i - \hat{y}_i)^2$
bisa cari successor based on turunan
partial dgn bobot $\langle b_0, b_1, \dots, b_n \rangle$

Least square (LS)

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

$$\hat{y} = b_0 + b_1 x_1 + \dots + b_n x_n$$

$$\langle b_0, b_1, \dots, b_n \rangle$$

iterative

\rightarrow buat koefisien gradien (ga iteratif)

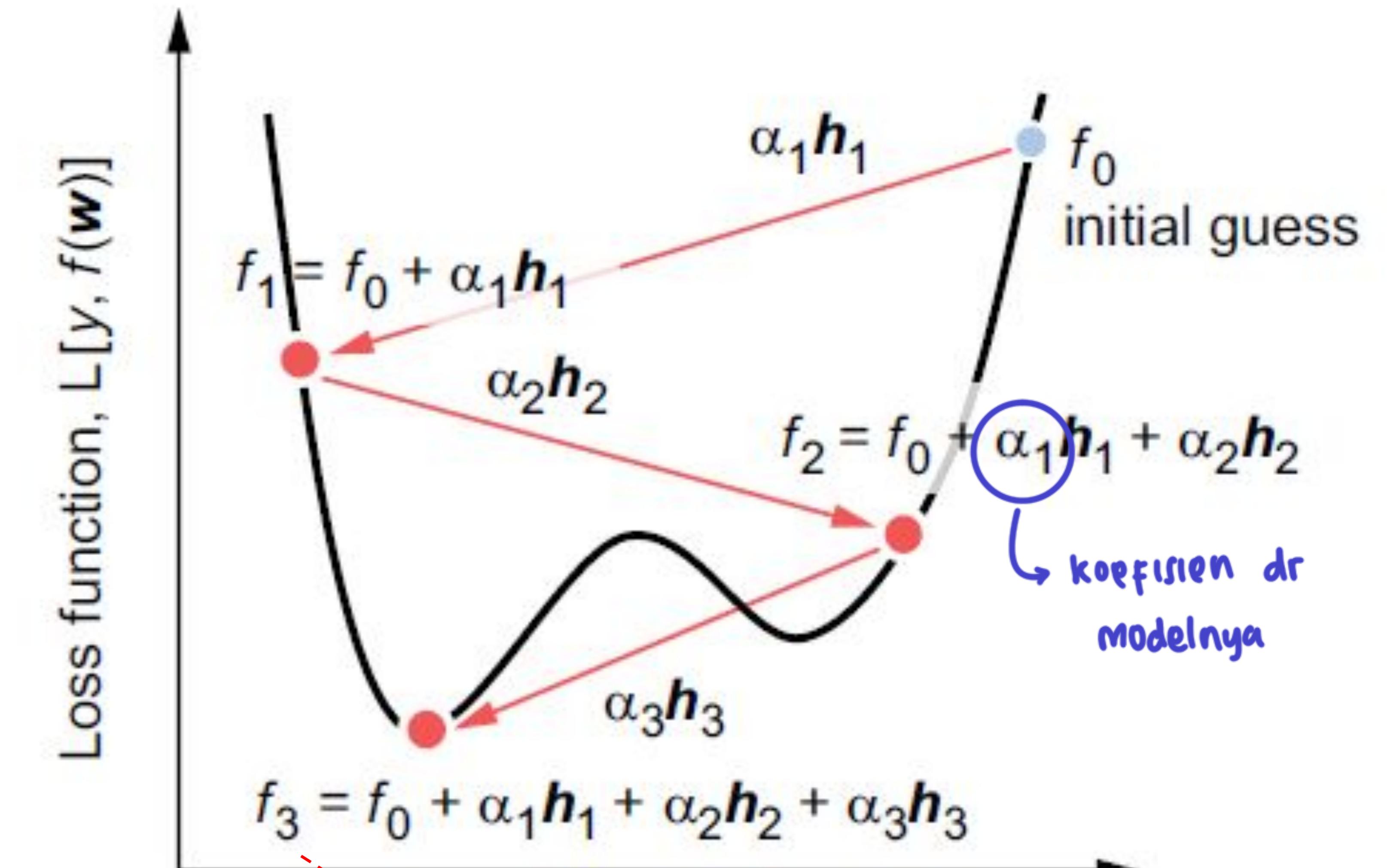
Gradient Boosting

Each iteration additively updates the current model with the new approximate weak gradient estimate (the regression tree, h_t).

$$\text{new model} = \text{old model} + (\text{step length}) * (\text{update direction})$$

$$F_{t+1}(x) = F_t(x) + \alpha_t \cdot h_t(x)$$

$$F_{t+1}(x) = F_0(x) + \alpha_1 \cdot h_1(x) + \alpha_2 \cdot h_2(x) + \dots + \alpha_{t-1} \cdot h_{t-1}(x) + \alpha_t \cdot h_t(x)$$



Set of all possible candidate models, f

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

Gradient Boosting: Residual

- residual(y, \hat{y})=residual($y, h(x)$)= $y-h(x)$
- Squared loss function (SSE) to model h :

$$f_{\text{loss}}(y, \hat{y}) = (y - h(x))^2 / 2$$

$$\text{gradient}(y, \hat{y}) = \partial f_{\text{loss}} / \partial h(y, h(x)) = -(y - h(x))$$

$$\text{gradient}(\text{true, predicted}) = \frac{\partial f_{\text{loss}}}{\partial h}(y, h(x)) = -(y - h(x))$$

- Residual is the negative gradient:

$$\text{residual}(y, \hat{y}) = -\partial f_{\text{loss}} / \partial h = y - h(x) \# \text{ for squared loss}$$

$$\begin{aligned} f_{\text{loss}} &= \frac{1}{2} (y - h)^2 \\ \frac{\partial f_{\text{loss}}}{\partial h} &= \frac{1}{2} \cdot 2 \cdot (y - h) (-1) \\ -\frac{\partial f_{\text{loss}}}{\partial h} &= y - h \end{aligned}$$

residual

Gradient Boosting: Basic Algorithm

initialize: $F = f_0$, some constant value

for $t = 1$ to T :

1. compute the negative residuals for each example, $r_i^t = -\frac{\partial L}{\partial F}(x_i)$
2. fit a weak decision tree regressor $h_t(\mathbf{x})$ using the training set $(x_i, r_i)_{i=1}^n$
3. compute the step length (α_t) using line search
4. update the model: $F_t = F + \alpha_t \cdot h_t(\mathbf{x})$

Gradient Boosting vs

Gradient Descent

- instead of using the **negative gradient**, we use an **approximate gradient** trained on the negative residuals
- instead of checking for convergence, the algorithm **terminates** after a finite, maximum number of iterations T.

