



Lecture 19: Ensemble Learning

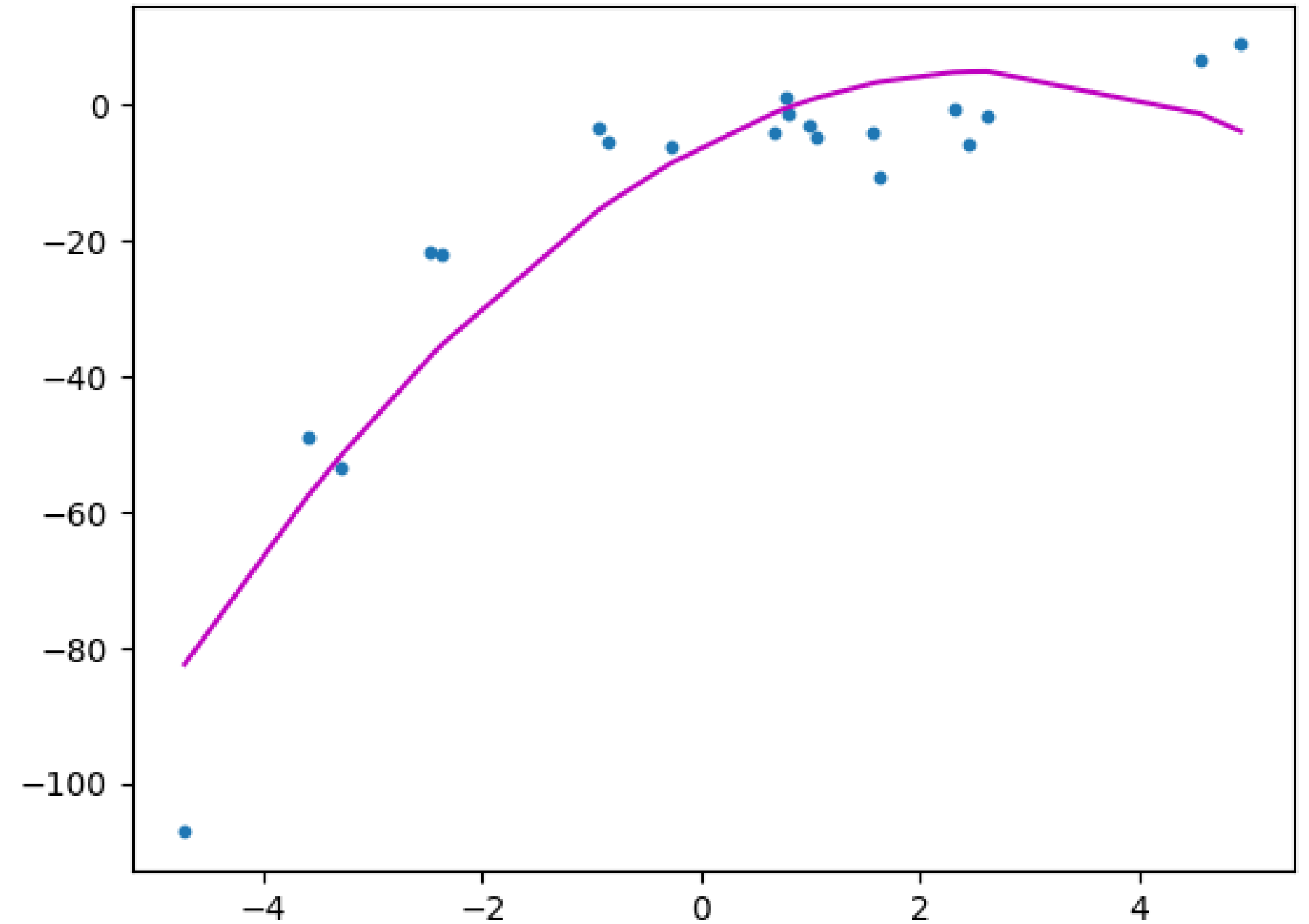
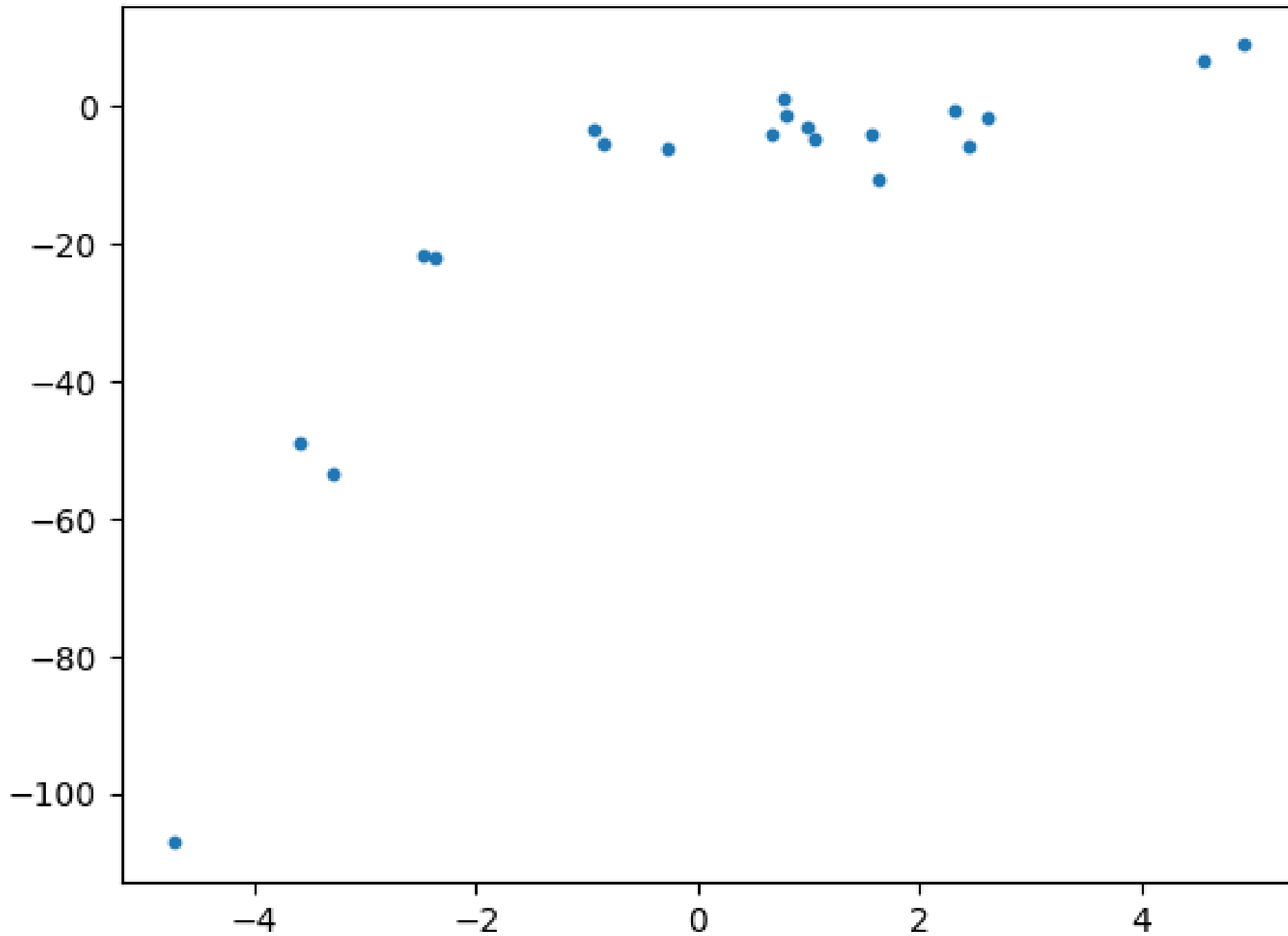
Recap

- Decision Tree
- Gini Impurity
- DT Pruning

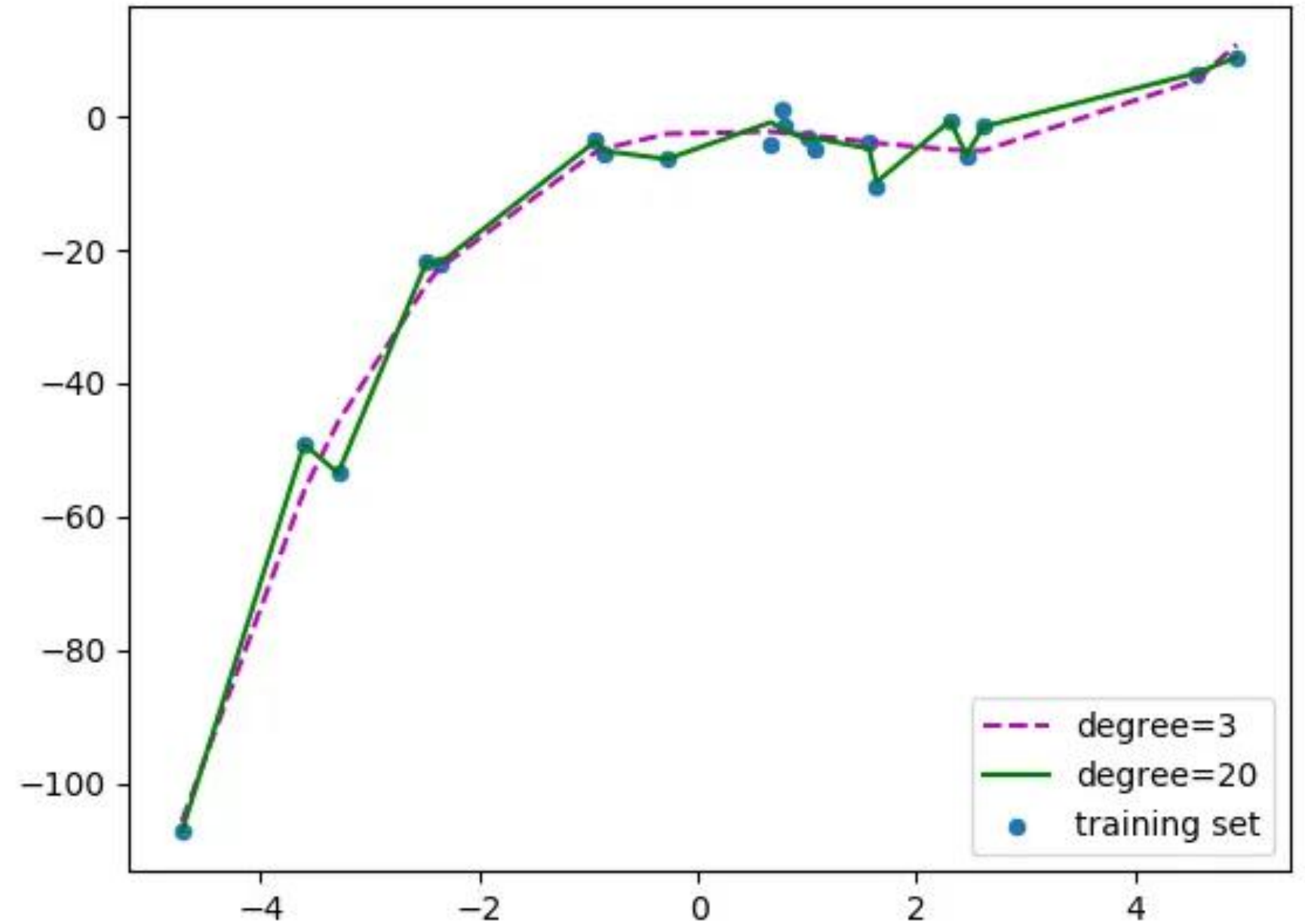
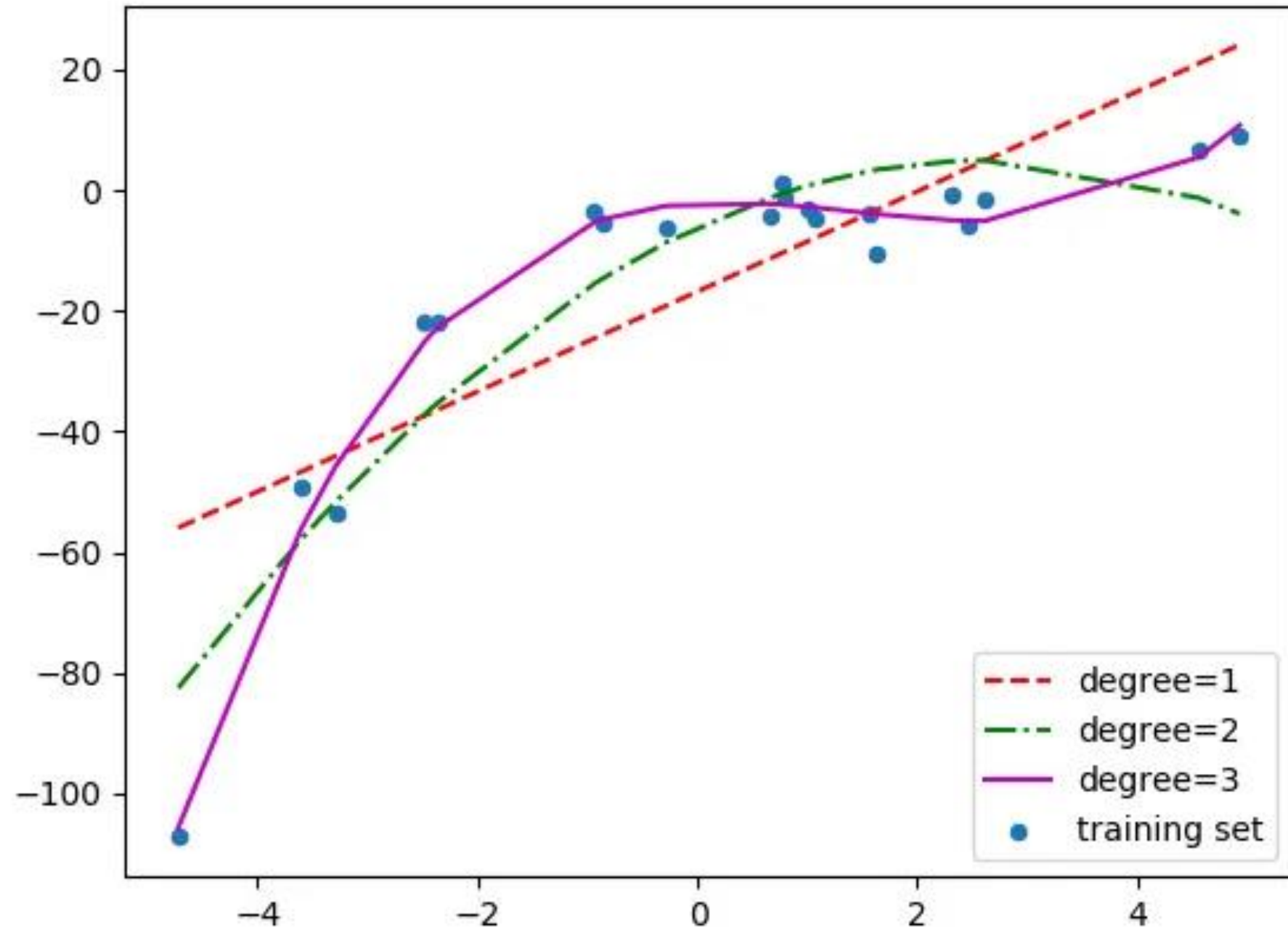


