



Report di Progetto

ISW2 A.A. 2021/2022

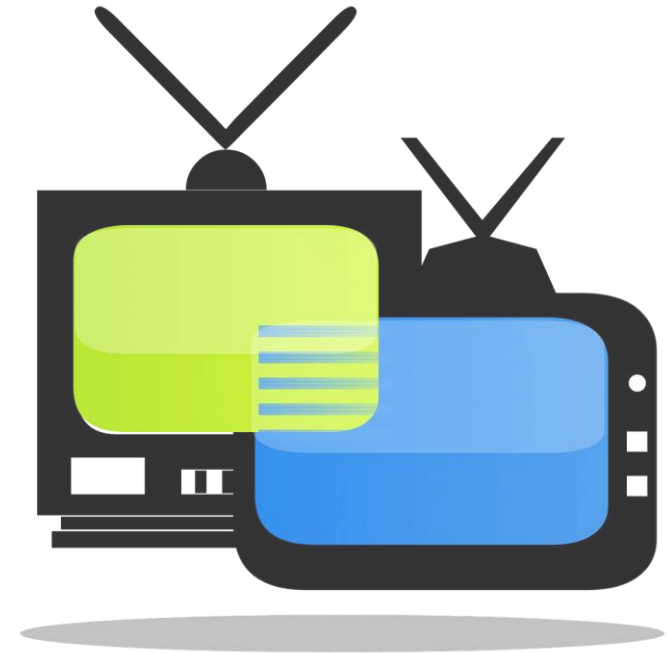
Machine Learning for Software Engineering
Deliverable

GIANMARCO BENCIVENNI - 0285847

Contenuto della presentazione

In questa presentazione :

- ▶ [Gli obiettivi dello studio](#)
- ▶ [Architettura dell'applicazione](#)
- ▶ [L'importanza dei dati e delle loro caratteristiche](#)
- ▶ [Data retrieval – merging Jira and Git](#)
- ▶ [Scelta e calcolo delle metriche](#)
- ▶ [Labeling](#)
- ▶ [Proportion increment](#)
- ▶ [Addestramento : walk-forward](#)
- ▶ [Tecniche di data processing : sampling e feature selection](#)
- ▶ [Evaluation : risultati](#)
- ▶ [Links e Informazioni](#)



Gli Obiettivi dello studio

► What

Addestrare un insieme di classificatori affinché siano in grado di **identificare**, in un certo istante del ciclo di vita di un certo progetto software, quali siano gli **artefatti software** che **più probabilmente sono difettosi** (*buggy/defective*).

Confronto e valutazione delle prestazioni di tali predittori, prendendo in considerazione i valori assunti dalle metriche di valutazione presentate durante il corso.

Lo studio prevede l'applicazione pratica dei concetti di *Machine Learning for Software Engineering*, visti a lezione del Prof. Falessi, sui due progetti open-source:

- **Apache Bookkeeper**
- **Apache Storm**



Gli Obiettivi dello studio



► Why

Un predittore correttamente addestrato può rivelarsi estremamente utile ai fini della **riduzione dei costi** (in termini di **tempo** e di **denaro!**) **nelle attività di software testing e di revisione del codice**, permettendo ai teams di test e sviluppo di disporre di un **approccio mirato**, onde **evitare testing e ricerca estensivi**.



► How

Il conseguimento degli obiettivi prefissati si snoda in alcuni **passaggi** base, che andremo ad affrontare in maggiore dettaglio nel corso della presentazione.

- **Creazione del Dataset** (*data retrieval, software metrics, proportion, labeling*)
- **Pre-processing** (*sampling, feature selection, cost sensitive analysis*)
- **Training** (*walk forward*).
- **Evaluation** (*precision, recall, auc, kappa, accuracy*)



Architettura dell'Applicazione 1/2

- ▶ L'applicazione sviluppata per questo studio è basata sul design pattern architetturale **BCE**.
- ▶ Principali classi di **BOUNDARY**:
 - ▶ **GitRepositoryManager**
Interfacciamento ai log delle git repositories dei 2 progetti Apache.
Utilizzo della libreria jGit.
 - ▶ **JiraTicketManager**
Interfacciamento al server di **Jira**, tramite **API Rest**, come visto in laboratorio.
Reperimento di tutte le informazioni relative agli issues (di tipo **bug**) che sono stati **chiusi**.
 - ▶ **DatasetBuilder**
Inizializzazione e popolamento dei file csv contenenti i dataset.
A partire da una `MultiKeyMap (key: (filename, version), value: metrics)`.

