

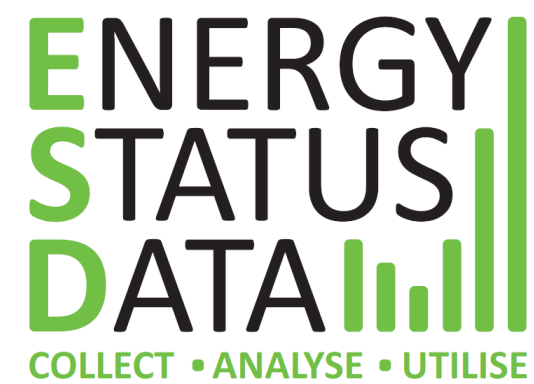
# Navigating Complex Machine Learning Challenges in Streaming Data

ECML Tutorial 2024

Heitor Murilo Gomes<sup>1\*</sup>, Marco Heyden<sup>2</sup>  
Maroua Bahri<sup>3,4</sup>

<https://heymarco.github.io/ecml24-streamingchallenges/>

\* Corresponding author: heitor.gomes@vuw.ac.nz



[1] Victoria University of Wellington, New Zealand, [2] KIT, Germany, [3] INRIA Paris, France,  
[4] Sorbonne Université, France

<https://capymoia.org/>

# Classification algorithms

# Hoeffding Tree\*

\* Also known as Very Fast Decision Tree (VFDT)

**Goal:** Grow a decision tree incrementally

This means that after every new training instance,  
the tree may grow

**Key question:** When should a split happen?

**Hypothesis:** A small sample is often enough to choose a near  
optimal split decision

# *Hoeffding Bound*

It is a statistical inequality that provides a theoretical guarantee on the convergence of sample averages to the true mean with a high probability

In other words, the **Hoeffding Bound** helps in determining whether **the observed differences in the attributes' merit (purity) are statistically significant** or merely due to random variation

# *Hoeffding Bound*

When should we split a node?

Let  $X_1$  and  $X_2$  be the top 2 most informative attributes\*

Is  $X_1$  a stable option?

Hoeffding bound, split on  $X_1$  if  $G(X_1) - G(X_2) > \epsilon$

Where  $G(*)$  is a purity measure  
(e.g. Gini index, Information gain)

\* The top attributes to split, the ones that will cause the splits to be “purer”

# Hoeffding Bound

When should we split a node?

Let  $X_1$  and  $X_2$  be the top 2 most informative attributes\*

$$\epsilon = \sqrt{\frac{R^2 \ln 1/\delta}{2n}}$$

Is  $X_1$  a stable option?

Hoeffding bound, split on  $X_1$  if  $G(X_1) - G(X_2) > \epsilon$

$R$  = Range of observed random variable

$\delta$  = The desired probability of the estimate not being within  $\epsilon$  of its expected value

$n$  = Number of observed instances

Where  $G(*)$  is a purity measure (e.g. Gini index, Information gain)

\* The top attributes to split, the ones that will cause the splits to be “purer”

# Hoeffding Tree

## wrap-up

- $\epsilon$  decreases with  $n$  (or the more instances observed)
- HT builds a tree that converges to the tree built by a batch learner given sufficiently large data
- A *grace period* can be used to avoid “splitting too fast”
- There are better options w.r.t. theoretical guarantees (See McDiarmid Trees\*), but HTs still works well in practice

\* Rutkowski, L., Pietruczuk, L., Duda, P., & Jaworski, M. (2012). Decision trees for mining data streams based on the McDiarmid's bound. *IEEE Transactions on Knowledge and Data Engineering*.

# Other Streaming Decision Tree algorithms

- The **Extremely Fast Decision Tree (EFDT) [1]** algorithm improves upon the Hoeffding Tree by using a more relaxed criterion for splitting nodes, allowing it to grow the tree more quickly.
  - EFDT can revisit and revise earlier splits if better splits are found later, which can lead to subtree pruning and potentially sudden drops in accuracy.
- **PLASTIC [2]** avoids EFDT's subtree pruning by restructuring the tree rather than pruning it when revising splits. Streaming Data session, 11:00am, Thursday
  - This allows PLASTIC to maintain accuracy by rearranging the tree structure without discarding information

[1] Manapragada, Chaitanya, Geoffrey I. Webb, and Mahsa Salehi. Extremely fast decision tree. ACM SIGKDD International, 2018

[2] M. Heyden, H. M. Gomes, E. Fouché, B. Pfahringer, and K. Bohm. Leveraging Plasticity in Incremental Decision Trees. ECML-PKDD, [To Appear] 2024



