

# Ensemble Learning

Concepts and Algorithms of Artificial Intelligence (WS2022)

**Lecturer:** Sharwin Rezagholi

**Authors:** Stefan Lackner, Bernhard Knapp, Sharwin Rezagholi



## Clustering

k-means  
Hierarchical clustering  
DB-scan

## Regression

KNN regression  
Regression trees  
Linear regression  
Multiple regression  
Ridge and Lasso regression  
Neural networks

## Classification

KNN classification  
Classification trees  
Ensembles & boosting  
Random Forest  
Logistic regression  
Naive Bayes  
Support vector machines  
Neural networks

## Supervised learning

# AI

## Data handling

EDA  
Data cleaning  
Feature selection  
Class balancing  
etc

## Non-supervised learning

## Dimensionality reduction

PCA / SVD  
tSNE  
MDS

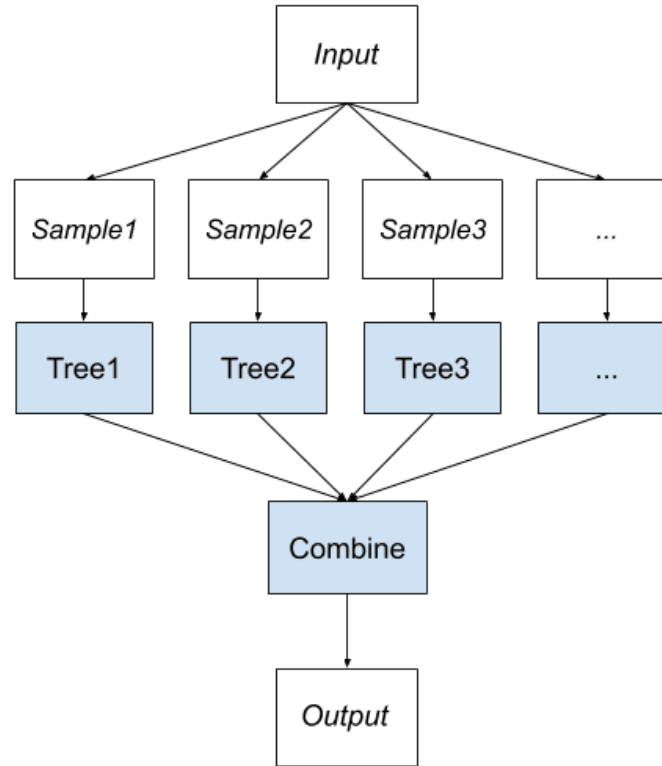
## Reinforcement learning

Covered in a separate lecture.

# Contents

- **Motivation and Demotivation**
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- Simple Voting
- Bagging and Pasting
- Random Patches and Random Subspaces
- 2 Popular Ensemble Techniques (Overview)
- Stacking
- Recap & Exercises

# Ensemble Learning: Basic Idea



# Motivation

- Ensemble learning is a process of **combining multiple models** (classifiers or regressors) to solve an ML problem.
- Ensembles (especially random forests) and boosting **are powerful** and deliver competitive results.
- Ensembles are **conceptually easy to understand**.
- Given sufficient computational power, ensemble learning is **easy to implement**.
- The bagging method of ensemble learning is **easy to parallelize**.
- There are many ensemble methods with different aims and properties.

# Demotivation

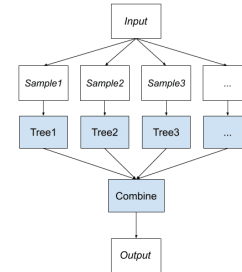
- Ensemble methods **increase training and test time.**
- **Performance increases level out.** In some domains (e.g. computer vision, natural language) best results are obtained by deep learning, usually without ensemble methods.
- Ensemble methods turn interpretable base-learners (more on terminology later) into **black box models**. Additional analytical methods are needed to extract „the meaning“ of an ensemble.

# Contents

- Motivation and Demotivation
- **Ensemble Learning (Basic Idea)**
- Bias-Variance Decomposition
- Simple Voting
- Bagging and Pasting
- Random Patches and Random Subspaces
- 2 Popular Ensemble Techniques (Overview)
- Stacking
- Recap & Exercises

# Ensemble Learning

- The **combination of different models** to obtain a final result.
- “Different models” may refer to
  - **different algorithms** (e.g. trees, SVMs, ...),
  - **same algorithm trained differently** (e.g. trees trained on different subsets of the data, or on differently weighted data).
- The **output of all models** is used to obtain a final result.
- Single models are called **base-learners**.
- **Applicable to many problems:** Regression, classification, and clustering (consensus-clustering).



