



CAS ETH Machine Learning in Finance and Insurance

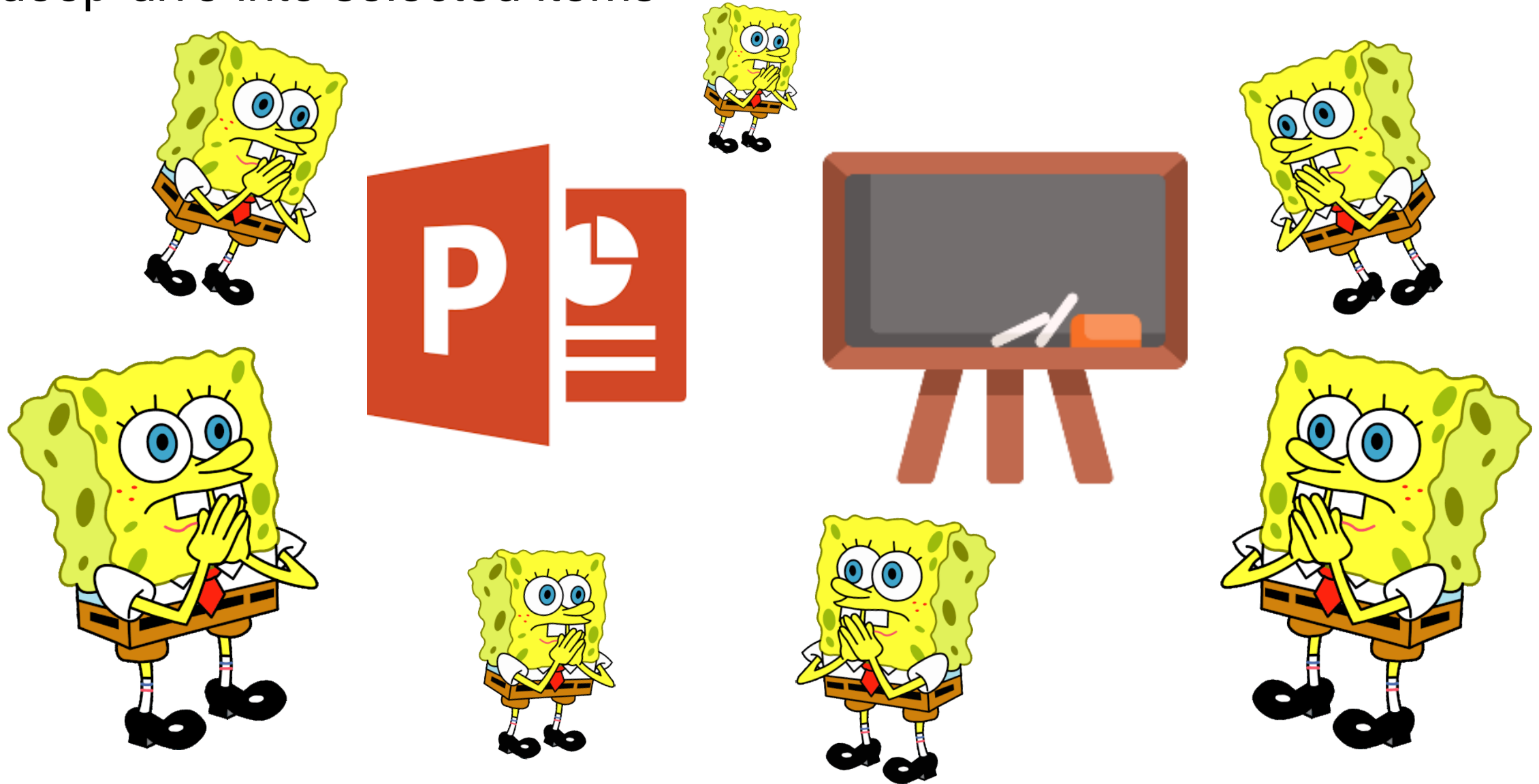
BLOCK I. Introduction to Machine Learning. Lecture 7.

Dr. A. Ferrario, ETH Zurich and UZH

On May 17th...

- 1 We introduced the decision trees, namely, an “interpretable” method to learn (binary) decisions from data
- 2 We discussed the famous Classification And Regression Tree algorithm by Breiman et al. (1984) in some detail and learned how to prune trees
- 3 We showed that decision trees are prone to overfitting and show high variance

We will use slides to introduce our topics and the blackboard to deep-dive into selected items



Machine Learning Methods (Part 7)

- Bagging and Random Forests
- Boosting

Decision trees are easy to train and “interpretable”. However...

...**they are prone to overfitting** (large trees are needed to capture complex, nonlinear decision boundaries) and **tend to show high variance**

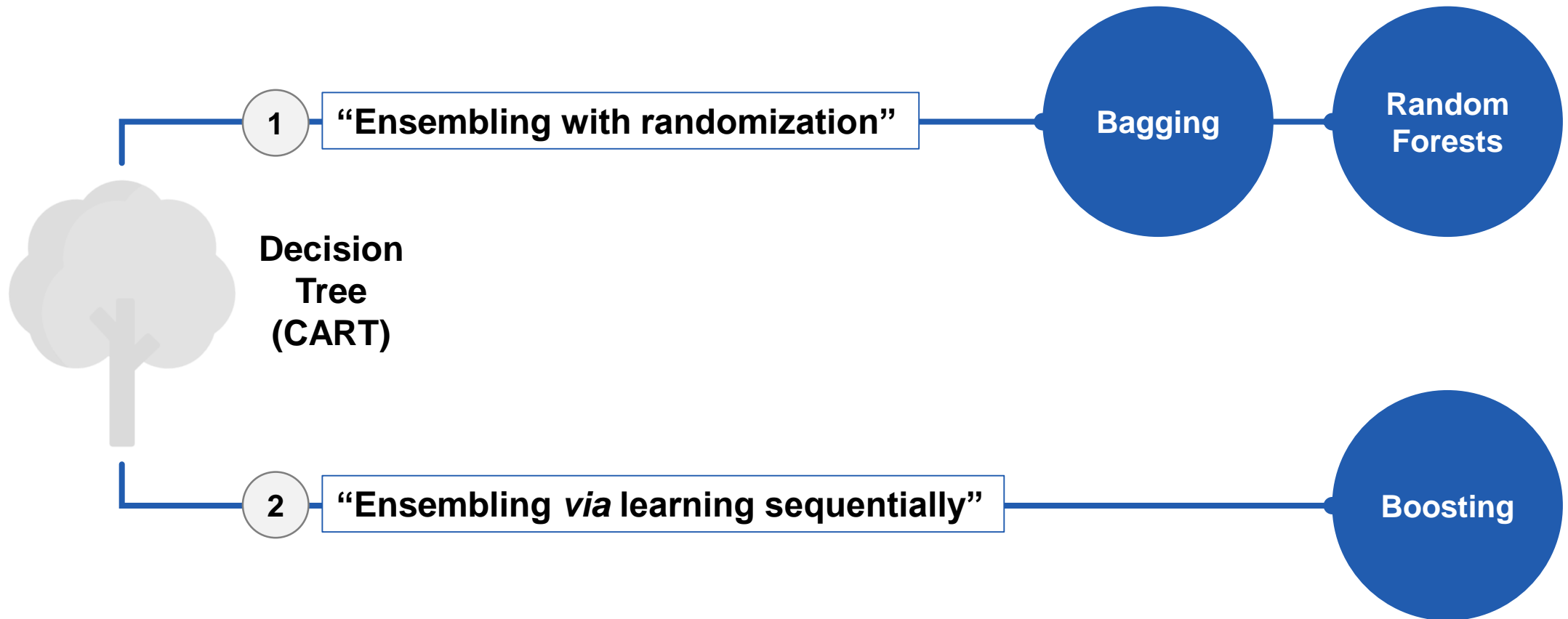
IDEA: “Wisdom of the crowd”

If we **appropriately combine many trees**, we may end up with a better model



Ensembling

The plan for today: Ensembling



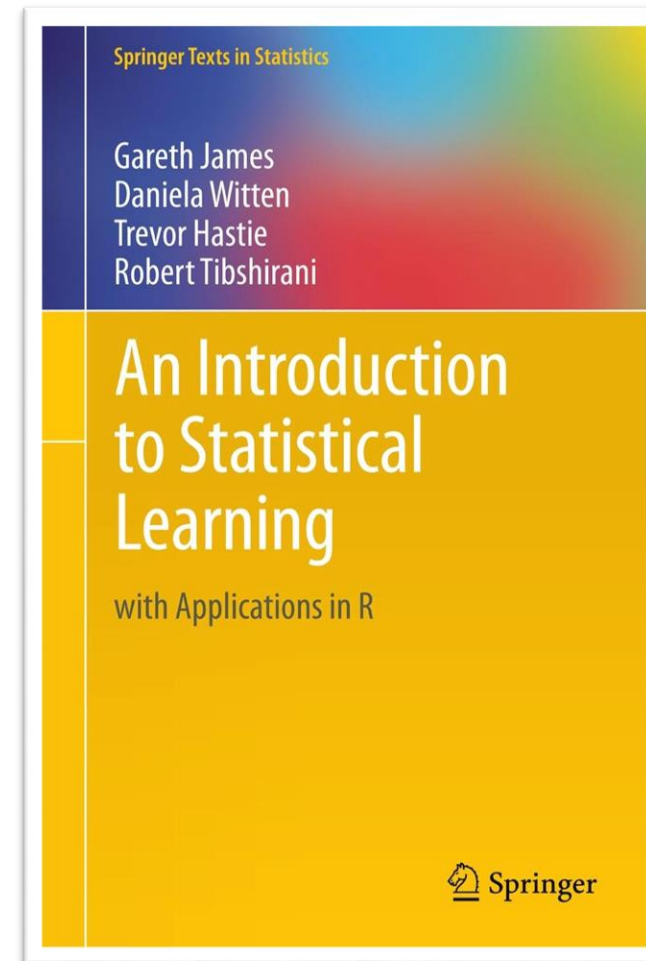
“Ensembling with randomization”

Bagging and Random Forests

Bagging = Bootstrap aggregating

A simple method to reduce the variance of decision trees

“[Bagging] is a general-purpose procedure for reducing the variance of a statistical learning method; we introduce it here because it is particularly useful and frequently used in the context of decision trees” (pag. 316)



James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning. New York: Springer

Bagging = Bootstrap aggregating

Bagging consists of two steps

We **bootstrap the original training set** as generally do not have access to multiple training datasets.

I

We generate B bootstrapped training datasets, **by randomly sampling with replacement B times from the original training data set.**

II

We train B decision trees on the B bootstrapped training datasets and we **aggregate their predictions.** The B trees are grown deep and are not pruned (high variance).

Formulae for regression and classification

Bagging = Bootstrap aggregating

Bagging consists of two steps

1

We **bootstrap the original training set** as generally do not have access to multiple training datasets.

I

We generate B bootstrapped training datasets, **by randomly sampling with replacement B times from the original training data set**.

II

We train B decision trees on the B bootstrapped training datasets and we **aggregate their predictions**. The B trees are grown deep and are not pruned (high variance).

(Bootstrap = random sampling with replacement)

	X	Y		X	Y		X	Y		X	Y
1			1			1			1		
2			2			2			N		
3			3			1			3		
...				
N-1			2			1			N		
N			N-1			N			N		
Original			Bootstrap 1			Bootstrap 2			Bootstrap 3		

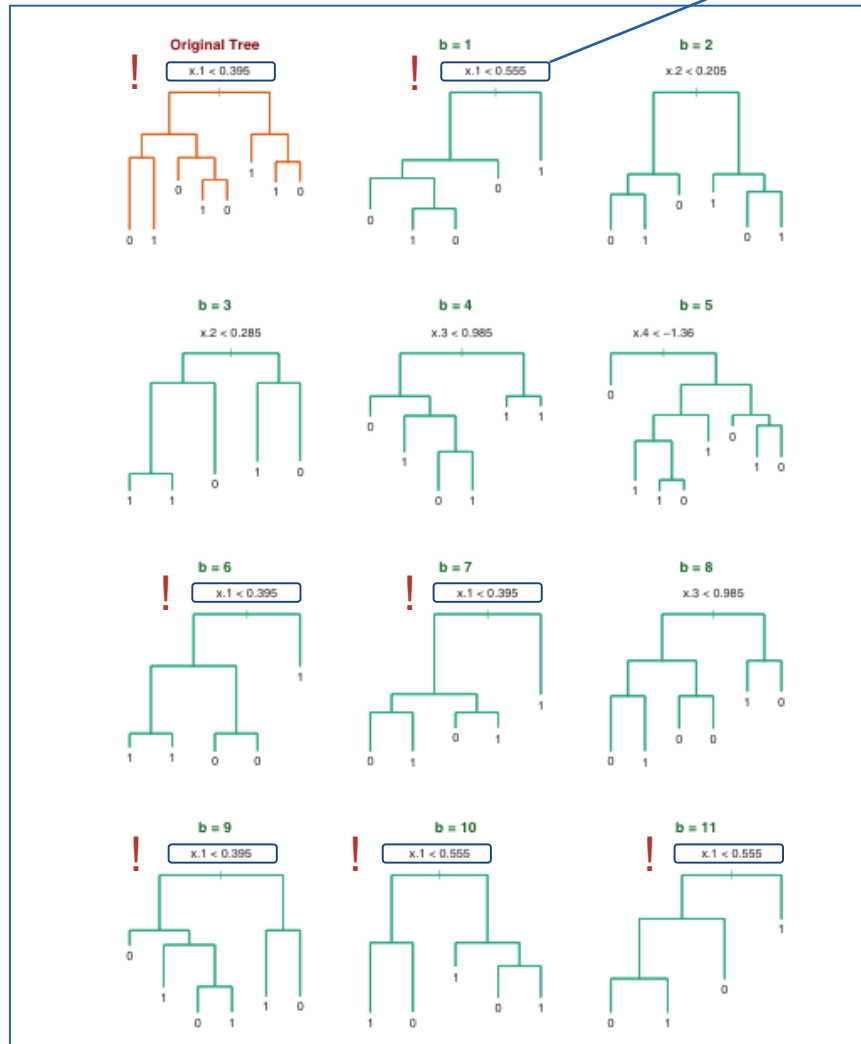
Formulae for regression and classification

Bagging = Bootstrap aggregating

Bagging consists of two steps

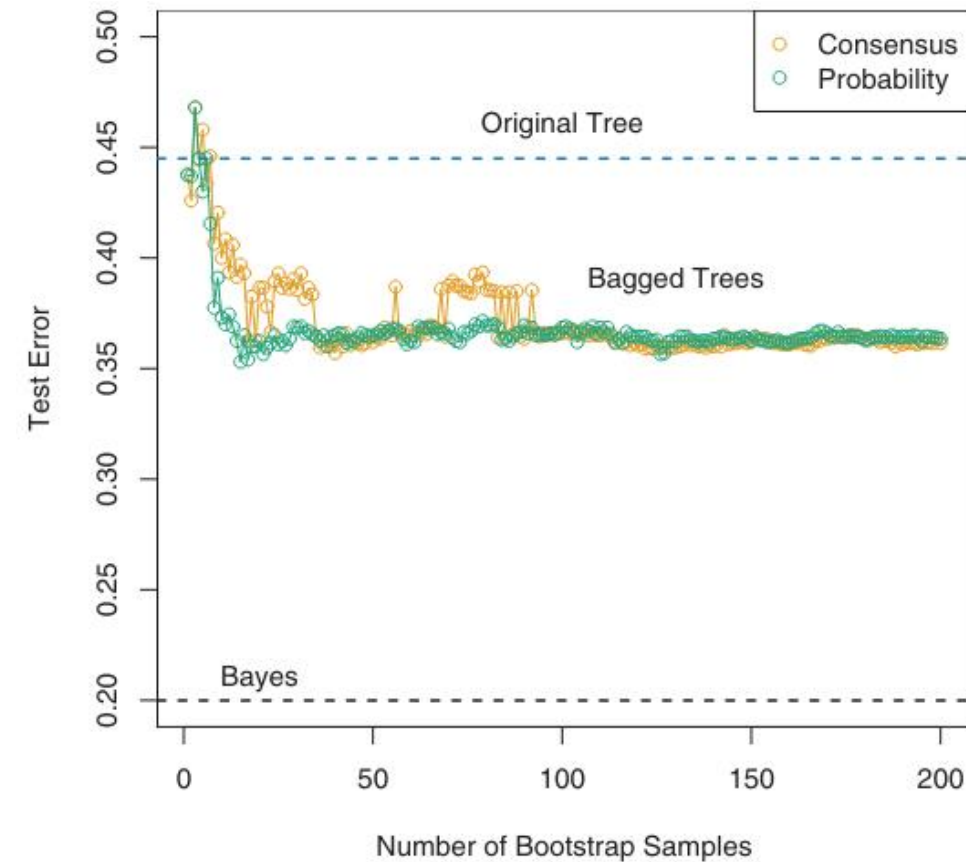
1

Source: Figure 8.9 in Hastie et al. (2009)*



Hey! 6 out of 11 bagged tree have a root using x_1 ...