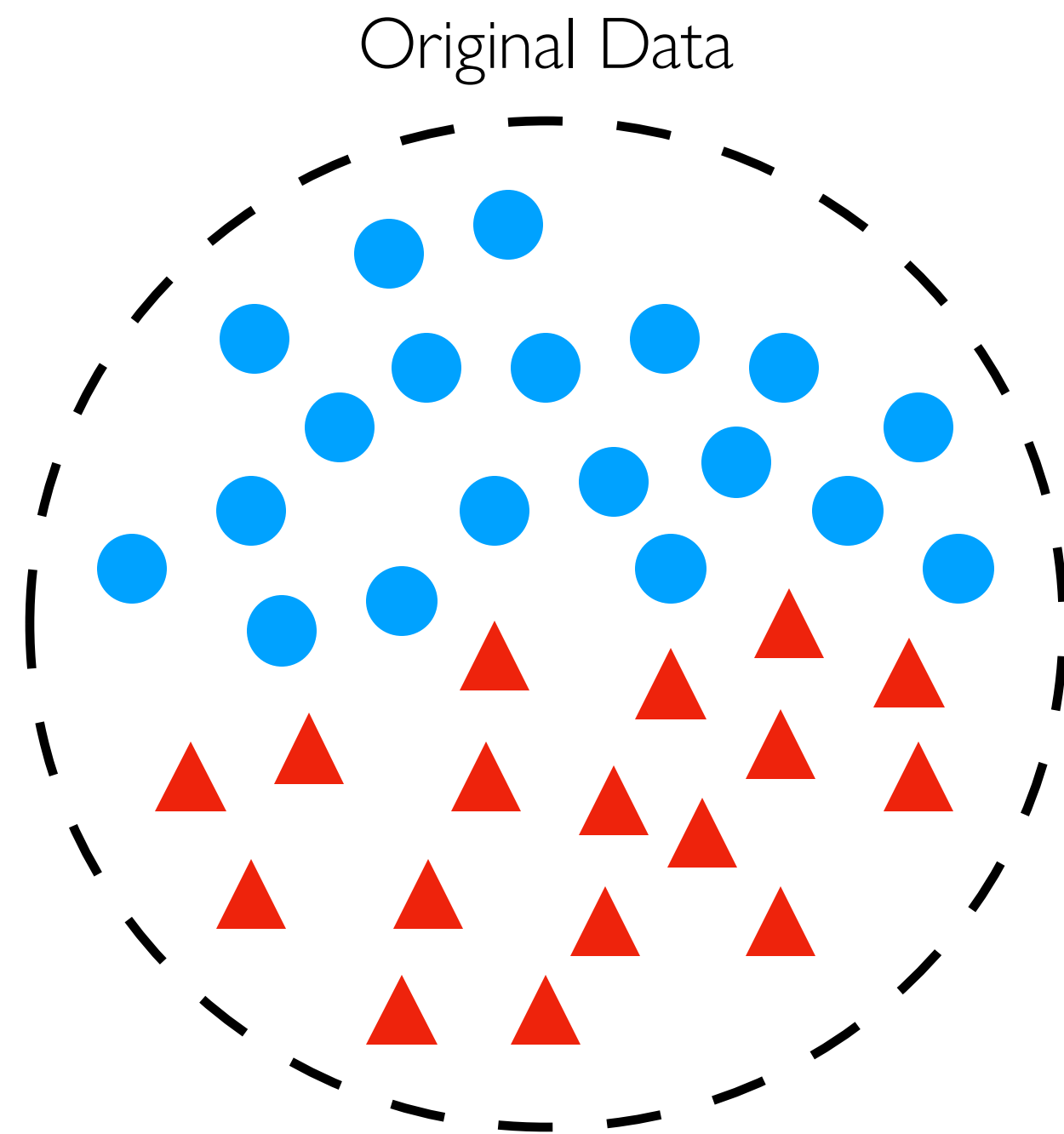
# Bagging

## Bootstrap Aggregating

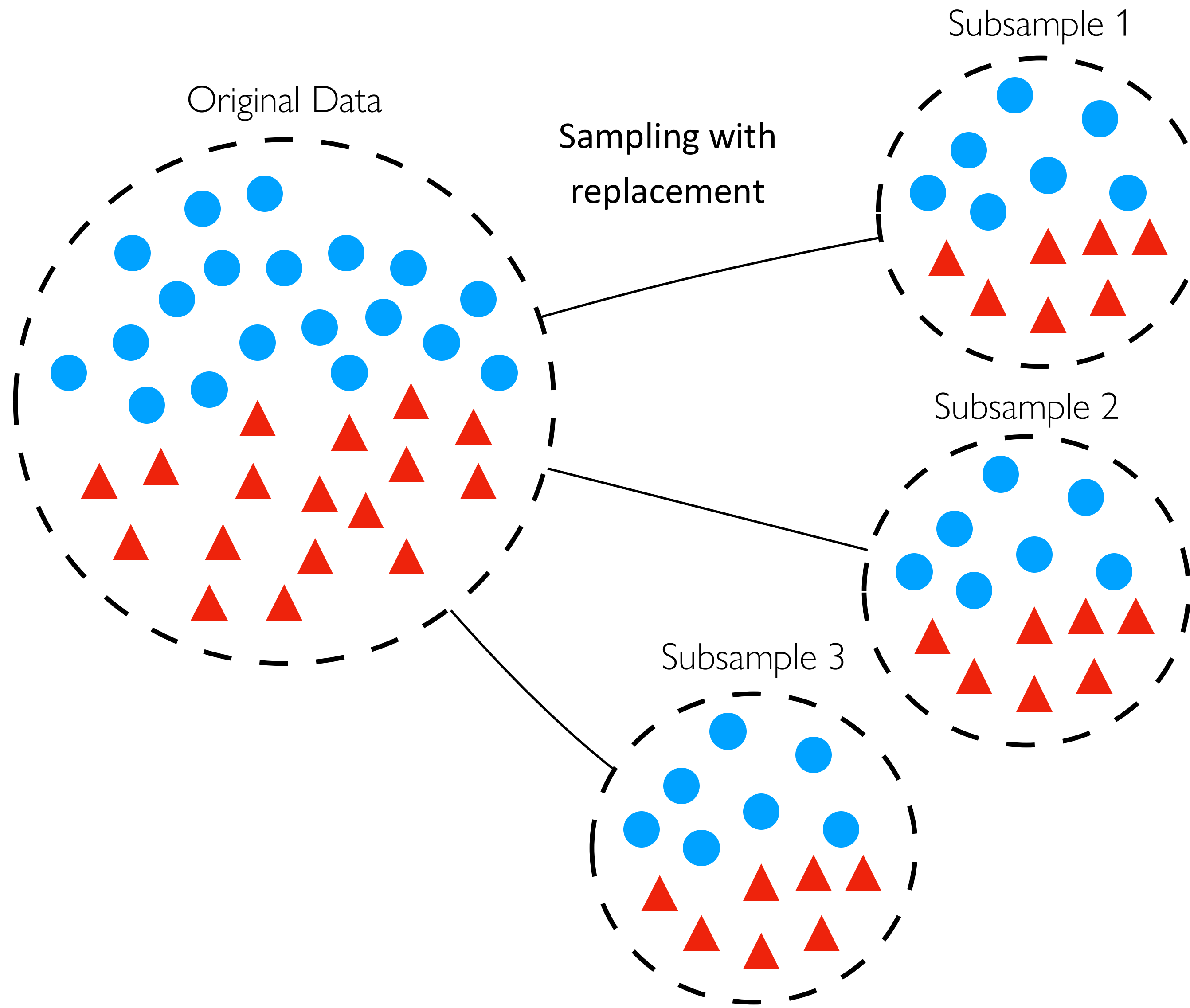
**Bagging** trains each model of the ensemble with a **bootstrap sample** from the original dataset.

Every bootstrap contains each original sample **K** times, where  **$\Pr(K=k)$**  follows a binomial distribution.

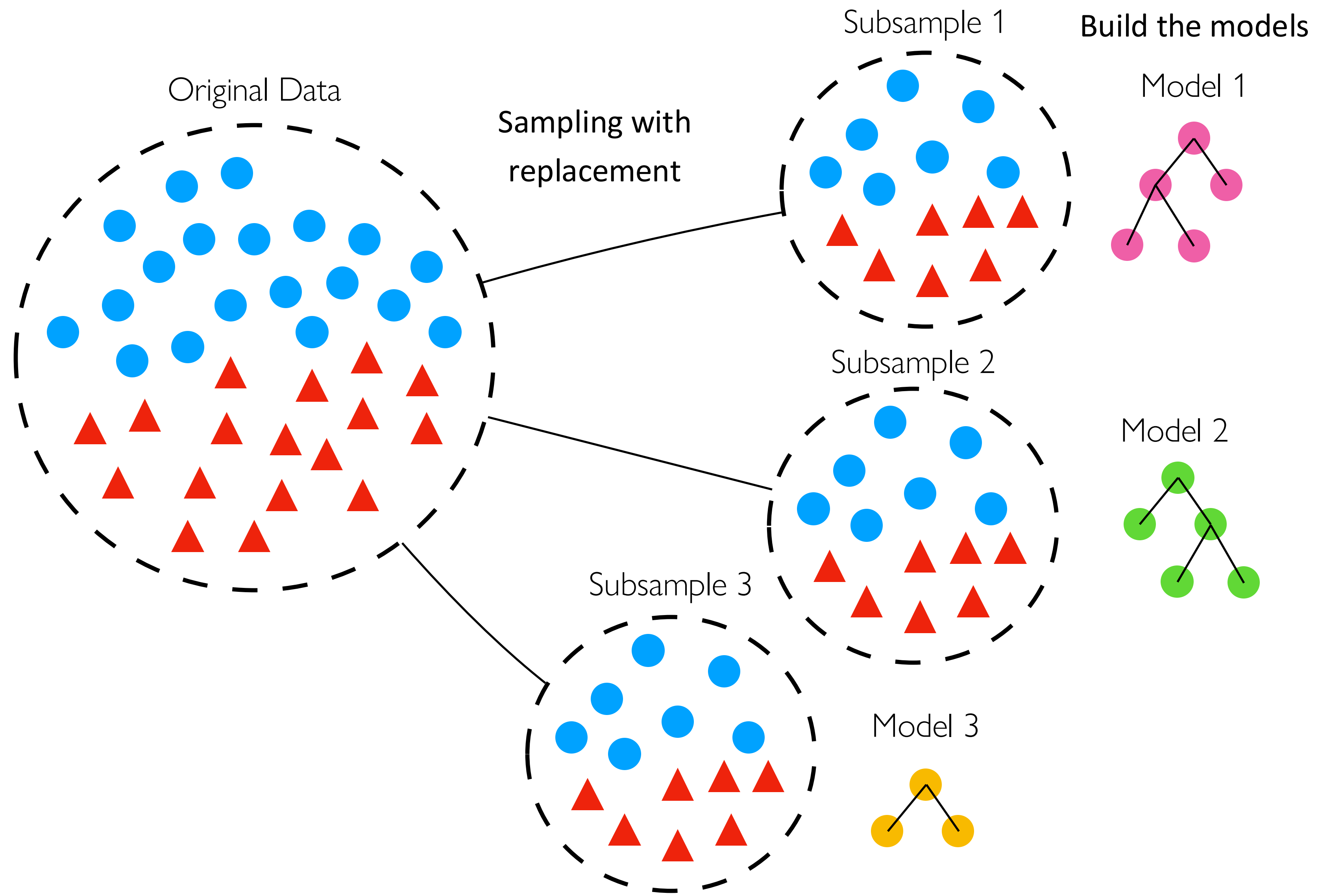
# Bagging



# Bagging



# Bagging



# Bagging

On average for each subsample:

~64% of the instances are from the original dataset

~37% are repeated instances

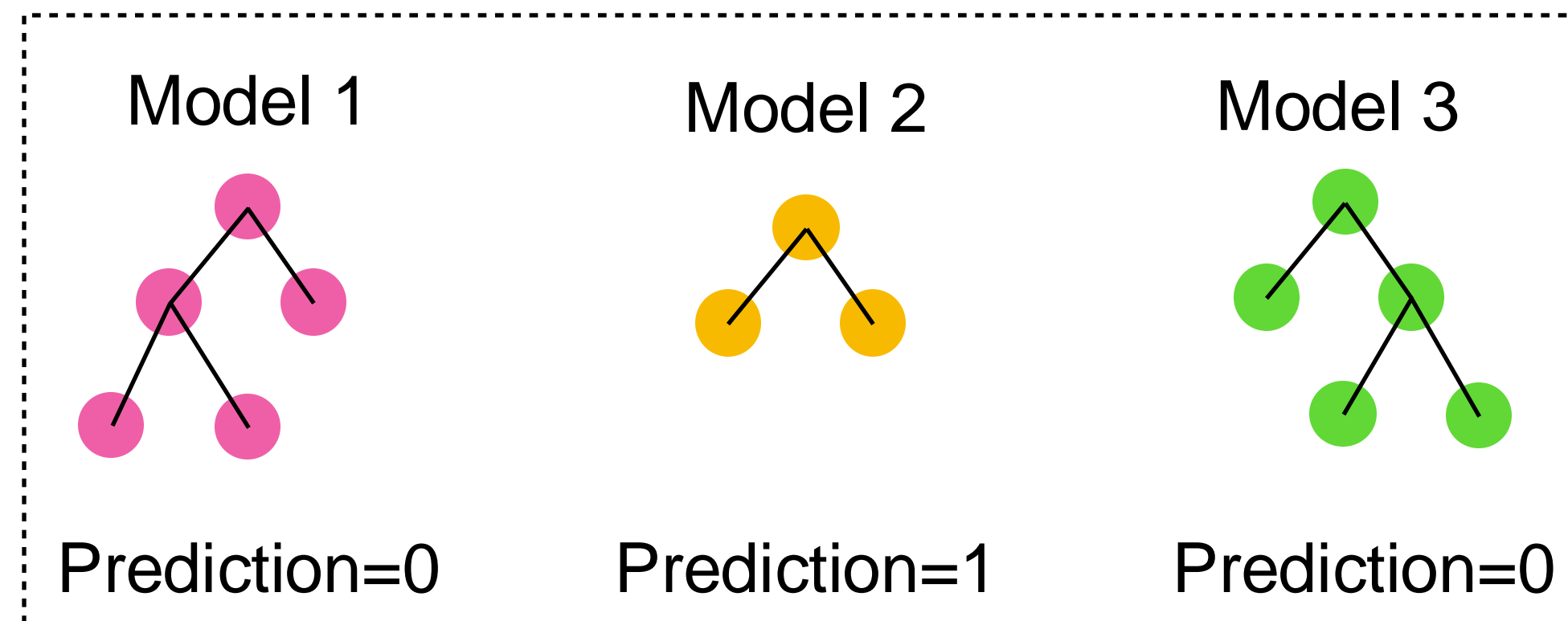
~37% of the original instances are not present\*

\* Out-Of-Bag (OOB)

# Bagging

The **predictions** of each learner are **aggregated** using majority vote to obtain the final prediction.

Prediction for a given instance X...



Ensemble  
Prediction=0

# Online Bagging

- We cannot apply Bagging directly to data streams...
- Unfeasible to store all data before creating each bootstrap subsample

We need to build the subsamples online

# Online Bagging

- Given a dataset with **N** samples
- In Bagging, every bootstrap contains each original sample **K** times, where **Pr(K=k)** follows a binomial distribution
- Oza and Russel found out that for large **N**, the binomial distribution tends to a **Poisson(1)** distribution
- Online Bagging instead of sampling with replacement, gives each example a weight according to **Poisson(1)** distribution



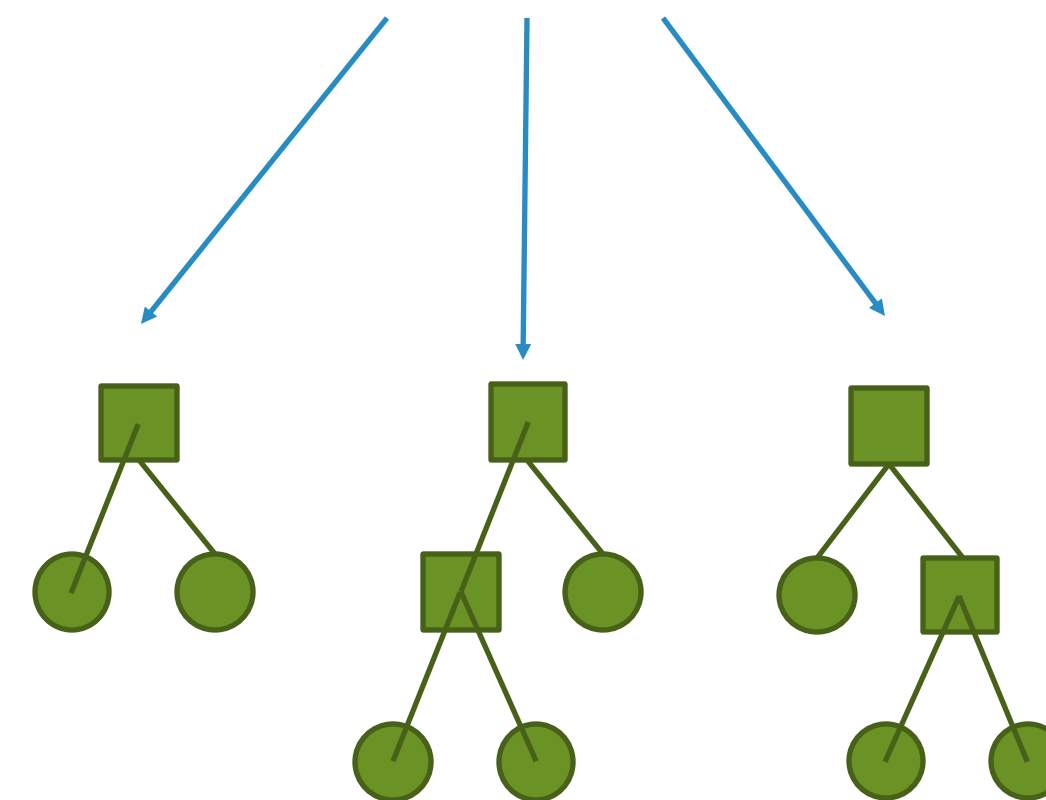
# Online Bagging

```
 $k \leftarrow \text{Poisson}(\lambda=1)$   
if  $k > 0$  then  
   $l \leftarrow \text{FindLeaf}(t, x)$   
   $\text{UpdateLeafCounts}(l, x, k)$ 
```

**Practical effect:** train learners with different subsets of instances.

stream ...  $(x^t, y^t)$  ...