Underfitting, Overfitting, Bias Variance trade off

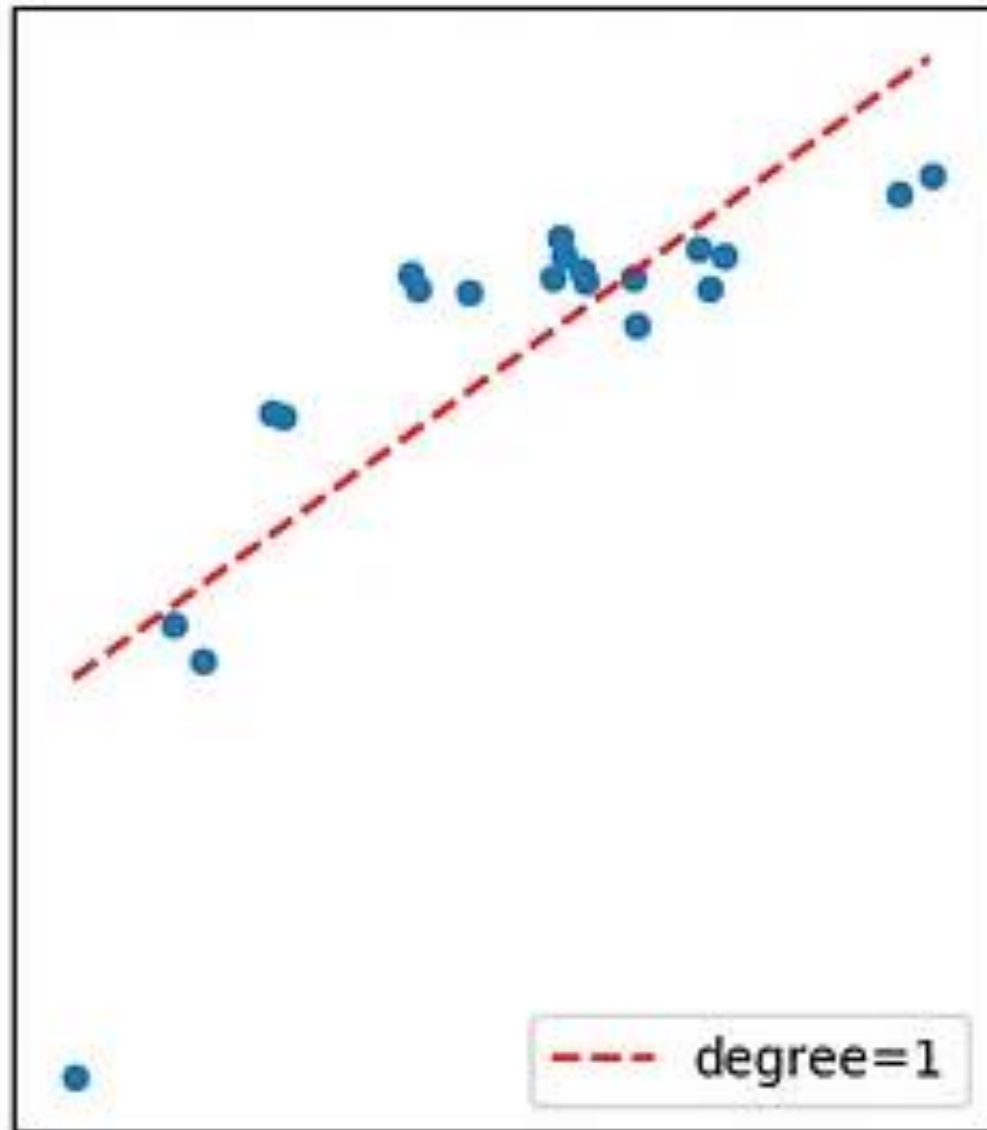
Fitting a curve



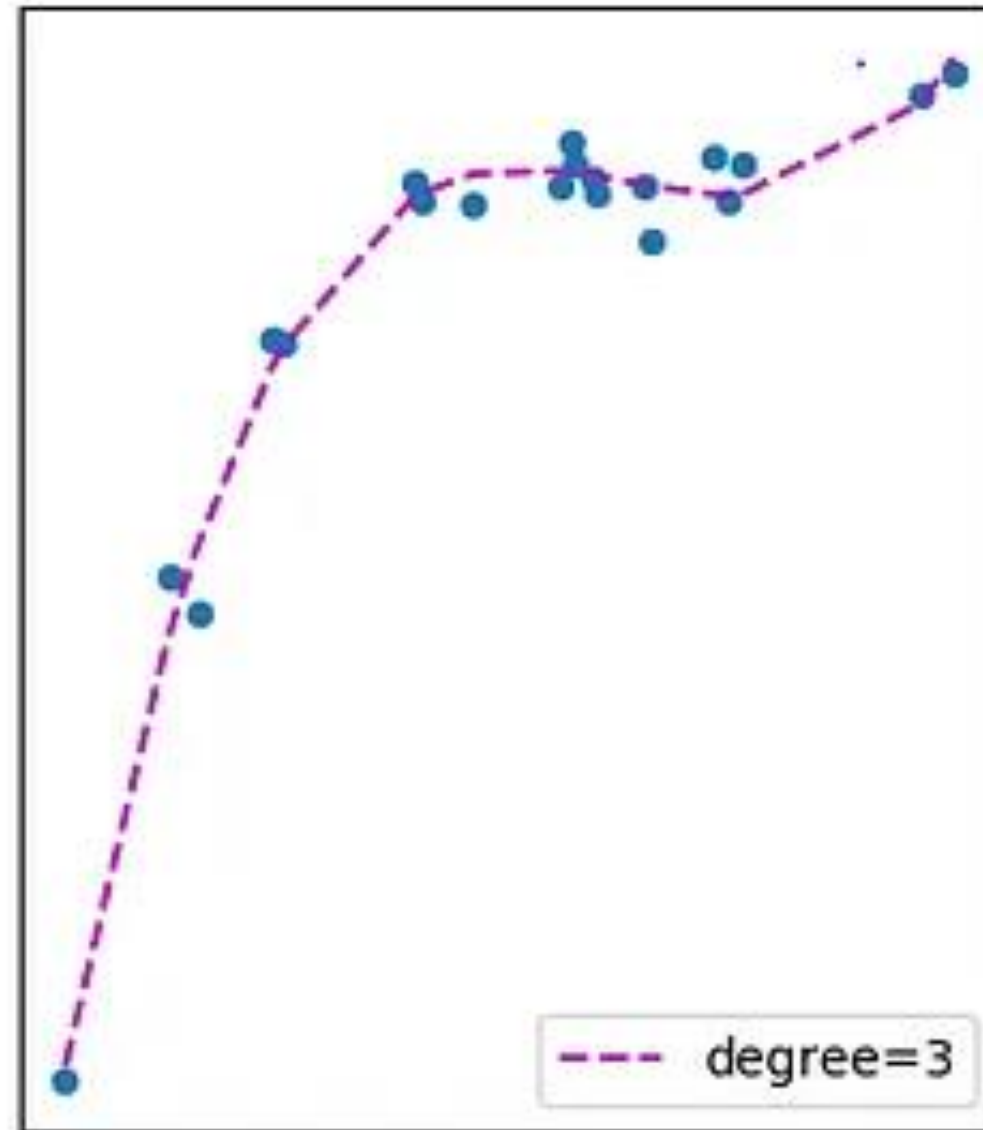
How many degrees is good?