Gradient Boosting: Training

```
from scipy.optimize import minimize_scalar
from sklearn.tree import DecisionTreeRegressor

def fit_gradient_boosting(X, y, n_estimators=10):
    n_samples, n_features = X.shape
    n_estimators = 10
    estimators = []
    F = np.full((n_samples, ), 0.0)

    for t in range(n_estimators):
        residuals = y - F
        h = DecisionTreeRegressor(max_depth=1)
        h.fit(X, residuals)

        hreg = h.predict(X)
        loss = lambda a: np.linalg.norm(
            y - (F + a * hreg))**2
        step = minimize_scalar(loss, method='golden')
        a = step.x

        F += a * hreg
        estimators.append((a, h))

    return estimators
```

Fits weak regression tree (h_t) to the examples and residuals

Gets dimensions of the data set

Initializes an empty ensemble

Predicts the ensemble on the training set

Computes residuals as negative gradients of the squared loss

Gets predictions of the weak learner, h_t

Sets up the line search problem

Finds the best step length using the golden section search

Updates the ensemble predictions

Updates the ensemble

Gradient Boosting: Inference

```
def predict_gradient_boosting(X, estimators):
    pred = np.zeros( (X.shape[0], ) )           ← Initializes all the predictions to 0

    for a, h in estimators:
        pred += a * h.predict(X)                ← Aggregates individual predictions from each regressor

    y = np.sign(pred)                         ← Converts weighted predictions to -1/1 labels
    return y
```

Gradient Boosting: Han et al. (2022)

Algorithm: Gradient Tree Boosting for Regression.

Input:

- D , a set of n training tuples $\{(x_1, y_1), \dots, (x_n, y_n)\}$, where x_i is the attribute vector of the i th training tuple and y_i is its true target output value;
- k , the number of rounds (one base regression model is generated per round);
- a differential loss function $\text{Loss} = \sum_{i=1}^n L(y_i, F(x_i))$.

Output:

A composite regression model $F(x)$.

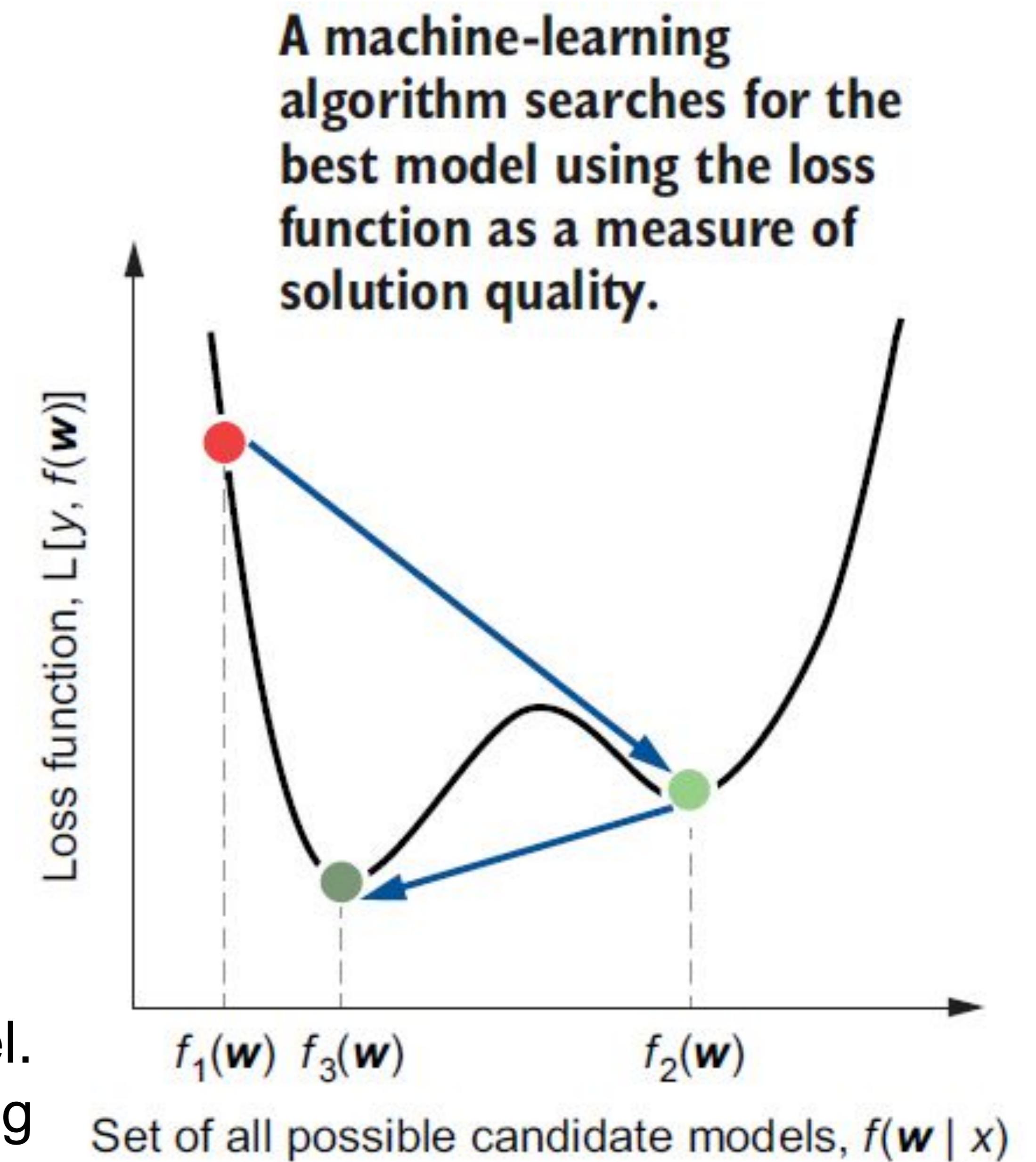
Method:

- (1) initialize the regression model $F(x) = \frac{\sum_{i=1}^n y_i}{n}$; Outputs a constant (i.e., the average output of all training tuples).
 - (2) **for** $t = 1$ to k **do** // construct a new weak learner $M_t(x)$ for each round:
 - (3) **for** $i = 1$ to n //each training tuple:
 - (4) calculate $\hat{y}_i = F(x_i)$; //predicted value by the current model $F(x)$
 - (5) calculate the negative gradient $r_i = -\frac{\partial L(y_i, \hat{y}_i)}{\partial \hat{y}_i}$;
 - (6) **endfor**
 - (7) fit a regression tree model $M_t(x)$ for the training set $\{(x_1, r_1), \dots, (x_n, r_n)\}$; the negative gradient r_i is treated as the targeted output value of the i th training tuple.
 - (8) update the composite regression model $F(x) \leftarrow F(x) + M_t(x)$.
 - (9) **endfor**
- $L(y_i, F(x_i)) = \frac{1}{2}(y_i - \hat{y}_i)^2$
- $r_i = y_i - \hat{y}_i$

Loss Function: Additional Slides

- Loss functions explicitly measure the fit of a model on a data set. Most often, we measure loss with respect to the true labels, by quantifying the error between the predicted and true labels.
- Given a loss function, training is the search for the optimal model that minimizes the loss, as illustrated in this figure.

An optimization procedure for finding the best model. Machine-learning algorithms search for the best model among all possible candidate models. The optimization procedure sequentially identifies increasingly better models f_1 , f_2 , and the final model, f_3 .



Gradient-based Optimization: Additional Slides

Many methods attempt to **use the gradient** of the landscape to find an optimum value. The gradient of the objective function is a vector ∇f that gives the magnitude and direction of the steepest slope.

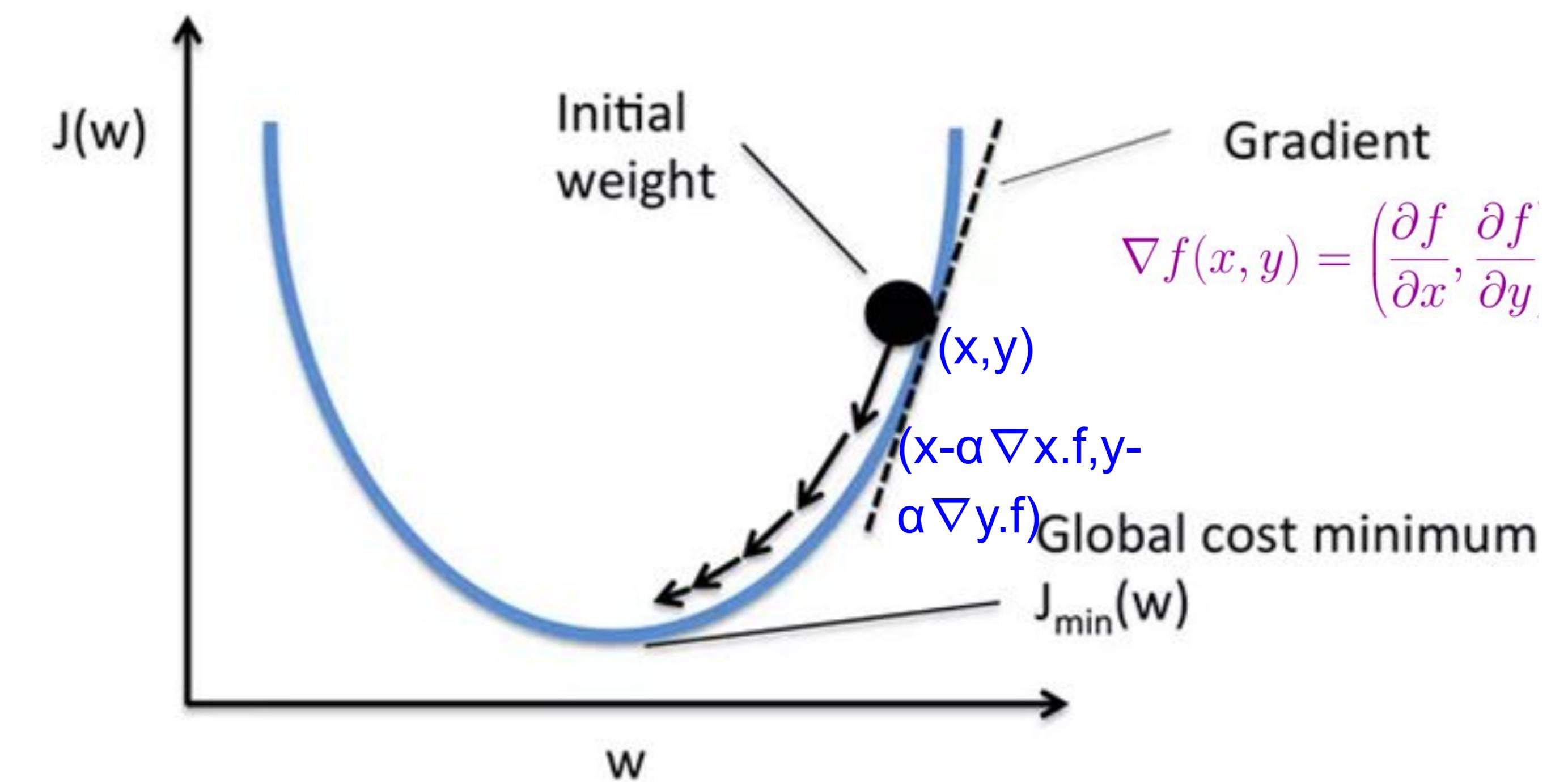
Given objective function $f(x,y)$ to be minimized, and **current state is (x,y)** , gradient is a vector:

$$\nabla f(x, y) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

Use gradient to set successor state:

$$x \leftarrow x - \alpha \nabla_x f$$

$$y \leftarrow y - \alpha \nabla_y f$$



<https://medium.com/data-science-group-iitr/loss-functions-and-optimization-algorithms-demystified-bb92daff331c>

The basic problem is that if α is too small, too many steps are needed; if α is too large, the search could overshoot the maximum.