$k$  “weight” train



# Subsamples

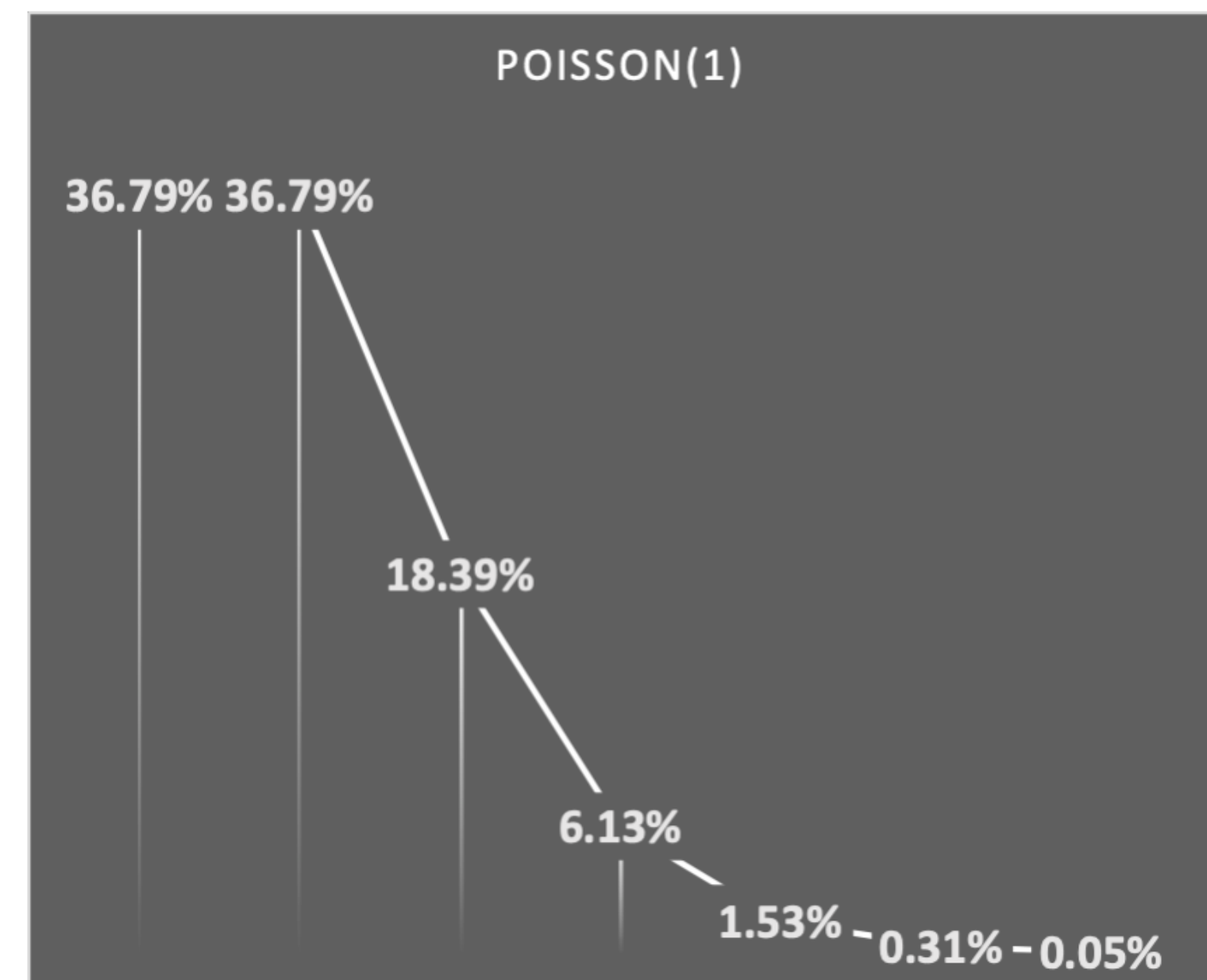
## Batch bagging

~64% from the original dataset

~37% are repeated

~37% are not present

## Online bagging



$k = 0$  1 2 3 4 5 6 ...

# Adaptive Random Forest (ARF)

Streaming version of the original Random Forest by Breiman

Uses a variation of the Hoeffding Tree

## Main differences:

Online bagging, base learner & detectors

## Overview:

1. Online bagging
2. Random subset of features
3. Drift detector for each tree

Hulten, G., Spencer, L., & Domingos, P. (2001). Mining time-changing data streams. In *ACM SIGKDD*.

Breiman, L. (2001). Random forests. *Machine learning*.

Gomes, H. M., Bifet, A., Read, J., ..., T. (2017). Adaptive random forests for evolving data stream classification. *Machine Learning*.

# ARF: Detect and Adapt

- One **Warning** and one **Drift** detector per base model
- Relies on the **Adaptive WINdow** (ADWIN) algorithm for detection (other algorithms could be used)
- **Background learners** are started once a warning is detected, their subspace of features may not correspond to the subspace of features used by the “*foreground*” learner.
- Once a drift is detected, the **background learner replaces** the “*foreground*” learner.

# Boosting and Gradient Boosting

- XGBoost and CatBoost are popular **batch gradient boosting** methods
- One key challenge when adapting such algorithms other than a stream setting includes concept drift recovery

# Boosting on Streams

- [2001] **OzaBoost** [1] uses weights from a Poisson(1) distribution to train multiple times using a given instance (similar to Online Bagging)
- [2012] **Online Smooth Boost** [2] is analogous to batch SmoothBoost, thus it uses a smooth distribution for weight assignment
- [2020] **Gradient boosted AXGB** [3] use Mini-batch trained XGBoost as its base learners and adjusts the booster when concept drifts are detected by ADWIN
- [2024] **Streaming Gradient Boosted Trees (SGBT)** [4] performs **better than** bagging-based stream learners

[1] N. Oza and S. Russel "Online bagging and boosting" *Artificial Intelligence and Statistics*, 2001

[2] Chen, Shang-Tse, Hsuan-Tien Lin, and Chi-Jen Lu. "An online boosting algorithm with theoretical justifications." *International Conference on International Conference on Machine Learning*. 2012.

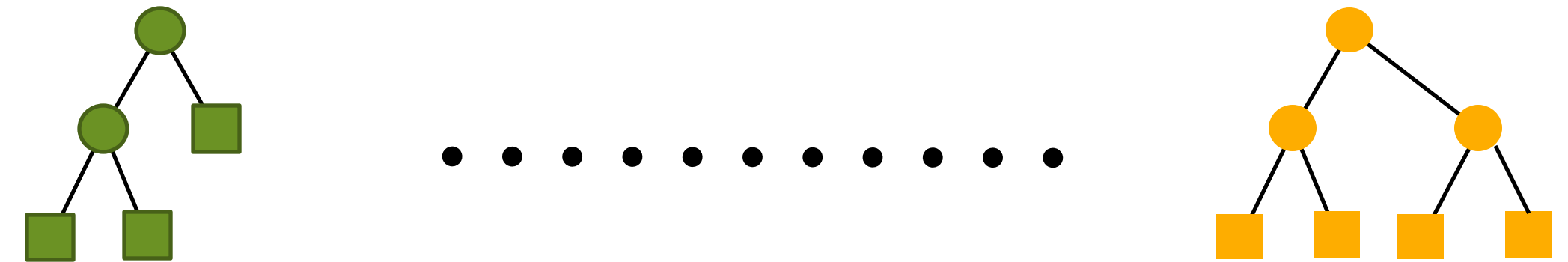
[3] Montiel, J., Mitchell, R., Frank, E., Pfahringer, B., Abdessalem, T., Bifet, A.: Adaptive xgboost for evolving data streams. In: 2020 IJCNN

[4] Gunasekara, N., Pfahringer, B., Gomes, H., & Bifet, A. (2024). Gradient boosted trees for evolving data streams. *Machine Learning*, 113(5), 3325-3352.

# Regression algorithms

# Adaptive Random Forest Regression

- Similar to ARF for classification
- builds **regression trees**
- for **prediction**, uses **mean of predictions** (by each tree)

