



# Online Machine Learning

---

## Ensemble Methods

Antonio Carta

antonio.cart@unipi.it

- **Basics:** what is an ensemble and why we need them in OML
- **Offline Ensembling** Methods
- **Online Ensembling** Methods

# General Concepts

*In 1907, [Sir Francis Galton](#) asked 787 villagers to guess the weight of an ox. None of them got the right answer, but when Galton averaged their guesses, he arrived at a near perfect estimate*

# What is an Ensemble?



- A set of models
- Combined together
- With a mechanisms to compute the output

## Questions:

- How to train each model (e.g. what data to use)
- How to compute the output

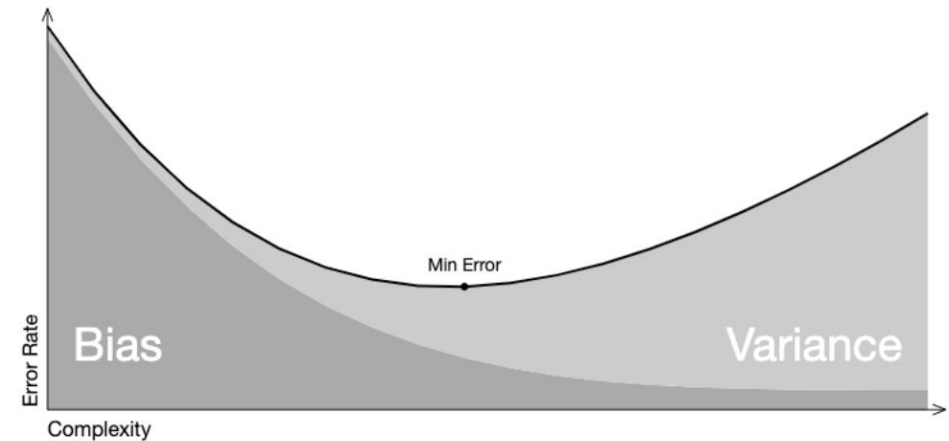
Many techniques we see today are offline methods. Why are we studying them?

- Ensembles are a popular technique in OML
- In the presence of **drifts**, models will be retrained, forgetting old concepts
- In the presence of **recurrent concepts**, this is a problem
- We can fix it by having a **separate model for each concept**
- Then, we can **select the correct model** at each step (easier said than done)

# Bias-Variance Tradeoff



- **Bias:**
  - **Simple model** → ignores relevant information, error due to **bias is high**
  - **Complex models** → error due to bias decreases (**low bias**)
- **Variance:**
  - Simple model → less error due to variance
  - Complex model → error due to variance increases.
- **Tradeoff:** The ideal model finds the optimal bias-variance tradeoff



# Ensemble Definition



*“An ensemble can be described as a **composition** of multiple **weak learners** to form one with (expected) higher predictive performance (strong learner), such that a weak learner is loosely defined as a learner that performs slightly better than random guessing”*

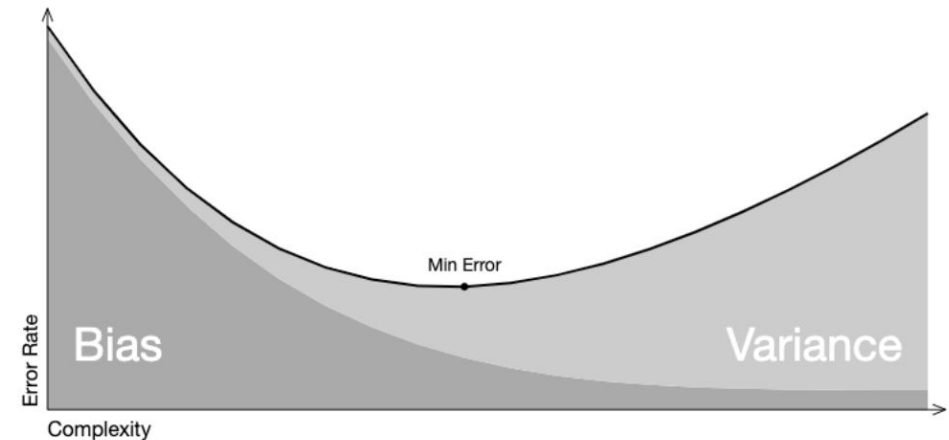
*Freund and Schapire, 1997*



# Weak Learners and Ensembling



- **Weak learners** have high-bias/low-variance
- An **ensemble** of weak learners improves the performance
- How do we choose and train the weak learner?
- **KEY IDEA:** The bias is decreased if the errors are decorrelated. We have to ensure the weak learners are different.