Architettura dell'Applicazione 2/2

► Principali classi **CONTROLLER**:

- **IssueLifeCycleManager**
Controllore per la logica di business legata alla **costruzione del dataset**.
 - Retrieval dei dati grezzi
 - Elaborazione delle software metrics di input
 - Labeling e scrittura delle entries sul file csv.
- **ClassifierModel**
Controllore per la logica di business legata all'**addestramento dei predittori** (Ibk, Random Forest, Naive Bayes, Cost Sensitive) e alla loro **valutazione**.

► Principali classi **ENTITY**:

- **IssueObject**
Rappresenta il singolo **issue**, che è definito principalmente attraverso gli attributi ID, IV, OV, AVs, FV, **commits[]**
- **CommitObject**
Rappresenta la singola **commit**, che è definita attraverso gli attributi ID, referenceIssue, **files[]**
- **FileObject**
Rappresenta il singolo **file**, che è definito attraverso gli attributi **filepath**, **version**, (**metrics**), **buggyness**. Il file object è l'entità che rappresenta una entry del dataset, una volta popolato in tutti i suoi attributi.

L'importanza dei dati e delle loro caratteristiche 1/2

► In questo studio, si è scelto un **approccio within-project** per la raccolta dei dati.

► L'**affidabilità di un predittore** dipende in larga parte dalla **qualità del dataset!**



I. **«Size matters»**

Tanto più è ampio (e vario) il dataset, tanto più accurato sarà il predittore.
I due progetti Apache hanno LOC alta, N° release alto, N° classi alto
... possibilità di **generare dataset di grandi dimensioni!**

II. **Maturazione dei dati**

Il N° release è sufficientemente elevato da poter assumere un **buon grado di maturità** per entrambi i progetti
Requisito fondamentale affinché un dataset contenga il **meno rumore (noise/garbage)** possibile.

L'importanza dei dati e delle loro caratteristiche 2/2

III. «Garbage-in... Garbage-out!»

Un dataset troppo giovane (immaturo) è un dataset che è **molto probabilmente ricco di garbage/rumore**. Addestrare un predittore con dati rumorosi produrrebbe ovviamente dei **risultati** altrettanto **inaccurati**.



IV. Fenomeno dello Snoring

Una classe contenente *sleeping defects* è **snoring**, poiché produce **rumore nel dataset!**
Per **ridurre il fenomeno dello snoring**:



1. **Generare le etichette** (label) **considerando tutte le revisioni disponibili**
2. ...**Usando per l'addestramento** solamente la porzione di dataset relativa alla **prima metà delle release**, poiché i **dati** sono **più maturi** e meno affetti da snoring.

Data retrieval – Merging Jira and Git 1/2

Due componenti fondamentali per la fase di costruzione del dataset sono:

- ▶ **Jira** è l'**Issue Tracking System** utilizzato per entrambi i progetti Apache.
Jira permette di **risalire a tutti gli issues (di tipo bug) chiusi**, indicando:
 - ▶ La Opening e la Fixed Versions (**OV, FV**) per i tickets associati ai Bugs.
 - ▶ Un **ID univoco associato all'issue/ticket**.
 - ▶ Le Affected Versions (**AV**), ovvero le versioni del progetto software affette dal Bug associato al ticket.
- ▶ Retrieval degli issues tramite **API Rest**.
- ▶ Il **codice sviluppato** è derivato dagli esempi forniti a laboratorio dal Prof. Falessi.
Si discriminano:
 - ▶ **Issues aventi Affected Versions dichiarate in Jira**
 - ▶ **Issues non aventi Affected Versions dichiarate in Jira**

Data retrieval – Merging Jira and Git 2/2

- ▶ **Git** è il **Version Control System** utilizzato per entrambi i progetti Apache.