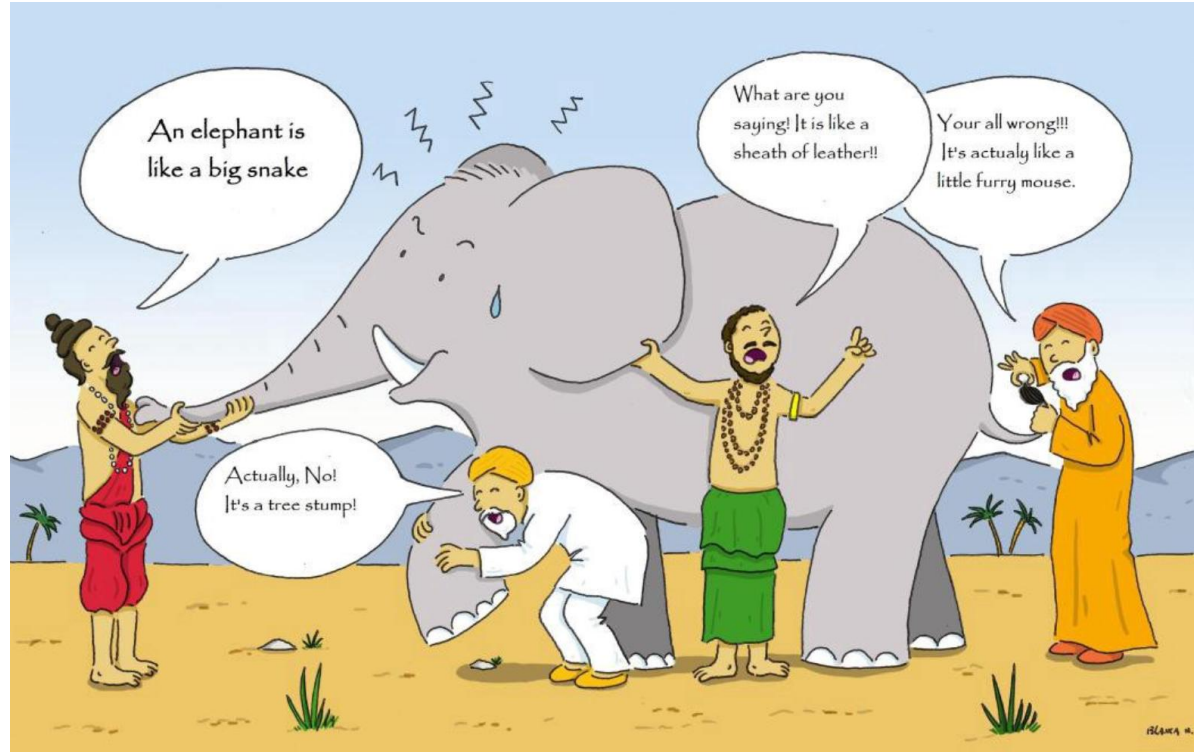


# Weak Base Learners vs. Strong Ensembles

- In ensemble learning **base-learners do not need to be strong.**
- A **strong model** is a **high-accuracy-classifier/regressor.**
- **Weak learners may only be slightly better than chance** (e.g. only reaching 0.51 accuracy).
- **Combining many weak learners can lead to a strong ensemble** if base learners are sufficiently independent.

# Many weak learners together can do a great job

(... while one single learner might achieve less than the sum of many weak learners ...)



[<https://medium.com/ml-research-lab/ensemble-learning-the-heart-of-machine-learning-b4f59a5f9777>]

# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- **Bias-Variance Decomposition**
- Simple Voting
- Bagging and Pasting
- Random Patches and Random Subspaces
- 2 Popular Ensemble Techniques (Overview)
- Stacking
- Recap & Exercises

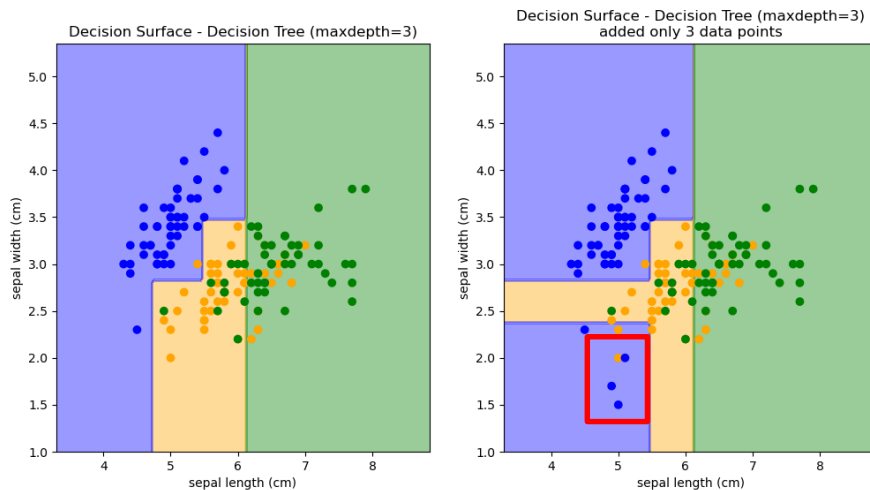
# Diversity of Classifiers

Diversity among base learner can be reached by:

(1) using **different algorithms**,

(2) using the same **high-variance algorithm** on different training sets (e.g. resampling).

Recall: A high-variance model changes much if the data changes little (example: decision trees).



High variance:  
3 new datapoints lead to a very different tree.

# Bias-Variance Decomposition

- Consider a regression.
- We make errors when using the fitted model for prediction.
- If we use a sufficiently large test set we can approximate the expected mean squared error (MSE) of the model.
- The expected MSE on test data can be **additively decomposed into 3 quantities**:
  - (1) The **variance of the model**,
  - (2) the **squared bias**,
  - (3) the **variance of the error**.

# Bias-Variance Decomposition

Error = Variance + Bias + Noise

$$\mathbb{E}\left(y_i - \hat{f}(x_i)\right)^2 = \text{Var}\left(\hat{f}(x_i)\right) + \text{bias}\left(\hat{f}(x_i)\right)^2 + \text{Var}(\epsilon)$$

$(x_i, y_i)$ , a tuple from the test set

$\mathbb{E}\left(y_i - \hat{f}(x_i)\right)^2$ , expected MSE on test data if model is trained on different training sets

$\text{Var}\left(\hat{f}(x_i)\right)$ , variance of the model (how different the model is when using different training sets)