- **Diversity**: induce diversity among learners
- **Combination**: combine the predictions
- **Adaptation (for OML)**: adapt to evolving data

## **Advantages:**

- High Predictive performance
- Flexibility

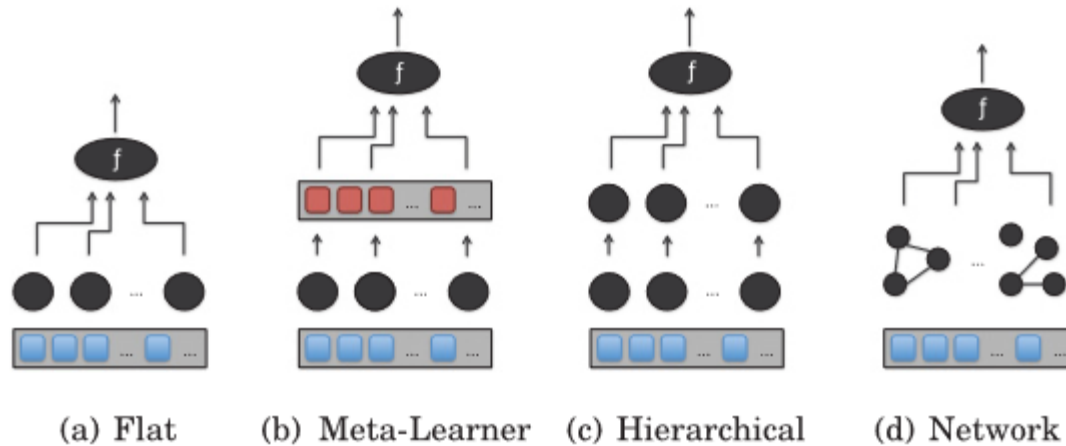
## **Disadvantages:**

- Computational resources

- Ensembles work best when the base models are different and have decorrelated errors
- **horizontal partitioning**: training classifiers in different chunks of data
- **vertical partitioning**: training with different subsets of features

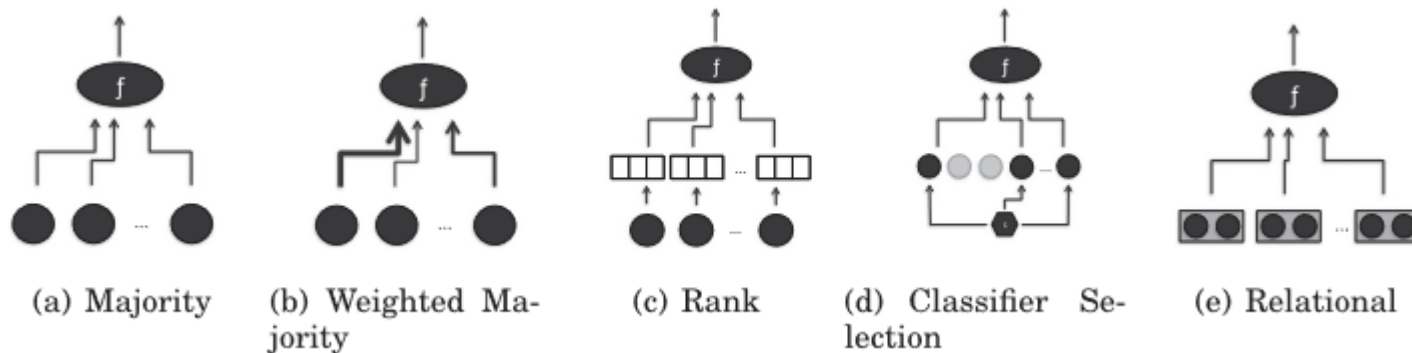
# Ensemble – Combination

- **Flat:** base models trained on input data, decision fusion via simple combination
- **Meta-learner:** the combination function is a model itself (meta-learned)
- **Hierarchical:** structured organization of the base learners
- **Network:** base models are nodes in a graph. The graph structure informs the combination scheme



# Ensemble – Voting Schemes

- **Majority vote:** every classifier has the same weight
- **Weighted majority:** gives a different weight to each classifier
- **Classifier selection:** uses a dynamic criterion to select the best classifier for the current sample



# Ensemble – Adaptation



- How many base models?
- **Fixed:** the number of base learners cannot grow
- **Dynamic:** add classifiers on the fly
  - We still need to limit the memory growth somehow

# Offline Ensembles

# Bagging – Bootstrap Aggregating



**IDEA:** Increase diversity by training on different subsets of the data

**Bootstrap:** sample with replacement

- An example of horizontal partitioning

**Training:**

- Sample  $m$  iid datasets of dimension  $n$  from the original data (with replacement)
- Train  $m$  models independently, one for each bootstrapped set

**Inference:** Combine outputs (average or majority voting)

**Properties:**

- Bagging reduces the variance
- Base models should have high variance to encourage diversity