Source: Figure 8.10 in Hastie et al. (2009)

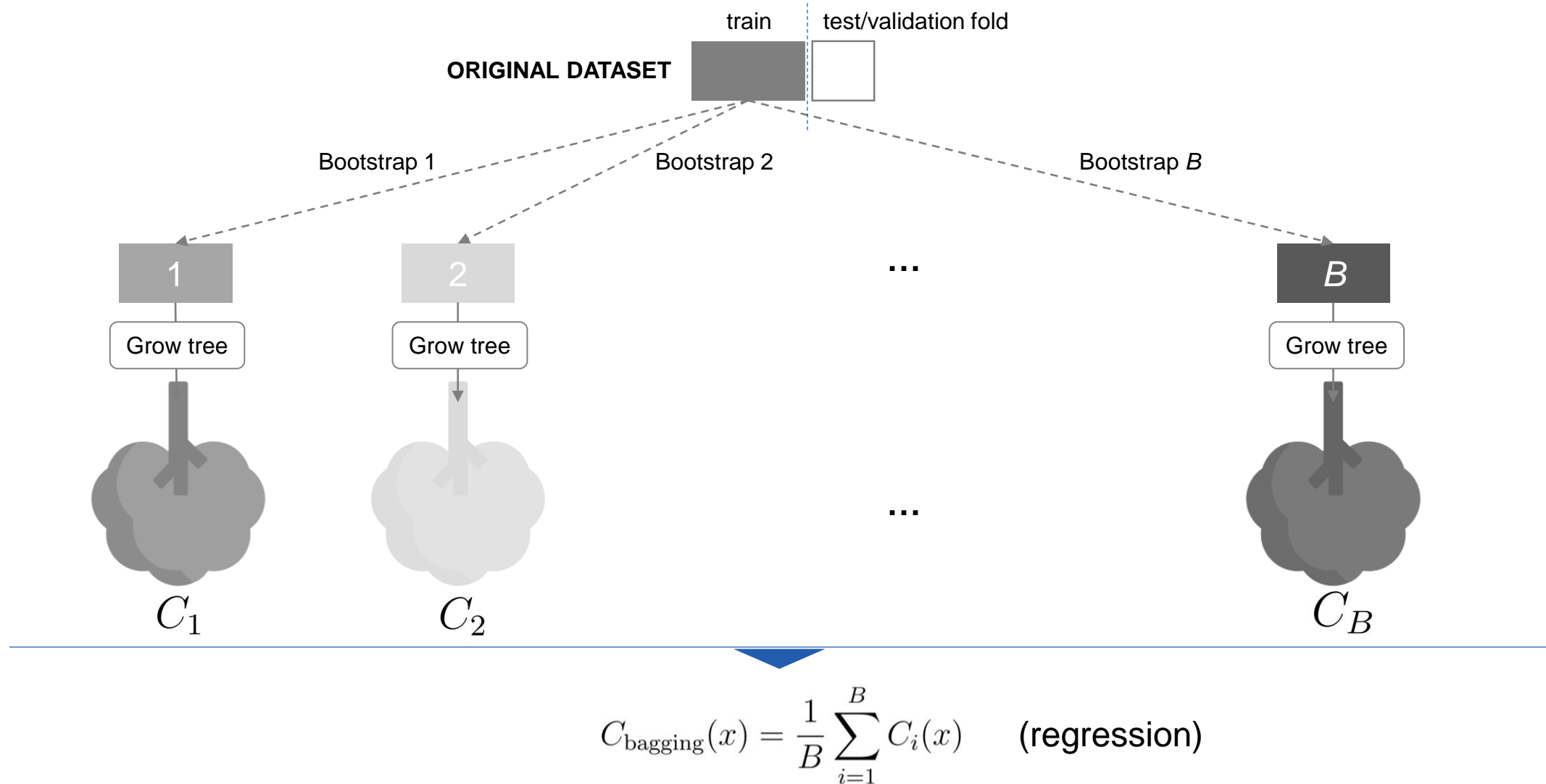


*Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction* (Vol. 2, pp. 1-758). New York: Springer.

Bagging: A visual summary

1

Choose how many bootstrapped training datasets: B



Bagging: some results from Breiman's original article

Classification trees on medium-sized datasets (in 1996)

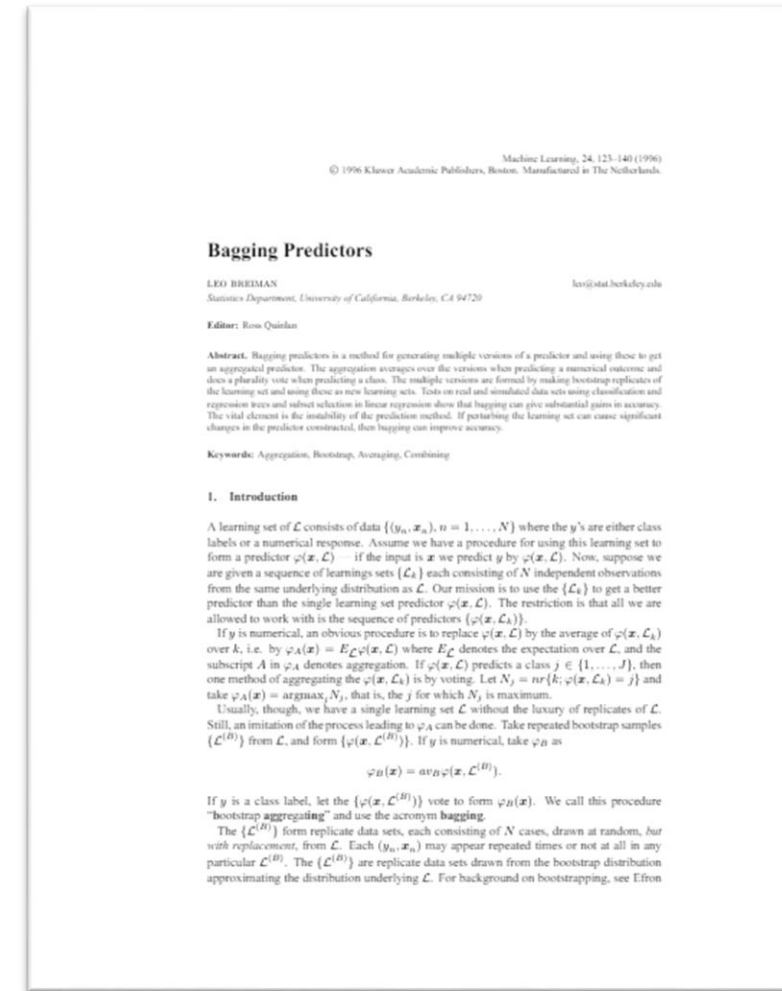
Data Set	# Samples	# Variables	# Classes
waveform	300	21	3
heart	1395	16	2
breast cancer	699	9	2
ionosphere	351	34	2
diabetes	768	8	2
glass	214	9	6
soybean	683	35	19

Table 2. Misclassification Rates (%)

Data Set	\bar{e}_S	\bar{e}_B	Decrease
waveform	29.1	19.3	34%
heart	4.9	2.8	43%
breast cancer	5.9	3.7	37%
ionosphere	11.2	7.9	29%
diabetes	25.3	23.9	6%
glass	30.4	23.6	22%
soybean	8.6	6.8	21%

Error
classification
tree

Error bagging
ensemble

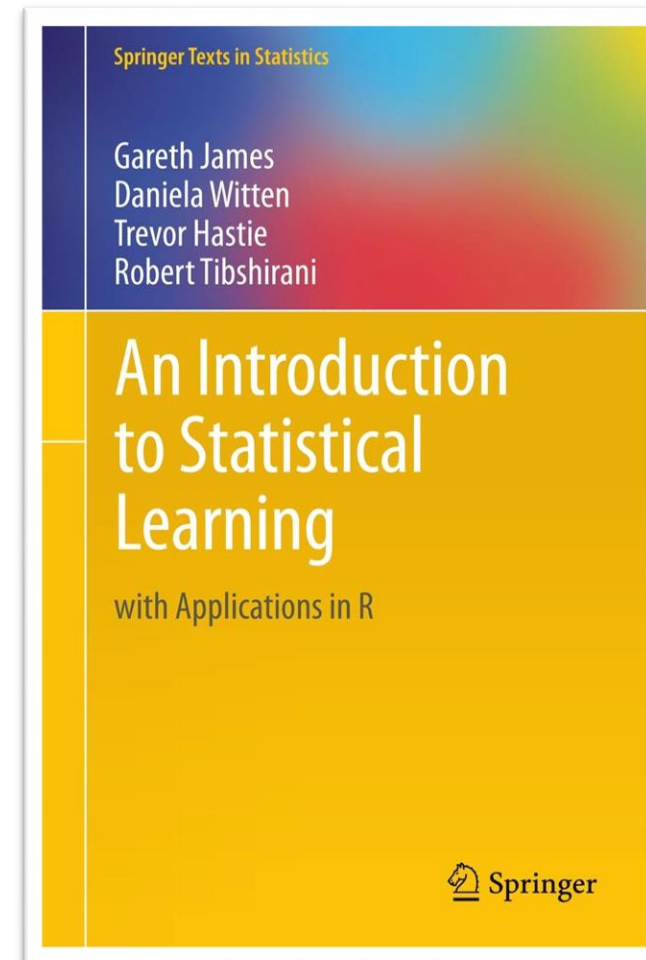


Breiman, L. (1996). Bagging predictors. *Machine learning*, 24, 123-140.

An important limitation in bagging

The ensembles of trees in bagging tend to be highly correlated

“Suppose there is one very strong predictor in the data set, along with a number of other moderately strong predictors. Then in a collection of bagged trees, most or all of the trees will use this strong predictor in the top split. Consequently, all of the bagged trees will look quite similar to each other. Hence the predictions from the bagged trees will be highly correlated. Unfortunately, averaging many highly correlated quantities does not lead to as large of a reduction in variance as averaging many uncorrelated quantities” (pag. 320)



James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning. New York: Springer

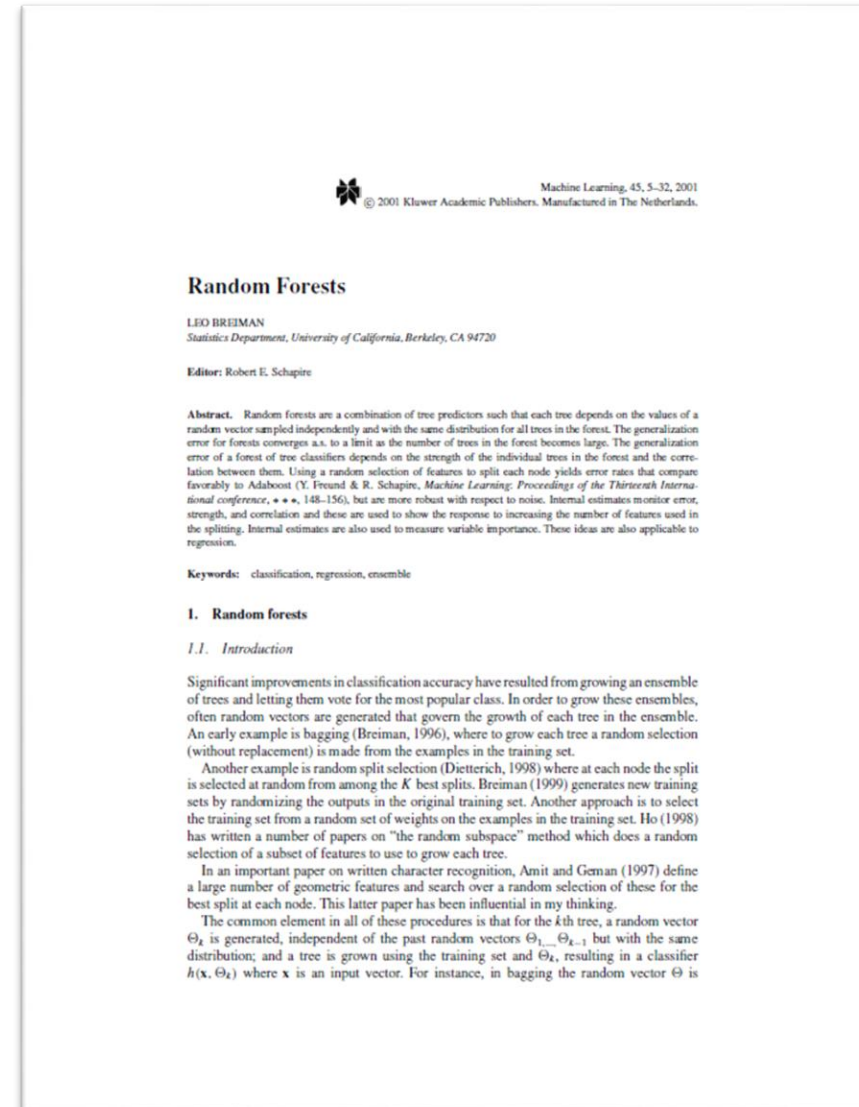
A solution: The random forests algorithm

Yet another breakthrough result from Leo Breiman

1

“Random forests are a combination of tree predictors such that each tree depends on the values of a random vector sampled independently and with the same distribution for all trees in the forest” (pag. 5, Breiman 2001)

“The forests studied here consist of using randomly selected inputs or combinations of inputs at each node to grow each tree” (pag. 10, Breiman 2001)



Breiman, L. Random Forests. Machine Learning 45, 5–32 (2001).

The randomness in the random forest algorithm: overview

1

1

Randomness 1/2 (same as in bagging)

Random forests = ensemble of trees grown by resampling (with replacement) training data, and selecting a random subset of $1 \leq m \leq p$ predictors at each step, for each tree.

2

Randomness 2/2 (new!)

By adding the randomness in the ensemble, the result is to decrease the variance, i.e., dependence on training data, of the forest. This is in contrast with individual trees, which typically **exhibit high variance** and tend to **overfit data**.

3

Standard choices: (1) $m=p$ (bagging), (2) $m = \sqrt{p}$

Where do we inject randomness into the forest?