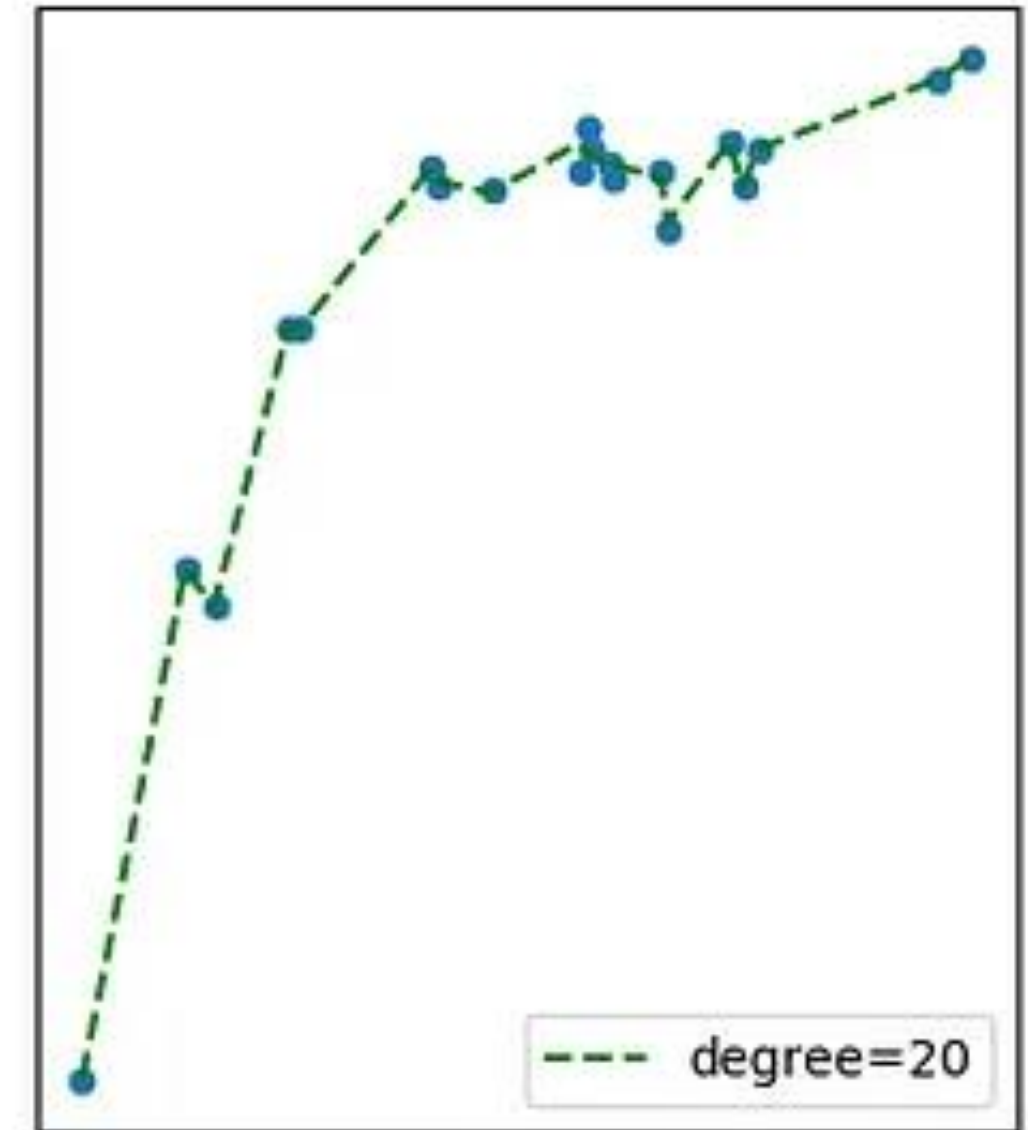
Underfitting & overfitting



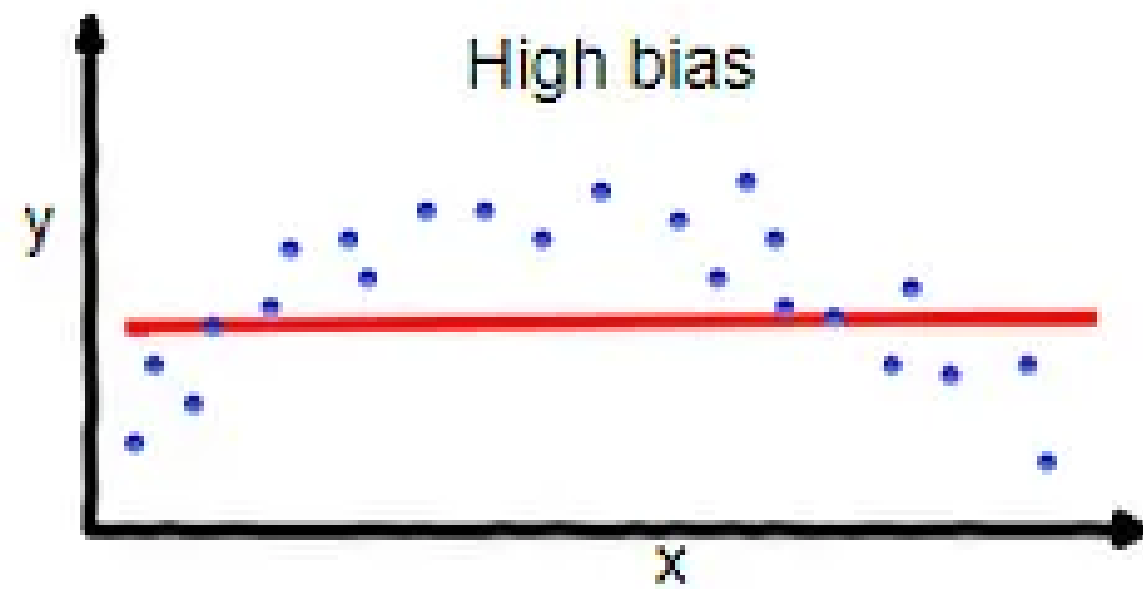
Underfit
High Bias
Low Variance



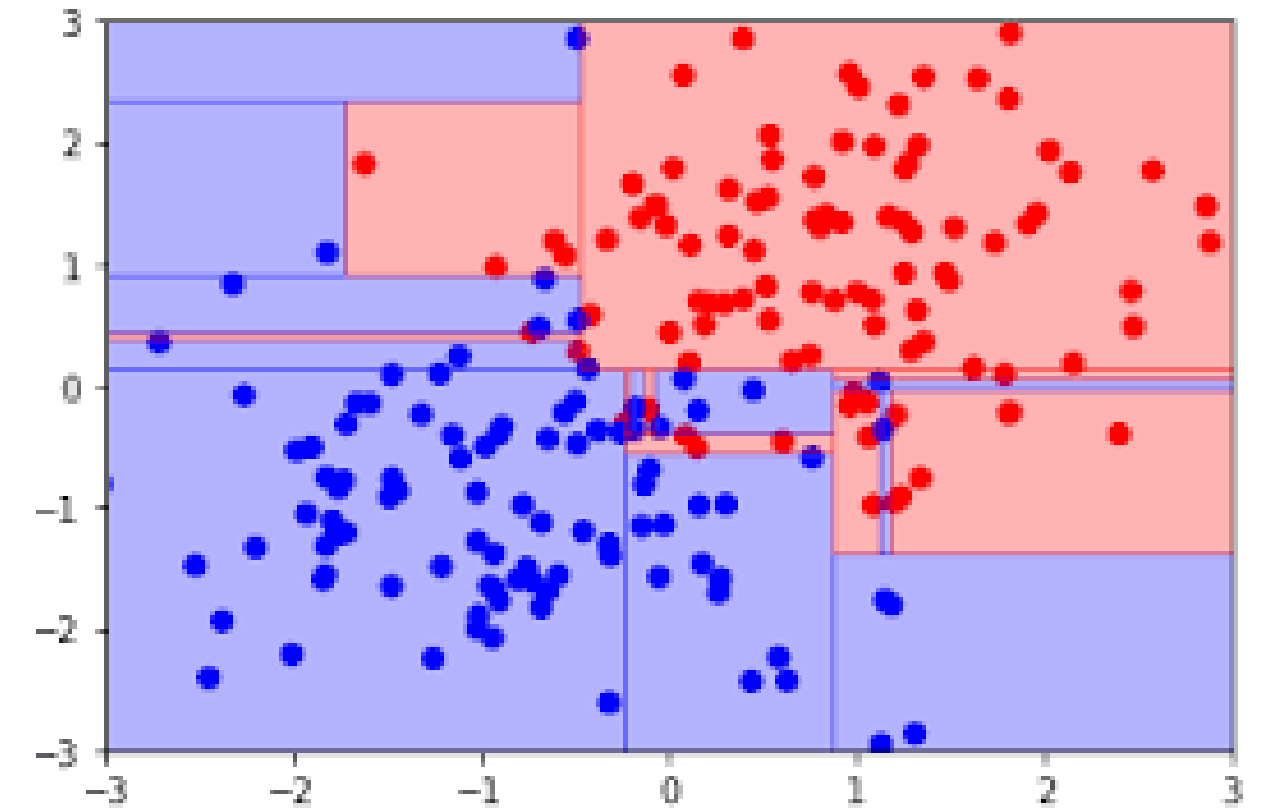
Correct Fit
Low Bias
Low Variance



Overfit
Low Bias
High Variance



Add some new
data points,
remove some.
What happens in
each case?



Think Linear Regression

Bias occurs when an algo has *limited flexibility* to learn the true signal from a dataset.

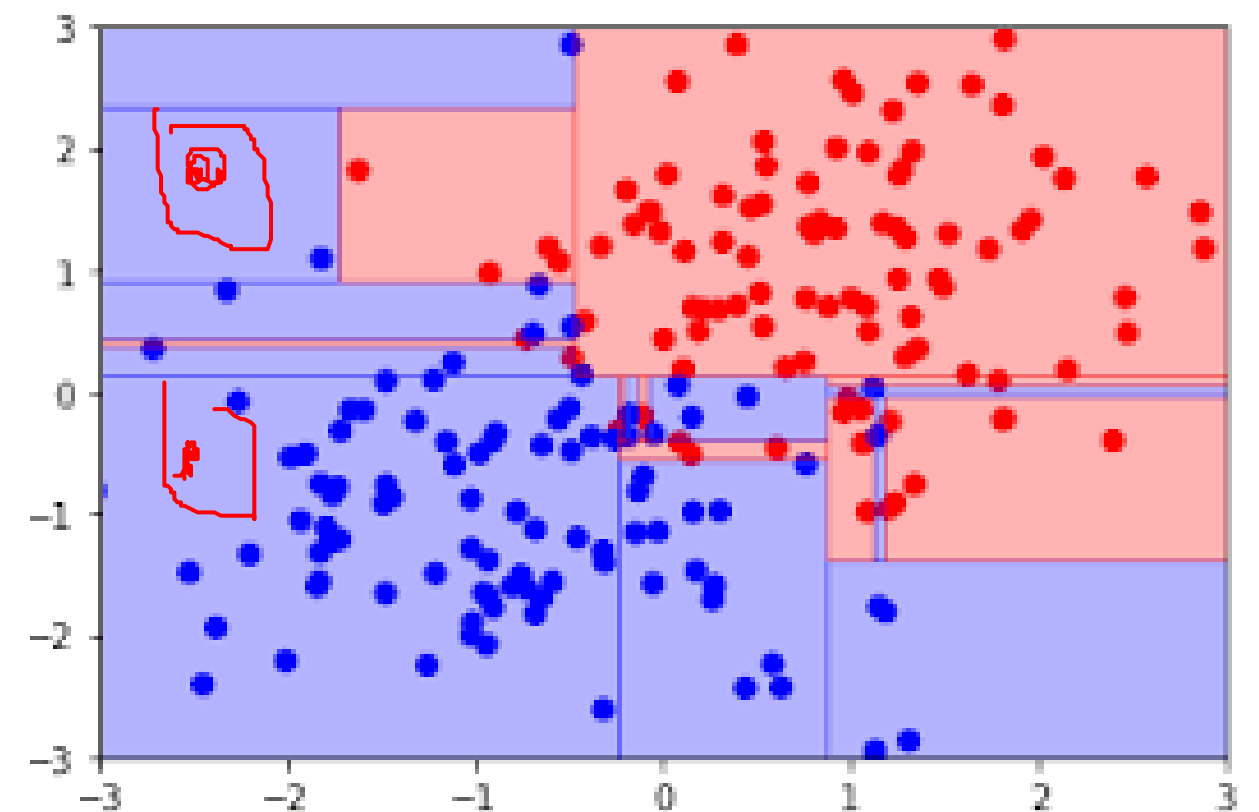
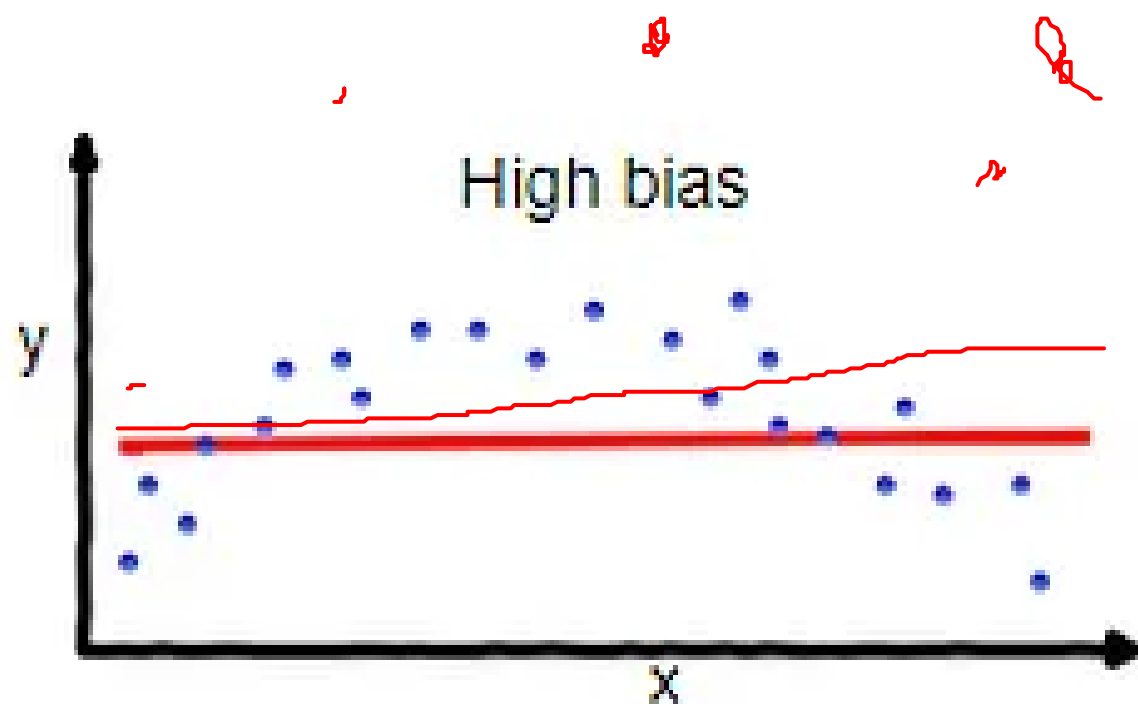
Think Decision Tree

Variance refers to an algo's *sensitivity* to specific sets of training data.

- 5 different training sets (imagine bootstrapping)
- Same algorithm trained on 5 data sets

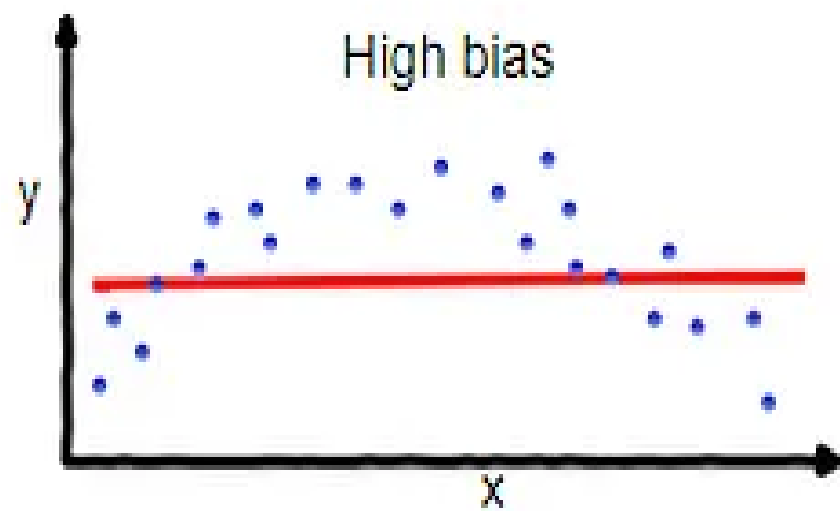
High bias, low variance
algorithms train models that
are consistent, but inaccurate
on average.

High variance, low bias
algorithms train models that
are accurate *on average*, but
inconsistent.



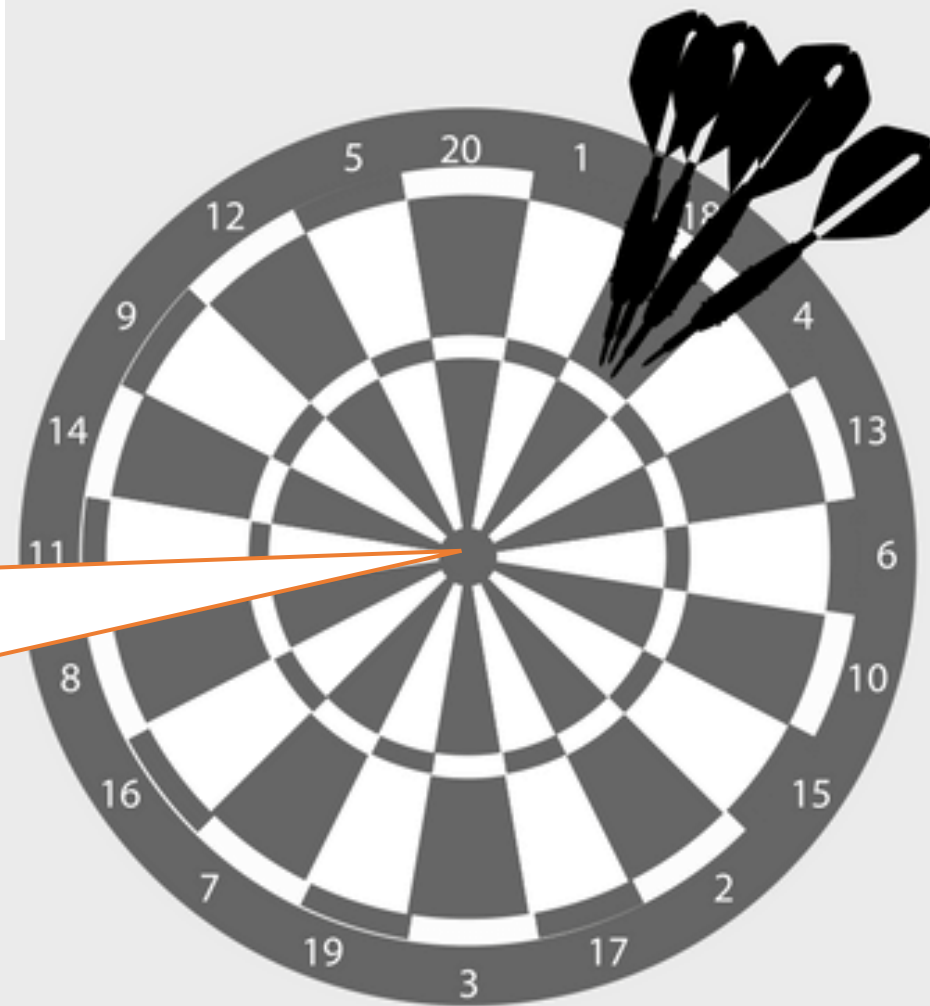
High bias, low variance algorithms train models that are consistent, but inaccurate *on average*.