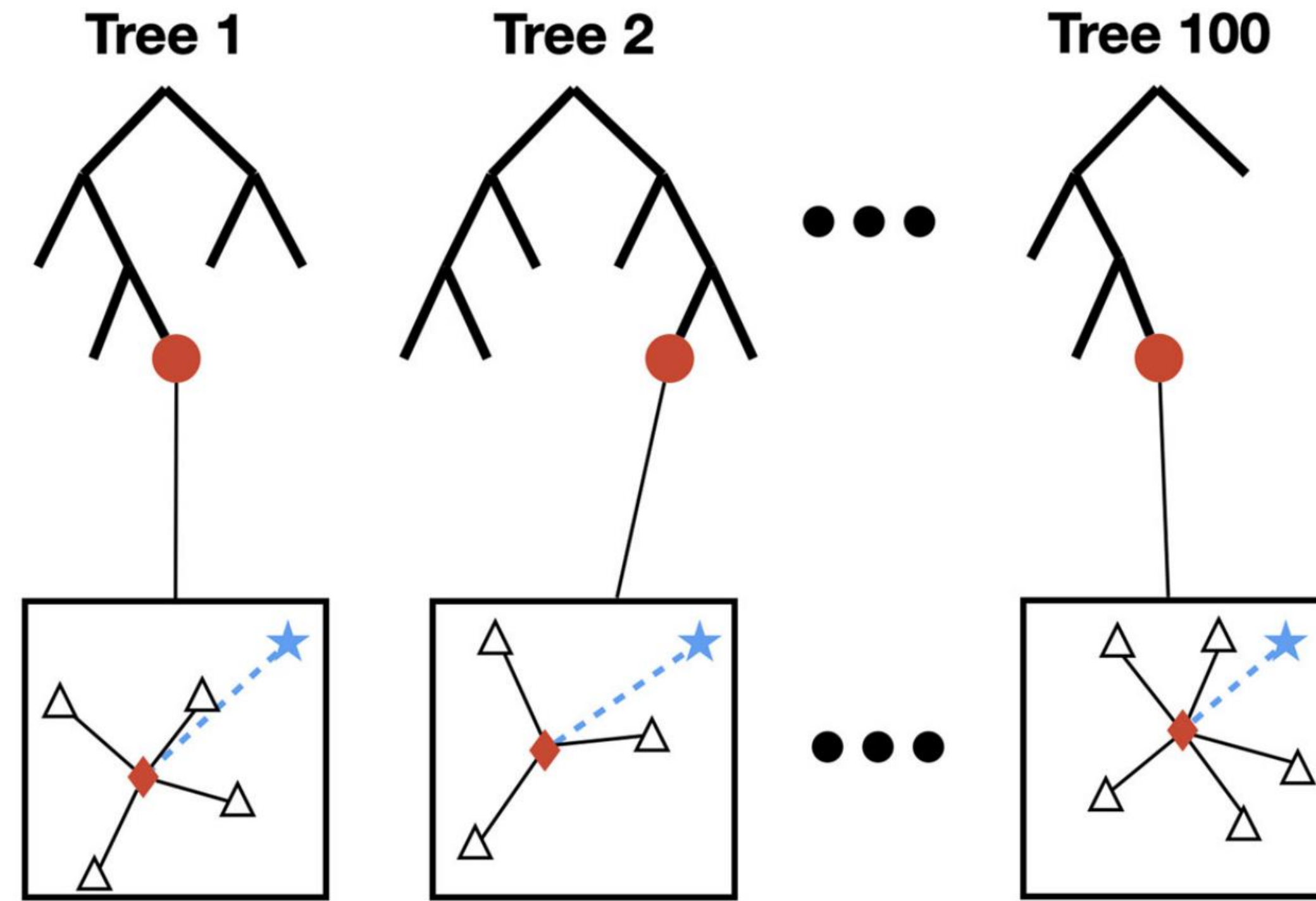




# Self-Optimizing k-Nearest Leaves (SOKNL)

- **Extends** Adaptive Random Forest Regression
- Generates a **representative data point (centroid)** in each leaf by **compressing** information from all instances in that leaf
- During **prediction**, calculates **distances** between **input instance** and **centroids** for **relevant leaves**
- Uses **only k leaves** with **smallest distances** for **prediction**
- **Dynamically tuning k** values based on **historical information**

# Self-Optimizing **k**-Nearest **L**eaves (SO**k**N**L**)



# Prediction Intervals

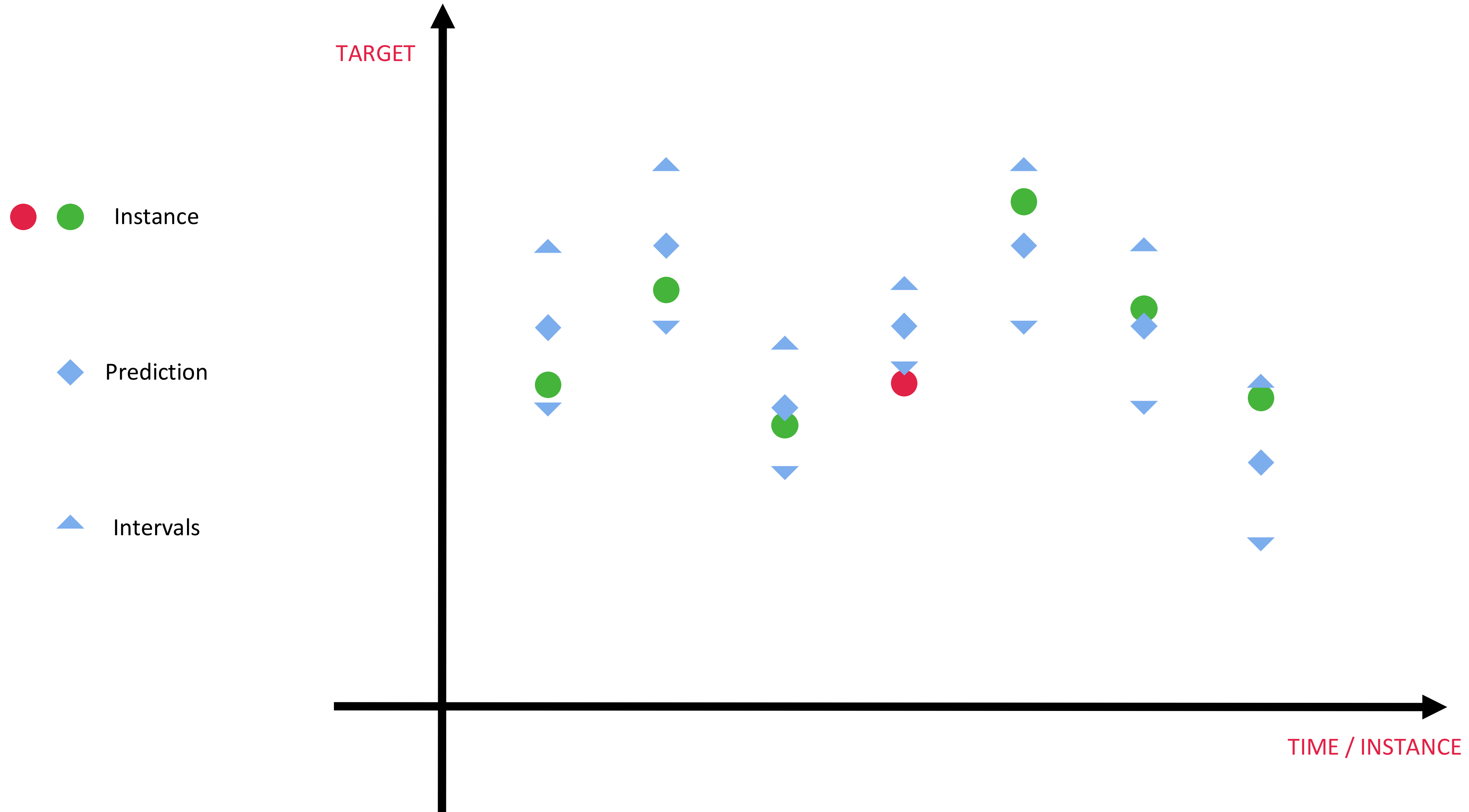
# Prediction Intervals

Prediction Intervals (PIs) are very useful improve our confidence in predictions yield in regression tasks