1

`sklearn.ensemble.RandomForestClassifier`

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='sqrt', max_leaf_nodes=None,  
min_impurity_split=None, bootstrap=True, oob_score=False, warm_start=False, n_jobs=None, random_state=None,  
verbose=0, ccp_alpha=0.0):
```

`sklearn.ensemble.RandomForestRegressor`

```
class sklearn.ensemble.RandomForestRegressor(n_estimators=100, *, criterion='squared_error', max_depth=None,  
min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features=1.0, max_leaf_nodes=None,  
min_impurity_split=None, bootstrap=True, oob_score=False, warm_start=False, n_jobs=None, random_state=None,  
verbose=0, ccp_alpha=0.0):
```

`max_features : {"sqrt", "log2", None}, int or float, default="sqrt"`

The number of features to consider when looking for the best split:

- If int, then consider `max_features` features at each split.
- If float, then `max_features` is a fraction and `max(1, int(max_features * n_features_in_))` features are considered at each split.
- If "sqrt", then `max_features=sqrt(n_features)`.
- If "log2", then `max_features=log2(n_features)`.
- If None, then `max_features=n_features`.

the randomness in the ensemble, the variance decreases, i.e.,

of the forest.
al trees, which
and tend to

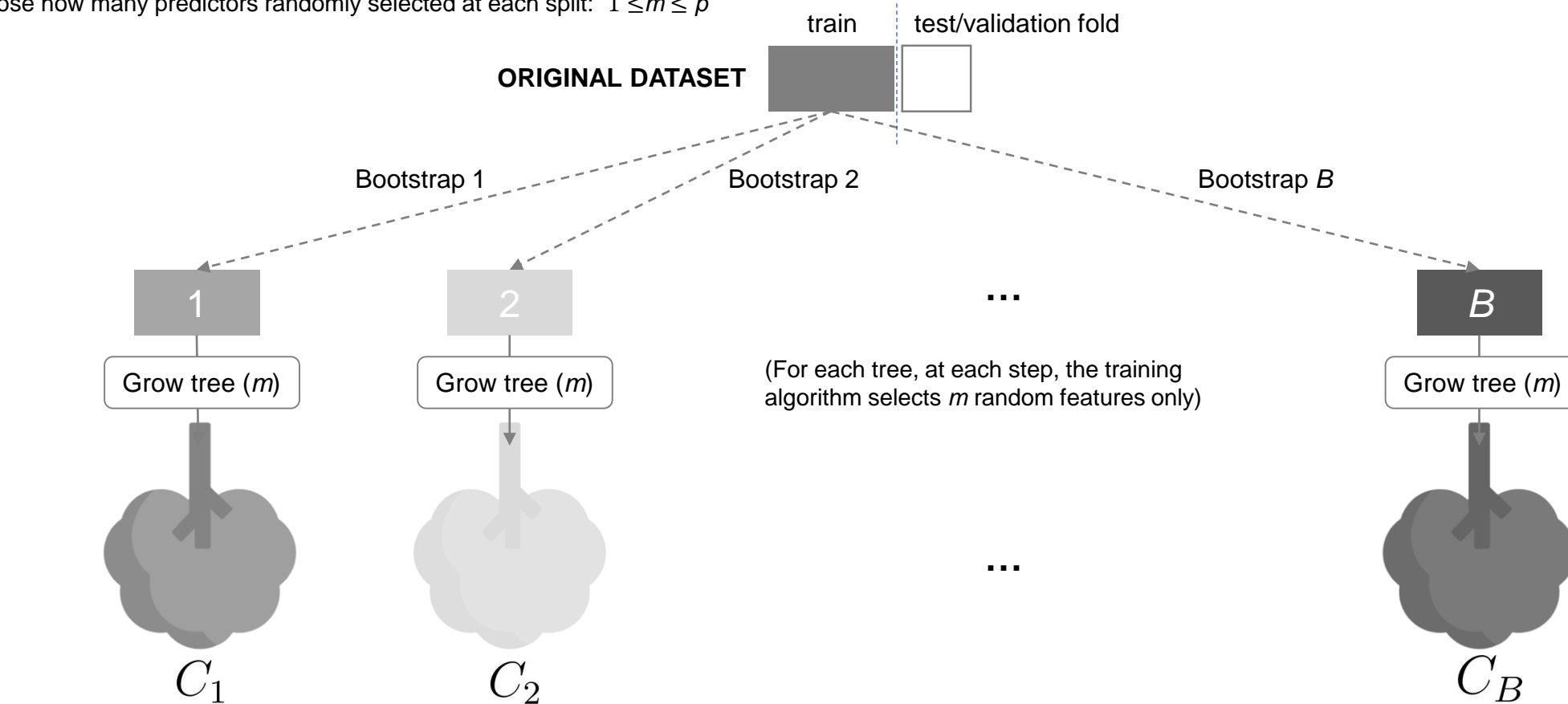
bagging), (2) $m = \sqrt{p}$

Random forests: A visual summary

1

Choose how many bootstrapped training datasets: B

Choose how many predictors randomly selected at each split: $1 \leq m \leq p$



$$C_{\text{RF}}(x) = \frac{1}{B} \sum_{i=1}^B C_i(x) \quad (\text{regression})$$

Random forests: An example from James et al. (2013), chapter 8

Gene expression dataset

1

“We applied random forests to a high-dimensional biological data set consisting of expression measurements of 4,718 genes measured on tissue samples from 349 patients. There are around 20,000 genes in humans, and individual genes have different levels of activity, or expression, on particular cells, tissues, and biological conditions. In this data set, each of the patient samples has a qualitative label with 15 different levels: either normal or 1 of 14 different types of cancer. Our goal was to use random forests to predict cancer type based on the 500 genes that have the largest variance in the training set.” (pag. 321, James et al. (2013))

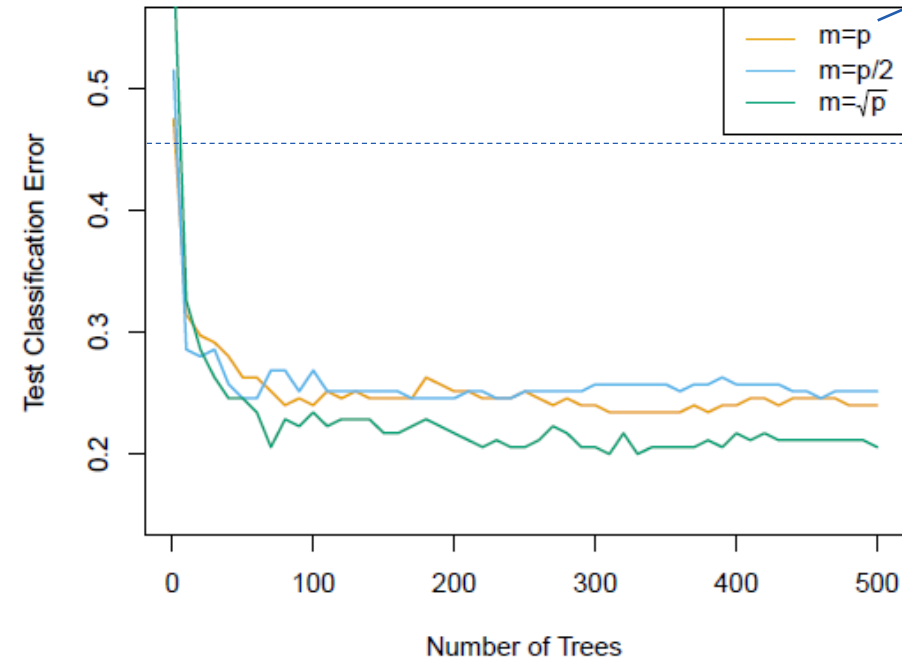


FIGURE 8.10. Results from random forests for the 15-class gene expression data set with $p = 500$ predictors. The test error is displayed as a function of the number of trees. Each colored line corresponds to a different value of m , the number of predictors available for splitting at each interior tree node. Random forests ($m < p$) lead to a slight improvement over bagging ($m = p$). A single classification tree has an error rate of 45.7%.

“Ensembling *via* learning sequentially”