High variance, low bias algorithms train models that are accurate *on average*, but inconsistent.

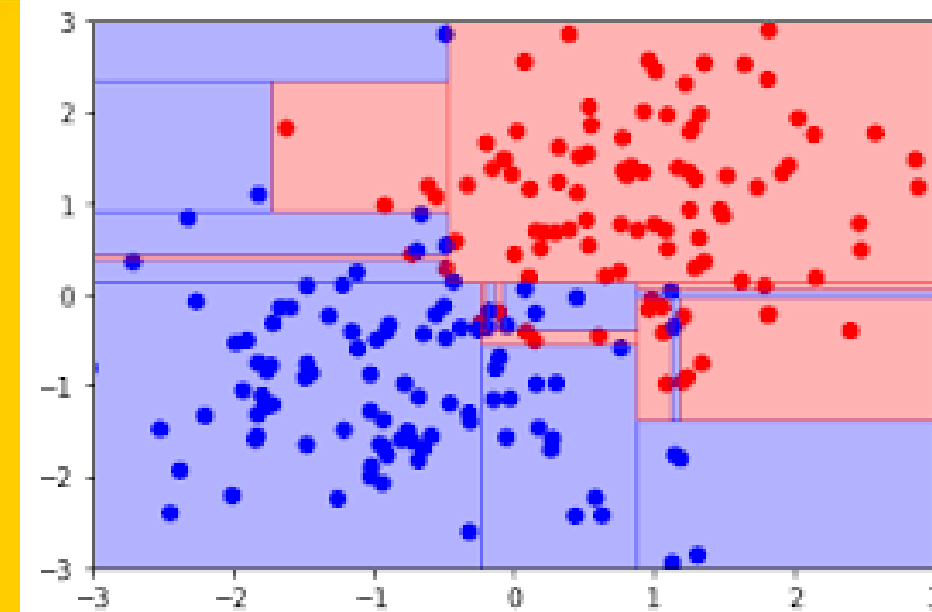
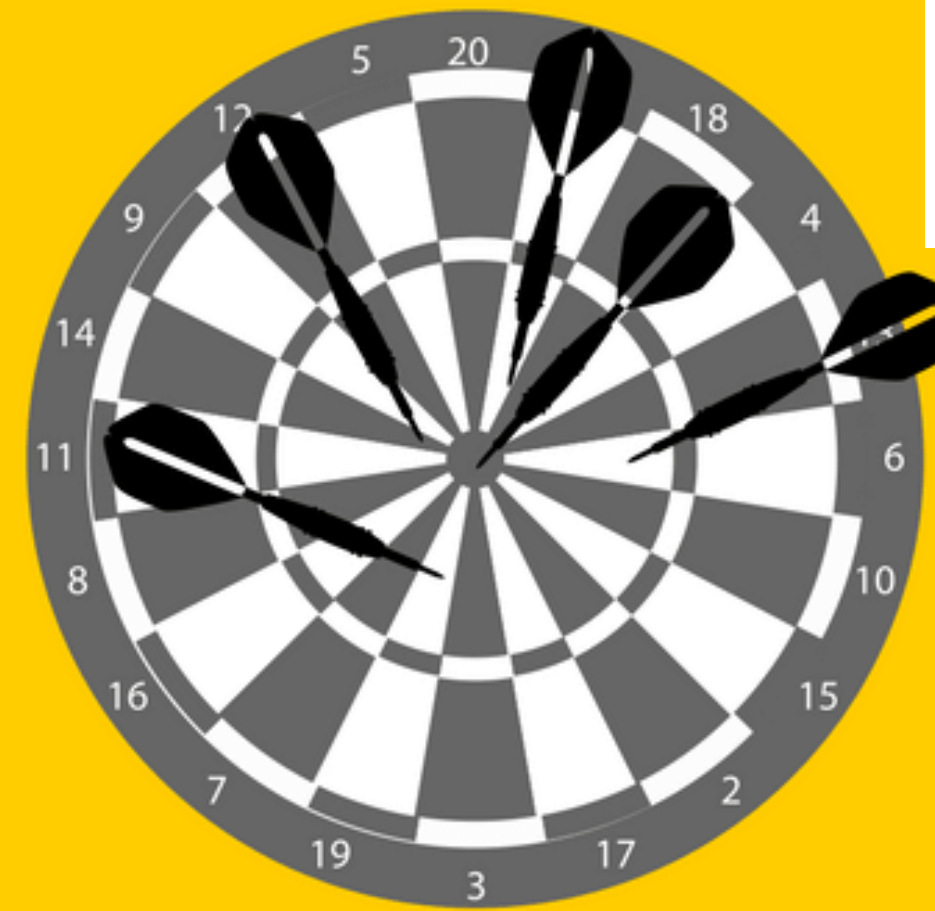


Bulls eye means best model & consistent model

High Bias
Low Variance



High Variance
Low Bias



But why is there a tradeoff?

Low variance algos tend to be **less complex**, with simple or rigid underlying structure.

- e.g. Regression
- e.g. Naive Bayes
- *Linear algos*
- *Parametric algos*

Low bias algos tend to be **more complex**, with flexible underlying structure.

- e.g. Decision trees
- e.g. Nearest neighbors
- *Non-linear algos*
- *Non-parametric algos*

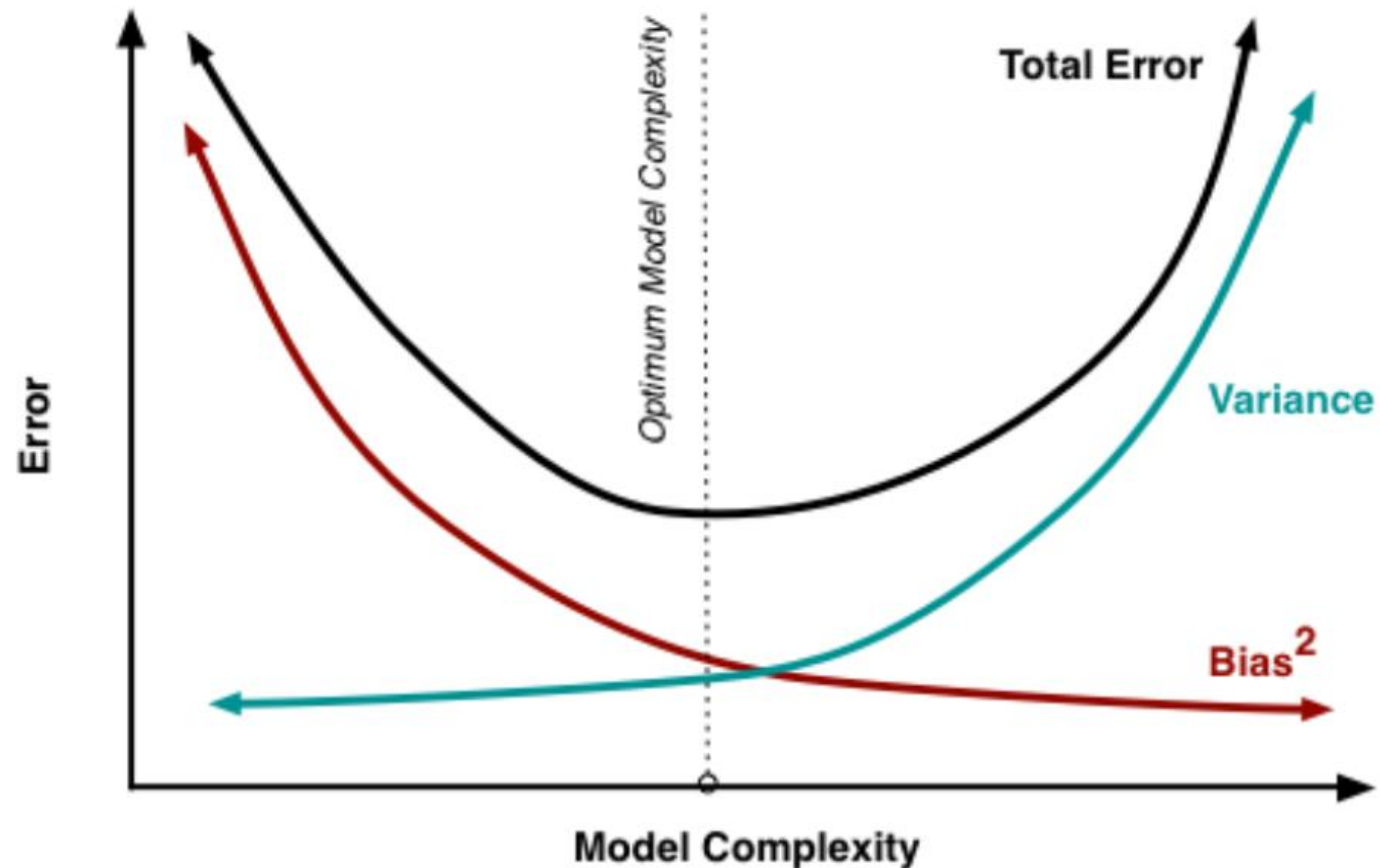
Within each algo family, there's a tradeoff too...

For example, regression can be **regularized** to further reduce complexity.

For example, decision trees can be **pruned** to reduce complexity.

Bias Variance Plot

$$\text{Total Error} = \text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$





Wisdom of the Crowds

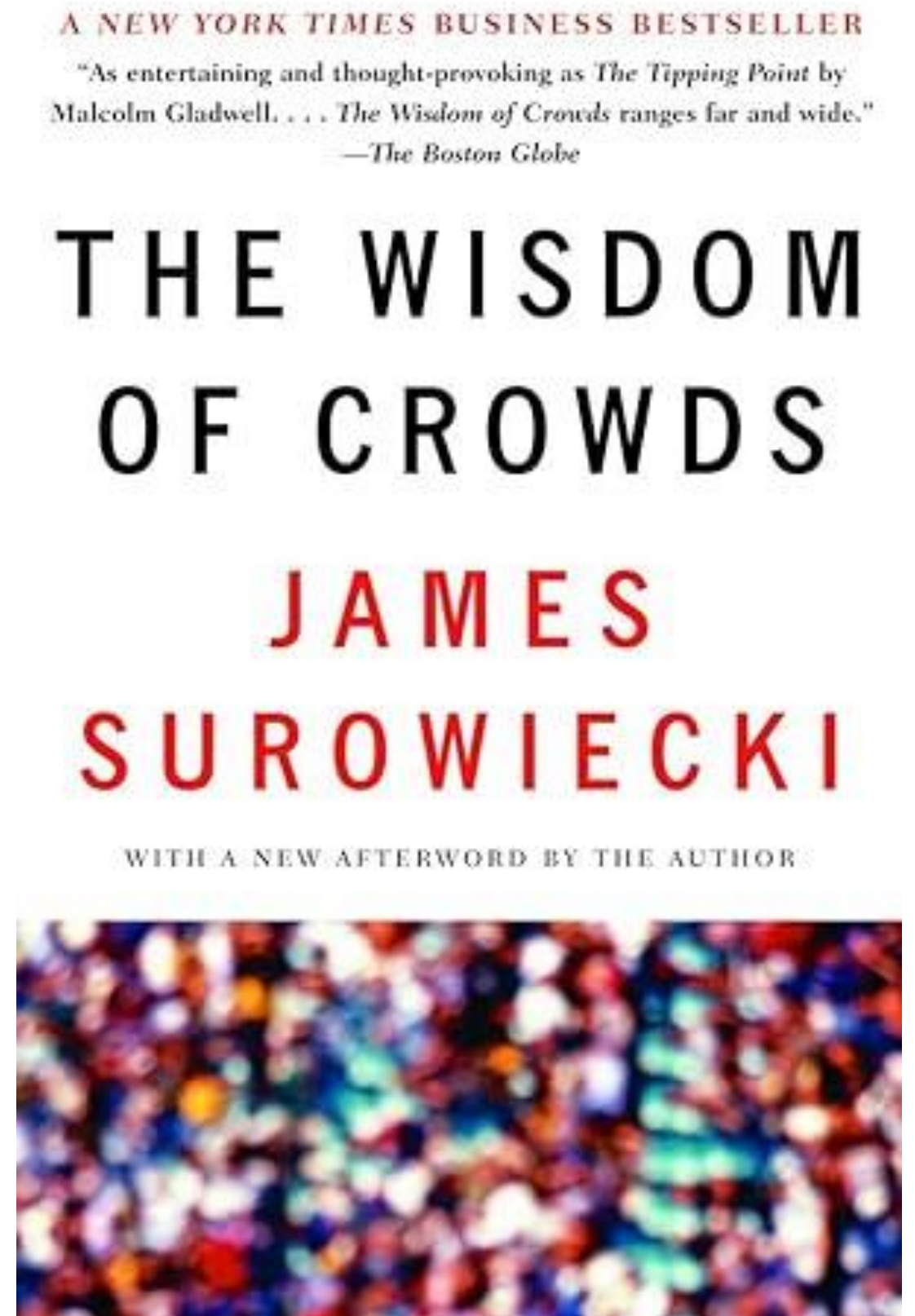
A game of guessing?