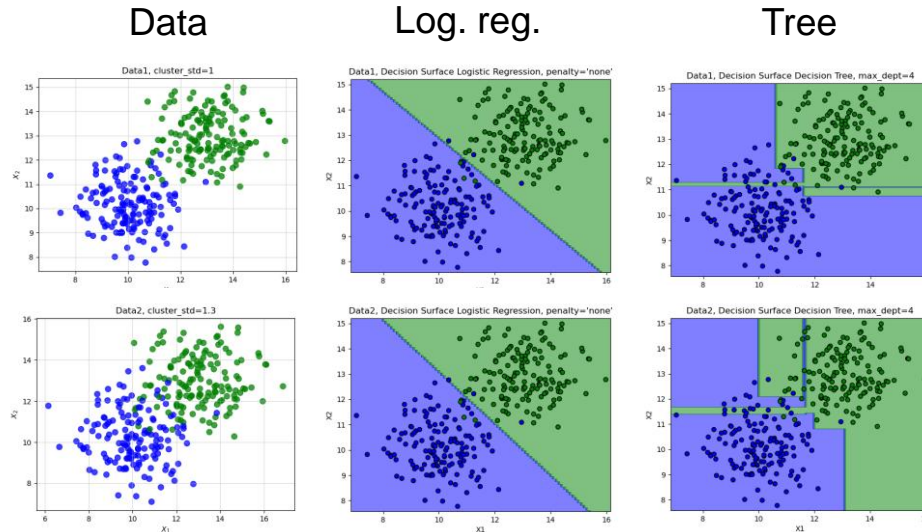
$\text{bias}\left(\hat{f}(x_i)\right)$ , error due to choice of model class (example: linear model for non-linear data-generating process)

$\text{Var}(\epsilon)$ , irreducible variance of the error term

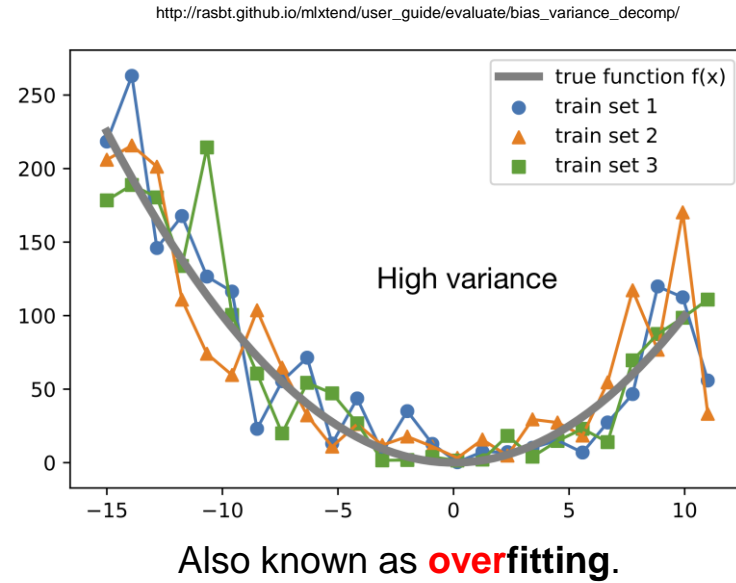
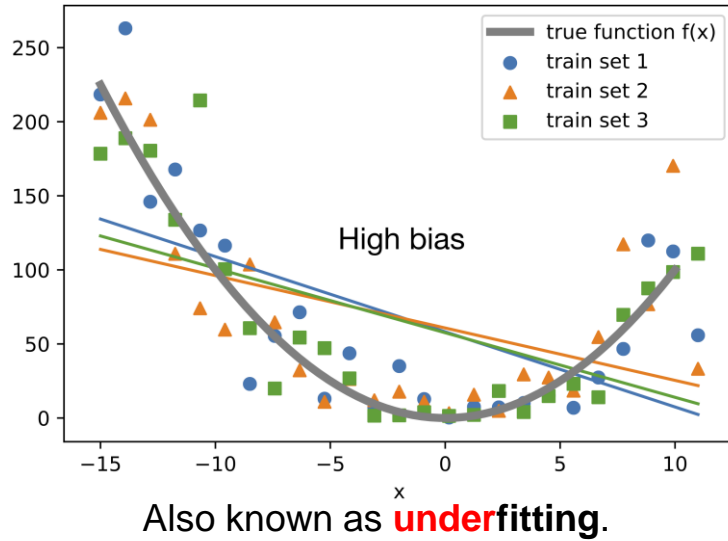
*We cannot do anything about the **irreducible** part but we can handle **bias** and **variance**.*

# High vs. Low-Variance Models

- Decision tree:  $\mathbb{E}(y_i - \hat{f}(x_i))^2 = \text{Var}(\hat{f}(x_i)) + \text{bias}(\hat{f}(x_i))^2 + \text{Var}(\epsilon)$
- Logistic regression:  $\mathbb{E}(y_i - \hat{f}(x_i))^2 = \text{Var}(\hat{f}(x_i)) + \text{bias}(\hat{f}(x_i))^2 + \text{Var}(\epsilon)$



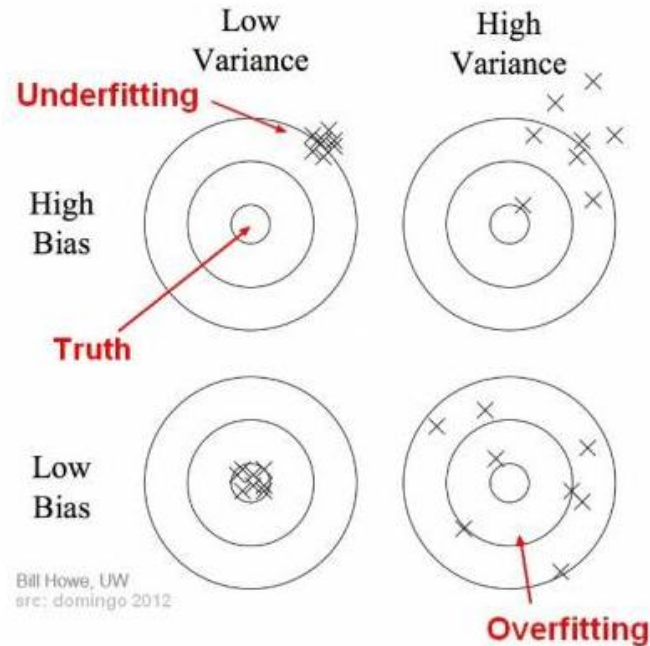
# High vs. Low-Variance Models



- The **bias** stems from **erroneous assumptions in the learning algorithm**. High bias can cause an algorithm to miss the relevant relations between features and target outputs (**underfitting**).
- The **variance** stems from **sensitivity to small changes** in the training data. High variance may result from an algorithm modeling the noise in the training data (**overfitting**).