Eksplorasi Mandiri

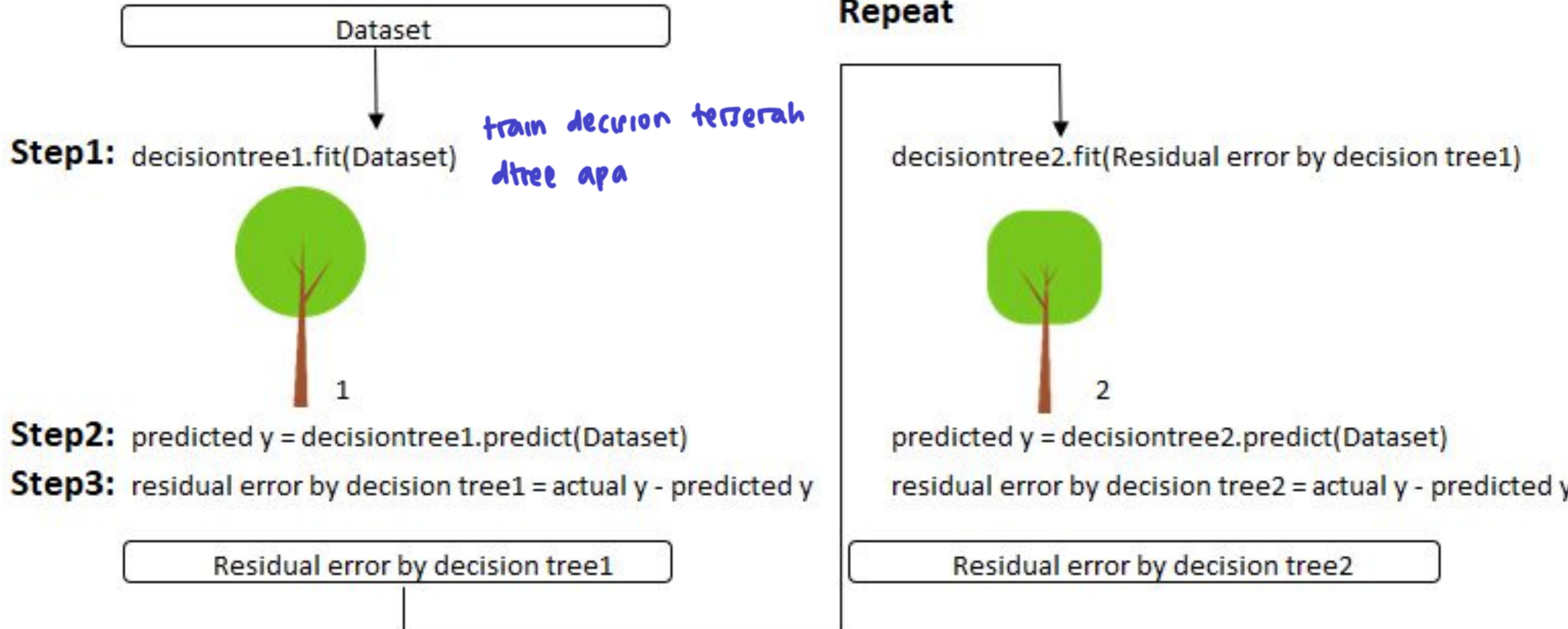
Kerjain !!!

sp tau ditanya di ujian

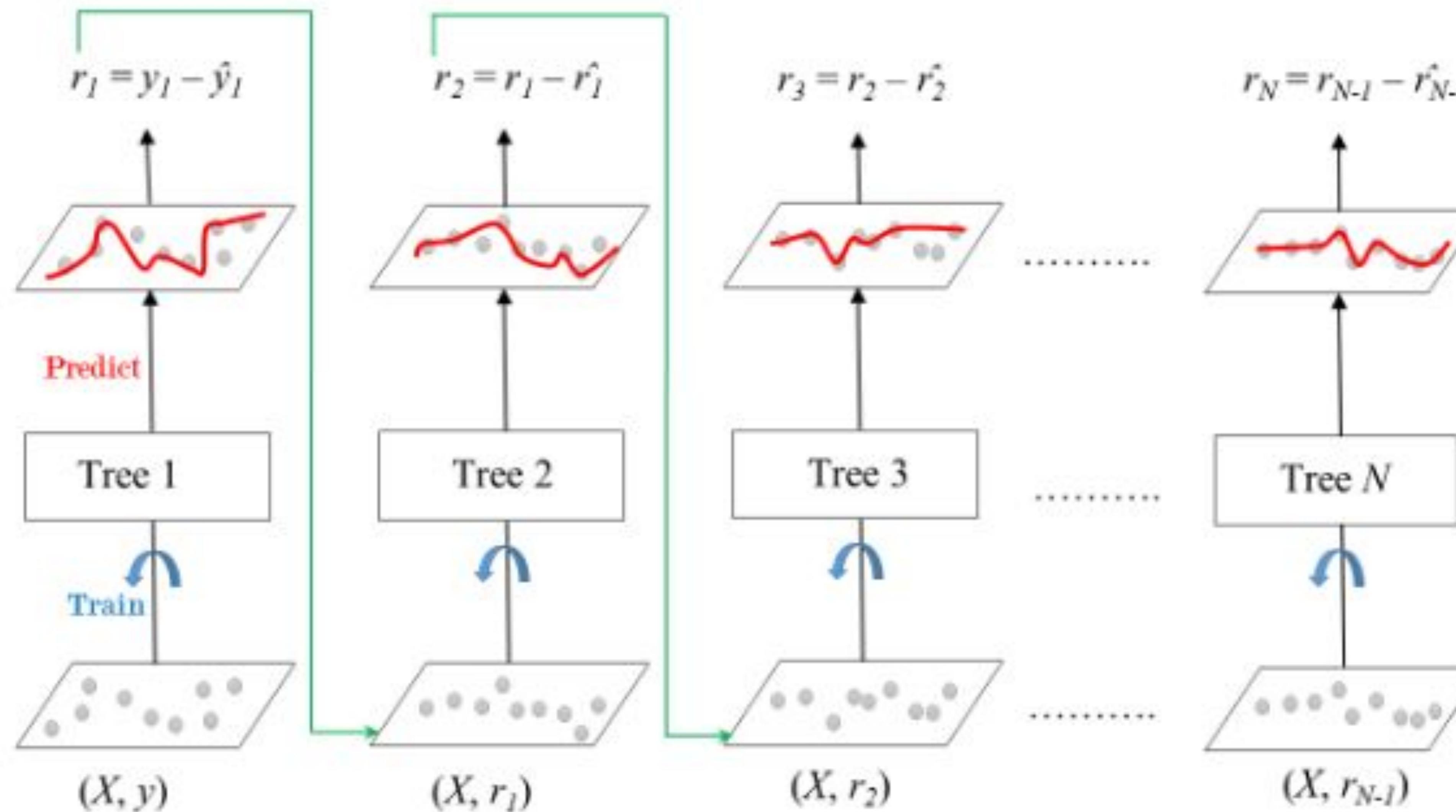
- Eksplorasi library `RandomForestClassifier`, `BaggingClassifier`, `AdaBoostClassifier`, `LightGBM` dari library `sklearn.ensemble`
- Ikuti hands-on pada subbab **2.5 Case study: Breast cancer diagnosis** buku teks Kunapuli, G. (2023). *Ensemble methods for machine learning*. Simon and Schuster