- How many candies in the jar?
- How to get the most correct answer?
- No answer is correct, some answers are useful
- Especially if averaged over a group



The wisdom of crowds

- “Under the right circumstances, groups are remarkably intelligent, and are often smarter than the smartest people in them”
- Using a Group for better decision
 - Diverse background
 - Independent decision by each individual
 - Good method for aggregation



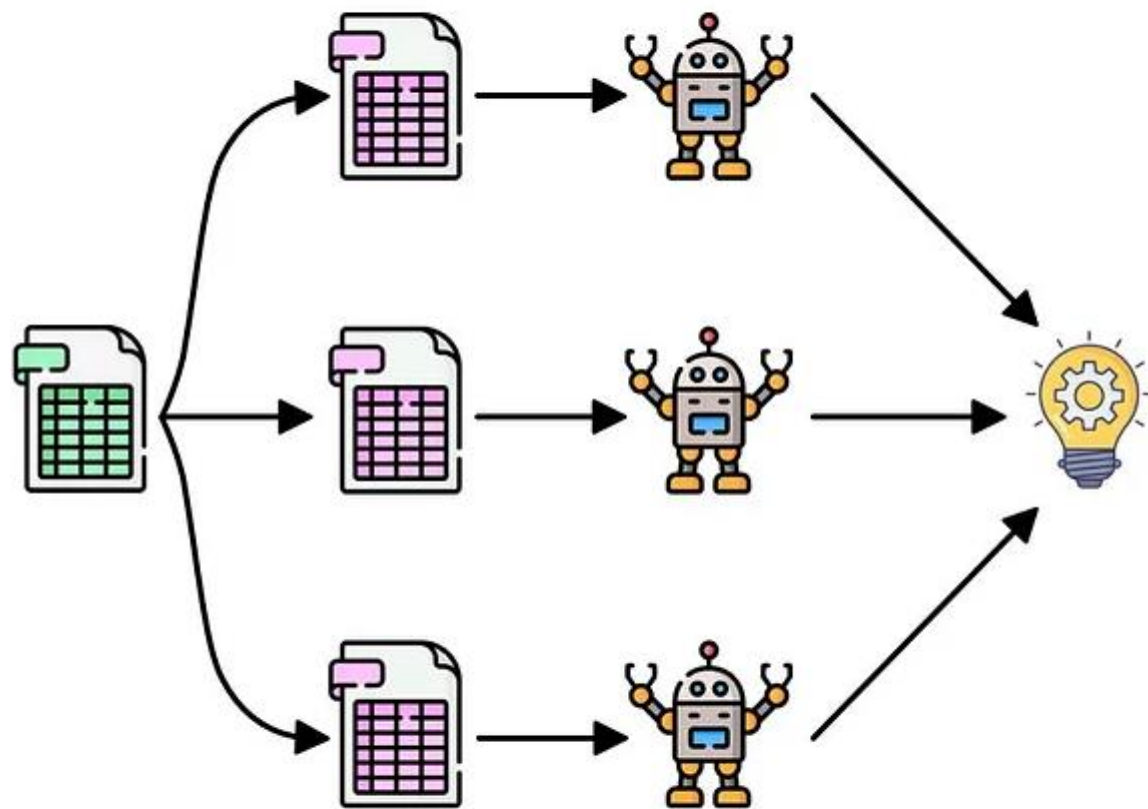


Ensemble Learning

Ensemble Learning

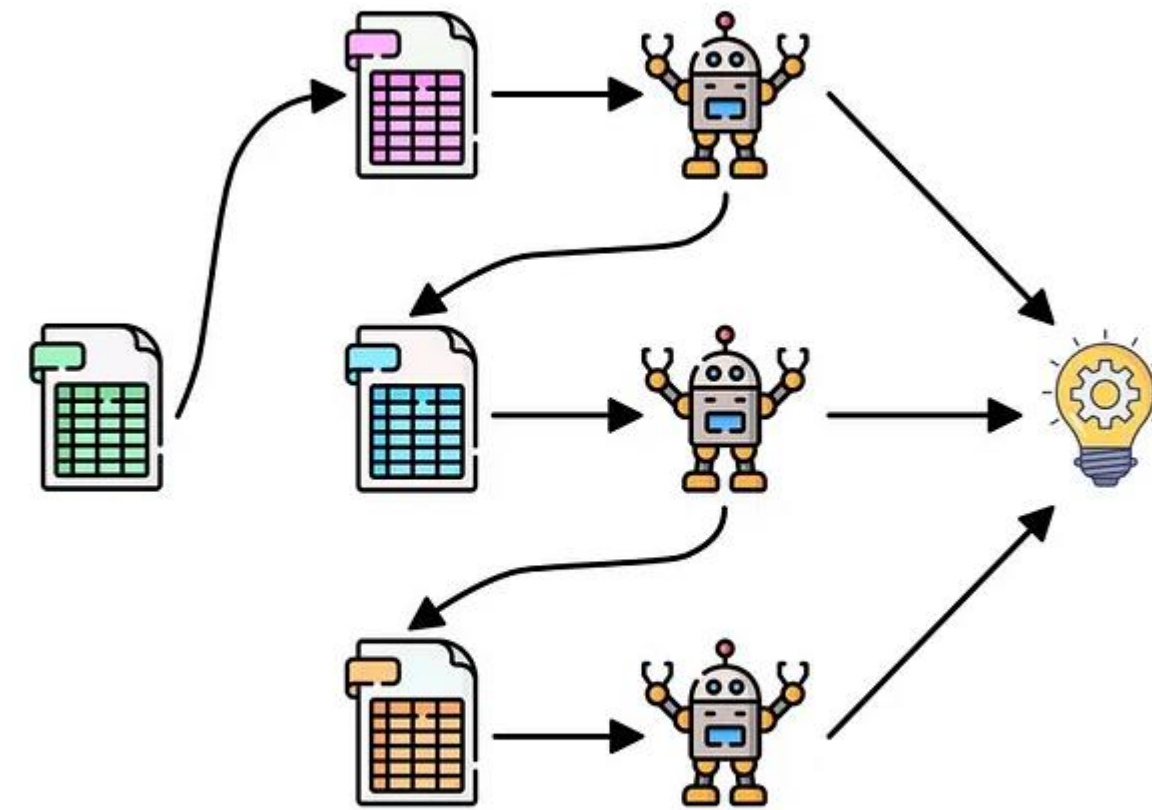
- Multiple ML model used together for prediction

Bagging



Parallel

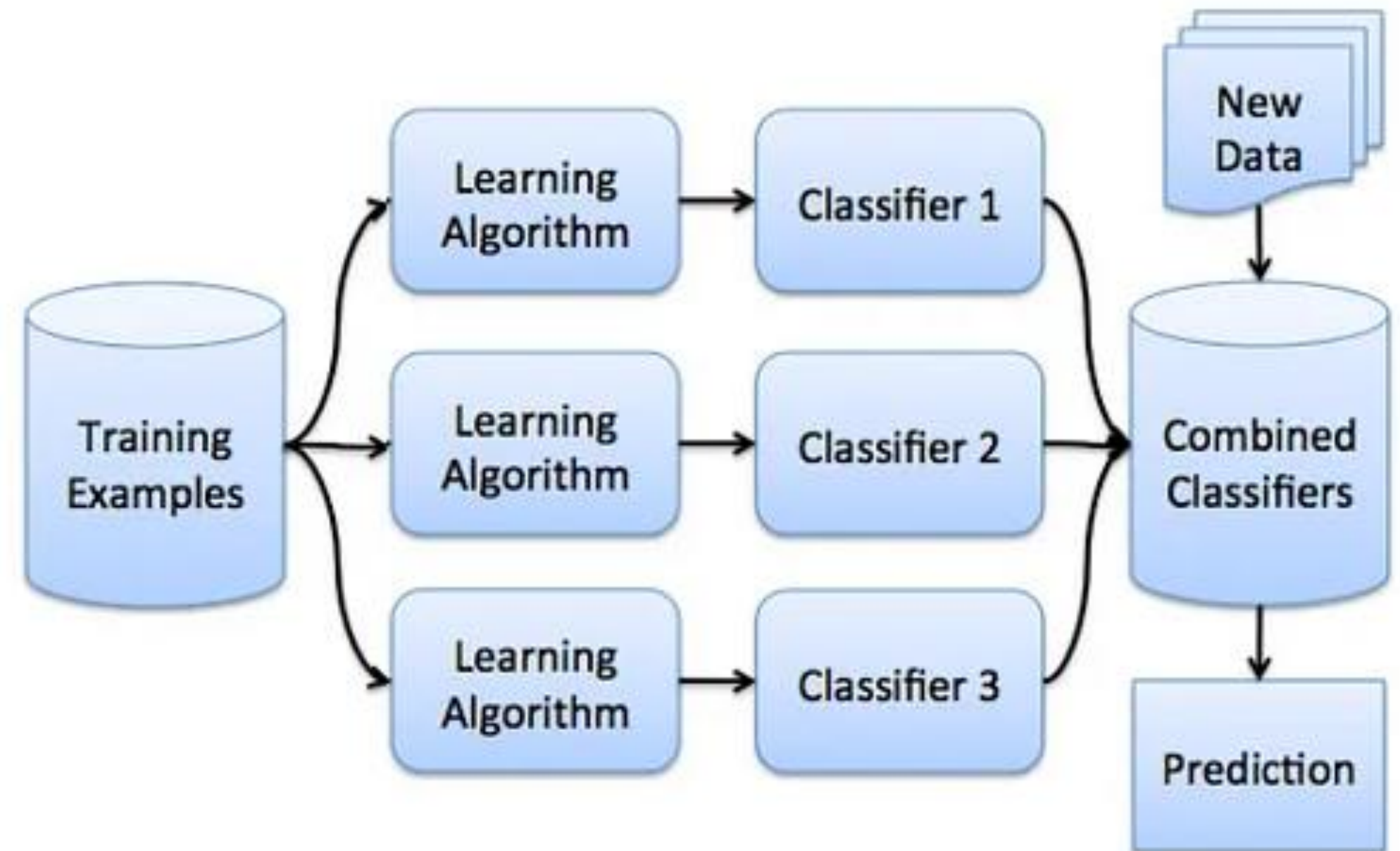
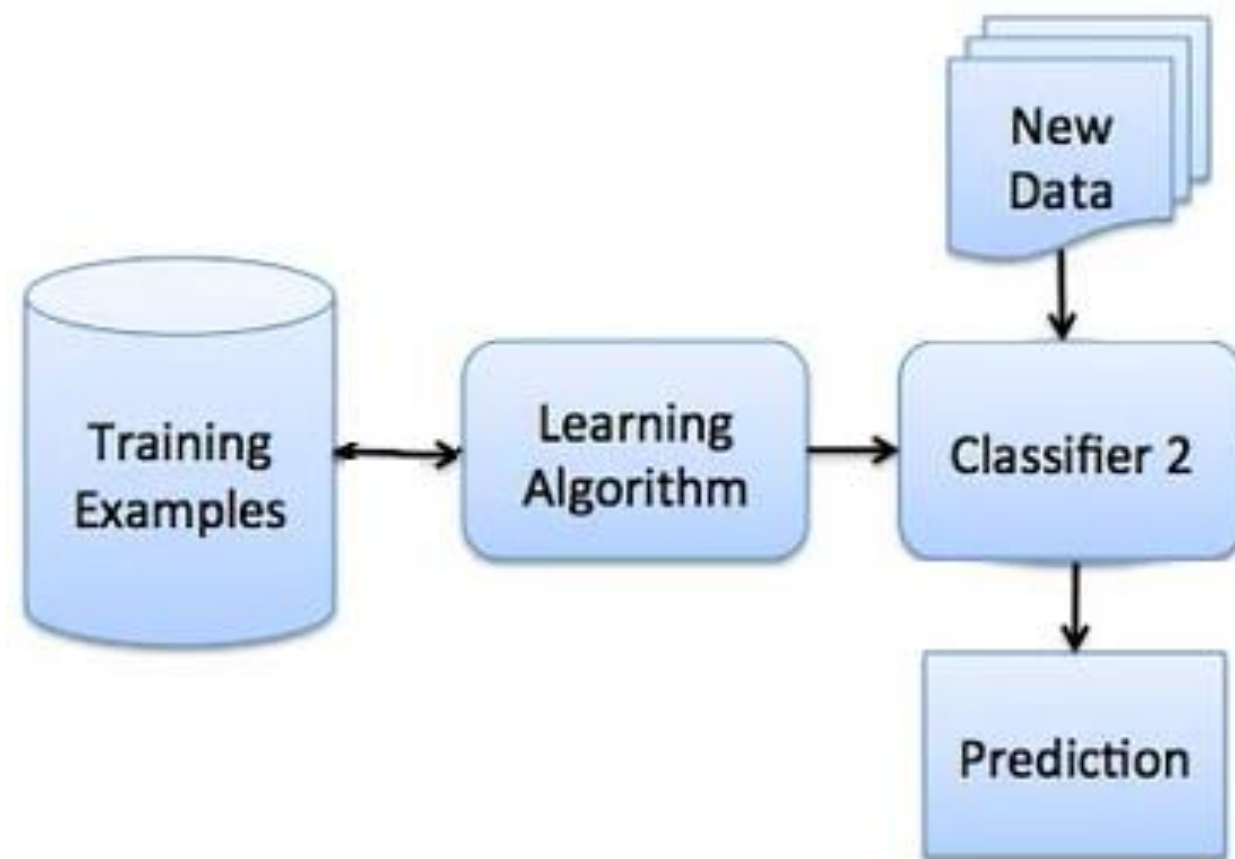
Boosting