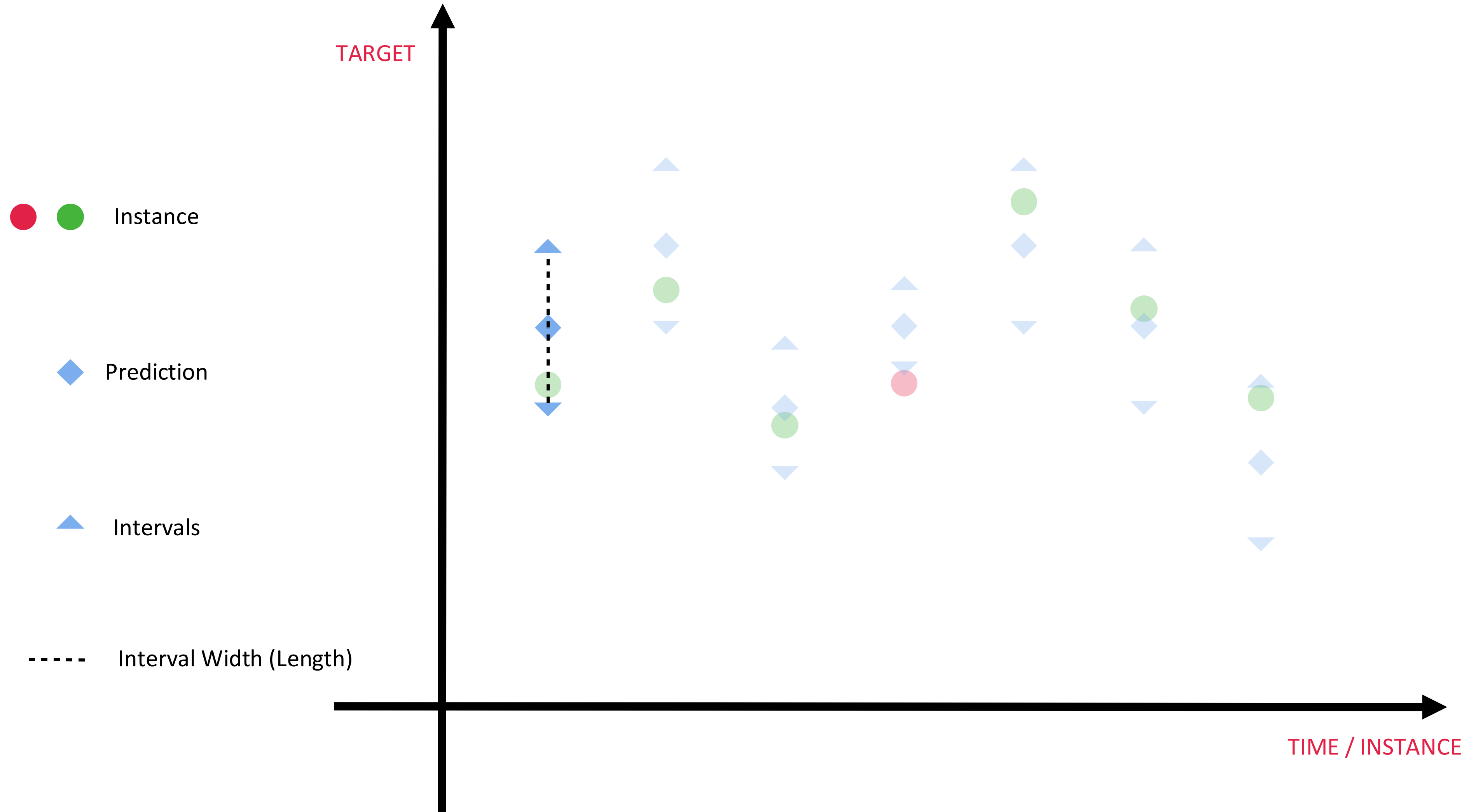
## **Challenge**

Traditional PI methods were not designed to adapt to evolving streams (i.e. those with concept drift)

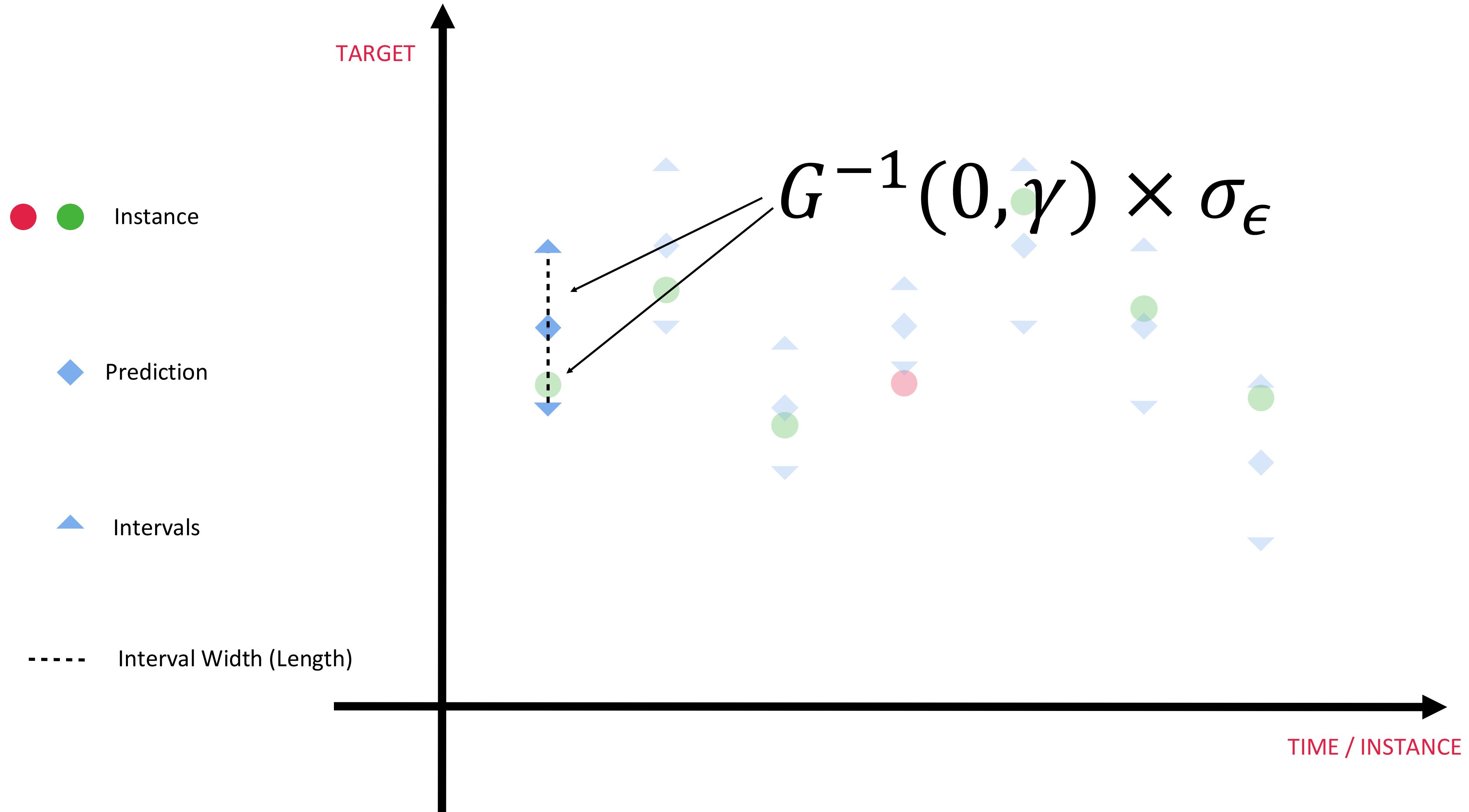
# Prediction Interval



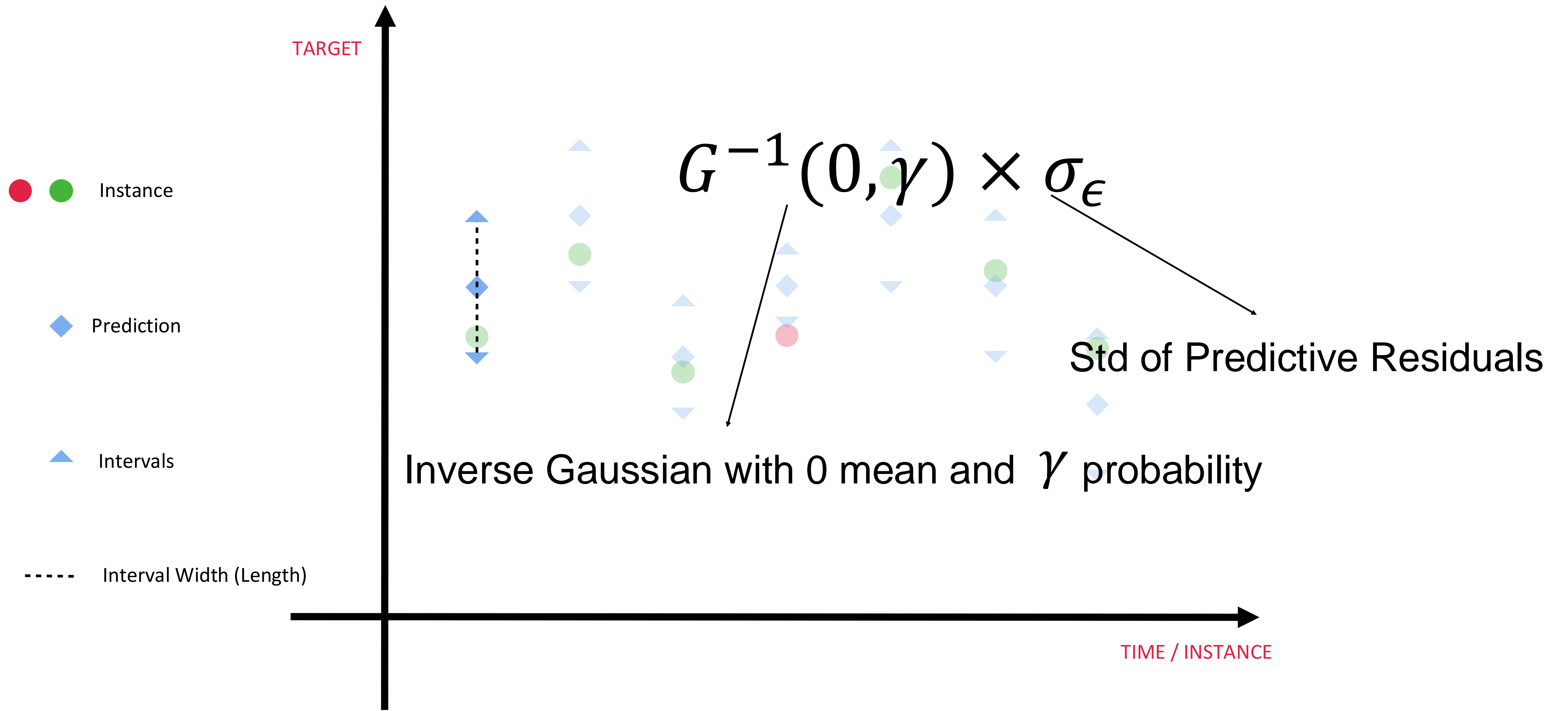
# Prediction Interval



# Mean and Variance Estimation (MVE)



# Mean and Variance Estimation (MVE)





## Mean and Variance Estimation (MVE)

$$\Pr\left(y \in \left[\hat{y} - G^{-1}(0,\gamma) \times \sigma_\epsilon, \hat{y} + G^{-1}(0,\gamma) \times \sigma_\epsilon\right]\right) \approx \gamma$$

$$\text{PI}_{\text{MVE}} \in \left(\hat{y} - G^{-1}(0,\gamma) \times \sigma_\epsilon, \hat{y} + G^{-1}(0,\gamma) \times \sigma_\epsilon\right)$$