Original Training Set

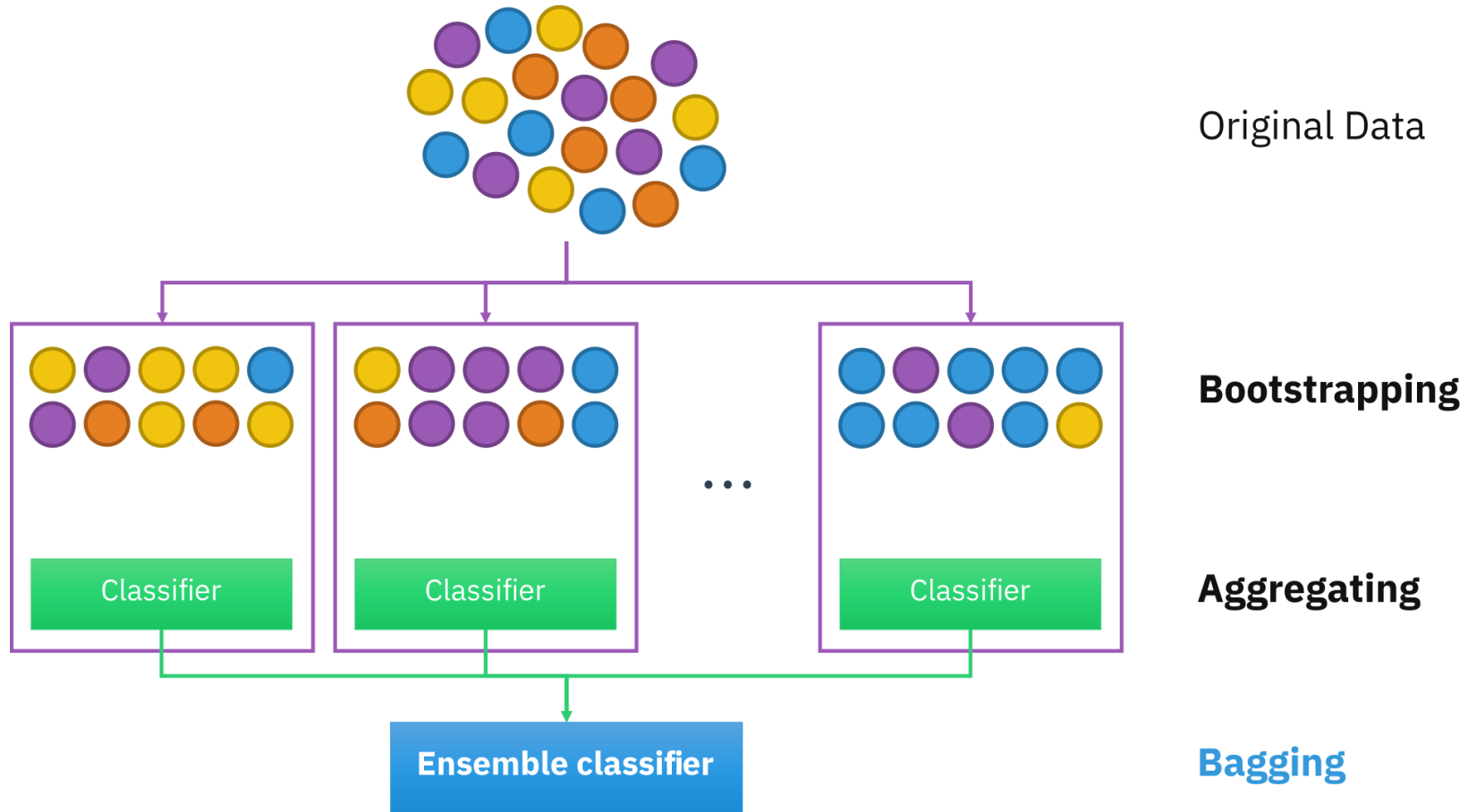


Possible Bootstrapped Sets

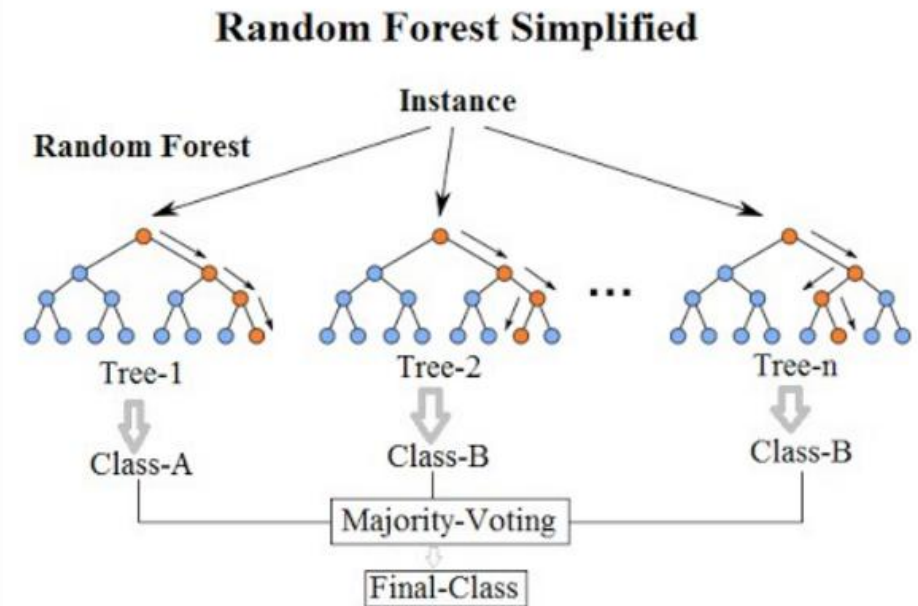




# Bagging

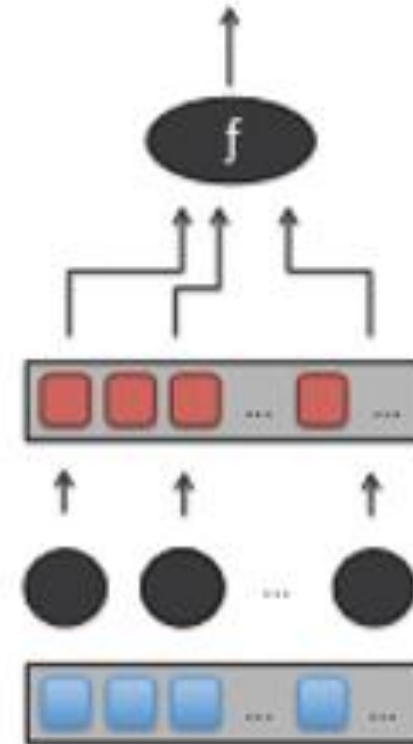


- Popular example of bagging
- **Train**  $m$  decision trees, each one on their own bootstrapped dataset
- **Weak learners:**
  - Limit the tree depth
  - Vertical partitioning: train each decision tree on a subset of features



# Stacking

- aka meta-learner combination scheme
- Train base models
- Fit a meta-learner to combine the outputs
  - Often a very simple model, such as a linear combination of the outputs
  - Example:  $y_{ensemble} = \sum_i \alpha_i y_{base}^i$