Gradient Boosting: Use Residual Error

decision start → depth nya 1

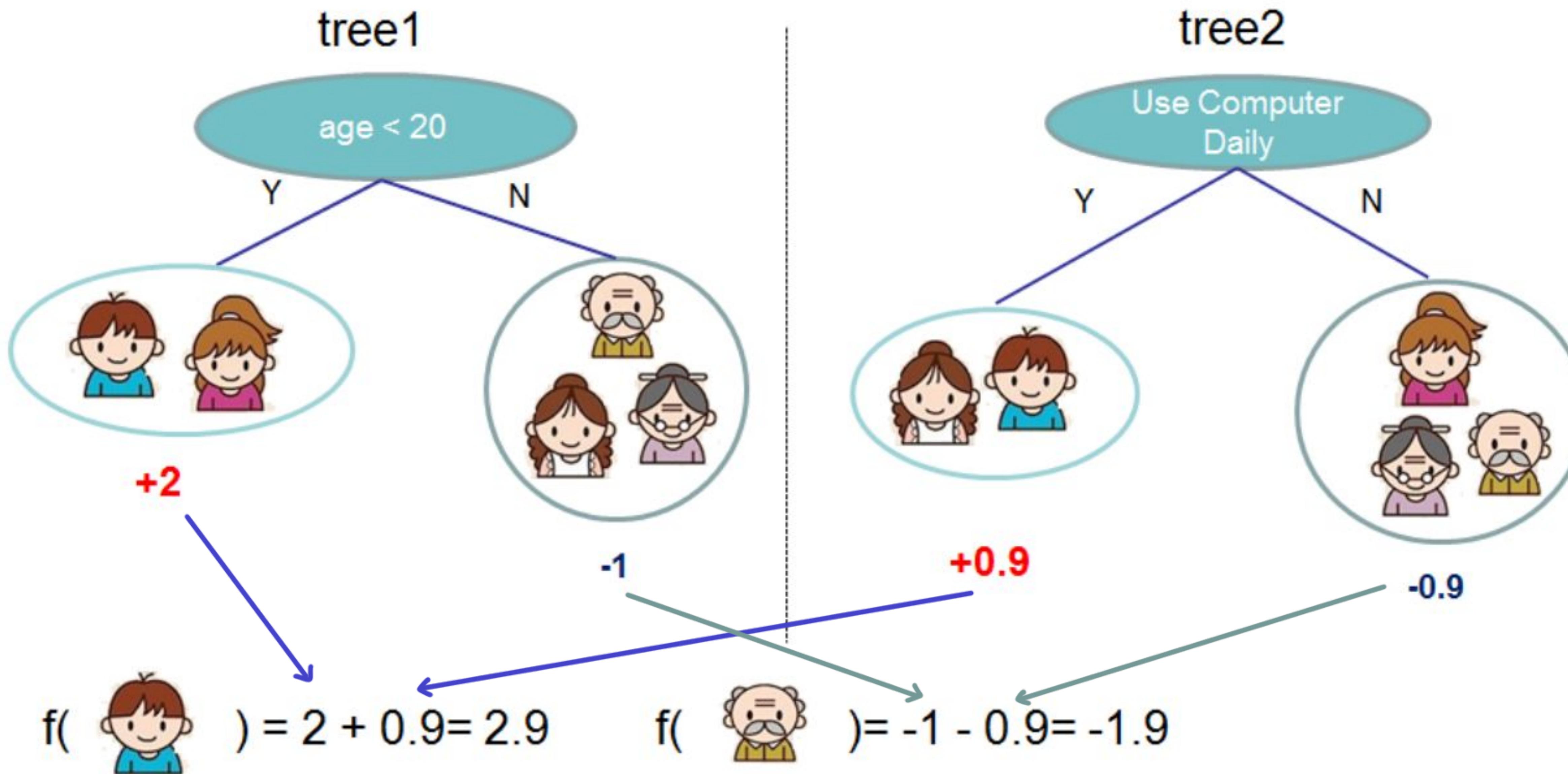


Gradient Boosting: Use Residual Error



$$y(\text{pred}) = y_1 + (\eta * r_1) + (\eta * r_2) + \dots + (\eta * r_N)$$