$\gamma$  : Confidence Level / Significance Level

# Adaptive Prediction Interval (AdaPI)

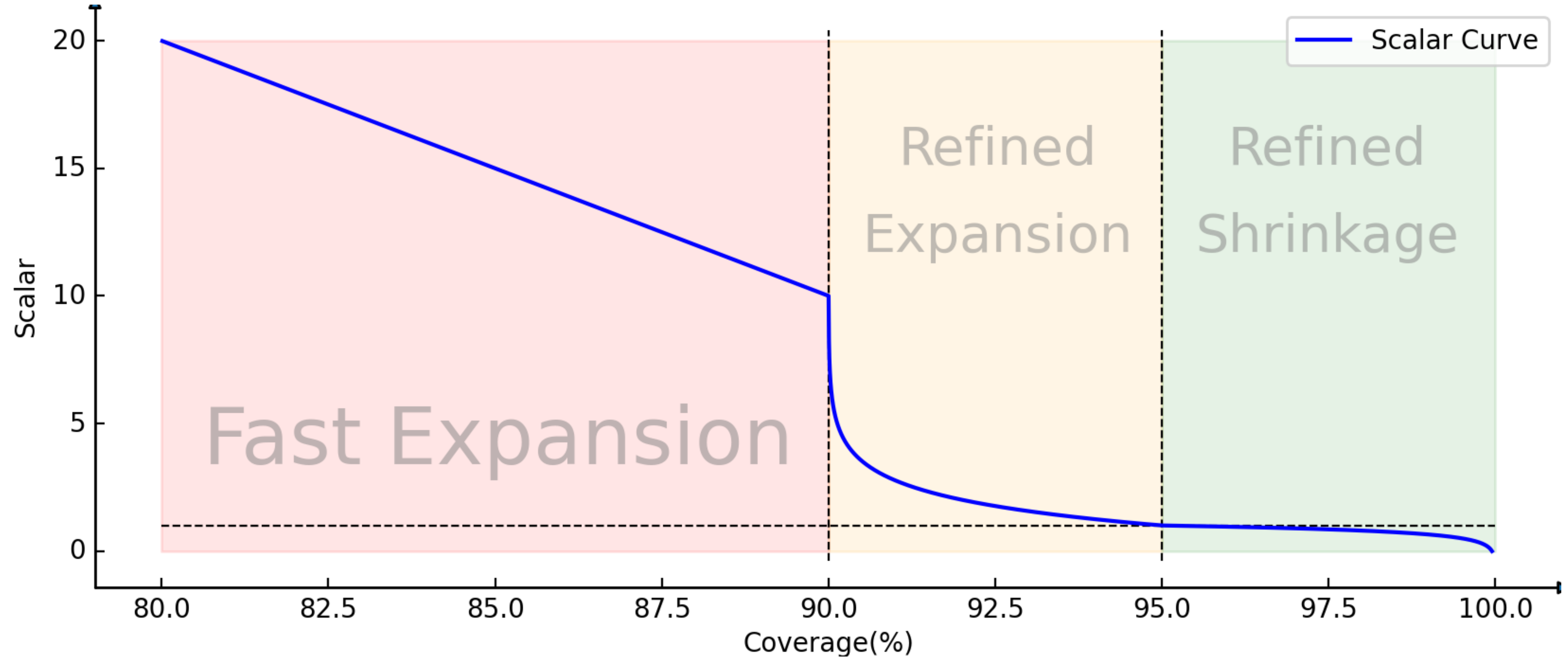
$$\Pr\left(y \in \left[\hat{y} - \mathcal{S} \times G^{-1}(0, \gamma) \times \sigma_\epsilon, \hat{y} + \mathcal{S} \times G^{-1}(0, \gamma) \times \sigma_\epsilon\right]\right) \approx \gamma$$

$$\text{PI}_{\text{AdaPI}} \in \left(\hat{y} - \mathcal{S} \times G^{-1}(0, \gamma) \times \sigma_\epsilon, \hat{y} + \mathcal{S} \times G^{-1}(0, \gamma) \times \sigma_\epsilon\right)$$