- Advantages: Can learn the combination scheme



# Online Ensembles

# Streaming Cross Validation



How do we split the stream if we want to train (and evaluate) an ensemble? We need to ensure diversity and robust training and evaluation.

## K-fold distributed cross-validation:

each sample is used for testing in one classifier selected randomly, and used for training on all the others

- Adaptation of offline cross-validation
- Good use of the data (only one 1/k samples are unused)
- high redundancy

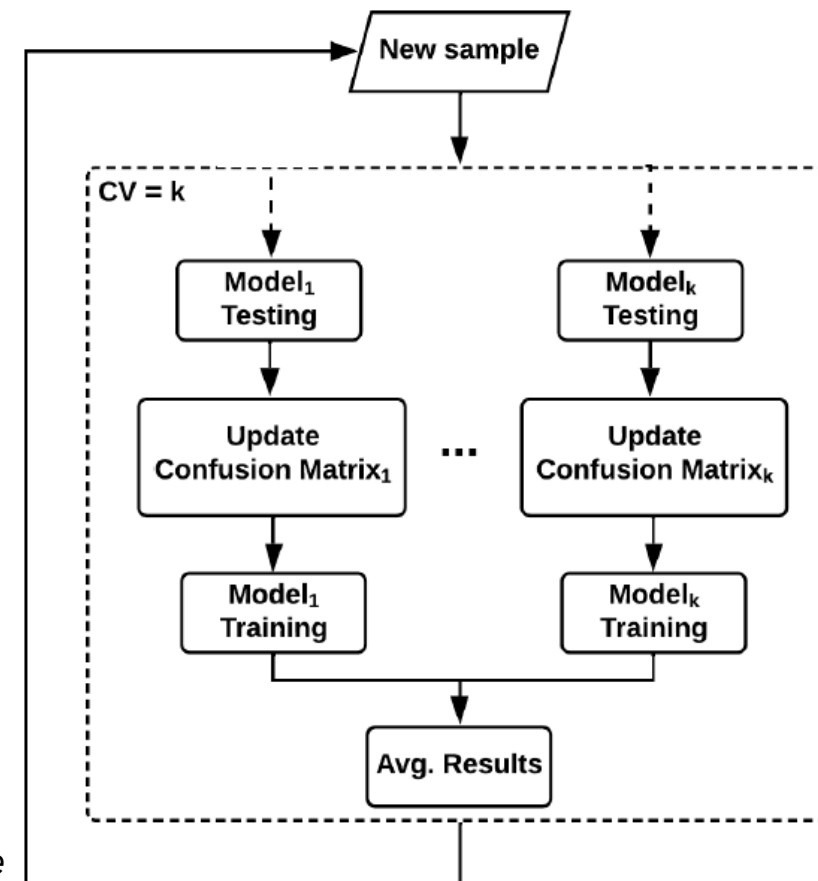
## K-fold distributed split-validation:

each sample is used for training in one classifier selected randomly, and for testing in all the other classifiers