Gradient Boosting Prediction: Example



<https://xgboost.readthedocs.io/en/latest/tutorials/model.html>

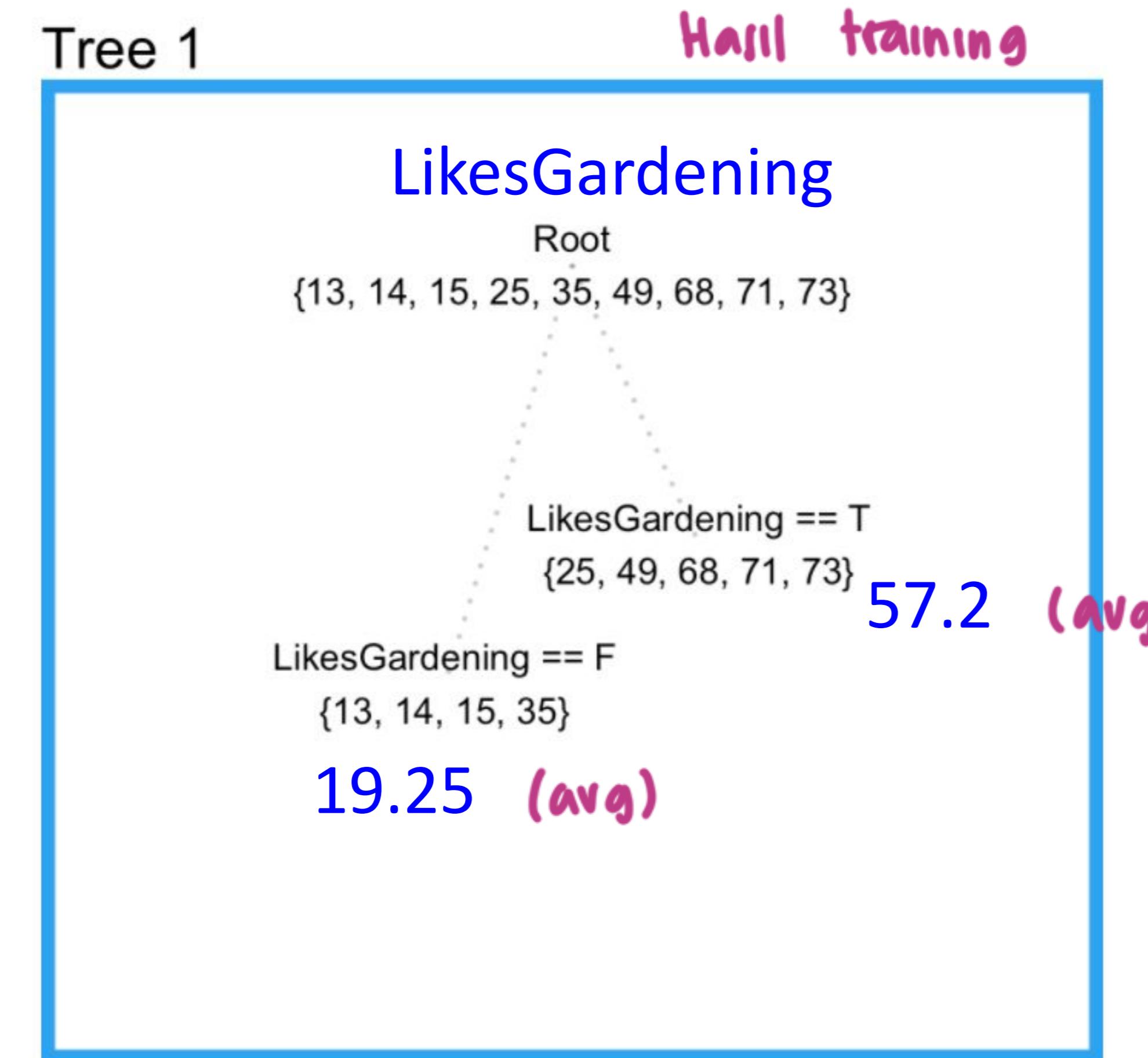
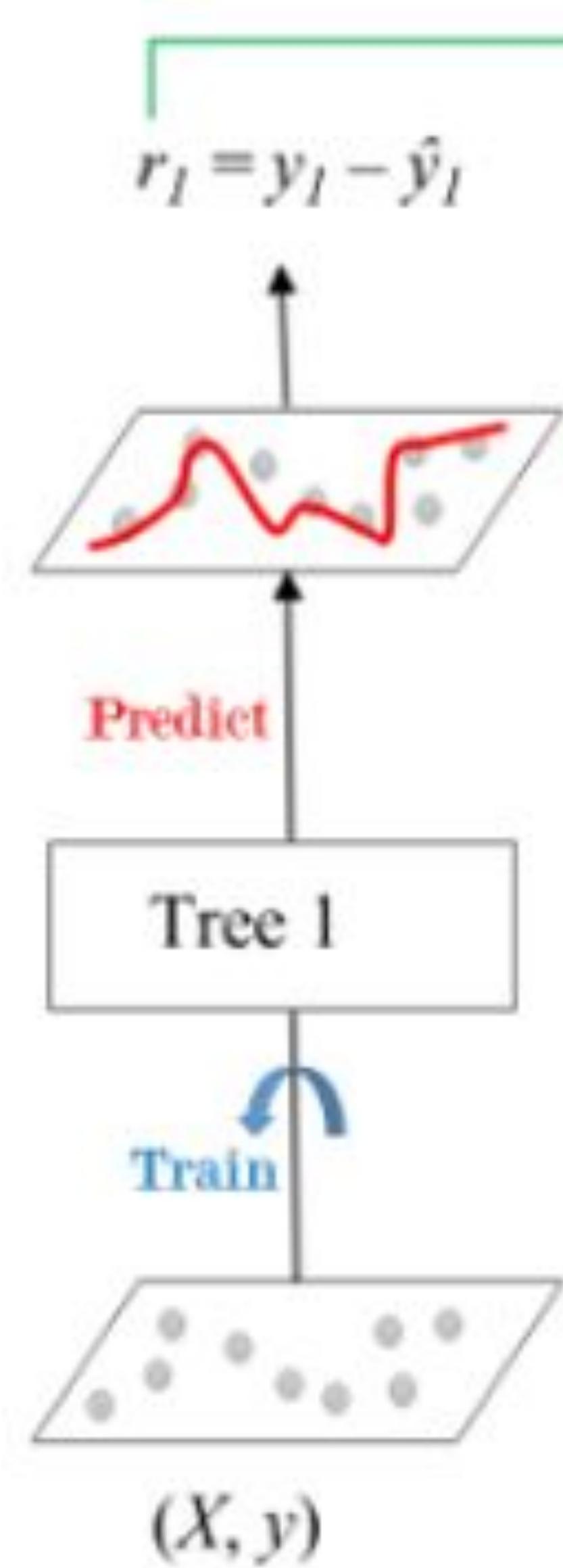
buat contoh $\alpha_1 h_1 + \alpha_2 h_2$