Sequential

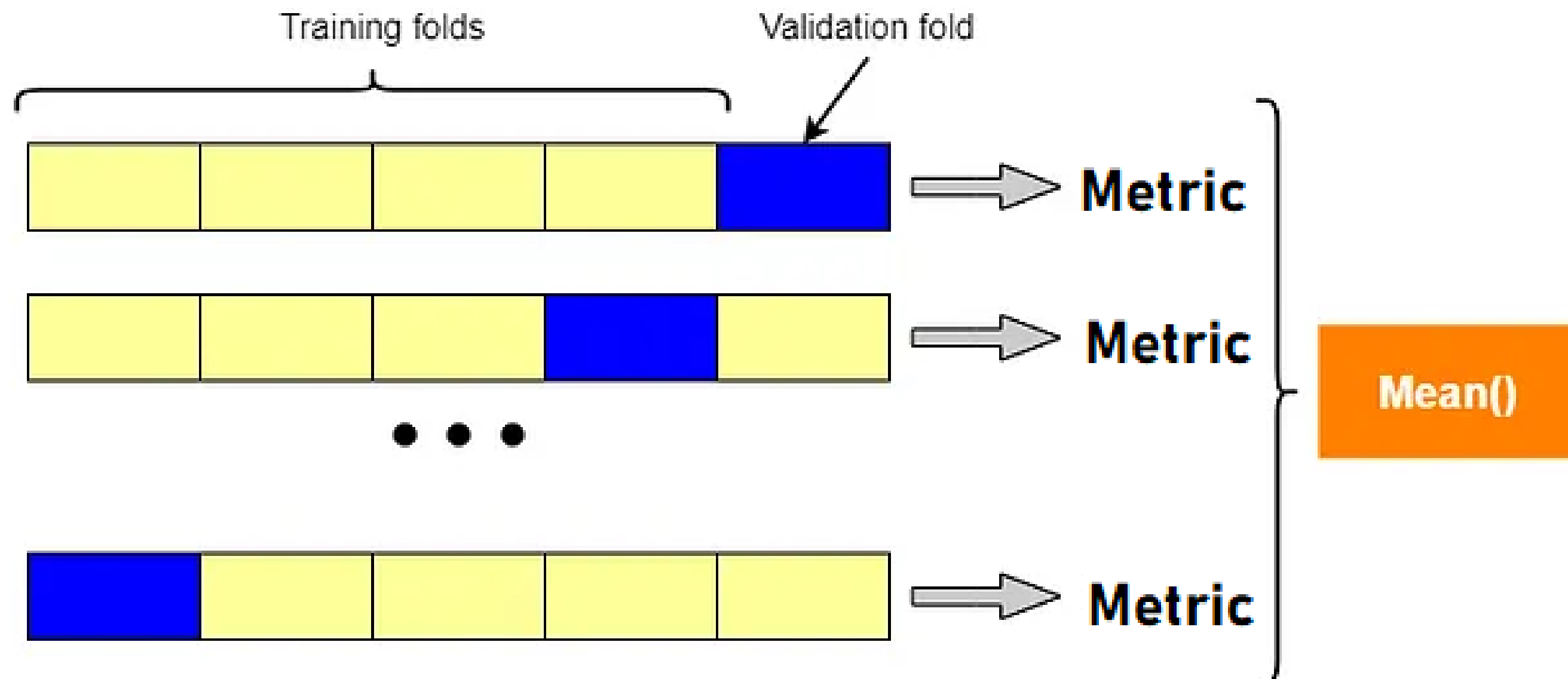
Majority Voting/Averaging models

- High variance from single model
- How about using multiple models?
 - Will it reduce variance?
 - Not if data is correlated



A second look at K-Fold CV

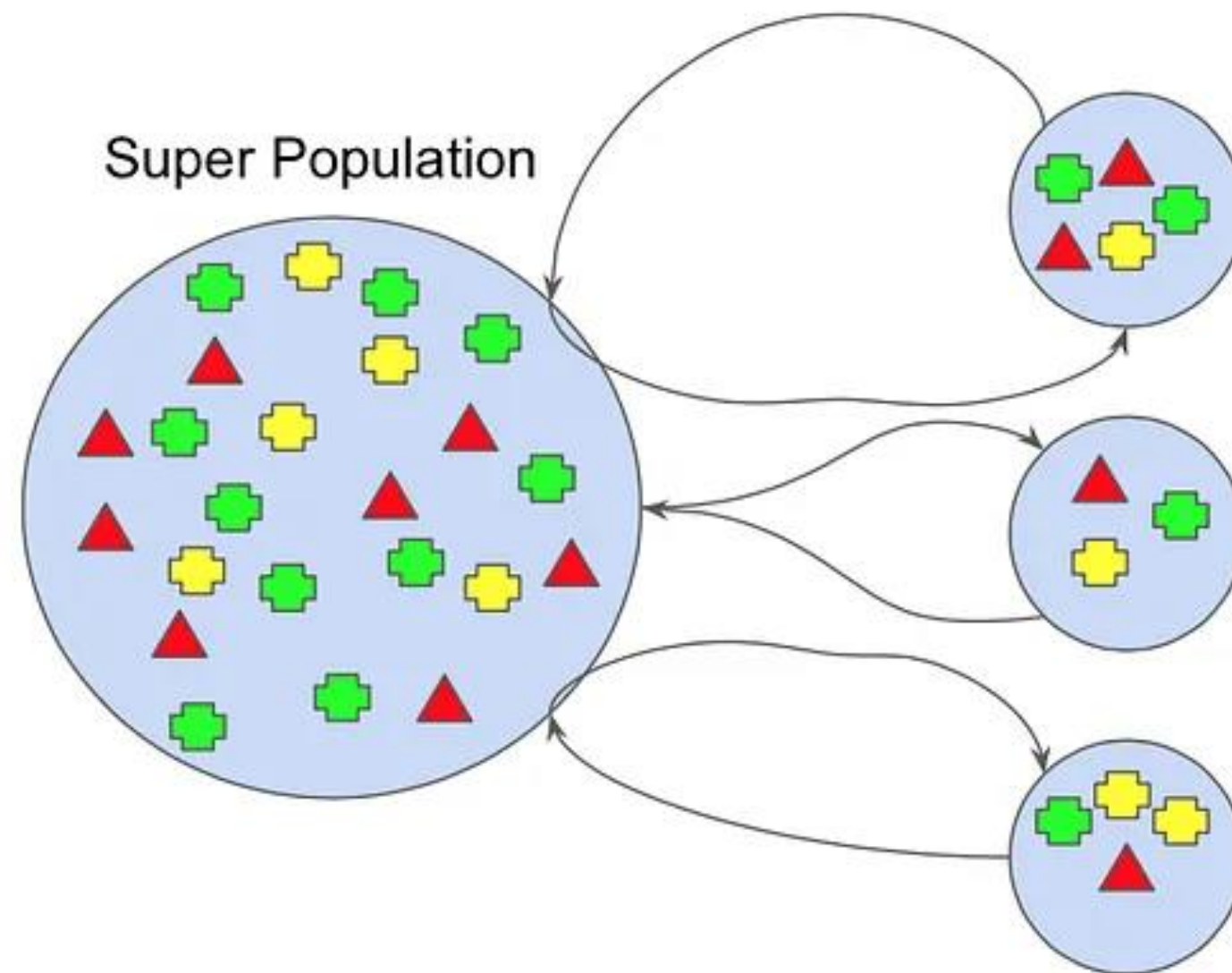
- Most data is repeated across folds
- IID from one record to next.
- Highly correlated from one fold to another



Solution: Bagging

- Bagging = Bootstrapping + Aggregation
- Bootstrapping: Sampling with Replacement

Data becomes uncorrelated



Sample Population 1

Std(a)

Sample Population 2

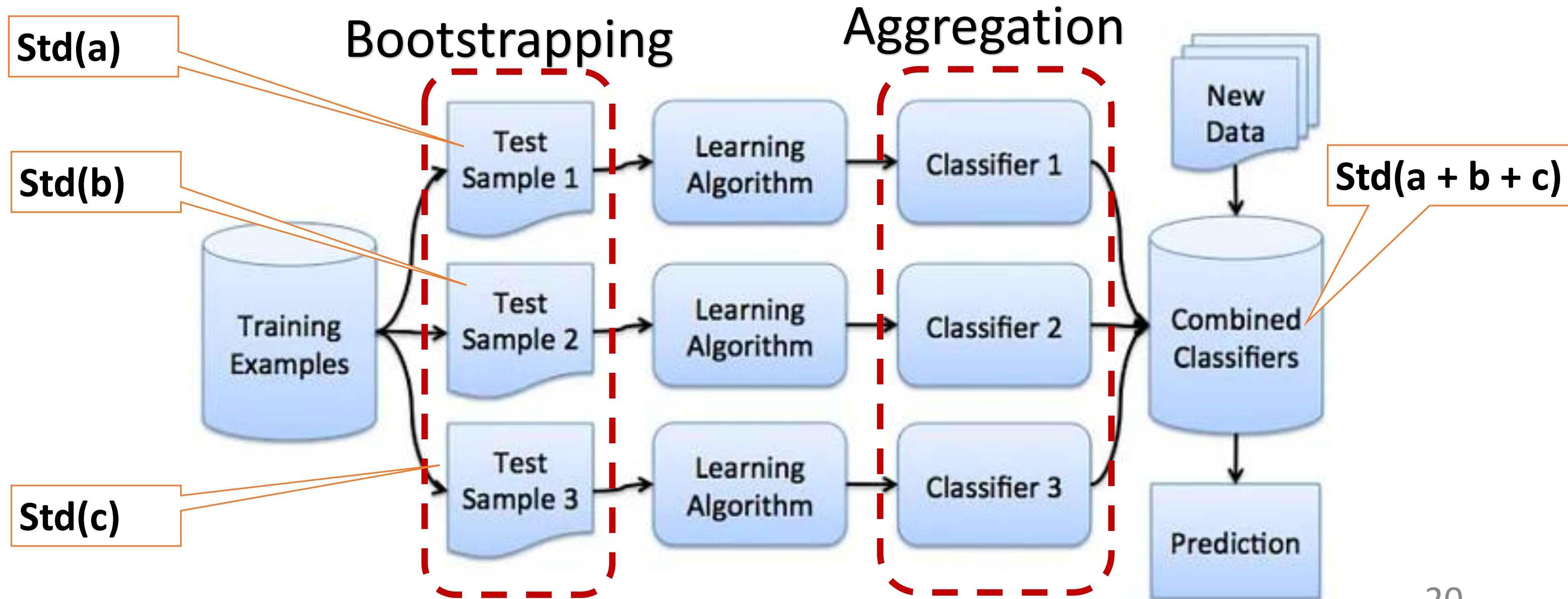
Std(b)

Sample Population 3

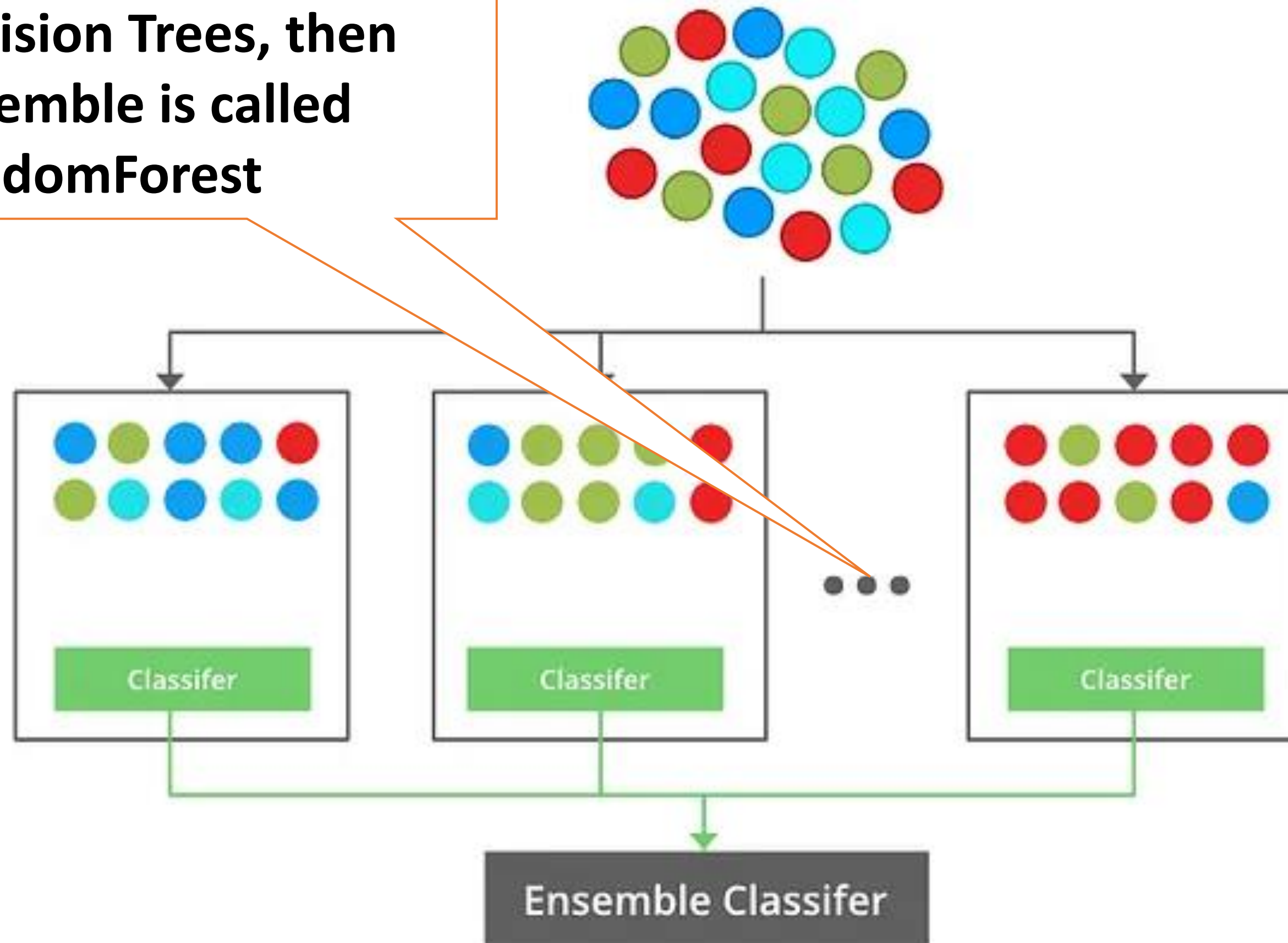
Std(c)

Solution: Bagging (Contd.)