# High vs. Low-Variance Models



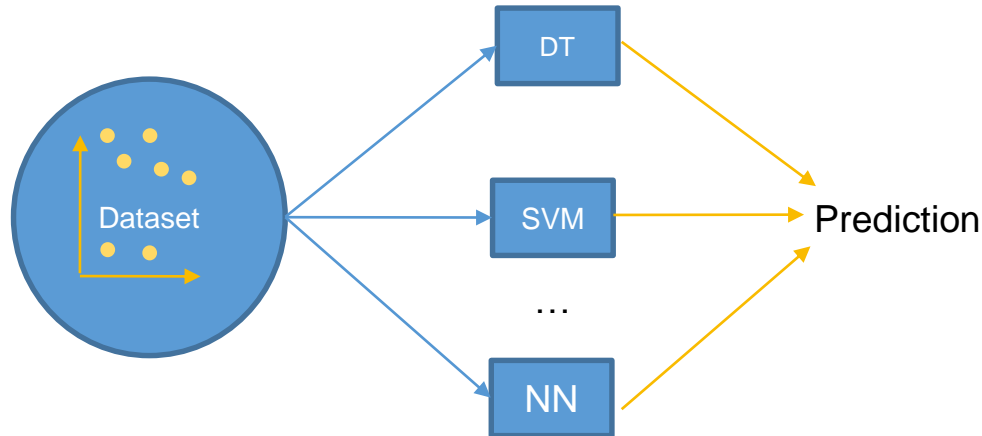
# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- **Simple Voting**
- Bagging and Pasting
- Random Patches and Random Subspaces
- 2 Popular Ensemble Techniques (Overview)
- Stacking
- Recap & Exercises

# Ensemble Method: Voting and Averaging

Combine the results of several classifiers by “letting them vote”.

- **Classification:**
  - **Hard voting:** Predict the plurality class.
  - **Soft voting:** Take the average of the estimated class probabilities.
- **Regression:**



# Simple Voting & Bias Reduction

- The **combination of different classifiers by voting** leads to a **reduction in bias** since different classifier algorithms are biased in different ways.
- The amount of reduction depends on the data and the difference in biases of single models.
- We assume here that each classifier is trained using the same training data. But we can do better (see next slides).

$$\mathbb{E}\left(y_i - \hat{f}(x_i)\right)^2 = \boxed{\text{Var}\left(\hat{f}(x_i)\right)} + \boxed{\text{bias}\left(\hat{f}(x_i)\right)^2} + \text{Var}(\epsilon)$$

↓ Voting  
ensemble

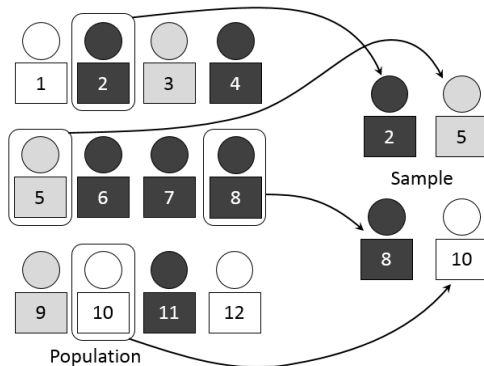
$$\mathbb{E}\left(y_i - \hat{f}(x_i)\right)^2 = \boxed{\text{Var}\left(\hat{f}(x_i)\right)} + \boxed{\text{bias}\left(\hat{f}(x_i)\right)^2} + \text{Var}(\epsilon)$$

# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- Simple Voting
- **Bagging and Pasting**
- Random Patches and Random Subspaces
- 2 Popular Ensembles Techniques (Overview)
- Stacking
- Recap & Exercises

# Bagging and Pasting

- We do not use all the training data at once.
- **Bagging** and **Pasting** are **resampling techniques**:
  - **Bagging (bootstrap aggregating)**: Repeatedly draw random samples from a training set **with replacement**.
  - **Pasting**: Draw random samples **without replacement** (each observation in the sample can be used once). Requires a sufficiently large dataset.
- Resampling techniques are used to **reduce variance**.



By Dan Kemler - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=36506020>

# Bagging and High-Variance Classifiers

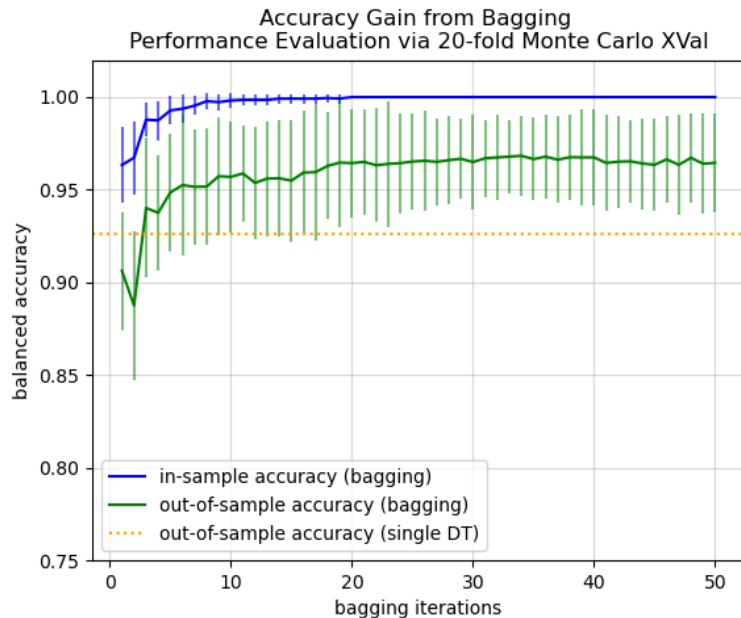
$$\mathbb{E}\left(y_i - \hat{f}(x_i)\right)^2 = \boxed{\text{Var}\left(\hat{f}(x_i)\right)} + \boxed{\text{bias}\left(\hat{f}(x_i)\right)^2} + \text{Var}(\epsilon)$$

↓ Bagging

$$\mathbb{E}\left(y_i - \hat{f}(x_i)\right)^2 = \boxed{\text{Var}\left(\hat{f}(x_i)\right)} + \boxed{\text{bias}\left(\hat{f}(x_i)\right)^2} + \text{Var}(\epsilon)$$

- Making your **base models as independent as possible** benefits the procedure.
- Increasing the variance of base-learners may **increase the generalization capability** of the ensemble.
- In practice one needs to **balance variance and the number of base-learners**.
- Low-variance models will not benefit (much) from bagging.