2

Boosting

“Ensembling with randomization”

Train deep trees on bootstrapped copies of training data (and randomly select a subset of predictors at each step for random forests), then aggregate their predictions



“Ensembling *via* learning sequentially”

Ensemble “**weak learners**” sequentially, and turn them into a “**strong learner**”...

“Ensembling with randomization”

Train deep trees on bootstrapped copies of training data (and randomly select a subset of predictors at each step for random forests), then aggregate their predictions



“Ensembling *via* learning sequentially”

Ensemble “**weak learners**” sequentially, and turn them into a “**strong learner**”...

I

Adaptive Boosting (AdaBoost)

II

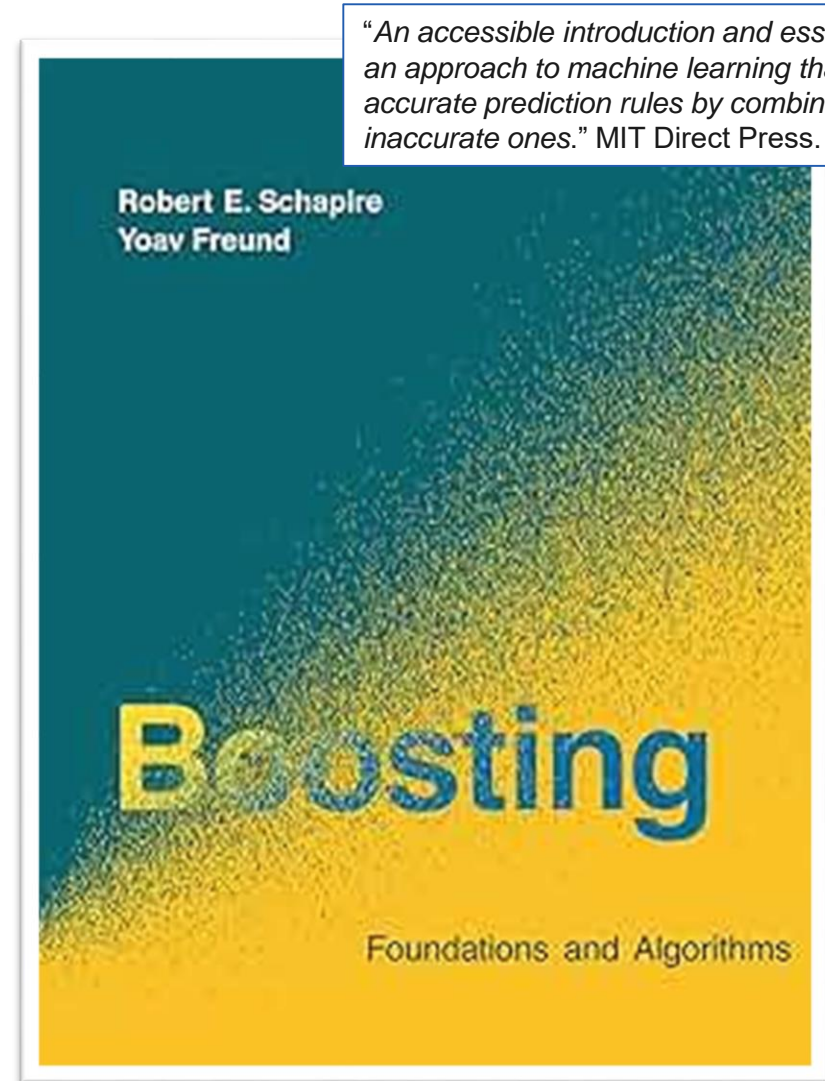
Gradient Boosting Machine

I Adaptive Boosting (AdaBoost) – from the works of Schapire and Freund

2

“In fact, Breiman (1996) (referring to a NIPS workshop) called AdaBoost with trees the “best off-the-shelf classifier in the world” (pag. 338) from

Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). The Annals of Statistics, 28 (2):337–407, 2000.



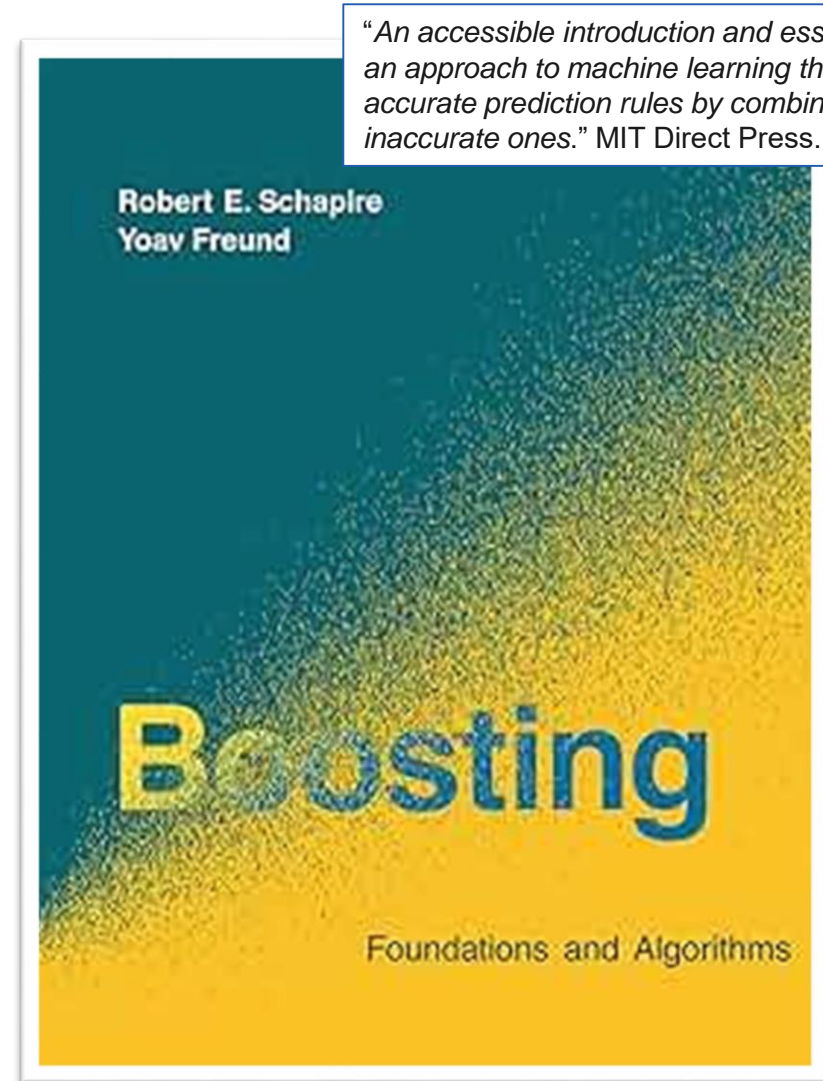
Schapire, R.E., Freund, Y. (2012). Boosting: Foundations and Algorithms. MIT Press.

Adaptive Boosting (AdaBoost) – from the works of Schapire and Freund

“An off-the-shelf” method is one that can be directly applied to the data without requiring a great deal of time-consuming data preprocessing or careful tuning of the learning procedure” (pag. 352, Hastie et al. (2013))

“In fact, Breiman (1996) (referring to a NIPS workshop) called AdaBoost with trees the “best off-the-shelf classifier in the world” (pag. 338) from

Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). The Annals of Statistics, 28 (2):337–407, 2000.



“An accessible introduction and essential reference for an approach to machine learning that creates highly accurate prediction rules by combining many weak and inaccurate ones.” MIT Direct Press.

Schapire, R.E., Freund, Y. (2012). Boosting: Foundations and Algorithms. MIT Press.

Adaptive Boosting (AdaBoost) – a brief definition from a funny preprint

“AdaBoost is a boosting algorithm that produces a single ensemble learner from a **sequential additive process** which involves **re-weighting of training data points** (therefore the name AdaBoost, or ‘Adaptive Boosting’ of weights) and **fitting of weak learners**” (pag. 1, emphasis ours)



Ferrario, Andrea and Hämmmerli, Roger, On Boosting: Theory and Applications (June 11, 2019). Available at SSRN: <https://ssrn.com/abstract=3402687> or <http://dx.doi.org/10.2139/ssrn.3402687>

Adaptive Boosting (AdaBoost) – a brief definition from a funny preprint

1

Sequential = information from step n is used at step $n+1$
Additive = the “strong learner” is obtained by aggregating the “weak learners” additively (linear combination)