- Aggregation = Combining classifiers



When all classifiers are
Decision Trees, then
ensemble is called
RandomForest

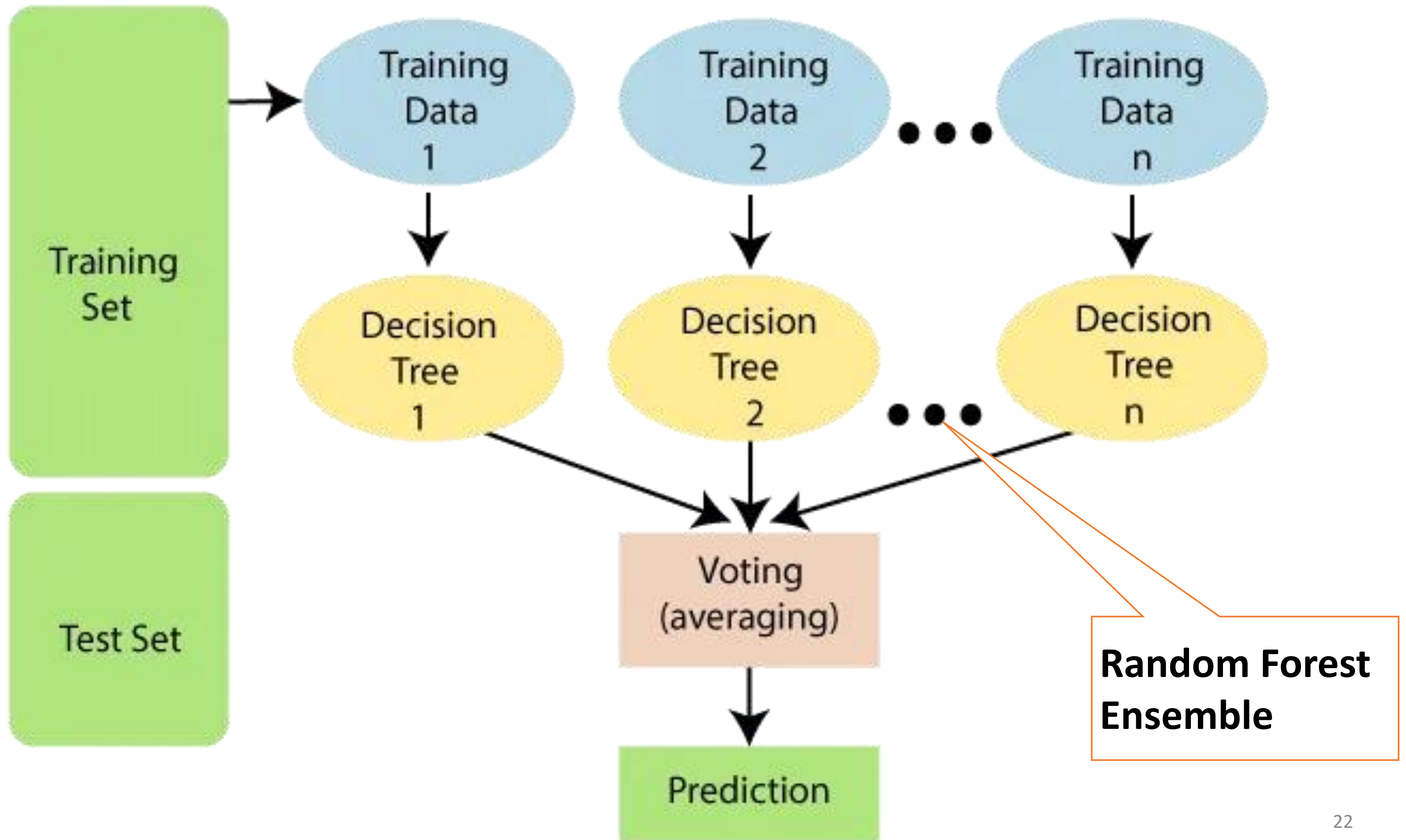


Original Data

Bootstrapping

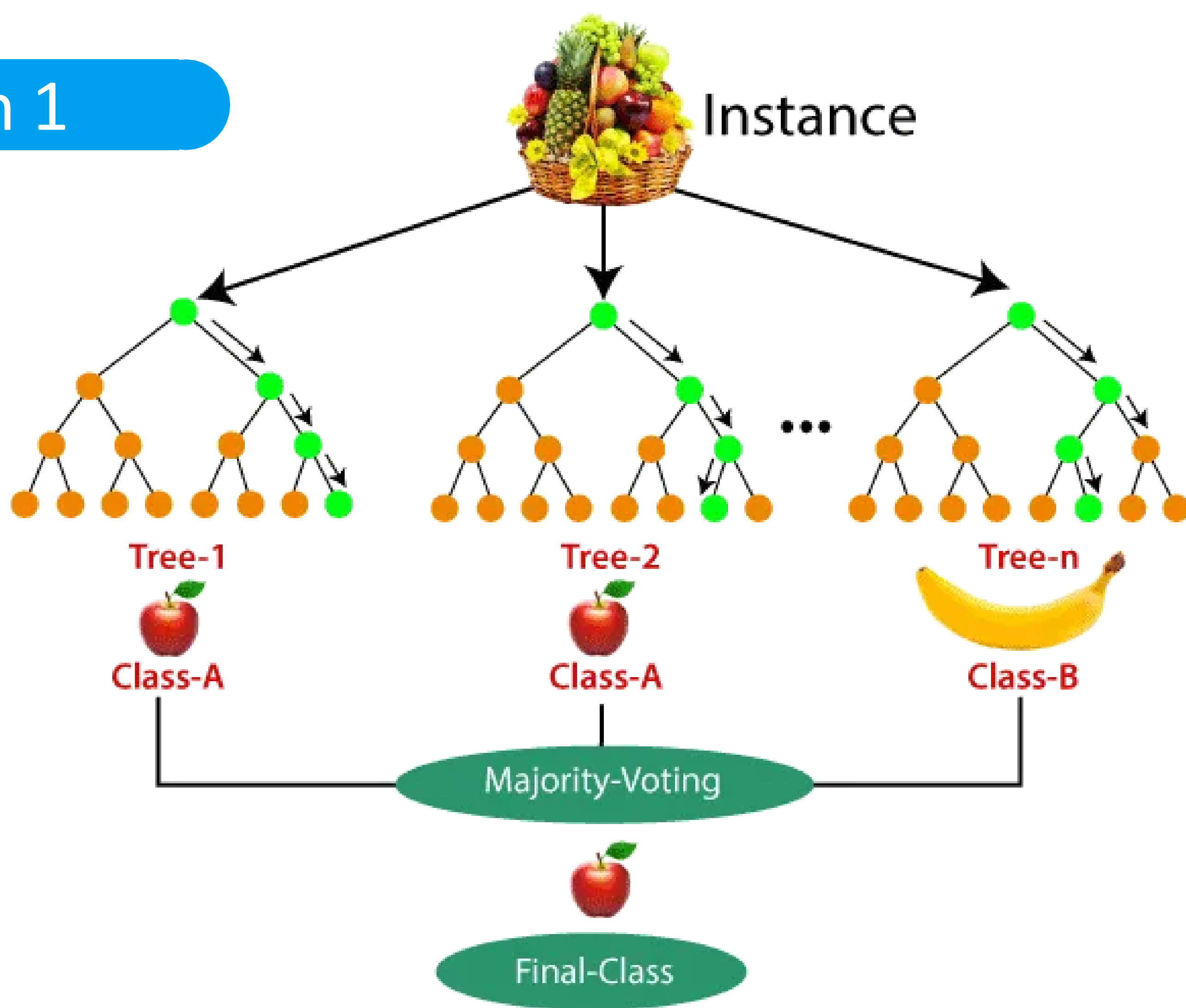
Aggregating

Bagging

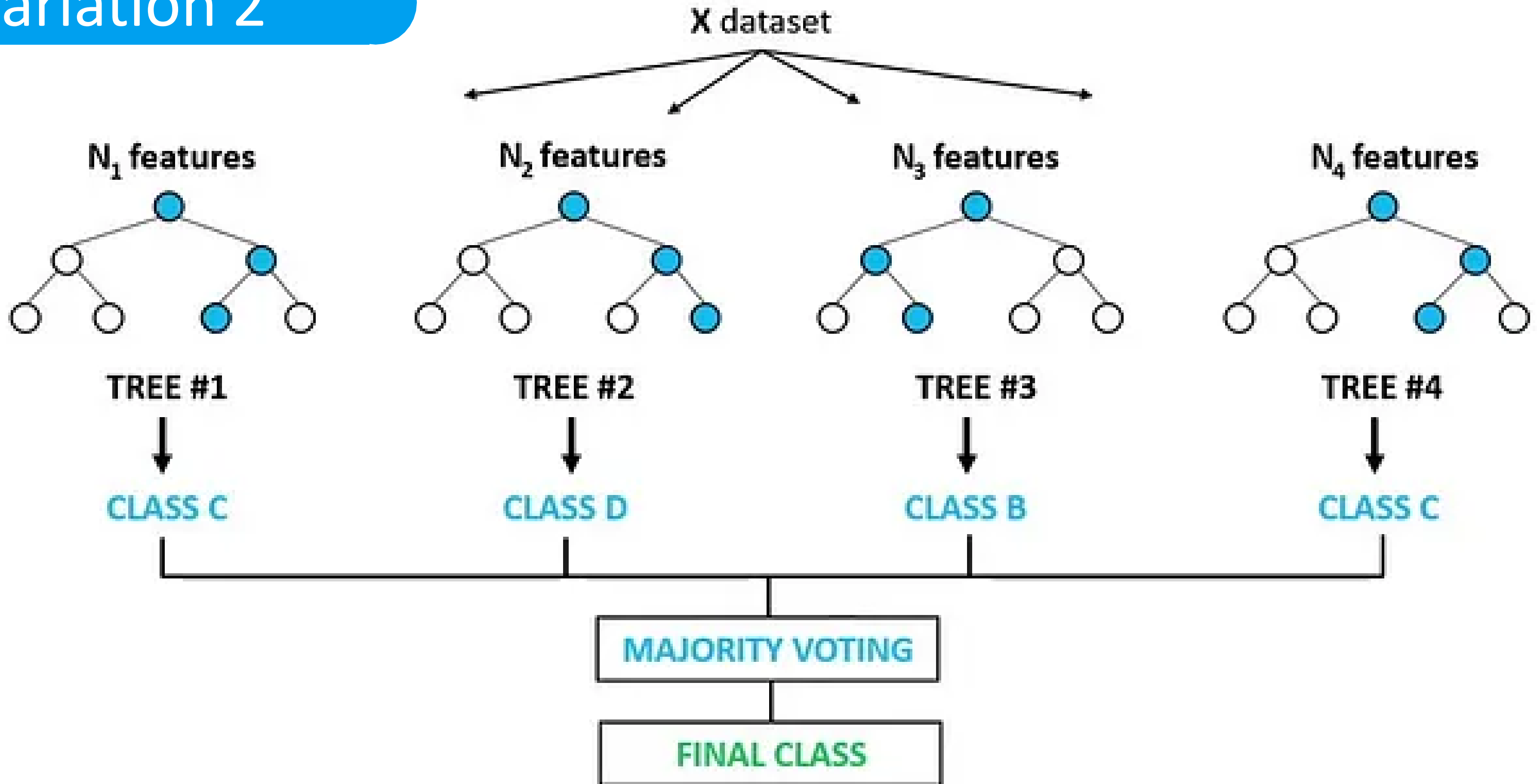




Variation 1



Variation 2



Random Forest hyperparameters

- Decision Tree hyper params
 - Criteria = gini/entropy
 - Max tree depth, Max leaf nodes
 - Min samples in split, Min samples in leaf
- Bagging hyperparams: Bootstrap Y/N
- Random Forest hyper params
 - Number of trees (n_estimators)
 - Max samples per tree
 - Max features per tree

Random Forest advantages

- Feature correlation does not matter
- Feature distribution does not matter
- No need to scale data
- Works well even when data is missing
- Overcomes the problem of overfitting
- Not very expensive
- Flexible and high accuracy