Algoritmo per il «log walk» e il calcolo delle software metrics

- ▶ Per ogni issue, **ricercare le commit che contengono l'issue ID** all'interno del **commit-message**, popolando così l'array di commits che è mantenuto come attributo per lo stato dell'IssueObject in questione.
- ▶ Problema : **incoerenza tra le versioni di Jira e git**
 - ▶ Assumo git come veritiero, essendo il versioning automatizzato
 - ▶ Scelgo come FV la versione relativa all'ultima commit che contiene l'issue-ID nel commit-message
- ▶ per ogni commit, acquisire i files toccati, filtrando in modo da ottenere esclusivamente files in **formato «.java»**.
- ▶ Per ogni file, istanziare un'entità **FileObject** definita da (filepath,version) e **calcolarne le metriche** (che specificheremo a breve)



Scelta e calcolo delle Metriche 1/2

Delle **metriche** viste a lezione, sono state scelte e calcolate le seguenti:

- ▶ **Number of Revisions** : numero di revisioni associate al file. Si è scelto in questo caso di registrarne il valore cumulativo, quindi non azzerando la metrica da una release alla successiva.
- ▶ **Age** : età del file, che in questo caso è stata misurata in settimane, a partire dalla creazione.
- ▶ **Churn** : somma, calcolata su tutte le revisioni, di LOCadded – LOCdeleted.
- ▶ **Number of Authors** : numero di autori che hanno lavorato sul file. Si è scelto in questo caso di registrarne il valore cumulativo, quindi non azzerando la metrica da una release alla successiva.
- ▶ **LOC** : valore medio, calcolato su tutte le revisioni di ogni release, delle linee di codice presenti nel file.
- ▶ **LOC touched** : somma, calcolata su tutte le revisioni, di LOCadded + LOCdeleted + LOCmodified.
- ▶ **Average LOC added** : valore medio delle linee di codice aggiunte al file, per release.
- ▶ **Max LOC added** : valore massimo delle linee di codice aggiunte al file, in per release.
- ▶ **Average Change Set Size** : valore medio di change set size nelle revisioni.
- ▶ **Max Change Set Size** : valore massimo del change set size nelle revisioni.

Scelta e calcolo delle Metriche 2/2

In aggiunta alle **metriche** viste in classe, ho deciso di **sperimentare un paio** metriche **aggiuntive** che ho pensato potessero essere delle caratteristiche **correlate alla buggyness di una classe**.

- ▶ **Number of Imports** : il numero di import nell'implementazione di una classe è legato al numero di moduli esterni che verranno utilizzati nel codice. Intuitivamente, quanto maggiori saranno i moduli da importare, tanto maggiore la complessità della classe, nonché la probabilità di iniettare difetti.
- ▶ **Number of Comments** : questa metrica ha secondo me un peso ancora maggiore rispetto alla precedente, e può portare un contributo importante alla classificazione soprattutto quando è accoppiata ad entries con un alto numero di LOC. L'assenza di commenti in una classe, infatti, la rende di difficile manutenibilità, e difficilmente estensibile senza introdurre difetti! Questo è vero soprattutto nelle classi con un alto livello di complessità (i.e. alto valore di LOC).

L'implementazione del calcolo delle metriche è consultabile nel codice sorgente della classe `GitRepositoryManager`, che si interfaccia direttamente ai logs delle commit dei due progetti. Si mostrano a seguire alcuni esempi di codice in snippets, relativi proprio al calcolo delle metriche.

```

public List<FileObject> computeMetricsAndAppendFile(CommitObject commitObject, List<FileObject>files, DiffEntry diff,

    int linesAdded = 0;
    int linesDeleted = 0;
    int linesReplaced = 0;
    RevCommit commit = commitObject.getCommit();

    if ( diff.getNewPath().endsWith( FILE_EXTENSION ) ){

        String filepath = diff.getNewPath();
        String fileText = getTextfromCommittedFile( commit, filepath );
        int fileAge = getFileAgeInWeeks( commit, filepath );

        for ( Edit edit : df.toFileHeader( diff ).toEditList() ) {

            if ( edit.getBeginA() < edit.getEndA() && edit.getBeginB() < edit.getEndB() ){
                linesReplaced += edit.getEndB() - edit.getBeginB();
            }
            if ( edit.getBeginA() < edit.getEndA() && edit.getBeginB() == edit.getEndB() ){
                linesDeleted += edit.getEndA() - edit.getBeginA();
            }
            if ( edit.getBeginA() == edit.getEndA() && edit.getBeginB() < edit.getEndB() ){
                linesAdded += edit.getEndB() - edit.getBeginB();
            }
        }

        int version = commitObject.getVersion();
        int loc = getLoc( fileText );
        int numImports = getNumImports(fileText);
        int numComments = getNumComments(fileText);
        files.add( new FileObject( filepath, version, fileAge, loc, linesAdded, linesDeleted, linesReplaced,
            changeSetSize, commitObject.getAuthorName(), numImports, numComments ) );
    }
}

```

Code View

13

Software Metrics

Le metriche per un dato file vengono calcolate sul posto, nel momento in cui si ricercano i cambiamenti apportati dalla commit che si sta analizzando.