# Increased Generalization Capability due to Bagging

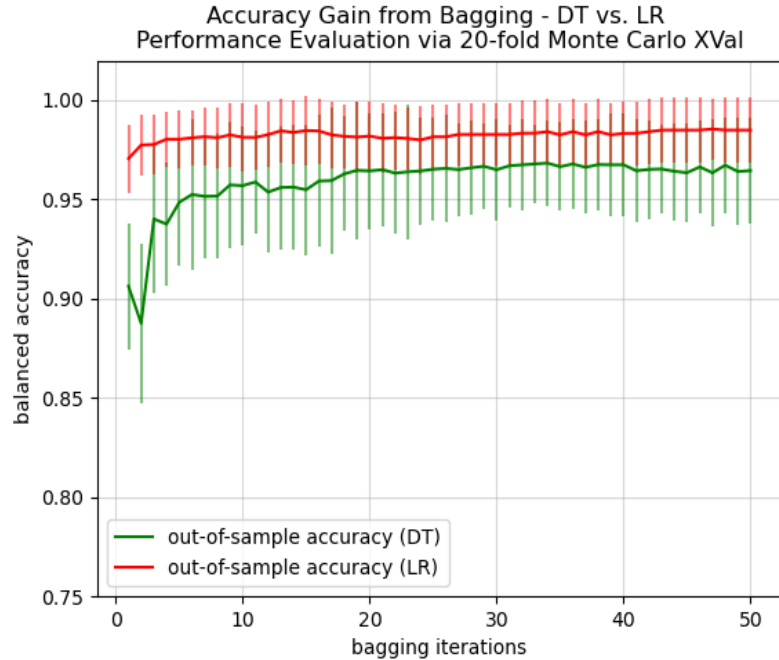


Note: The function `BaggingClassifier()` in SK-Learn performs soft voting if `predict_proba()` is available.

- **Gains from bagging arise quickly and taper out fast.**
- Training accuracy reaches 1.
- Test accuracy increased by 0.05 by bagging 20 base learners.
- No tuning was used for the single decision trees (default parameters are used).
- Unrestricted decision trees (which have very high variance) are used.



# Bagging: Decision Trees vs. Logistic Regressions



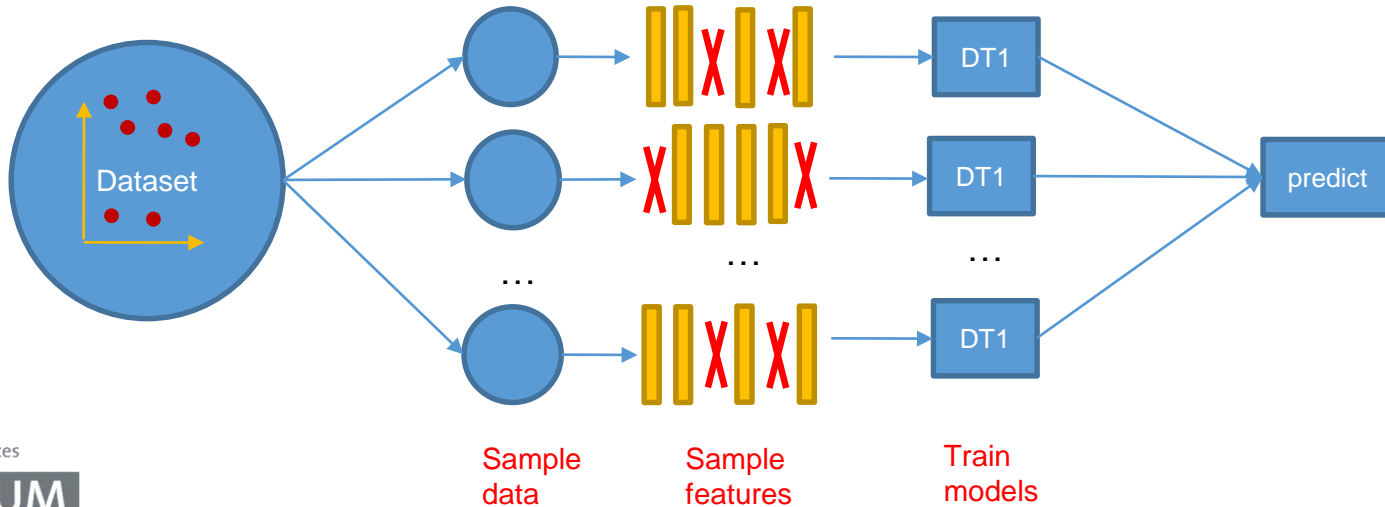
- Log. reg. performs better than bagged decision trees on this data.
- **Log. reg. is a low variance model**, hence it **does not gain much from bagging**.
- However: Be careful and try different things on your own datasets.

# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- Simple Voting
- Bagging and Pasting
- **Random Patches and Random Subspaces**
- 2 Popular Ensemble Techniques (Overview)
- Stacking
- Recap & Exercises

# Random Patches and Random Subspaces

- Bagging & Pasting reduce variance and work best when base learners are diverse.
- Idea: Increase diversity further by **sampling not only the data but also the features used in training base learners?**
- Doing this **in combination** with bagging is called **random patches**. Doing this **in isolation**, this is called **random subspaces**.