- Models trained on disjoint data
- Under utilization of data. Each model uses 1/k data for training

## K-fold distributed bootstrap-validation:

each sample is used for training in approximately 2/3 of the classifiers, with a separate weight in each classifier, and for testing in all the classifiers (we will see online bootstrap in the following slides)



# Accuracy-Weighted Ensemble



**IDEA:** train an ensemble of offline methods. Adapt the ensemble by removing the oldest models.

- Process the stream in chunks
- Whenever a new chunk is ready, train a new model to add to the ensemble
  - Remove the oldest one to keep a fixed number of models in the ensemble
- **LIMITATION:** We have to fix the window size



# Accuracy-Weighted Ensemble



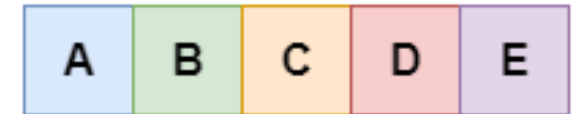
- **Weighting Scheme:** each classifier is weighted by the expected accuracy on the future data
  - Weight  $w_i = err_i - err_r$ 
    - $err_i$  error of classifier  $i$
    - $err_r$  error of random classifier
  - Ideally, estimated using a separate test set with recent data
  - Realistically estimated using the last chunk of data from the stream
- **Ensemble output:**  $f(x) = sign(\sum_i w_i C_i(x))$



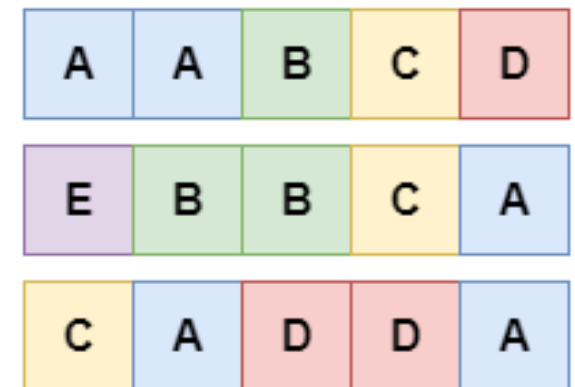
**IDEA:** We want to use bootstrap (remember: sampling with replacement). How can we do it online?

- **OFFLINE:** draw  $n$  random i.i.d. samples uniformly with replacement
- **ONLINE:** we give a discrete weight to the current sample in the stream
  - Equivalent to the number of copies of the current example that we would get in the offline setting
  - What distribution should we use?

Original Training Set



Possible Bootstrapped Sets





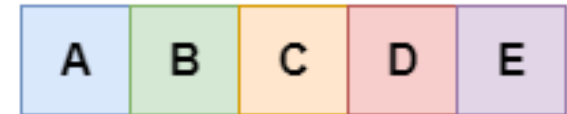
# Online Bagging – Binomial Distribution



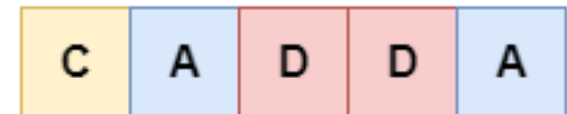
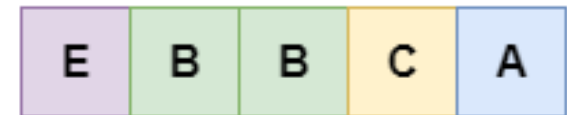
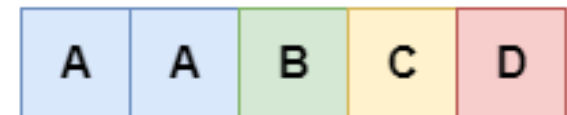
- $n$  = number of examples in a bootstrapped set
- $k$  = number of replicas for a specific example
- The number of replicas in a bootstrapped set follows a binomial distribution:
  - We perform  $n$  independent experiments where:
    - With probability  $p = \frac{1}{n}$  we pick the current example
    - With probability  $1 - p$  we pick another example

$$P(K = k) = \binom{n}{k} p^k (1 - p)^{n-k} = \binom{n}{k} \frac{1}{n^k} \left(1 - \frac{1}{n}\right)^{n-k}$$