XGBoost Example: Predict Person's Age

→ target (y)

PersonID	Age	LikesGardening	PlaysVideoGames	LikesHats
1	13	FALSE	TRUE	TRUE
2	14	FALSE	TRUE	FALSE
3	15	FALSE	TRUE	FALSE
4	25	TRUE	TRUE	TRUE
5	35	FALSE	TRUE	TRUE
6	49	TRUE	FALSE	FALSE
7	68	TRUE	TRUE	TRUE
8	71	TRUE	FALSE	FALSE
9	73	TRUE	FALSE	TRUE

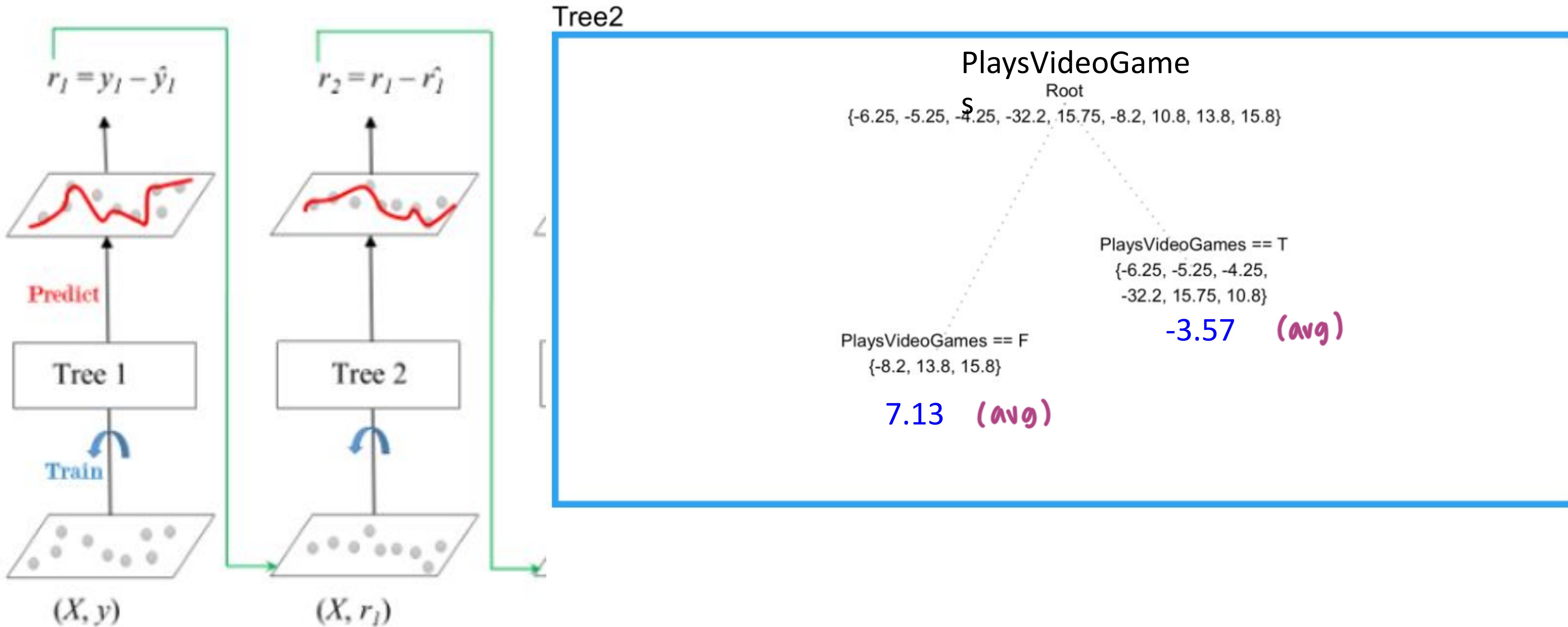
XGBoost: Initialization $F_0(x) = \text{Tree1}$, Predict, & r_1



age - Tree1 res next target
take ini, gapake age lg

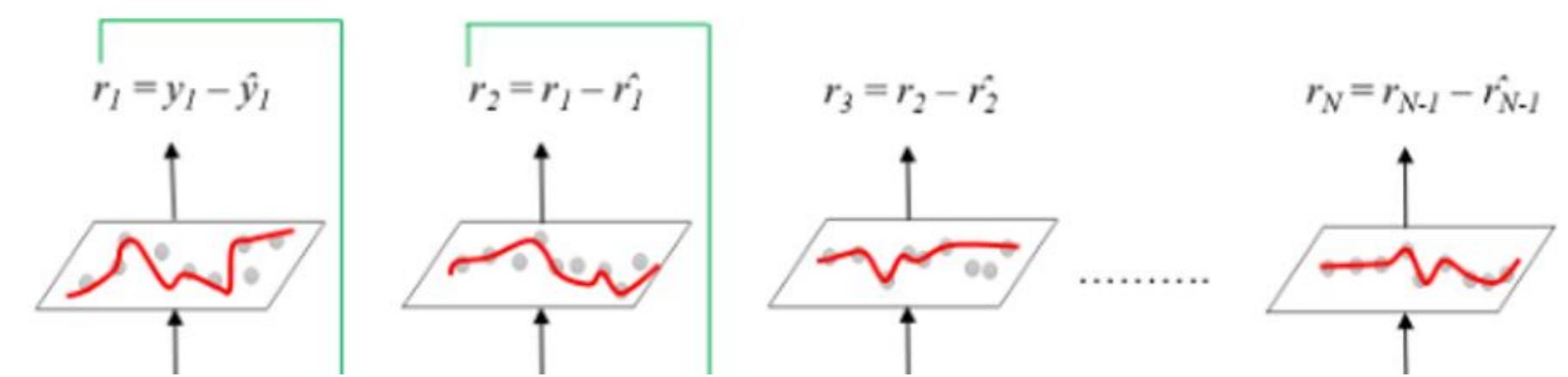
PersonID	Age	Tree1 Prediction	Tree1 Residual
1	13	19.25	-6.25
2	14	19.25	-5.25
3	15	19.25	-4.25
4	25	57.20	-32.20
5	35	19.25	15.75
6	49	57.20	-8.20
7	68	57.20	10.80
8	71	57.20	13.80
9	73	57.20	15.80