# Summary

- **Bagging and pasting:**
  - Resampling from data (observations).
- **Random Subspaces:**
  - Resampling features (columns).
  - Sometimes called “feature bagging”.
- **Random Patches:**
  - Resampling from data and features (observations and features).
  - Sometimes called “random subspace plus bagging”.

	A	B	C	D
1	7.9	3.8	6.4	2.0
2	5.8	2.7	5.1	1.9
3	5.1	3.8	1.6	0.3
4	4.4	2.9	1.4	0.2
5	5.1	3.8	1.6	0.2
6	5.4	3.4	1.7	0.2
7	7.0	3.2	4.7	1.4
8	6.7	3.3	5.7	2.1
9	6.7	3.1	4.4	1.4
10	5.5	2.5	4.0	1.3
11	5.1	3.5	1.4	0.3
12	5.5	4.2	1.4	0.2
13	4.8	3.4	1.6	0.2
14	5.4	3.0	4.5	1.5
15	6.1	2.6	5.6	1.4
16	6.7	2.5	5.8	1.8
17	6.3	2.8	5.1	1.5
18	4.6	3.6	1.0	0.2
19	6.0	2.9	4.5	1.5
20	6.1	3.0	4.6	1.4
21	4.7	3.2	1.3	0.2
22	5.8	2.8	5.1	2.4
23	5.9	3.0	4.2	1.5
24	6.5	3.2	5.1	2.0
25	6.1	2.8	4.7	1.2
26	5.8	2.6	4.0	1.2
27	6.3	2.7	4.9	1.8

Bagging

	A	B	C	D
1	7.9	3.8	6.4	2.0
2	5.8	2.7	5.1	1.9
3	5.1	3.8	1.5	0.3
4	4.4	2.9	1.4	0.2
5	5.1	3.8	1.6	0.2
6	5.4	3.4	1.7	0.2
7	7.0	3.2	4.7	1.4
8	6.7	3.3	5.7	2.1
9	6.7	3.1	4.4	1.4
10	5.5	2.5	4.0	1.3
11	5.1	3.5	1.4	0.3
12	5.5	4.2	1.4	0.2
13	4.8	3.4	1.6	0.2
14	5.4	3.0	4.5	1.5
15	6.1	2.6	5.6	1.4
16	6.7	2.5	5.8	1.8
17	6.3	2.8	5.1	1.5
18	4.6	3.6	1.0	0.2
19	6.0	2.9	4.5	1.5
20	6.1	3.0	4.6	1.4
21	4.7	3.2	1.3	0.2
22	5.8	2.8	5.1	2.4
23	5.9	3.0	4.2	1.5
24	6.5	3.2	5.1	2.0
25	6.1	2.8	4.7	1.2
26	5.8	2.6	4.0	1.2
27	6.3	2.7	4.9	1.8

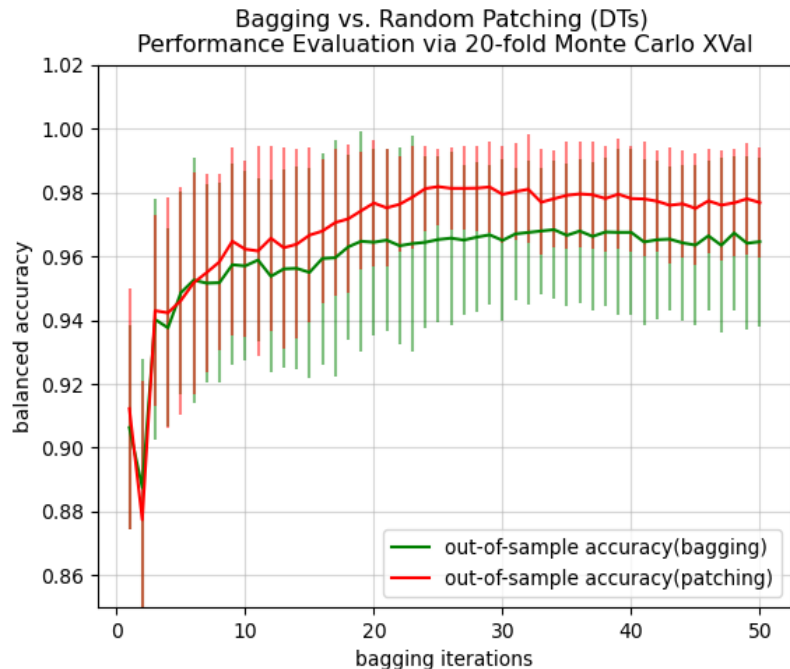
Random subspace

	A	B	C	D
1	7.9	3.8	6.4	2.0
2	5.8	2.7	5.1	1.9
3	5.1	3.8	1.6	0.2
4	4.4	2.9	1.4	0.2
5	5.1	3.8	1.7	0.2
6	5.4	3.4	1.7	0.2
7	7.0	3.2	4.7	1.4
8	6.7	3.3	5.7	2.1
9	6.7	3.1	4.4	1.4
10	5.5	2.5	4.0	1.3
11	5.1	3.5	1.4	0.3
12	5.5	4.2	1.4	0.2
13	4.8	3.4	1.6	0.2
14	5.4	3.0	4.5	1.5
15	6.1	2.6	5.6	1.4
16	6.7	2.5	5.8	1.8
17	6.3	2.8	5.1	1.5
18	4.6	3.6	1.0	0.2
19	6.0	2.9	4.5	1.5
20	6.1	3.0	4.6	1.4
21	4.7	3.2	1.3	0.2
22	5.8	2.8	5.1	2.4
23	5.9	3.0	4.2	1.5
24	6.5	3.2	5.1	2.0
25	6.1	2.8	4.7	1.2
26	5.8	2.6	4.0	1.2
27	6.3	2.7	4.9	1.8

Random patch

These illustrations are very simplified. Random choice is key in resampling.