Tramite un'istanza di **DiffEntry** è possibile risalire ai **cambiamenti** effettuati su di un file **rispetto** alla **parent commit**.

Nello snippet è possibile osservare parte della computazione delle varie metriche.

Ho voluto inserire il dettaglio relativo alla metrica «numero di commenti», nello **snippet sottostante**, poiché la reputo un'idea interessante.

```

public int getNumComments( String fileText ){
    int numComments = 0;
    try ( BufferedReader reader = new BufferedReader(new StringReader(fileText)) ) {
        for ( String line = reader.readLine(); line != null; line = reader.readLine() ) {
            if ( line.contains("//") ||
                line.contains("/*") ||
                line.contains("*/") ||
                !( line.endsWith(";") || line.endsWith("{") || line.endsWith("}") || line.equals("\n") || line.endsWith("\n")) ) {
                numComments ++;
            }
        }
    } catch ( IOException e ) {
        e.printStackTrace();
    }
    return numComments;
}

```

Labeling : IV, OV, AV, FV

È necessario etichettare i FileObjects (che saranno le **entry del dataset**) come **Buggy** o **Not Buggy**.

- ▶ **Ricapitolando**, ogni **File** è rappresentato da un **nome**, una **versione**, delle **metriche**, e un valore relativo alla **buggyness**
- ▶ Come visto più volte a lezione, **una classe, in una certa versione, deve essere classificata come buggy se la sua versione appartiene alle Affected Versions del ticket** di riferimento.
Le Affected Versions partono con la Injected Version e terminano con la versione precedente alla Fixed Version, che è la versione nel corso della quale il bug è stato fixato.
- ▶ **Problema** : c'è ancora una parte dei tickets per cui conosciamo solamente **OV** e **FV**!
- ▶ È necessario applicare un metodo per stimare le **Injected Version** degli issues per cui Jira non ne riporta il valore.
Come noto, esistono differenti metodi, ma quello che è stato utilizzato è il metodo **Proportion Incremental**.

Proportion - Increment

- ▶ La tecnica **Proportion** consente in generale di risalire alle Affected Versions (stimando il valore della Injected Version) laddove in Jira non siano disponibili.
- ▶ L'**assunzione** (dimostrata in diversi articoli scientifici che il Prof. Falessi ci ha proposto durante il corso) è che **i difetti** abbiano un **ciclo di vita stabile in termini della proporzione del numero di AV prima del fixing.**
 1. $P = (FV - IV) / (FV - OV)$
 2. $Predicted\ IV = FV - (FV - OV) \cdot P$
- ▶ La declinazione **Increment** di Proportion prevede, per ogni issue privo di IV/AV, di considerare il sottoinsieme degli issues (aventi IV/AV) che sono passati a «fixed» prima di esso, e calcolare il valore di **P** come la media delle proporzioni di questi ultimi.
A seguire, si stima il valore della IV dell'issue tramite la formula (2).
- ▶ Nella slide a seguire, si propone **l'implementazione di Proportion Increment** sviluppata per questo studio.

Code View

Gli snippets sono tratti dal codice sorgente della classe `IssueLifeCycleManager.java`

Nel corpo del metodo

`setAffectedAndInjectedVersionP` (snippet sottostante) viene effettuata la stima di IV e conseguentemente AV per ogni ticket che ne è sprovvisto, tramite la tecnica **Proportion Increment**.