XGBoost: Tree2 □ Combined Prediction



XGBoost: *defaultnya CART*
Classification & Regression Tree

$F_0(x) = \text{Tree1.predict} + \text{Tree2.predict}$



PersonID	Age	Tree1 Prediction	Tree1 Residual	Tree2 Prediction	Combined Prediction	Final Residual
1	13	19.25	-6.25	-3.57	15.68	2.68
2	14	19.25	-5.25	-3.57	15.68	1.68
3	15	19.25	-4.25	-3.57	15.68	0.68
4	25	57.20	-32.20	-3.57	53.63	28.63
5	35	19.25	15.75	-3.57	15.68	-19.32
6	49	57.20	-8.20	7.13	64.33	15.33
7	68	57.20	10.80	-3.57	53.63	-14.37
8	71	57.20	13.80	7.13	64.33	-6.67
9	73	57.20	15.80	7.13	64.33	-8.67

Questions ?

contoh soal

1. (Bobot 15) Diberikan prediksi 5 base-classifier (Decision Tree) dari model ensemble sebagai berikut. Jelaskanlah penentuan hasil prediksi model ensemble jika menggunakan skema:
I
(a) bagging; (b) Random Forest; (c) AdaBoost.
Jika diperlukan bobot model dalam perhitungan, gunakanlah bobot berikut ini terurut model 1..5 yaitu 0.15, 0.25, 0.3, 0.25, 0.05. Jawaban prediksi tanpa penjelasan caranya tidak akan dinilai.

Model 1	Model 2	Model 3	Model 4	Model 5
Positif	Positif	Negatif	Negatif	Positif

Jawab:

Bagging: argmax{positif: 3, negatif: 2} = positif

RandomForest: argmax{positif: 3, negatif: 2} = positif

AdaBoost: positif 0.15 => positif 0.40 => positif 0.40, negatif 0.3 => positif 0.4, negatif 0.55 => positif 0.45, negatif 0.55. Hasil: negatif

CPMK 1: Mampu melakukan training dalam pembangunan model pembelajaran mesin dengan berbagai teknik supervised learning.

2. (Bobot 10) Diberikan dataset dengan 10 instances 2 atribut. hasil pemodelan SVM

6. (Bobot 10) Terdefinisi model Random Forest 2 base-classifier (cls1 dan cls2) untuk klasifikasi 3 kelas c1-c2-c3. Pada proses klasifikasi suatu input data, cls 1 berhenti pada suatu daun dengan rasio instances pada daun tersebut adalah $c_1:c_2:c_3=3:0:2$, sedangkan cls2 berhenti pada daun dengan rasio instances $c_1:c_2:c_3=4:3:3$. Jelaskanlah hasil hard-classification dan soft-classification sebagai hasil inferensi Random Forest tersebut.

Jawab:

Hard classification:



- Cls1: c1 ($c_1:c_2:c_3=3:0:2$)
- Cls2: c1 ($c_1:c_2:c_3=4:3:3$)
- Hard classification (majority voting): c1

Soft classification:

- Cls1: [0.6,0,0.4]
- Cls2: [0.4,0.3,0.3]
- Soft classification: [0.5,0.15,0.35]

Table 5.1 Comparing AdaBoost, gradient descent, and gradient boosting

Algorithm	Loss function	Base-learning algorithm	Update direction $h_t(x)$	Step length α_t
AdaBoost for classification	Exponential	Classification with weighted examples	Weak classifier	Computed in closed form
Gradient boosting	User-specified	Regression with examples and residuals	Weak regressor	Line search

$$\alpha_t = \frac{1}{2} \log \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$$

Gradient Boosting: alpha value

```
from scipy.optimize import minimize_scalar
from sklearn.tree import DecisionTreeRegressor

def fit_gradient_boosting(X, y, n_estimators=10):
    n_samples, n_features = X.shape
    n_estimators = 10
    estimators = []
    F = np.full((n_samples, ), 0.0)

    for t in range(n_estimators):
        residuals = y - F
        h = DecisionTreeRegressor(max_depth=1)
        h.fit(X, residuals)
        hreg = h.predict(X)
        loss = lambda a: np.linalg.norm(
            y - (F + a * hreg))**2
        step = minimize_scalar(loss, method='golden')
        a = step.x

        F += a * hreg
        estimators.append((a, h))

    return estimators
```

Fits weak regression tree (h_t) to the examples and residuals

Updates the ensemble predictions

Gets dimensions of the data set

Initializes an empty ensemble

Predicts the ensemble on the training set

Computes residuals as negative gradients of the squared loss

Gets predictions of the weak learner, h_t

Sets up the line search problem

Finds the best step length using the golden section search

Updates the ensemble

In [31]:

```
1 from scipy.optimize import minimize_scalar
2 def f(x):
3     return 2*(x**2)+4*x-6
4 result = minimize_scalar(f)
5 x_min = result.x
6 x_min #x=-b/2a=-4/4=-1
```

Out[31]: -0.999999851900002