$\gamma$  : Confidence Level / Significance Level

$\mathcal{S}$  : Scalar

# Scalar for AdaPI



# Evaluation Metrics for PI

$$\text{Coverage} = \frac{1}{N} \sum_{i=1}^N I_i$$

$$\text{NMPIW} = \frac{\frac{1}{N} \sum_{i=1}^N (P_{u_i} - P_{l_i})}{R}$$

NMPIW : Normalized Mean Prediction Interval Width

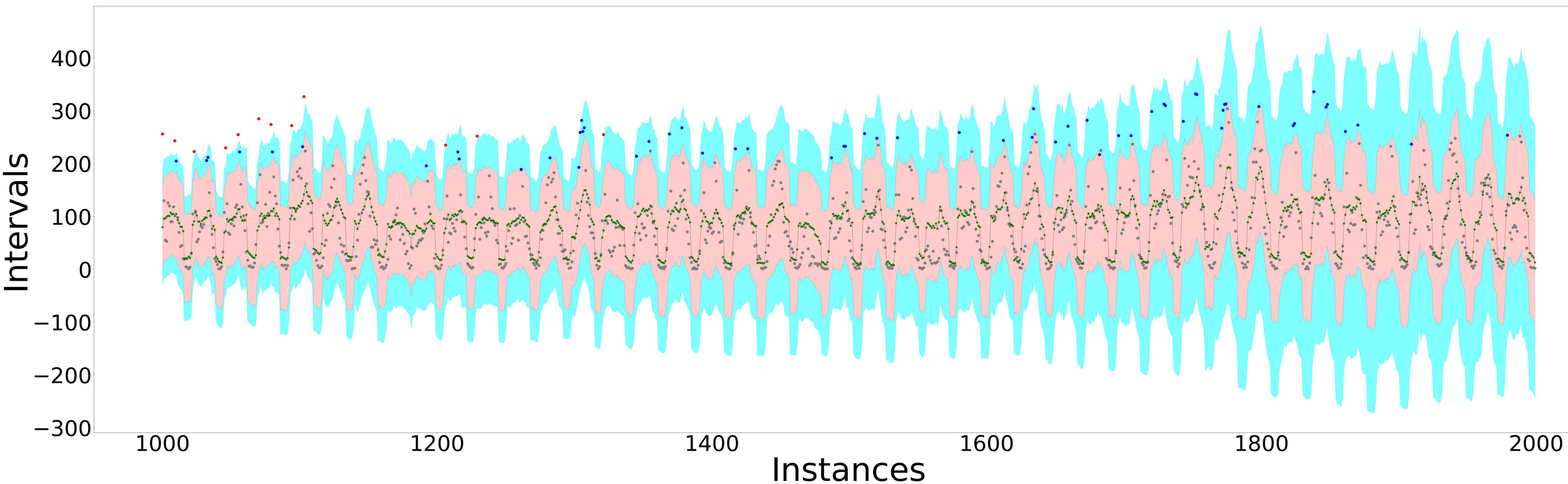
R : Range of Target Values

$P_u, P_l$  : Upper and Lower Bounds of Prediction Intervals



# Expansion Case

Datasets: Bike



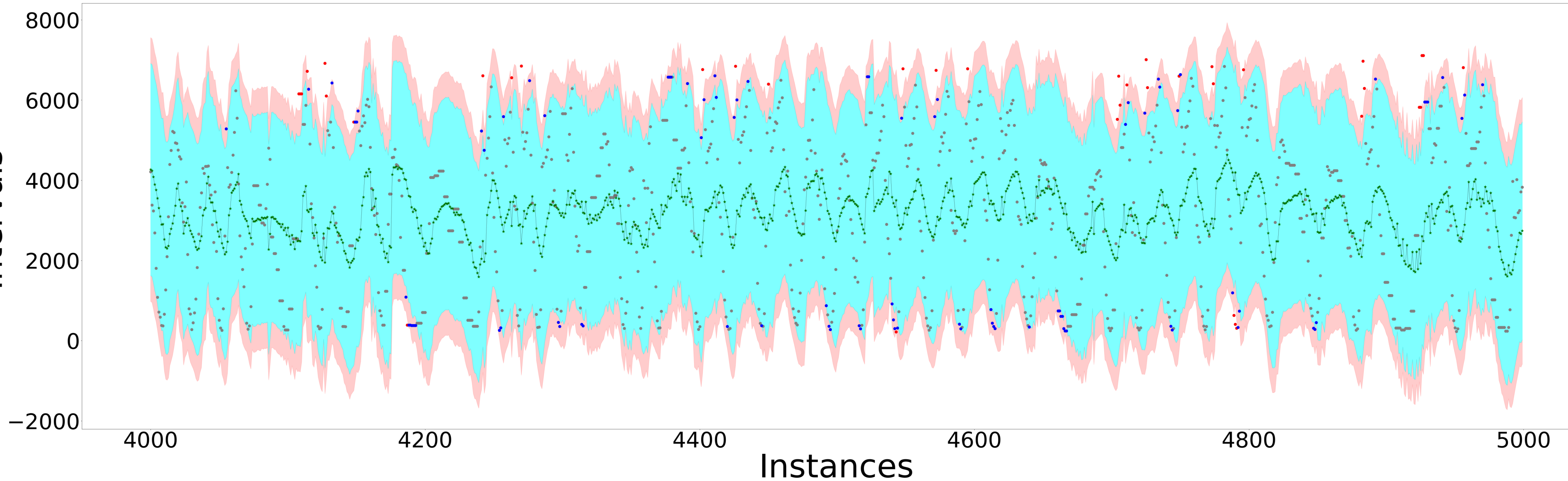
Blue area: **AdaPI**

Red area: **MVE**



# Shrinkage Case

Datasets: MetroTraffic

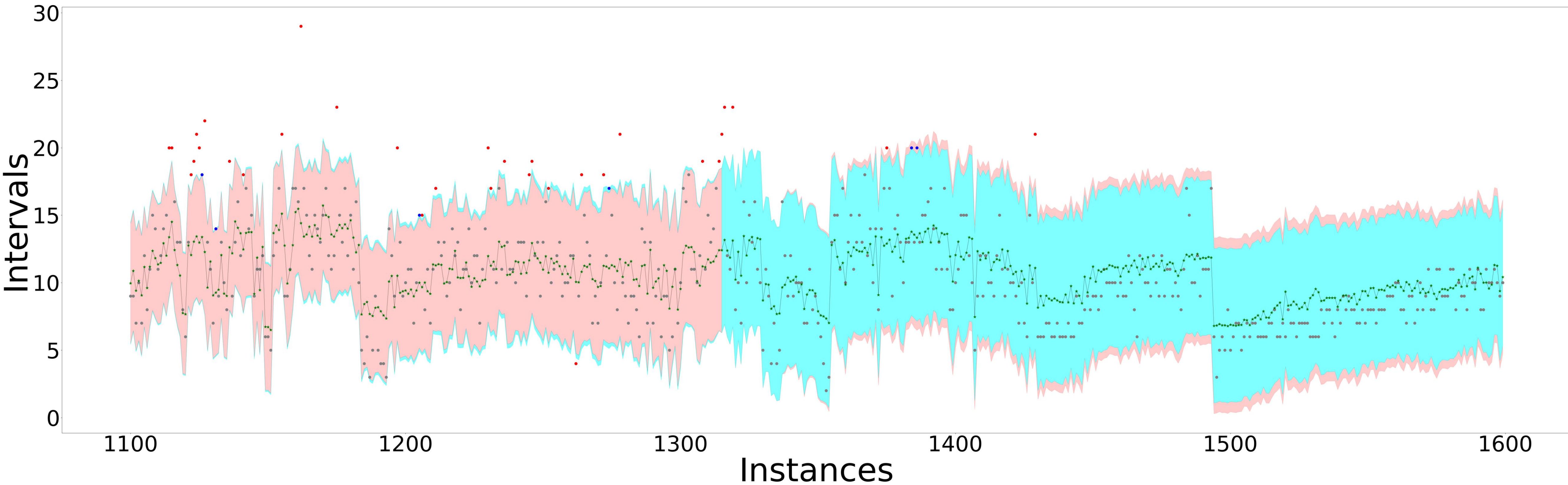


Blue area: **AdaPI**

Red area: **MVE**

# Switching Case

Datasets: Abalone



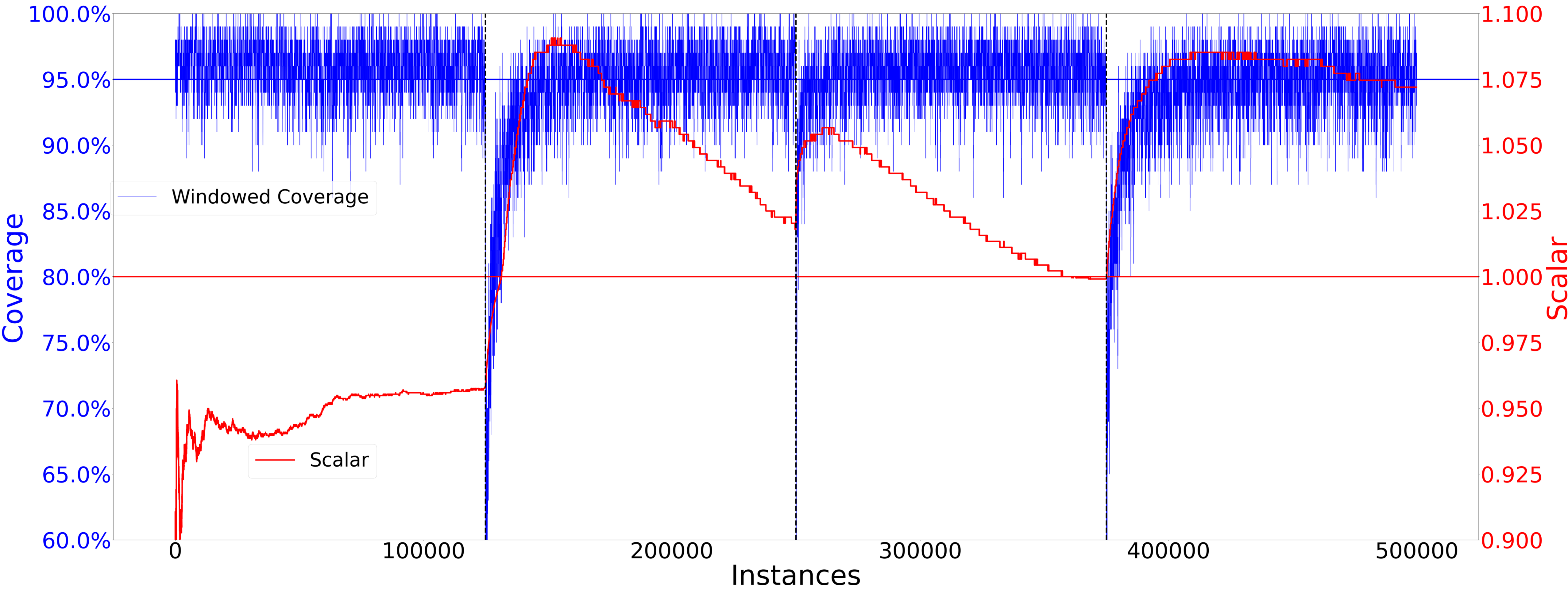
Blue area: **AdaPI**

Red area: **MVE**



# Dealing with Drifts

Datasets: HyperA





# Prediction Intervals Summary

Prediction Interval (PI) is essential for uncertainty quantification in regression tasks.

## **Challenges**

Traditional PI methods are not suitable for dynamic data streams.

## **Solution**

Mean and Variance Estimation (MVE); ADAPI\*

## **Evaluation**

Coverage

Normalised Mean Prediction Interval Width (NMPIW)

# Practical examples

02\_ECML2024\_supervised.ipynb