# Bagging vs. Random Patches



- Patching increases the diversity of base learners.
- Increased diversity leads to higher gains from bagging.
- Since the diversity of base learners is higher, gains from bagging taper out later (roughly at ensemble size 25).

# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- Simple Voting
- Bagging and Pasting
- Random Patches and Random Subspaces
- **2 Popular Ensemble Techniques (Overview)**
- Stacking
- Recap & Exercises

# Random Forests

- Random forests are **ensembles of unrestricted decision trees** obtained by a **variant of random patches**.
- Instead of sampling features after bootstrapping a sample, **features for spitting are resampled at every node** (for every split).
- The intensity of resampling increases from bagging, to random patches, to random forests.
- Random forests will be the focus of a dedicated lecture in this course.

# Boosting

Base learners are trained **sequentially** (not in parallel as in bagging).

Aim: Reduction of both, bias and variance.

Two approaches:

1. Critical datapoints are reweighted to emphasize “complicated” observations.
2. Single models are fitted to the residuals of the previous models.

Boosting will be the focus of a dedicated lecture in this course.



# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- Simple Voting
- Bagging and Pasting
- Random Patches and Random Subspaces
- 2 Popular Ensemble Techniques (Overview)
- **Stacking**
- Recap & Exercises

# Stacking

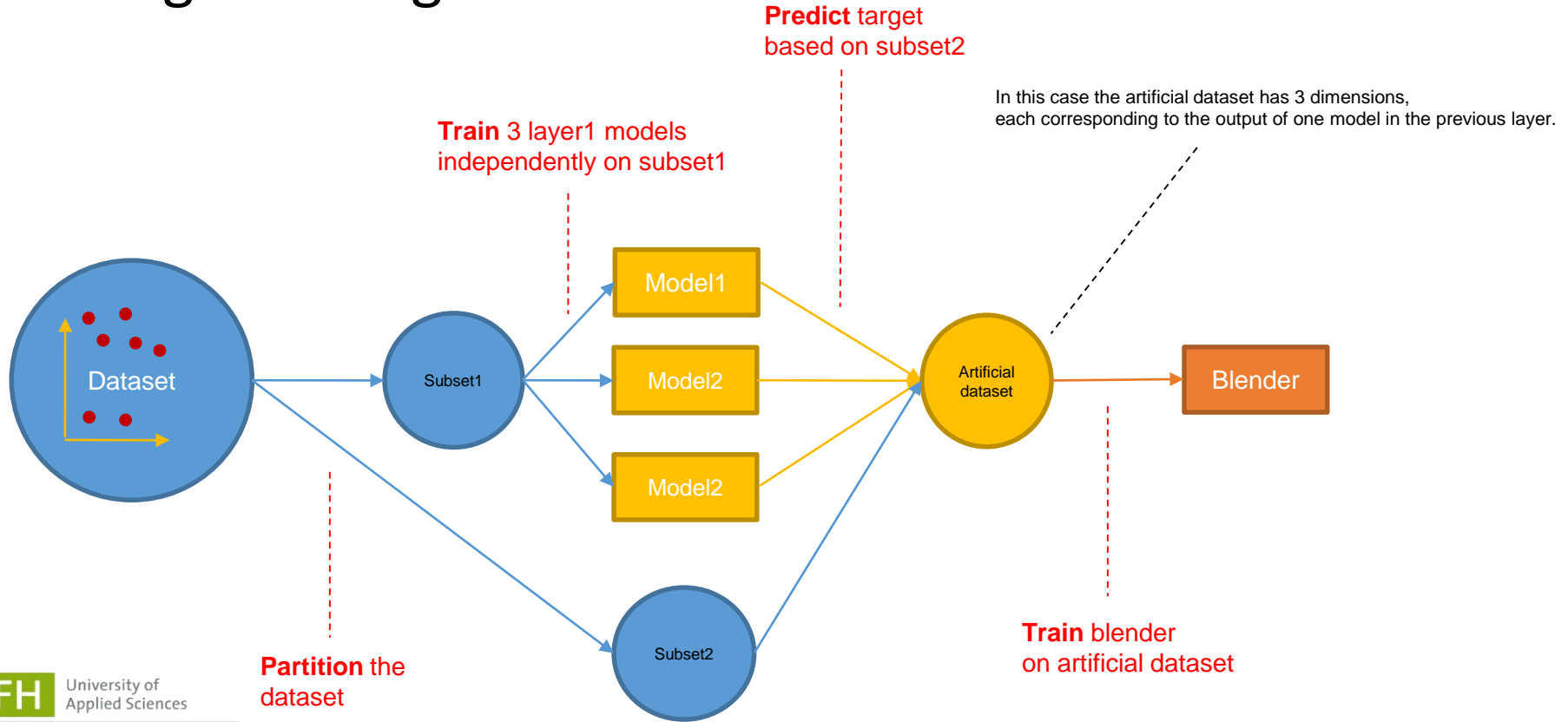
- Stacking is an **extension of the voting classifier/regressor** that **replaces plurality-vote/average by an additional classifier/regressor called the blender**.
- This is **vaguely similar to the intuition behind deep learning**: Layers of predictors are learned sequentially, each layer takes the predictions of the earlier layer as input.
- As in voting methods different models classes are combined.