Original Training Set



Possible Bootstrapped Sets



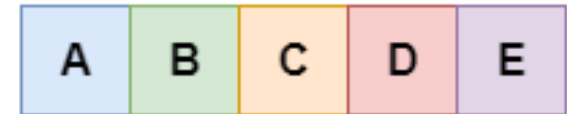
# Online Bagging – Binomial vs Poisson



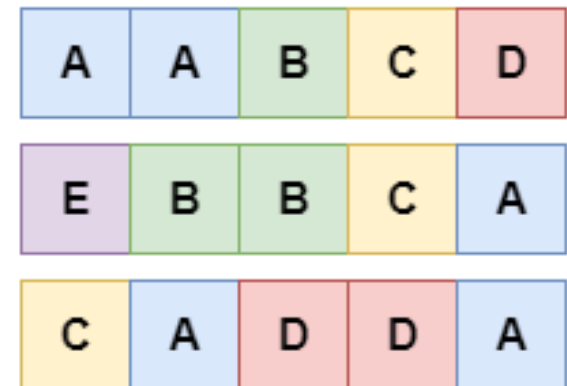
- For large  $n$ , the binomial can be approximated with a Poisson distribution
- **Online Bagging:** draw weights according to a Poisson with  $\lambda = 1$

$$\text{Poi}(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

Original Training Set



Possible Bootstrapped Sets



# Online Bagging – Pseudocode



ONLINE BAGGING(*Stream*,  $M$ )

Input: a stream of pairs  $(x, y)$ , parameter  $M =$  ensemble size

Output: a stream of predictions  $\hat{y}$  for each  $x$

```
1 initialize base models  $h_m$  for all  $m \in \{1, 2, \dots, M\}$ 
2 for each example  $(x, y)$  in Stream
3     do predict  $\hat{y} \leftarrow \arg \max_{y \in Y} \sum_{t=1}^T I(h_t(x) = y)$ 
4         for  $m = 1, 2, \dots, M$ 
5             do  $w \leftarrow \text{Poisson}(1)$ 
6             update  $h_m$  with example  $(x, y)$  and weight  $w$ 
```

Figure 7.2

Online Bagging for  $M$  models. The indicator function  $I(\text{condition})$  returns 1 if the condition is true, and 0 otherwise.

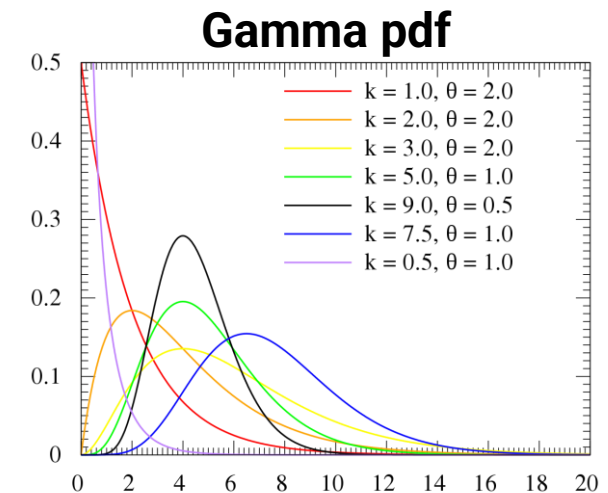
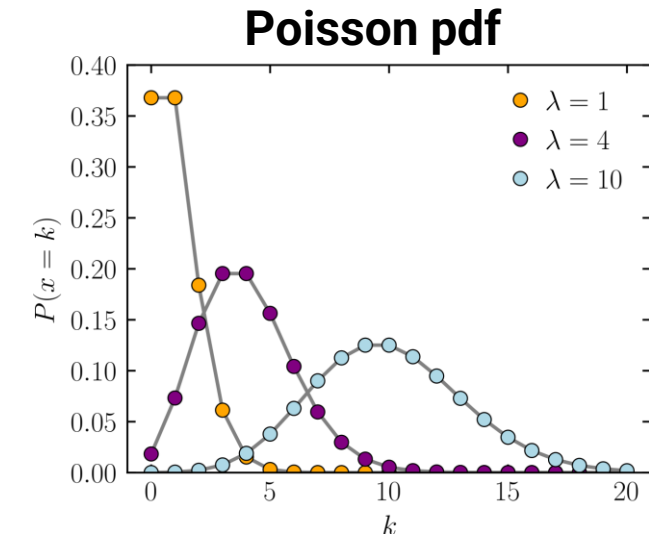
## IDEA: Bagging in the presence of Concept Drift

- **ADWIN Bagging:**
  - Bagging ensemble
  - ADWIN for concept drift detection
- **Keep  $m$  CD detectors (ADWIN), one for each base model**
  - REMEMBER: ADWIN looks at the stream of errors and consider a CD whenever the error decreases too much between the old and new window (according to a statistical test)
- **When a CD detector finds a CD:**
  - Remove the worst classifier («*replace the loser*» strategy)
  - Train a new one