“AdaBoost is a boosting algorithm that produces a single ensemble learner from a **sequential additive process** which involves **re-weighting of training data points** (therefore the name AdaBoost, or ‘Adaptive Boosting’ of weights) and **fitting of weak learners**” (pag. 1, emphasis ours)

3

Small trees, e.g., a “tree stump” (two leaves)

2

The training data set is **sequentially (re-)weighted** and weights are used to train each “weak learner” in the ensemble. The weights at step $n+1$ depend on the performance of the “weak learner” trained at step n .

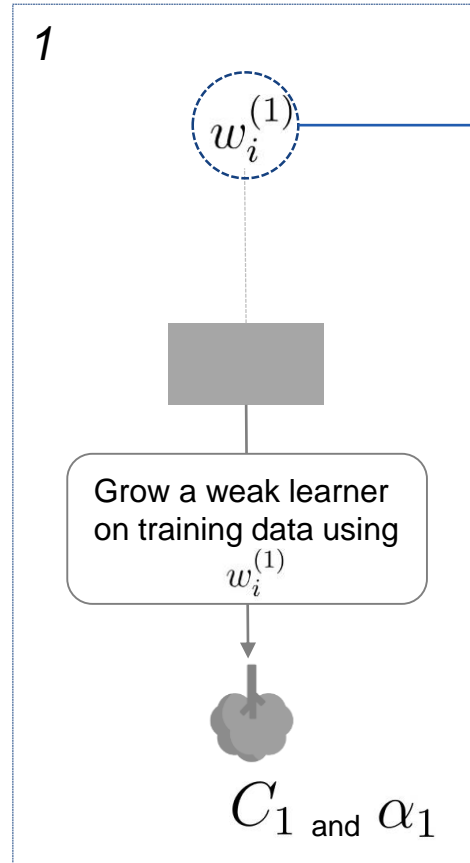


Ferrario, Andrea and Hämmmerli, Roger, On Boosting: Theory and Applications (June 11, 2019). Available at SSRN: <https://ssrn.com/abstract=3402687> or <http://dx.doi.org/10.2139/ssrn.3402687>

STEP

DATASET
WEIGHTS

SEQUENCE

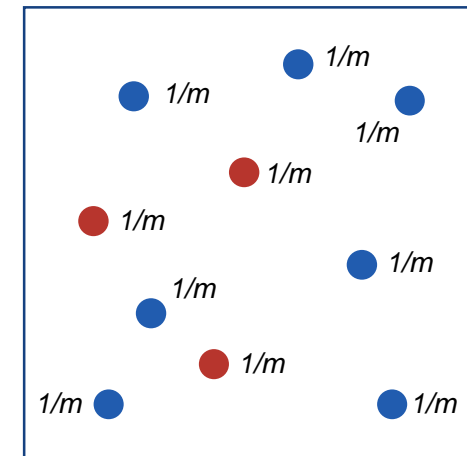


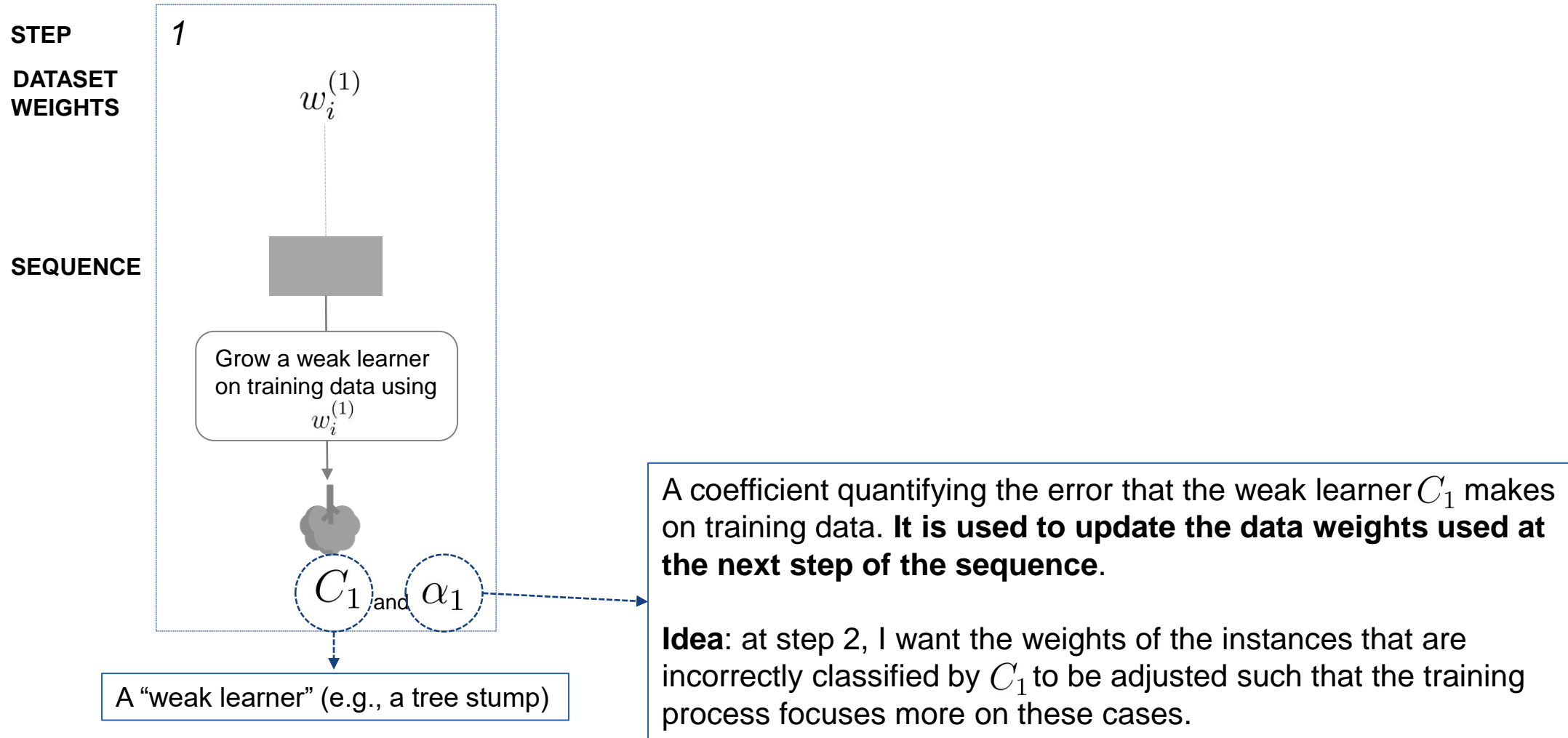
(classical choice in the AdaBoost articles and Schapire and Freund's monograph)

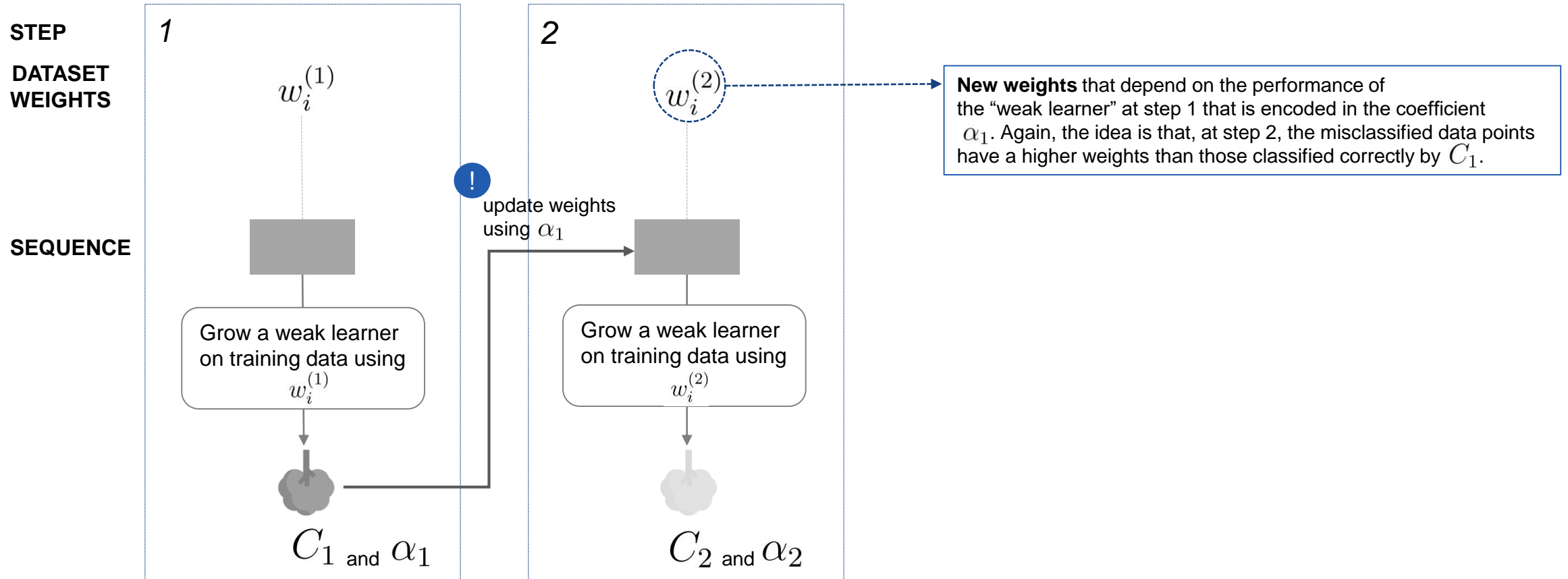
Training data

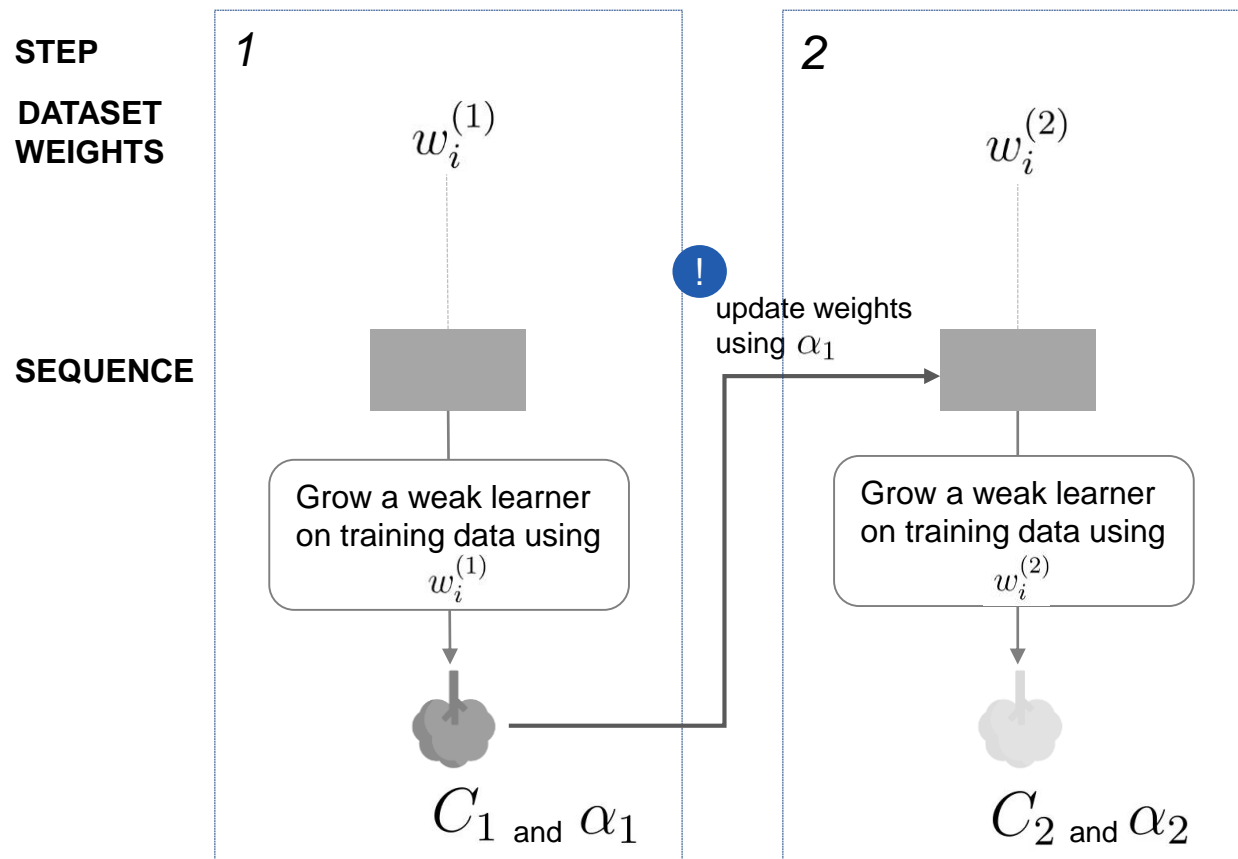


$$= \{(x_i, y_i) \in \mathcal{X} \times \{-1, +1\}, 1 \leq i \leq m\}$$

Initialized by hand: usually $1/m$ Example: Training data ($m=10$)

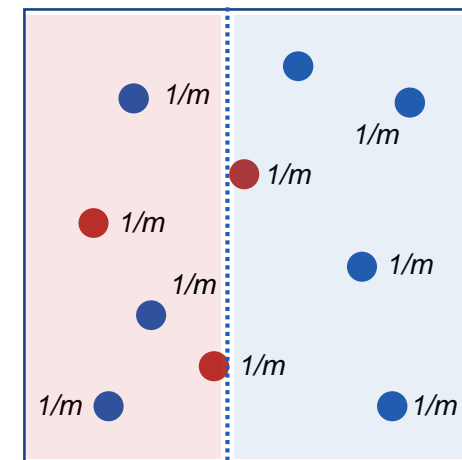




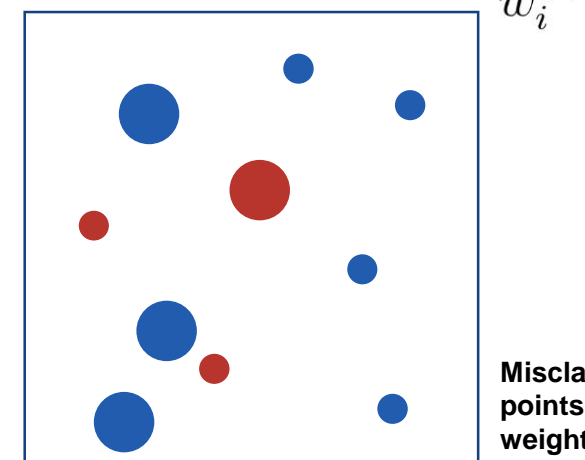


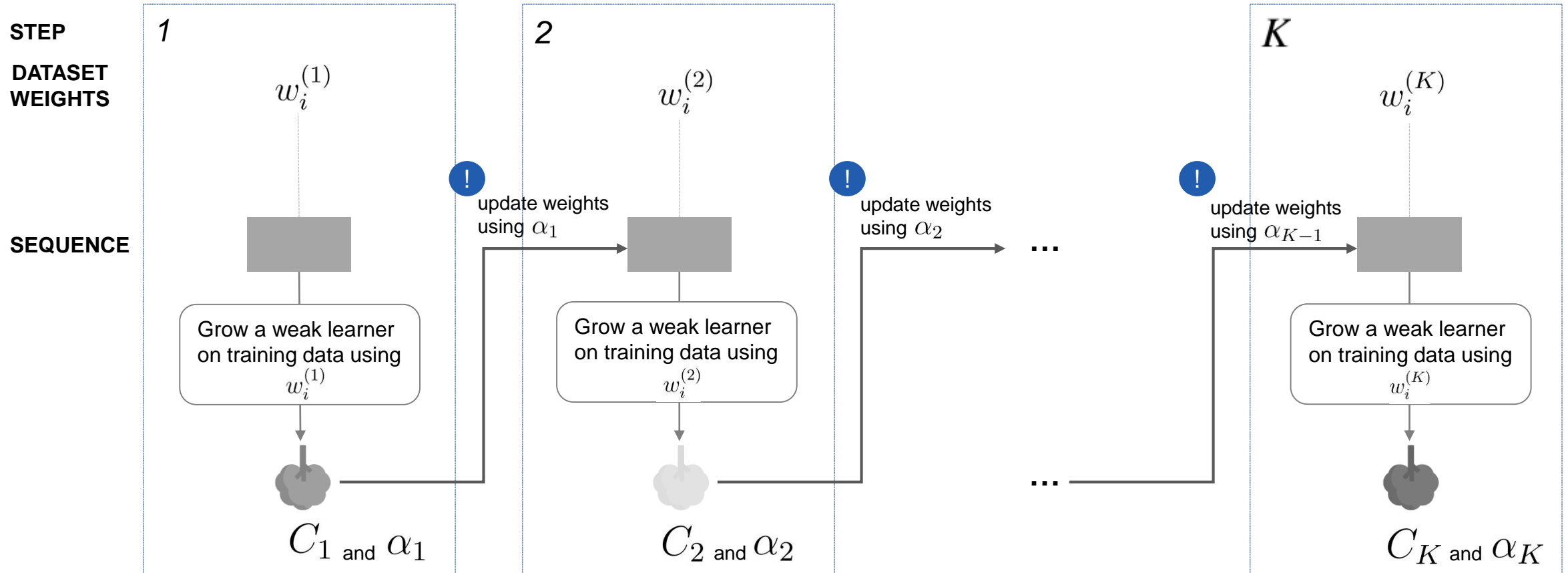
Example

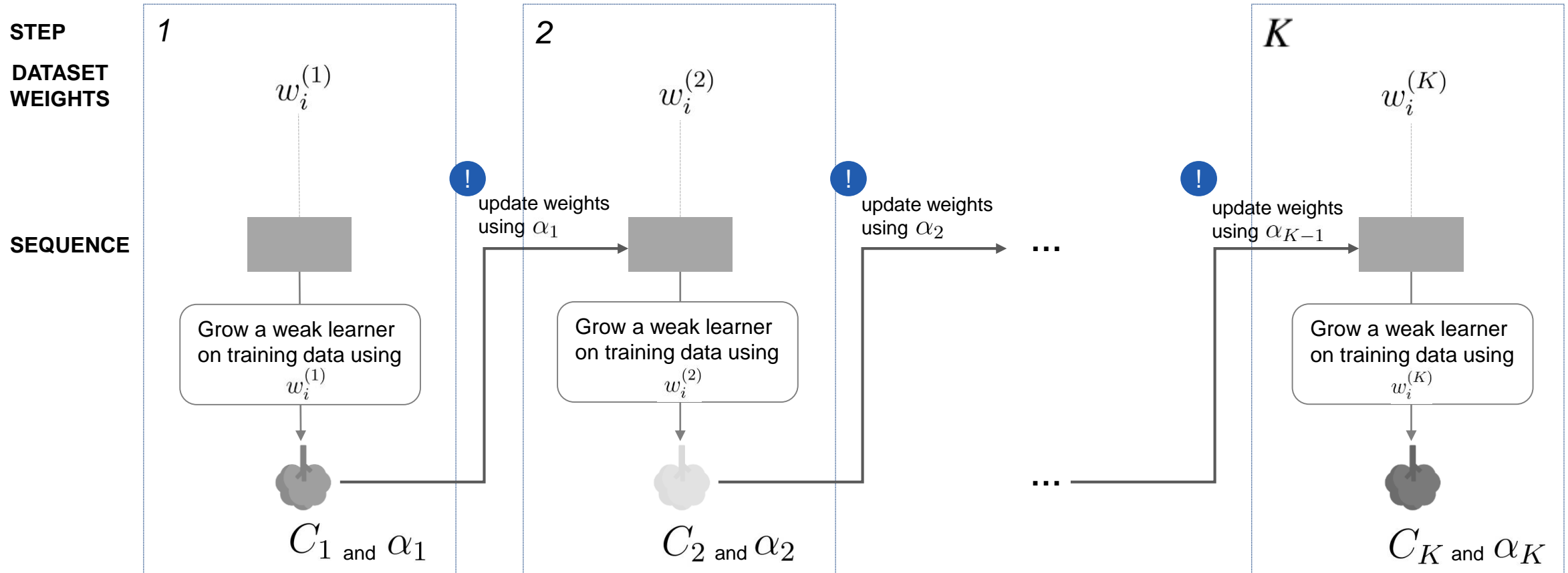
C_1 “weak learner” – a tree stump



update weights using α_1

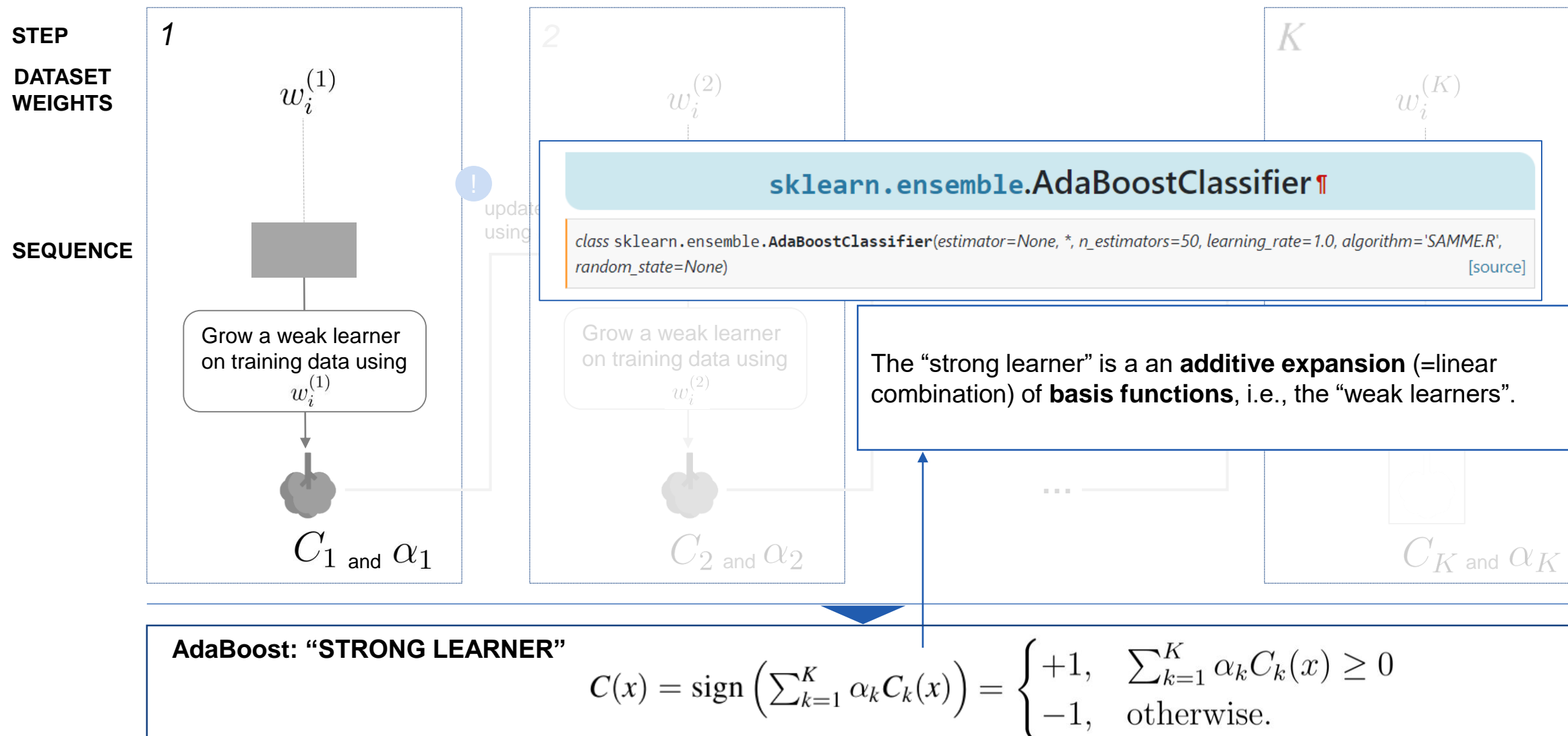






AdaBoost: “STRONG LEARNER”

$$C(x) = \text{sign} \left(\sum_{k=1}^K \alpha_k C_k(x) \right) = \begin{cases} +1, & \sum_{k=1}^K \alpha_k C_k(x) \geq 0 \\ -1, & \text{otherwise.} \end{cases}$$



The Gradient Boosting Machine

An important idea by Friedman

The idea behind *gradient* boosting is to compute a sequence of “weak learners” and combine them additively (“strong learner”) using the *gradient* of an **arbitrary (differentiable) loss function** chosen for the problem at hand and in a step-wise fashion.

The Annals of Statistics
2001, Vol. 29, No. 5, 1189–1232

1999 REITZ LECTURE

GREEDY FUNCTION APPROXIMATION: A GRADIENT BOOSTING MACHINE¹

BY JEROME H. FRIEDMAN

Stanford University

Function estimation/approximation is viewed from the perspective of numerical optimization in function space, rather than parameter space. A connection is made between stagewise additive expansions and steepest-descent minimization. A general gradient descent “boosting” paradigm is developed for additive expansions based on any fitting criterion. Specific algorithms are presented for least-squares, least absolute deviation, and Huber- M loss functions for regression, and multiclass logistic likelihood for classification. Special enhancements are derived for the