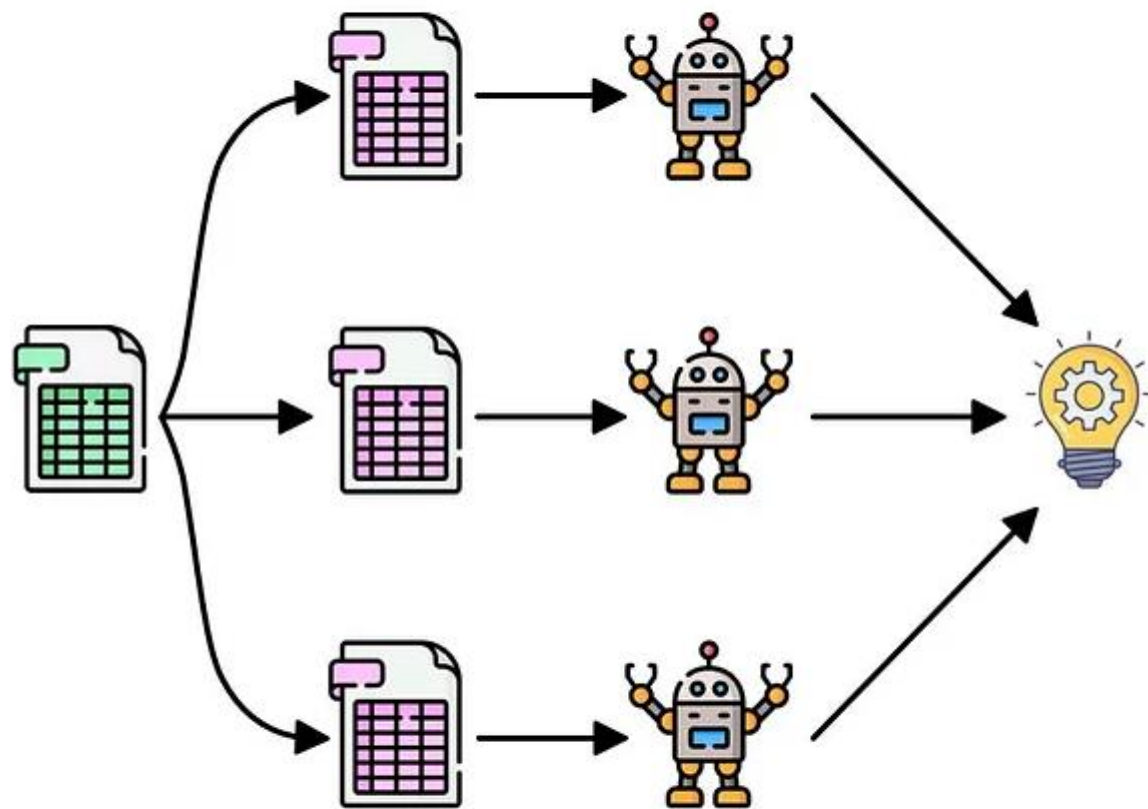


Boosting

Ensemble Learning

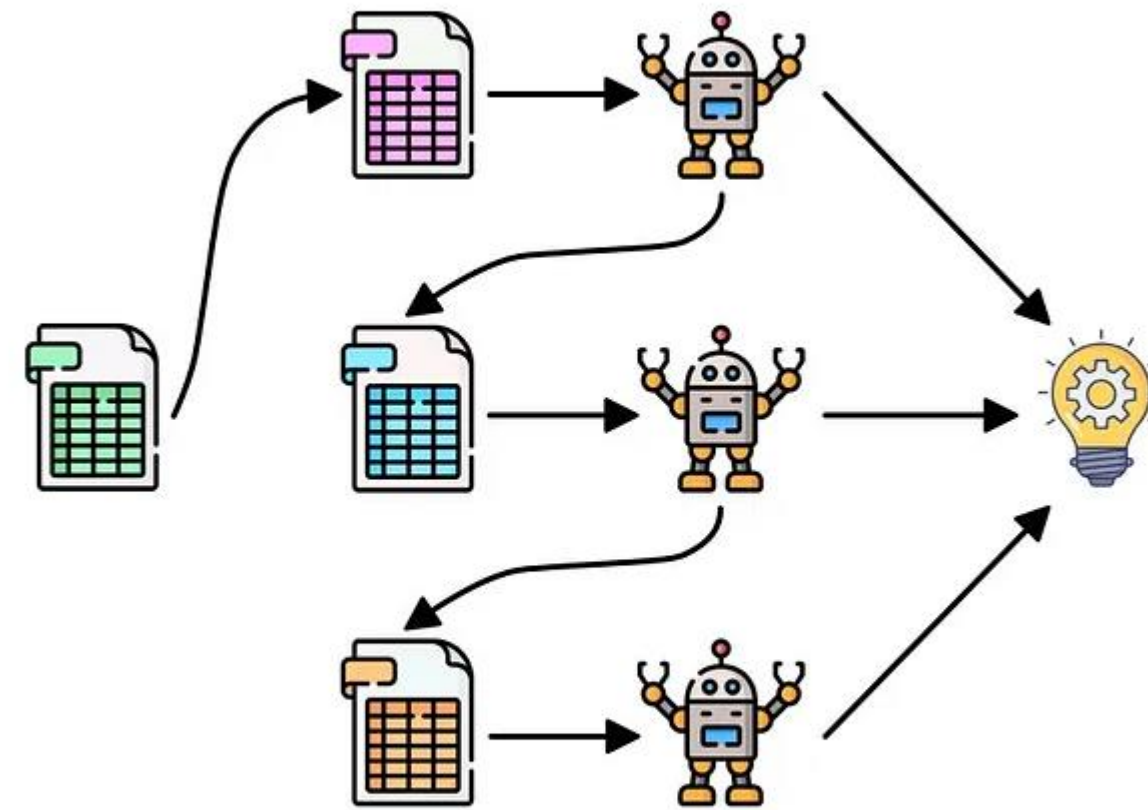
- Multiple ML model used together for prediction

Bagging



Parallel

Boosting



Sequential

BAGGING

BOOTSTRAP AGGREGATION

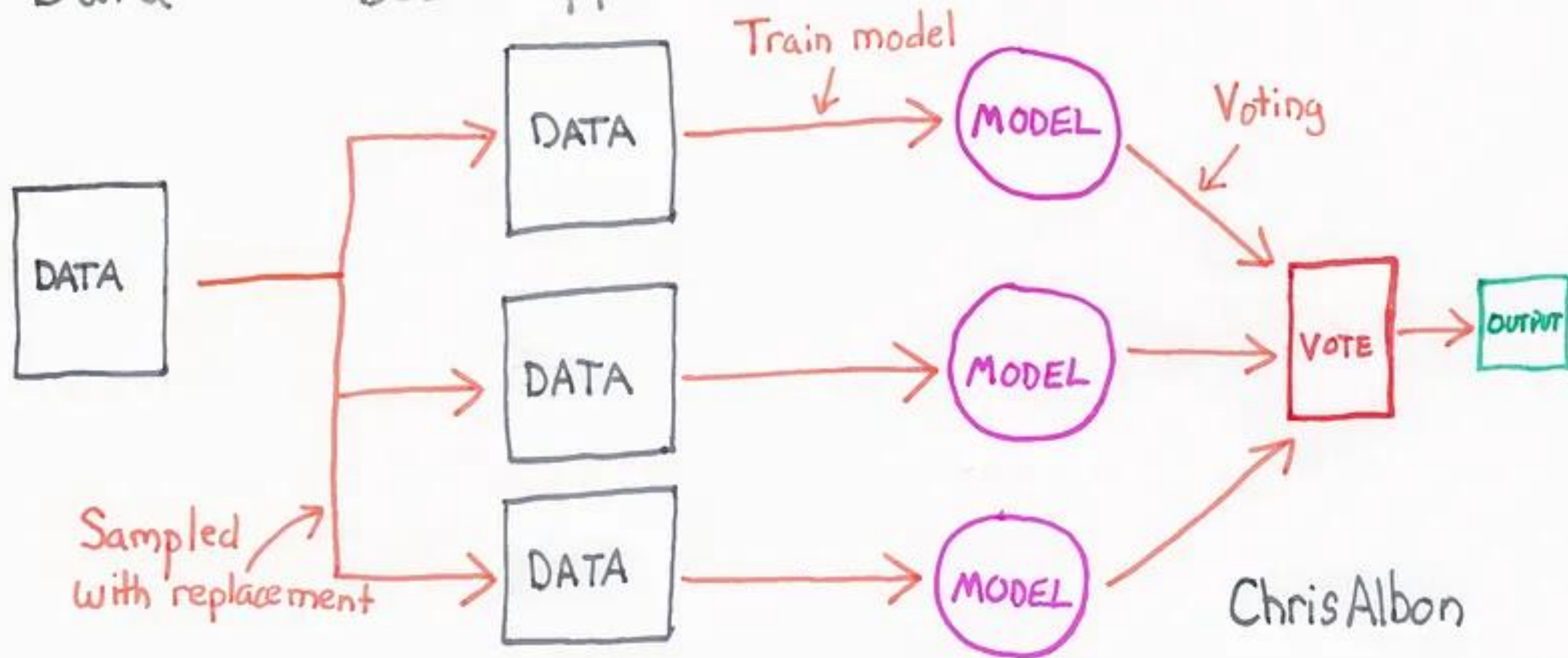
Data

Bootstrapped Data

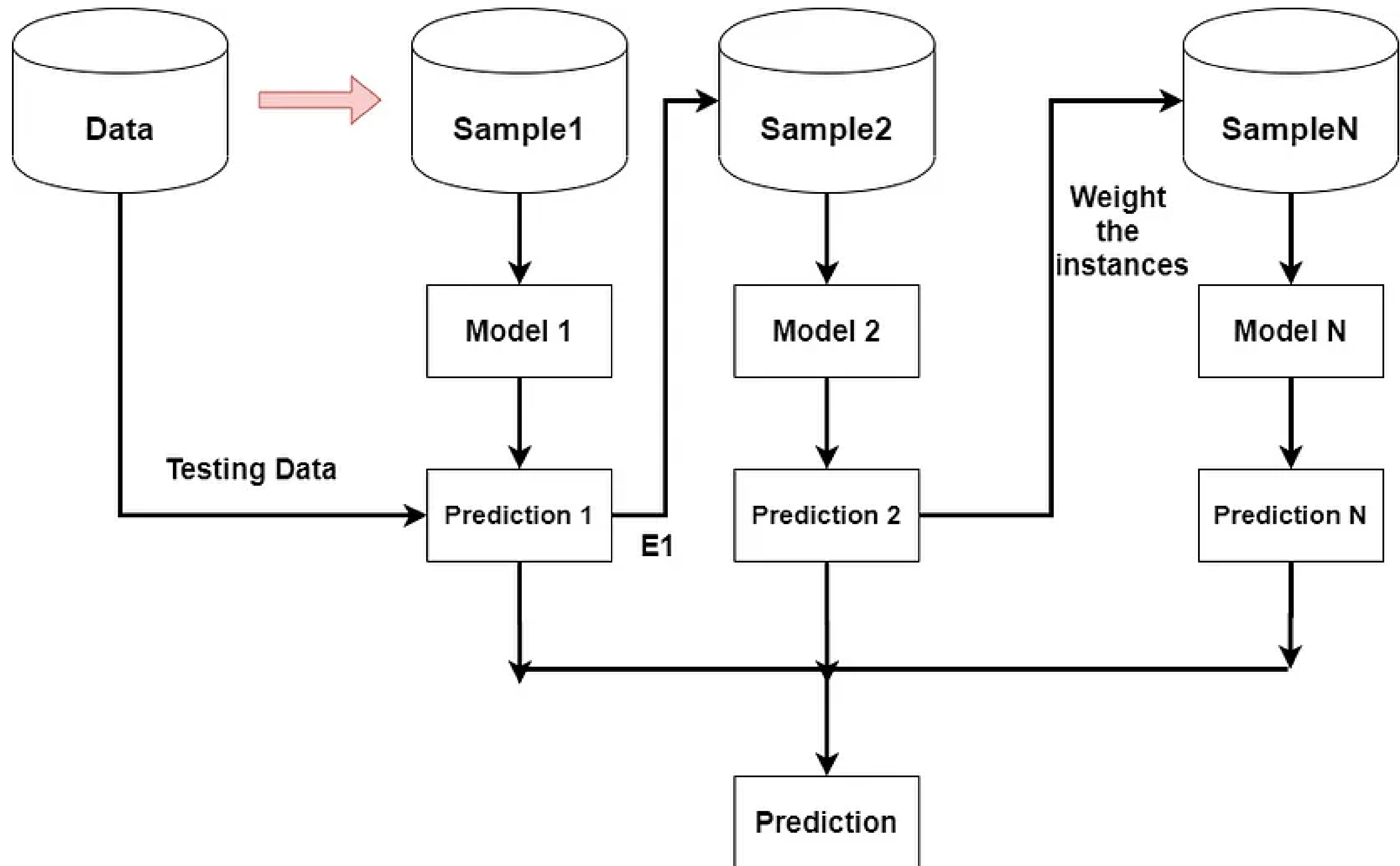
Models

Voting

Outcome



Chris Albon

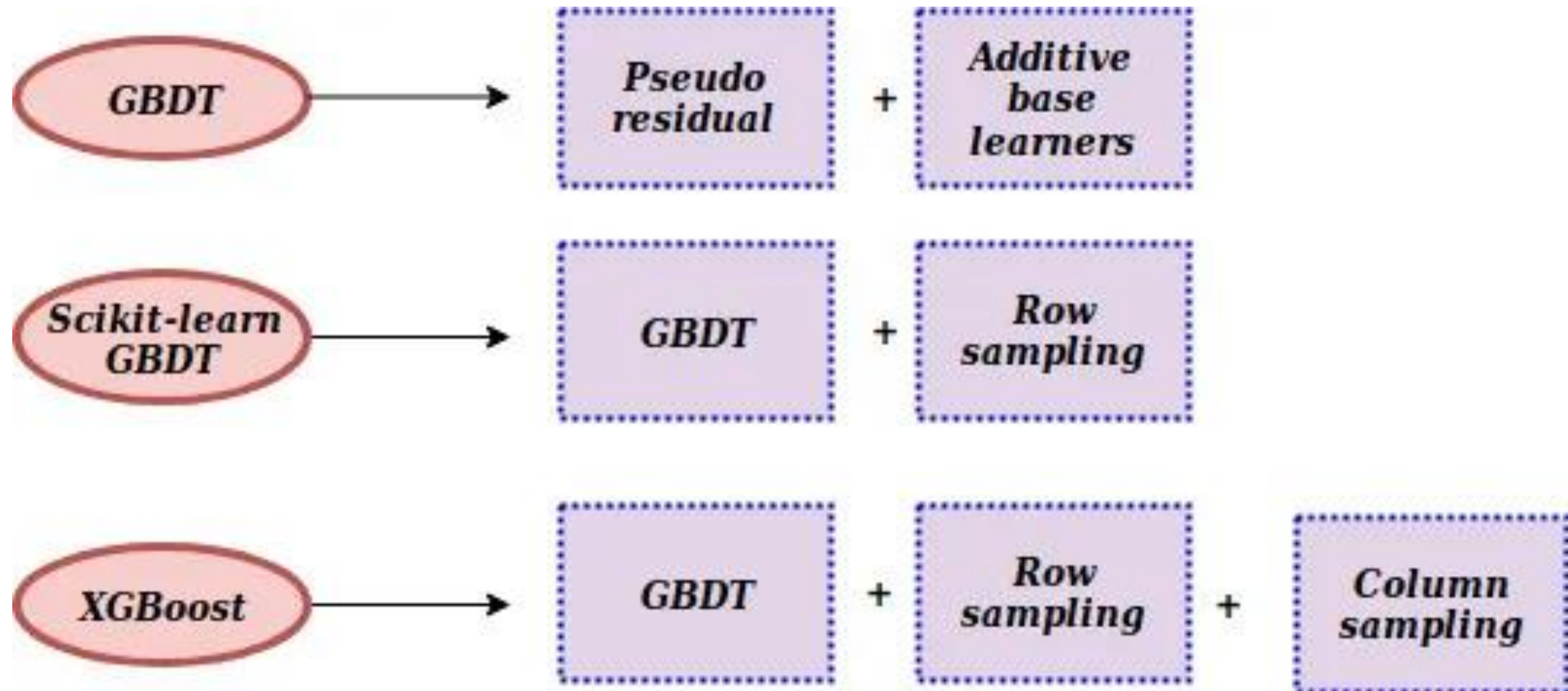


Boosting and Bias/Variance

- Boosting reduces Bias
- Hence low variance models are used as individual learner
 - E.g. Decision Tree with very less depth
- Combine high bias low variance models (weak learners) to make low bias low variance models

Boosting

- AdaBoost
- GradientBoost
- XGBoost



Helpful videos on boosting (Optional)

- AdaBoost
 - <https://www.youtube.com/watch?v=LsK-xG1cLYA>
- Gradient Boost (4 parts)
 - <https://www.youtube.com/watch?v=3CC4N4z3GJc>
- XGBoost (4 parts)
 - <https://www.youtube.com/watch?v=OtD8wVaFm6E>



QUESTIONS