# Leveraging Bagging



- Adding more randomness can improve the performance
- We can also change  $\lambda$  in a Poisson distribution to control the weights
  - Higher  $\lambda$  increases the mean and results in a flatter distribution
- Alternative distributions: Gamma  $\Gamma(k, \theta)$
- Example: **Leveraging Bagging**
  - Samples from Poisson with  $\lambda \geq 1$



Ref: MOA book (Ch. 7.4)

Image: Di Skbkekas - Opera propria, CC BY 3.0, <https://commons.wikimedia.org/w/index.php?curid=9447142>

Image: By Gamma\_distribution\_pdf.png: MarkSweep and Cburnettderivative work: Autopilot (talk) - Gamma\_distribution\_pdf.png, CC BY-SA 3.0,

<https://commons.wikimedia.org/w/index.php?curid=10734916>

# Leveraging Bagging



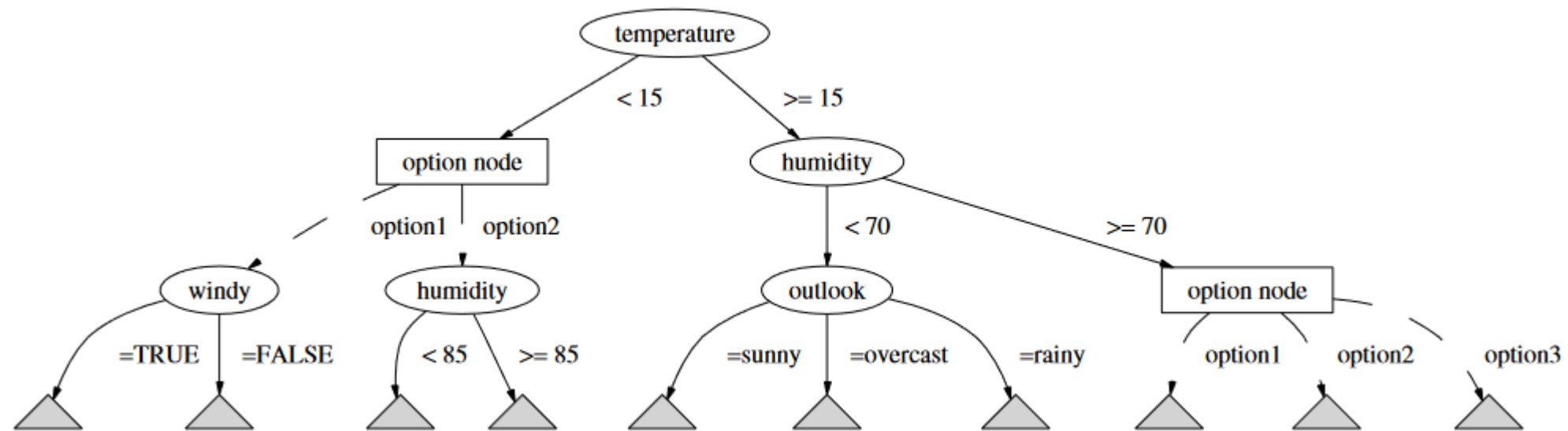
- Sample from **Poisson** with  $\lambda \geq 1$
- **ADWIN Bagging** with replace the loser strategy
- **Error-correcting output codes:**
  - In multi-class problems map set of classes into two sets  $\{0, 1\}$
  - Converts a multi-class problem into a binary classification problem
  - Each base model uses a different code
    - Increases diversity because each model is learning a different function

# Hoeffding Option Tree (HOT)



- Ensemble of trees in a single tree structure
- **Option Nodes** represent an ensemble of subtrees
- **INFERENCE:**
  - Each subtree is evaluated in parallel
  - Output from the leaves are aggregated (e.g. weighted voting)
- **TRAINING:**
  - If the splitting criterion for different attributes is similar and we have enough data (according to the Hoeffding bound), we split the tree into **several options**
  - Less instability
  - We can split earlier
  - REMINDER: in the VFDT, if attributes are close enough, we pick only the best.