particular case where the individual additive components are regression trees, and tools for interpreting such “TreeBoost” models are presented. Gradient boosting of regression trees produces competitive, highly robust, interpretable procedures for both regression and classification, especially appropriate for mining less than clean data. Connections between this approach and the boosting methods of Freund and Shapire and Friedman, Hastie and Tibshirani are discussed.

1. Function estimation. In the function estimation or “predictive learning” problem, one has a system consisting of a random “output” or “response” variable y and a set of random “input” or “explanatory” variables $\mathbf{x} = \{x_1, \dots, x_n\}$. Using a “training” sample $\{(y_i, \mathbf{x}_i)\}_{i=1}^N$ of known (y, \mathbf{x}) -values, the goal is to obtain an estimate or approximation $\hat{F}(\mathbf{x})$, of the function $F^*(\mathbf{x})$ mapping \mathbf{x} to y , that minimizes the expected value of some specified loss function $L(y, F(\mathbf{x}))$ over the joint distribution of all (y, \mathbf{x}) -values,

$$(1) \quad F^* = \arg \min_F E_{y, \mathbf{x}} L(y, F(\mathbf{x})) = \arg \min_F E_{\mathbf{x}} [E_y (L(y, F(\mathbf{x}))) \mid \mathbf{x}].$$

Frequently employed loss functions $L(y, F)$ include squared-error $(y - F)^2$ and absolute error $|y - F|$ for $y \in \mathbb{R}^1$ (regression) and negative binomial log-likelihood, $\log(1 + e^{-2yF})$, when $y \in \{-1, 1\}$ (classification).

A common procedure is to restrict $F(\mathbf{x})$ to be a member of a parameterized class of functions $F(\mathbf{x}; \mathbf{P})$, where $\mathbf{P} = \{P_1, P_2, \dots\}$ is a finite set of parameters whose joint values identify individual class members. In this article we focus

Received May 1999; revised April 2001.

¹Supported in part by CSIRO Mathematical and Information Science, Australia; Department of Energy Contract DE-AC03-76SF00515; and NSF Grant DMS-97-64431.

AMS 2000 subject classifications. 62-02, 62-07, 62-08, 62G08, 62H30, 68T10.

Key words and phrases. Function estimation, boosting, decision trees, robust nonparametric regression.

1189

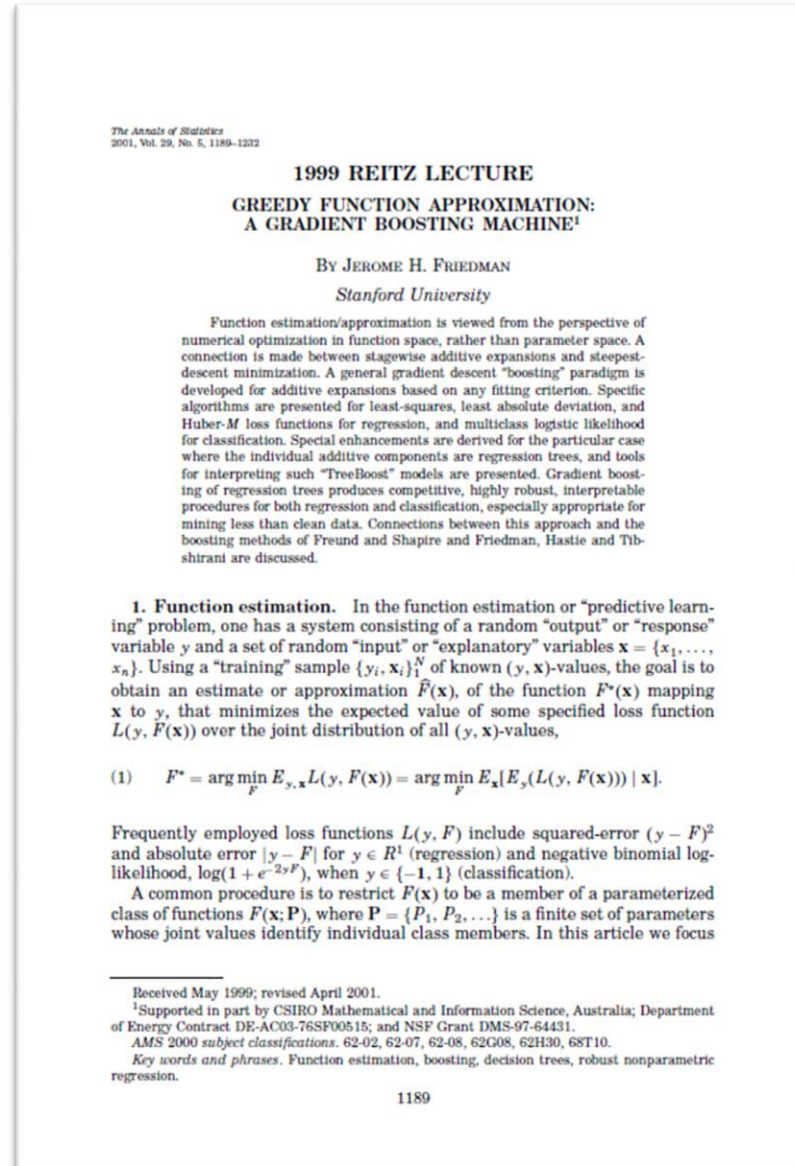
Jerome H. Friedman. "Greedy function approximation: A gradient boosting machine.." Ann. Statist. 29 (5) 1189 - 1232, October 2001.

The Gradient Boosting Machine

An important idea by Friedman

The idea behind *gradient* boosting is to compute a sequence of “weak learners” and combine them additively (“strong learner”) using the *gradient* of an **arbitrary (differentiable) loss function** chosen for the problem at hand and in a step-wise fashion.

Remark: this procedure is what AdaBoost also does...it took some time to researchers to prove this point.



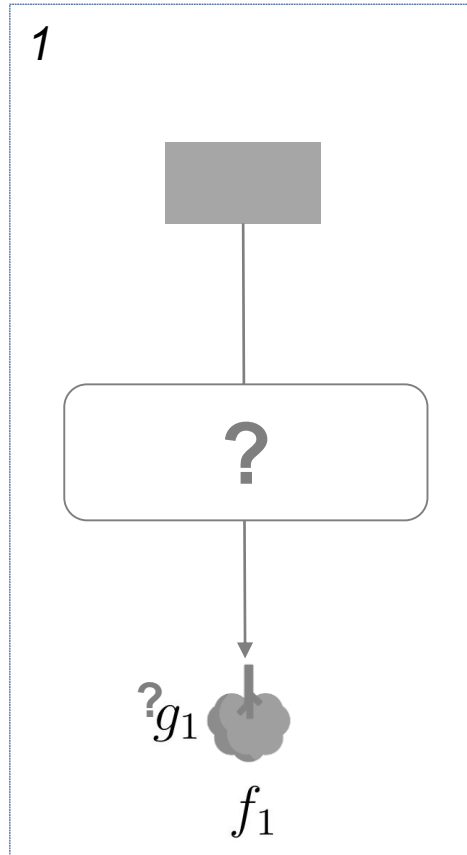