# Stacking

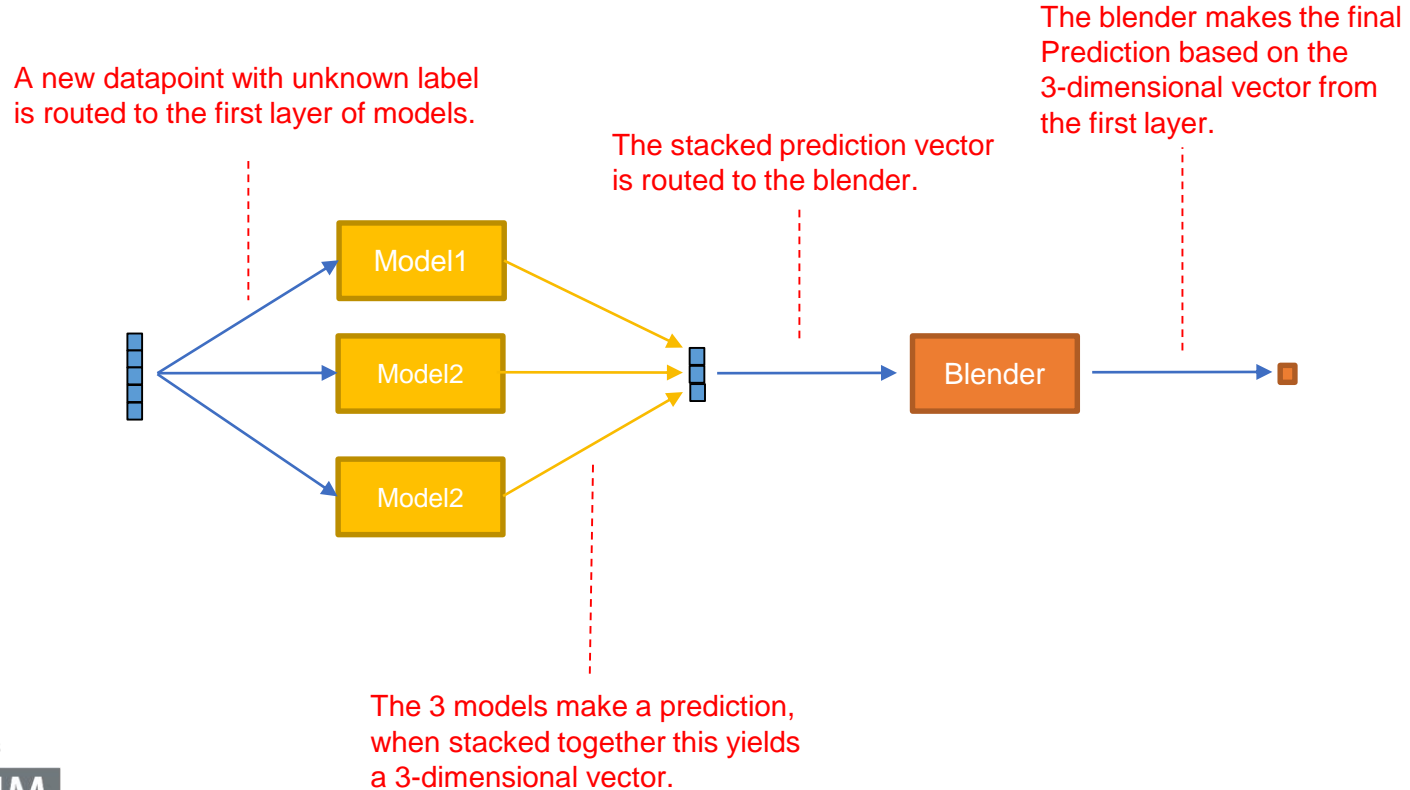
Procedure:

- 1) On a **dedicated data subset**, several **base learners are fitted** (perhaps using different algorithms).
  - 2) On a **second data subset**, the **blender** is trained. The blender is another model trained on the predictions of the other classifiers.
- 
- By partitioning the training data into  $m$  subsets, one can stack  $m$  layers.
  - **Partitioning is necessary** to ensure that inputs to subsequent layers are not overfit.

# Stacking: Training



# Stacking: Prediction



# Stacking: Regression vs. Classification

- In regression models the first layer contains models that output real-valued predictions.
- In classification models the first layer contains models that output **categorical values**.
- In classification using SK-Learn **use `predict_proba()`** in layer1.

# Parallel vs. Sequential Implementation

- Bagging, Pasting, Random Patches, and Random Subspaces can be implemented in **parallel**, therefore they are **easily scalable**.
- Boosting and Stacking are **sequential** in nature: **Scaling is computationally expensive**.

# Black-Box vs. White-Box

- Using ensembles, (simple) base learners can be combined into powerful models.
- The downside is that their **interpretability is lost** to a large extent. They become **black-box models**.
- Example: While a small decision tree is a white-box model, interpreting the behavior of an ensemble of decision tree



# Contents

- Motivation and Demotivation
- Ensemble Learning (Basic Idea)
- Bias-Variance Decomposition
- Simple Voting
- Bagging and Pasting
- Random Patches and Random Subspaces
- 2 Popular Ensemble Techniques (Overview)
- Stacking
- **Recap & Exercises**

# Recap

- **Ensemble learning is powerful.** It is the only way to obtain models for certain datasets.
- There is **no guarantee that ensembles will improve performance.** Experience shows that simple **voting** and **stacking** are less likely to increase performance than **ensembles based on resampling.**
- Ensembles that allow for parallel estimation are computationally efficient

# Assignment: Random patches.

a) Explain (and compare, for example in a table) a plurality-voting ensemble, bagging/pasting, random patches, random subspaces, stacking, and boosting.

Use simple self-made images or even hand drawings (of which you take a photo).

Use self-written explanations. Do not copy from the lecture slides or the internet (neither text nor images).

b) Implement a version of random patches without the use of a library. Use a base learner of your choice (e.g. a tree) and compare the performance with a single learner.

Use a dataset of your choice. You will need a dataset with a large number of observations and features.

# References

- Géron A. (2017): Hands-On Machine Learning with Scikit-Learn & Tensorflow. – O'Reilly.
- James G., Witten D., Hastie T., Tibshirani R. (2017): An introduction to Statistical Learning. – Springer.
- Kuhn M., Johnson K. (2016): Applied Predictive Modeling. – Springer.
- Berk R. (2016): Statistical Learning from a Regression Perspective. – Springer.
- [Bagging predictors | SpringerLink](#)
- [The random subspace method for constructing decision forests | IEEE Journals & Magazine | IEEE Xplore](#)