# Hoeffding Option Tree



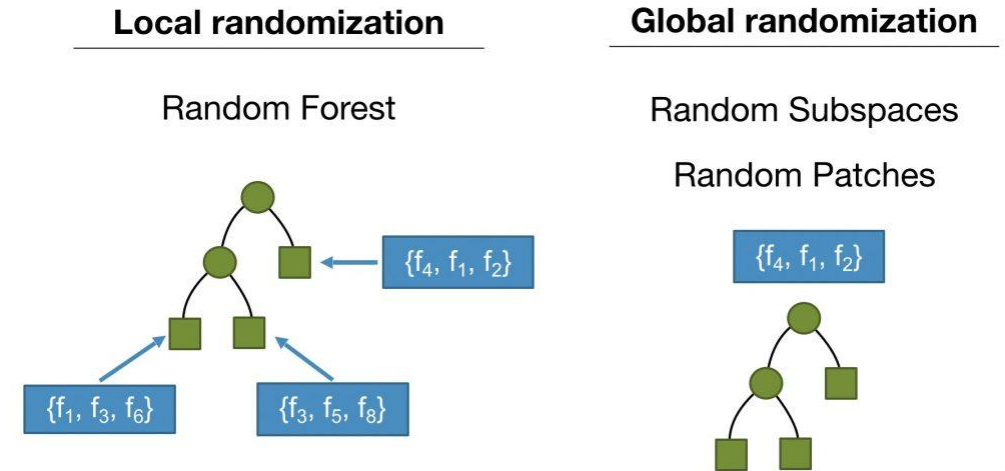


# Streaming Random Patches (SRP)



Adaptation of Random Forest for online learning

- **Adaptive Random Forest (ARF)** uses **local subspace randomization**: random subsets of features are set for each leaf to be considered for future node splits
- **SRP uses global subspace randomization**: each base model is trained on a randomly selected subset of features
- ARP and SRP are available in river



- **Recurrent Concepts:** in presence of concept drifts, previous concepts can reoccur
  - Seasonal trends (daily/weekly/yearly)
  - We don't want to learn from scratch every time
- **IDEA: keep two sets of classifiers:**
  - Library of available classifiers, currently inactive (past concepts)
  - Current ensemble of active models (current concepts)
  - **We need a policy to move the models from one set to the other**
    - in general, estimating the accuracy of a model on the current data requires much less samples than what we need for training

# Recurrent Concept Drift (RCD)



**IDEA:** keep a model for each concept, identify drifts, reuse models

- **DDM for CD Detection**

- REMINDER: DDM uses the errors to detect CD, assuming errors decrease over time for stationary data

- Train current model  $c_a$  and store training samples in buffer  $b_a$

- **At WARNING level**

- Start training a new model
- Train new classifier  $c_n$  and  $c_a$  (old classifier) and store samples in  $b_n$

- **IF DRIFT is confirmed**

- Compare  $b_n$  with all the stored buffers to check if it is a new/old concept

---

**Algorithm 1.** RCD algorithm
 

---

**Input:** ( $m$ ) Max buffer size, ( $\alpha$ ) Significance value, ( $e$ ) Ensemble size  
**Data:** ( $c_a$ ) Actual classifier, ( $b_a$ ) Actual buffer, ( $c_n$ ) New classifier, ( $S$ ) Data stream, ( $b_n$ ) New buffer  
**Result:**  $C$ : Classifiers list,  $B$ : Buffers list

```

1 begin
2    $C \leftarrow \{\text{Create}(c_a)\};$ 
3    $B \leftarrow \{\text{Create}(b_a)\};$ 
4   for each  $s \in S$  do
5      $level \leftarrow \text{DDM}(c_a, s);$ 
6     switch  $level$  bf do
7       case WARNING
8         if  $c_n = \text{null}$  bf then
9            $\text{Create}(c_n);$ 
10           $\text{Create}(b_n);$ 
11         end
12           $\text{SaveFIFO}(b_n, s, m);$ 
13           $\text{Train}(c_n, s);$ 
14         case DRIFT
15           if  $\text{StatTest}(C, B, b_n, \alpha)$  then
16              $(c_a, b_a) \leftarrow \text{Stored}(C, B, b_n);$ 
17           end
18         else
19            $\text{SaveFIFO}(C, c_n, e);$ 
20            $\text{SaveFIFO}(B, b_n, m);$ 
21            $c_a \leftarrow c_n;$ 
22            $b_a \leftarrow b_n;$ 
23            $c_n = b_n = \emptyset;$ 
24         otherwise
25           if  $c_n \neq \text{null}$  then
26              $c_n = b_n = \emptyset;$ 
27           end
28            $\text{SaveFIFO}(b_a, s, m);$ 
29         end
30       end
31        $\text{Train}(c_a, s);$ 
32     end
33 end
  
```

---

Init current model and buffer

WARNING: prepare a new model

DRIFT: check if current concept is already known

# Conclusion

- **Ensembles are one of the easiest and most reliable methods to improve the performance of an ML method**
  - As long as you can guarantee diversity
- **In the online world:**
  - Train ensembles on windows (AWE)
  - In nonstationary streams: each ensemble learns different concepts
  - Bagging -> online bagging

- Streaming Data Analytics Course - Emanuele Della Valle and Alessio Bernardo @ POLIMI
- MOA Book

- Lab with river
  - Concept drift
  - Online classification
  - Ensembles