Jerome H. Friedman. "Greedy function approximation: A gradient boosting machine.." Ann. Statist. 29 (5) 1189 - 1232, October 2001.

Gradient Boosting: A visual summary

Let us start with the very first step of the sequence

STEP

1



Key Inputs

Regression and classification

Training data $\text{[grey box]} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}, 1 \leq i \leq m\}$

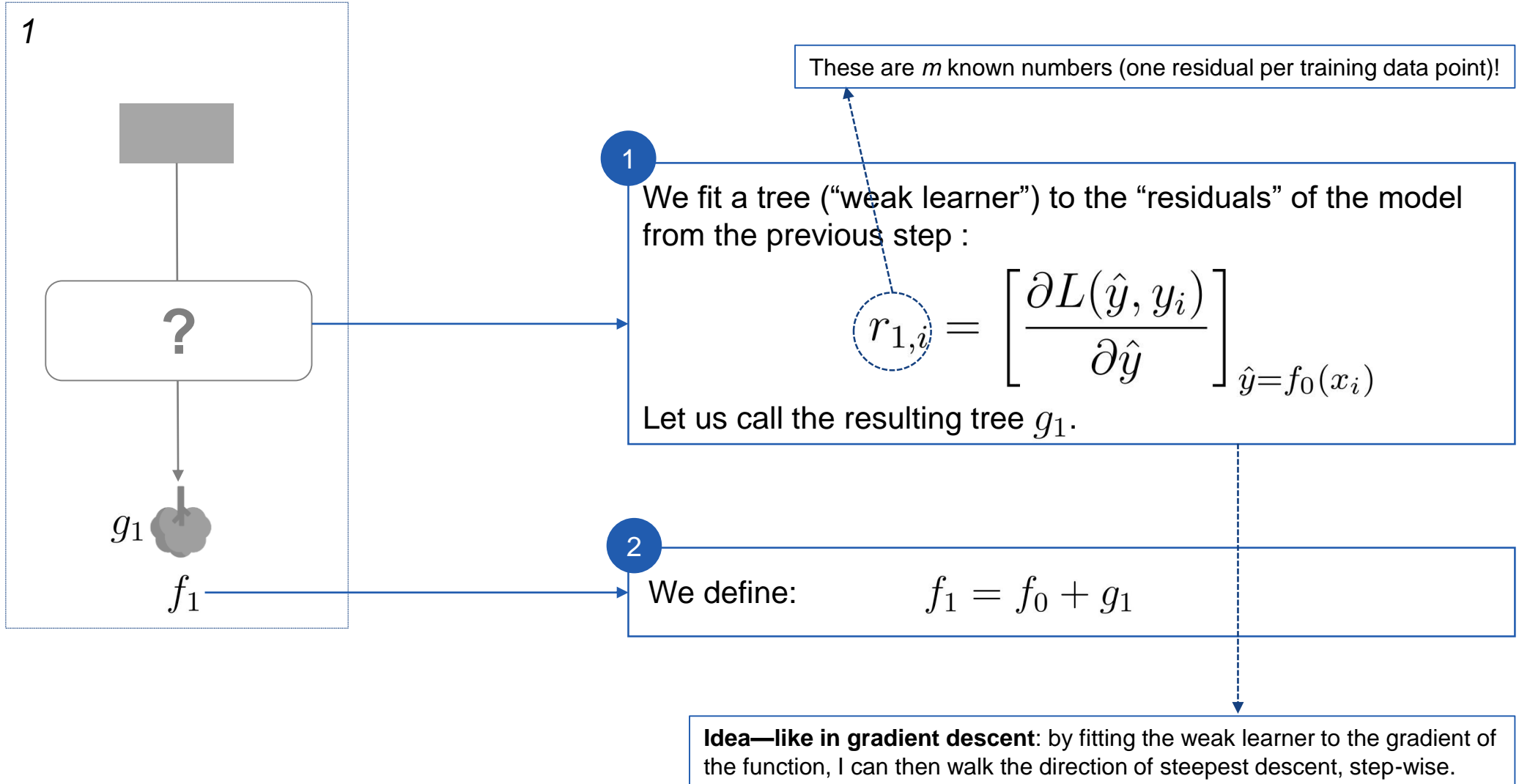
Loss function $L : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$
 $L(z, w) = (z - w)^2$ (regression)

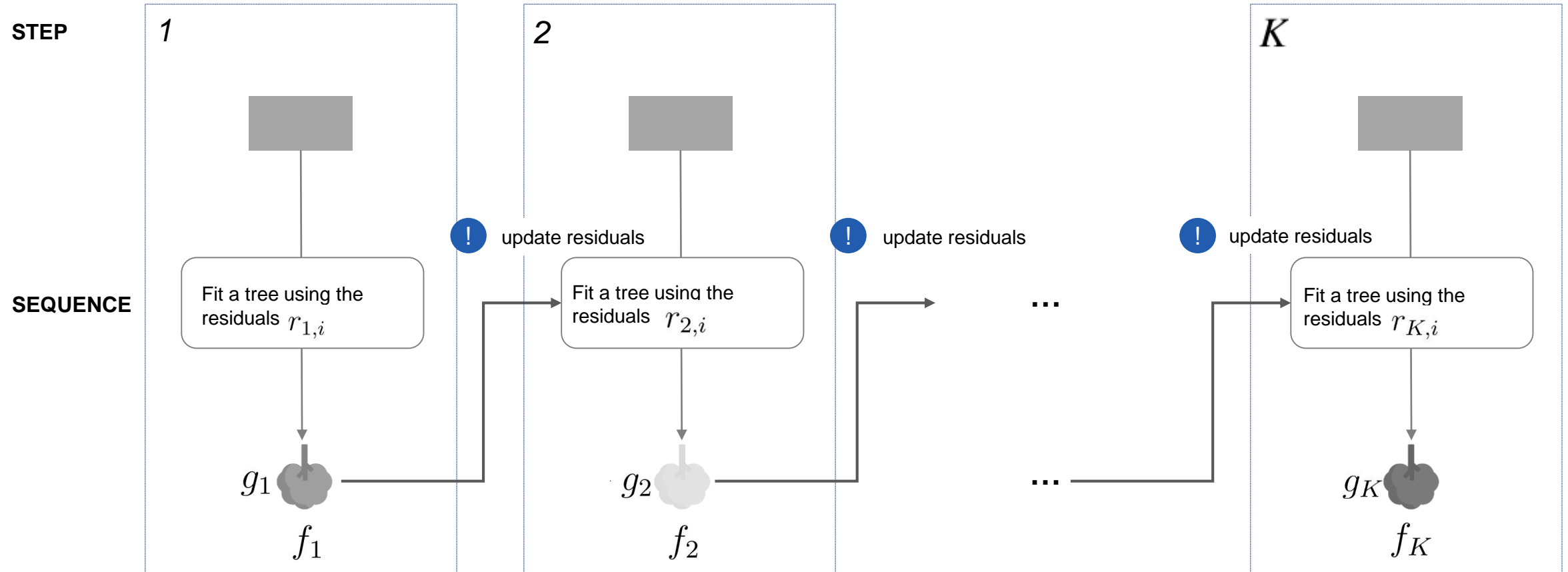
$f_0(x) = 0$ (convenient initialization – we need it later)

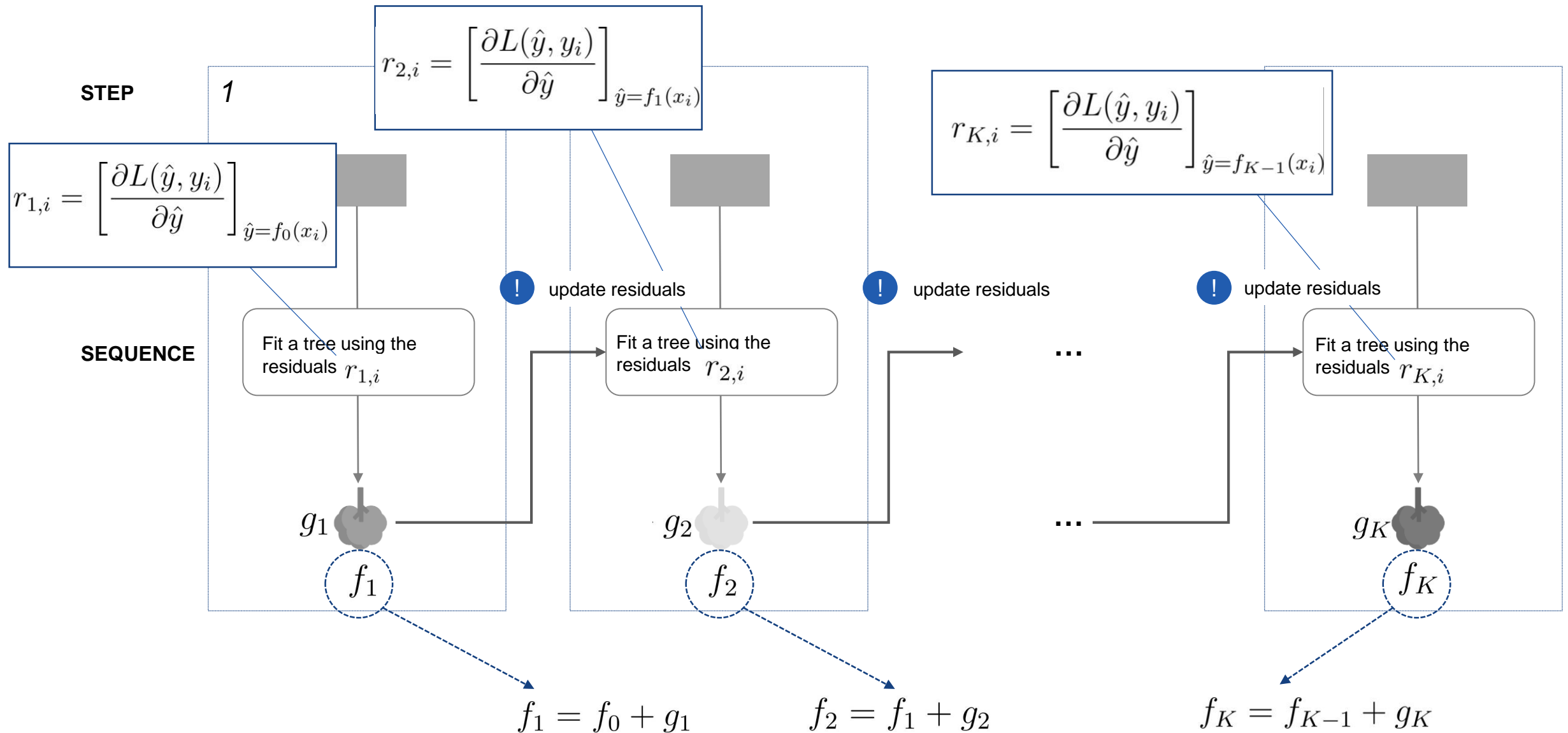
Gradient Boosting: A visual summary

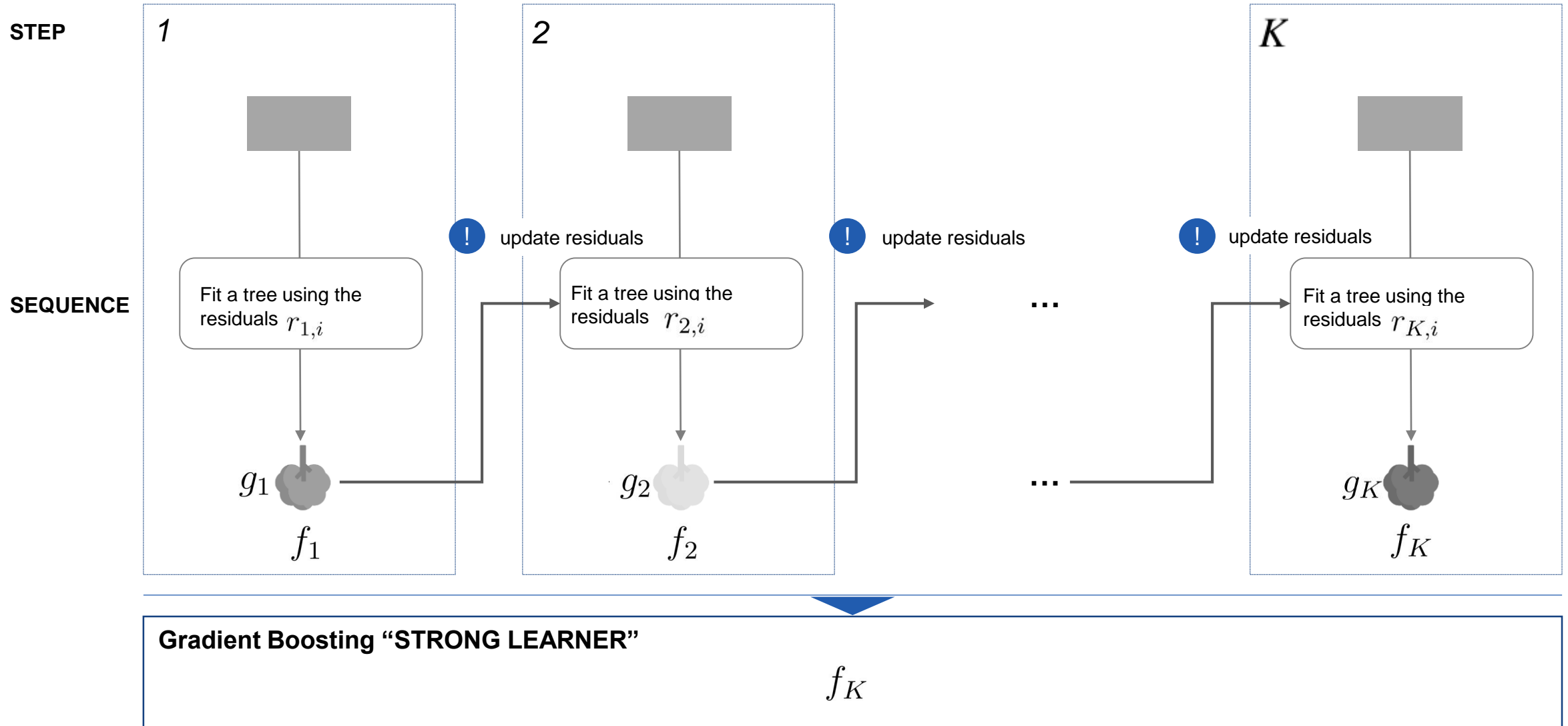
The core of the gradient boosting method: fitting a weak learner to the gradient of the loss function (also known as the “residuals”)

STEP



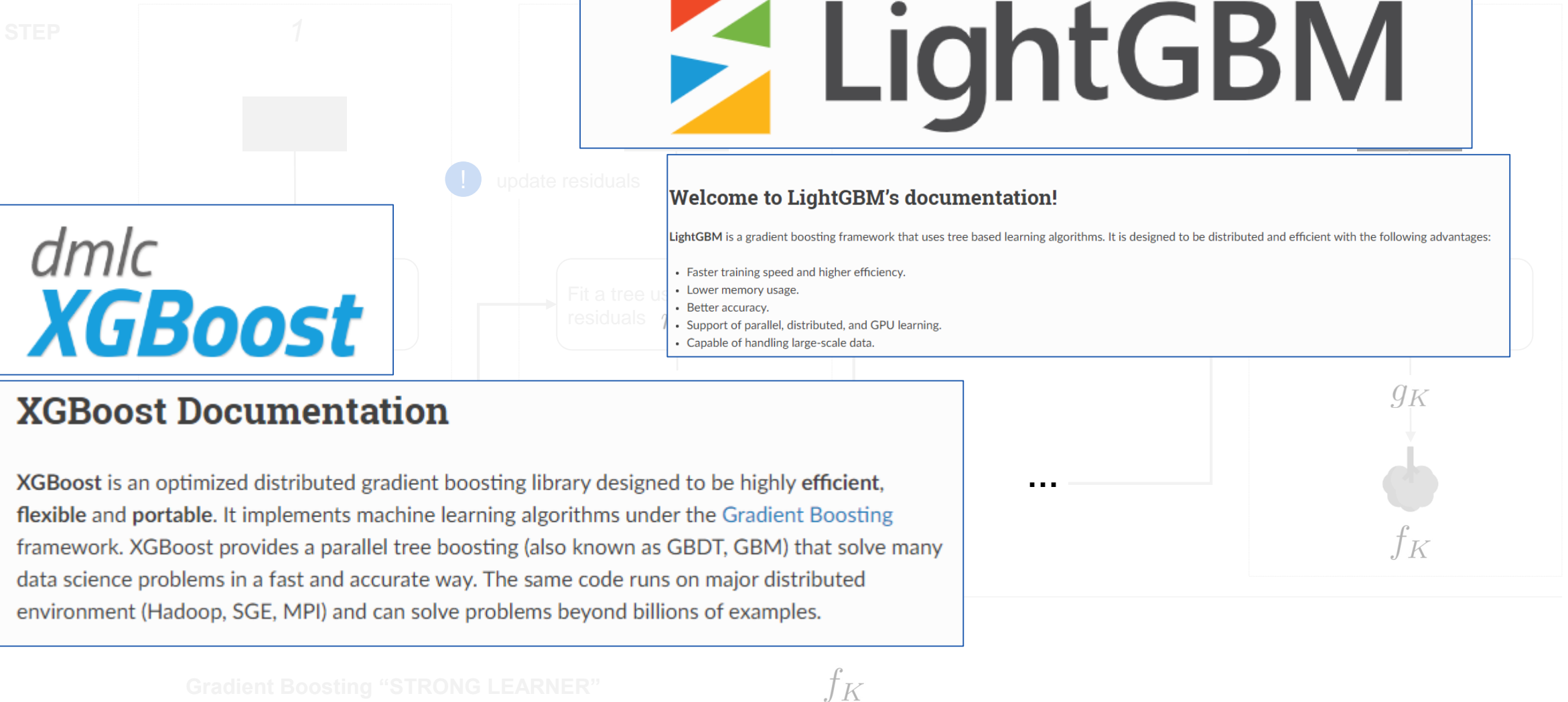






Two famous implementations of Gradient Boosting

Both are available in Python



Deep learning is not all you need: Tree-based ensembling shows high performance on tabular data



McElfresh, D., Khandagale, S., Valverde, J., Prasad C. V., Ramakrishnan, G., Goldblum, M., & White, C. (2024). When do neural nets outperform boosted trees on tabular data?. *Advances in Neural Information Processing Systems*, 36.



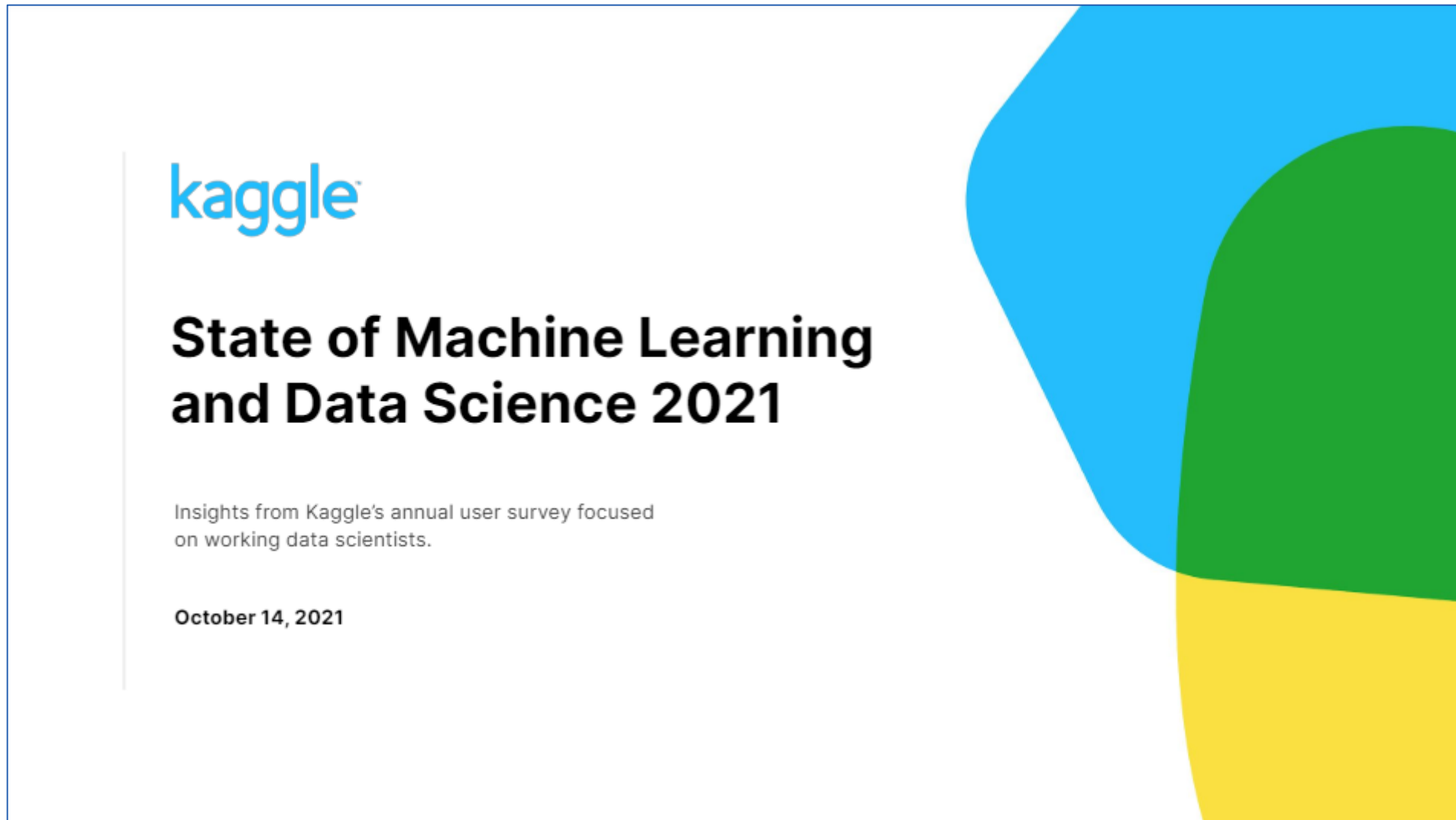
Grinsztajn, L., Oyallon, E., & Varoquaux, G. (2022). Why do tree-based models still outperform deep learning on typical tabular data?. *Advances in neural information processing systems*, 35, 507-520.



Shwartz-Ziv, R., & Armon, A. (2022). Tabular data: Deep learning is not all you need. *Information Fusion*, 81, 84-90.

Kaggle: tree-based ensemble methods are still a popular choice

The presentation can be found on Moodle



Kaggle: tree-based ensemble methods are still a popular choice

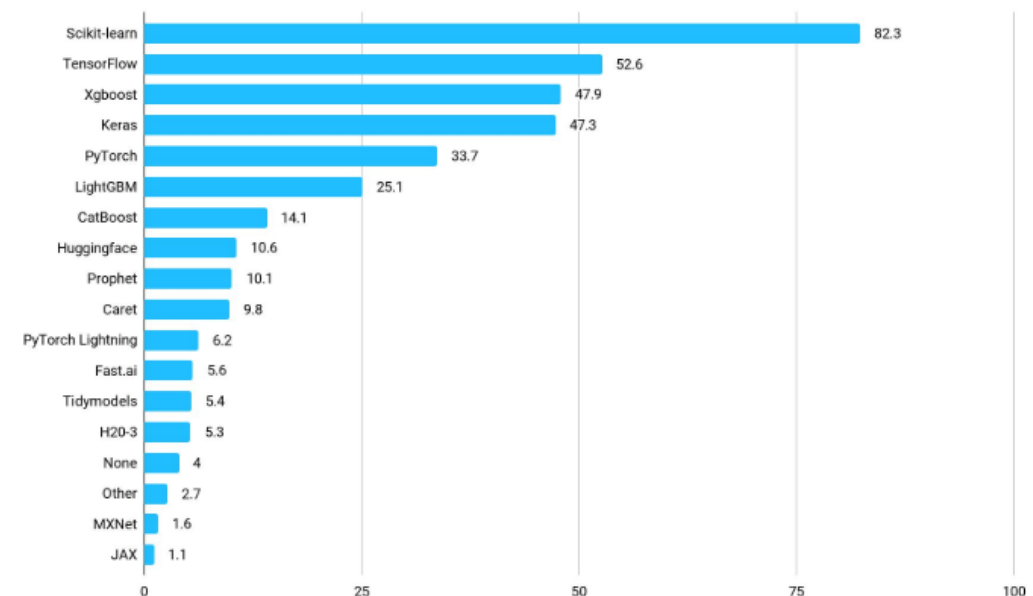
Machine Learning Frameworks

Python-based tools continue to dominate the machine learning frameworks.

Like last year, Scikit-learn, a swiss army knife applicable to most projects, is the top with over 80% of data scientists using it. TensorFlow and Keras, notably used in combination for deep learning, were each selected on about half of the data scientist surveys. Gradient boosting library xgboost is fourth, with about the same usage as 2020 and 2019.

The most popular of the new tools added to the survey this year is Huggingface reaching over 10%.

Machine Learning Framework Usage



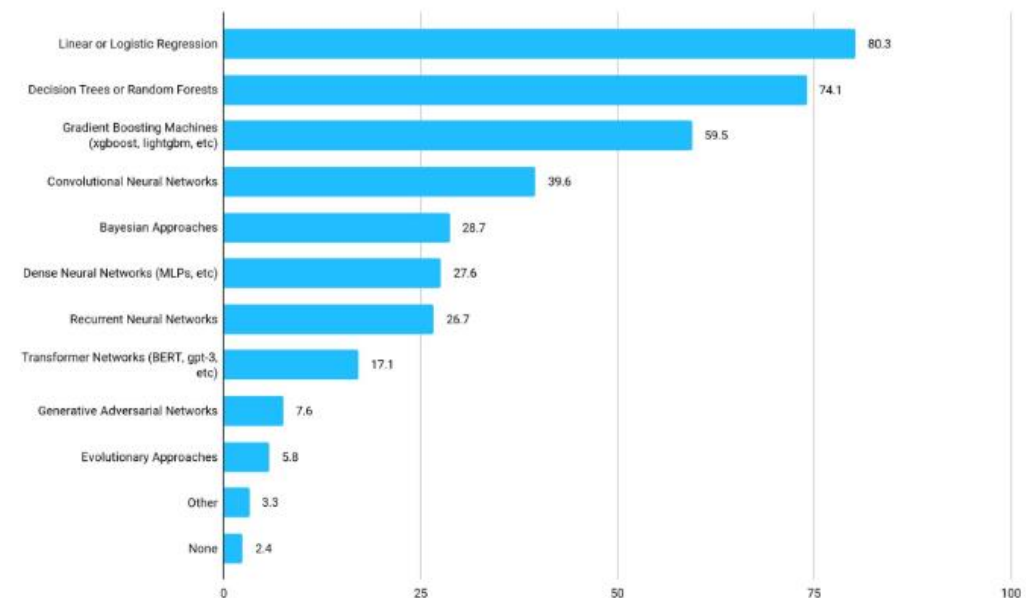
Kaggle: tree-based ensemble methods are still a popular choice

Methods & Algorithms

Like last year, the most commonly used algorithms were linear and logistic regression, followed closely by decision trees and random forests.

Of more complex methods, gradient boosting machines and convolutional neural networks were the most popular approaches.

Methods and Algorithms Usage



Technology

33

Kaggle | State of ML & Data Science 2021

Tree-based ensembling methods: Key takeaways

I

Ensembling allows taking care of the limitations that affect decision trees, namely, overfitting and high variance

II

Bagging and **random forest algorithms** inject randomness into learning by bootstrapping training data and de-correlating trees

III

Boosting methods, such as AdaBoost and gradient boosting machines, train “weak learners” sequentially and combine them additively to return a “strong learner”

IV

Research is actively investigating why **tree-based ensemble methods exhibit performance comparable with the one of deep learning methods** on tabular data. Tree-based methods are still actively used by data scientists across the world.

Lastly, a Google Colab notebook to test what we learned today

`ensemble.ipynb`

Reflection

*“Looking back on all the time we spent together...”**

*La Bouche, Be My Lover (1995) - <https://www.youtube.com/watch?v=ViP87WipSm0>

Please consider these questions and answer them: We will comment them together

1

What was your understanding of machine learning before starting Block I and how have our lectures together changed that?

2

What was the most engaging/surprising fact that you learned during Block I?

3

What is the topic related to machine learning that is still the most unclear/mysterious to you?

4

Can you specify a business problem you would like to tackle? Which machine learning method would you like to use?

5

How would you continue learning machine learning now?

Andrea's recommendations



1

You just started learning machine learning: now, increase time spent on self-study and improve your current understanding of methods and algorithms

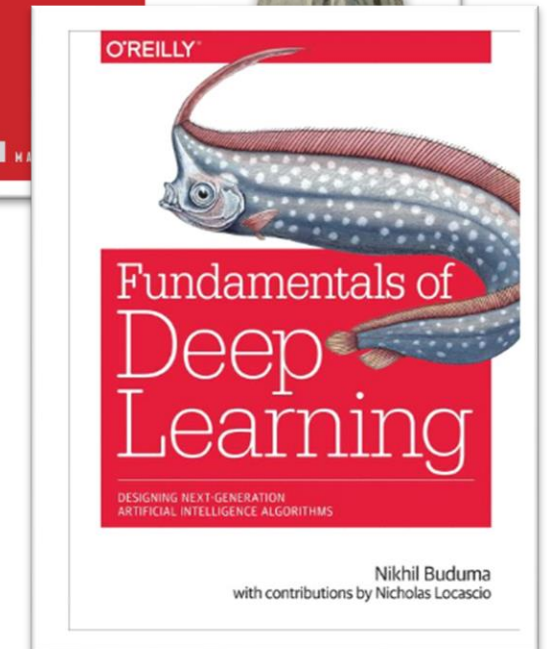
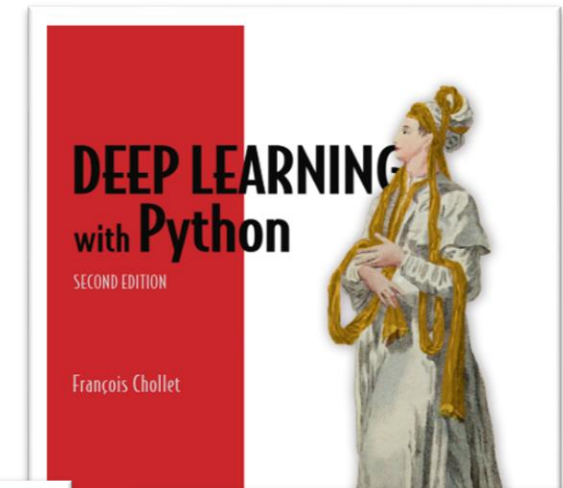
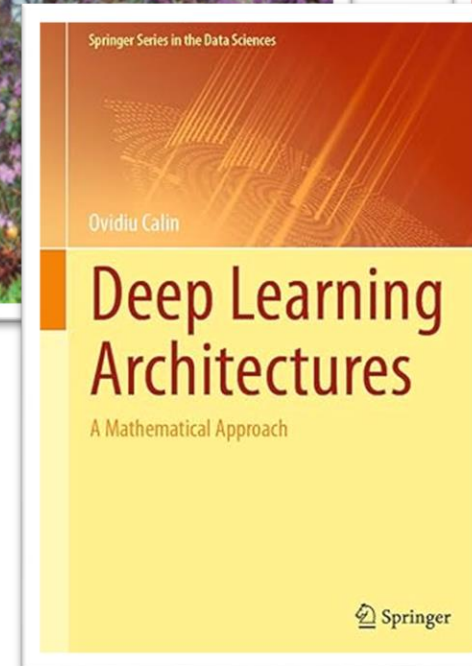
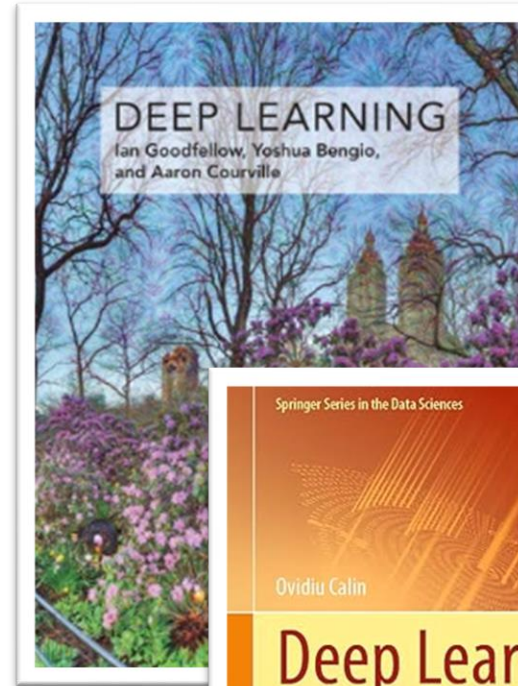
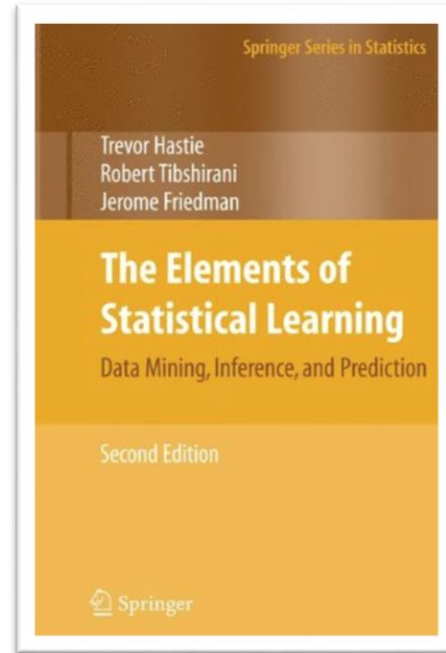
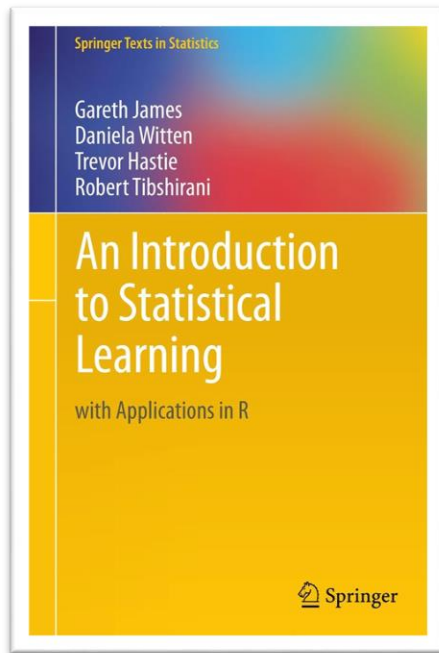
2

Identify a problem you would like to address, including a specific data modality such as text, images, or tabular data, and begin coding. Start simple!

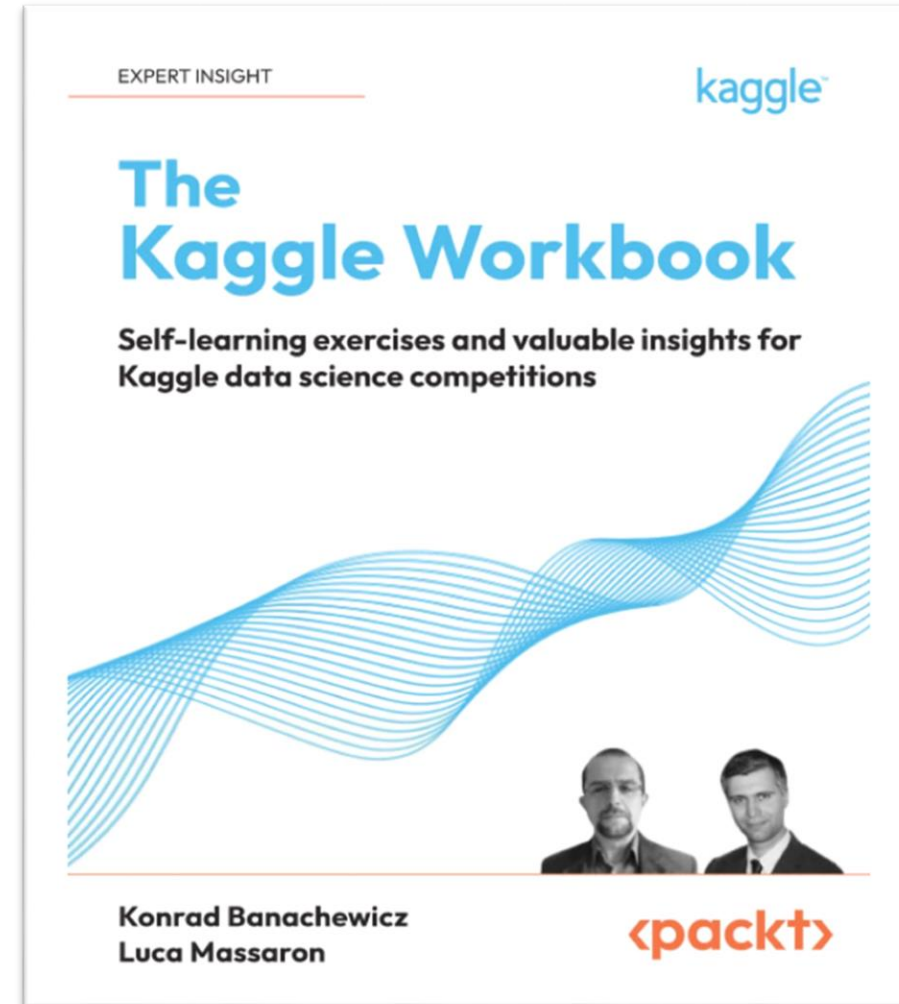
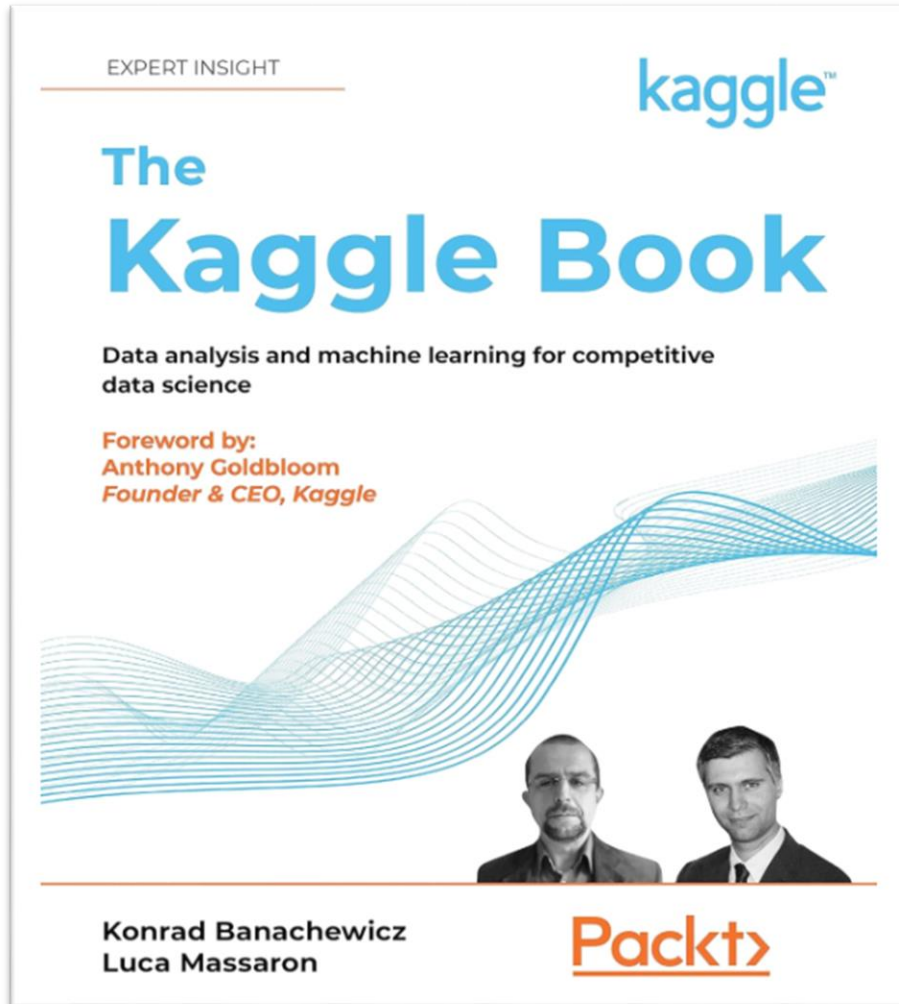
3

Most probably, many other people have tried to solve your problem already: **get connected, stay updated**

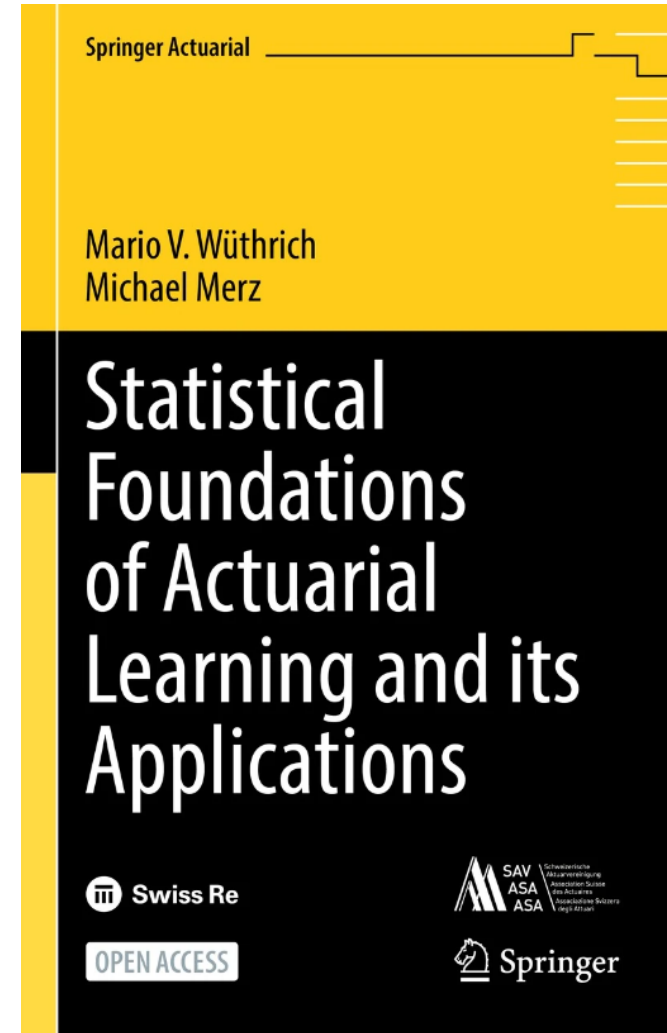
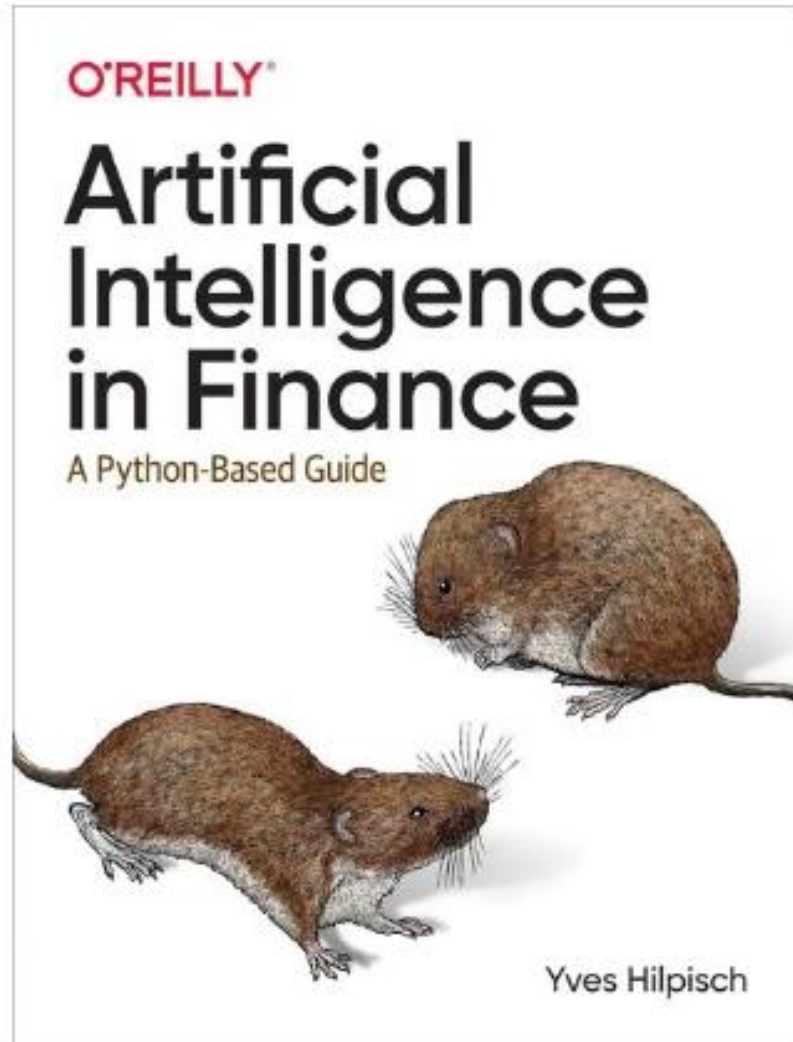
For those of you interested in machine learning, its models and algorithms



For those of you interested in competitive machine learning



For those of you interested in machine learning applications (finance and insurance)



After 7+1 lectures, 568 slides, 95 mini-exercises, and 7 Google Colab notebooks...

...this is the end of Block 1. Grazie per l'attenzione.