16

```
public int computeProportionIncremental( List<IssueObject> filteredIssues ){
    ArrayList<Double> proportions = new ArrayList<>();
    for ( IssueObject issue : filteredIssues ){
        if ( issue.getOv() != issue.getFv() ) {
            double fv = issue.getFv();
            double ov = issue.getOv();
            double iv = issue.getIv();
            double p = ( fv - iv ) / ( fv - ov );
            proportions.add( p );
        }
    }
    return (int) average( proportions );
}
```

```
public void setAffectedAndInjectedVersionsP(){
    for ( IssueObject targetIssue : issuesWithoutAffectedVersions ){
        int fv = targetIssue.getFv();
        int ov = targetIssue.getOv();
        List<IssueObject> filteredIssues = issuesWithAffectedVersions.stream()
            .filter(issue -> issue.getFv() <= fv)
            .collect(Collectors.toList());
        int p = computeProportionIncremental(new ArrayList<>(filteredIssues)); ★
        if ( fv == ov ) {
            // The following formula is an approximation of the correct one (the one in the else block).
            ★ targetIssue.setIv( ( fv - ( 1 * p ) ) );
            continue;
        } else{
            targetIssue.setIv( ( fv - ( ( fv - ov ) * p ) ) );
        }
        int minAVValue = targetIssue.getIv();
        int maxAVValue = ( targetIssue.getFv() - 1 );
        ArrayList<Integer> avs = new ArrayList<>( IntStream.rangeClosed(minAVValue, maxAVValue).boxed().collect(Collectors.toList()) );
        targetIssue.setAvs( avs );
    }
}
```

Ad ogni modo, come è possibile notare nella linea evidenziata in rosso, ho dovuto effettuare un'**approssimazione** nella stima di IV. Infatti, nei casi in cui il ticket per cui si vuole effettuare la previsione ha **FV=OV**, la formula per stimare il valore di IV «collapsa» in: $IV = FV$.

Pertanto, in questi casi (comunque una minoranza), ho deciso di sostituire, con un certo margine di errore, il termine $(FV-OV)$ con 1.

Addestramento – La Tecnica WalkForward



- ▶ Un **ingrediente fondamentale** di cui non abbiamo parlato è... il **tempo**!
- ▶ Il **dataset** costruito **ha un ordinamento per versione**.
Si tratta infatti di una **time-series**, in cui l'ordinamento dei data points è una componente necessaria ai fini del corretto addestramento del predittore.
- ▶ **Non posso utilizzare dati del futuro per predire il passato!**
- ▶ La **tecnica di addestramento walk-forward** è designata a risolvere il problema dell'addestramento su serie temporali. La logica di addestramento è la seguente :
 - ▶ Suddivisione del dataset in K batch ordinati (in questo caso, si suddivide per versioni)
 - ▶ Addestramento iterativo (in K-1 passaggi) prendendo alla j-esima iterazione i primi j batch come training set, ed il j+1-esimo batch come test set.
 - ▶ La valutazione finale del predittore si ottiene mediando sui risultati delle metriche di ogni iterazione.

Tecniche di data processing

Sampling e Feature Selection 1/2



- ▶ Il modulo dell'applicazione che implementa la logica relativa a **pre-processing dei dati, addestramento e valutazione dei predittori**, fa uso dell'API di **Weka**, una **libreria di Machine Learning** per applicazioni Java.
- ▶ **Sampling**
Spesso **le classi nel dataset sono sbilanciate in modo pronunciato**.
Il modello non avrà sufficienti esempi della classe minoritaria, e una volta addestrato non possiederà buone capacità di generalizzazione.
Per questo studio ho trovato necessario applicare il sampling, proprio perché il numero di **data points classificati come *buggy*** sono comunque in **forte minoranza** rispetto ai *not buggy*.
 - ▶ **Under-sampling** (API SpreadSubSample): diminuzione dei punti della classe maggioritaria
 - ▶ **Over-sampling** (API Resampling): aumento (per duplicazione) dei punti della classe minoritaria
 - ▶ **Smote** (API SMOTE): aumento dei punti della classe minoritaria attraverso la creazione di data points sintetici, a partire da quelli esistenti.

Tecniche di data processing

Sampling e Feature Selection 2/2



- ▶ **Feature Selection** : la tecnica della feature selection ha lo scopo di diminuire la complessità del modello (tagliando il numero di features di input si ha ovviamente un modello più semplice) e, teoricamente, di migliorare le performance del training rispetto al caso con l'intero set di features.
L'approccio specifico utilizzato per lo studio è denominato **best-first**, e si basa su un **algoritmo greedy**.
 - ▶ **Best-First** : Si inizia creando una suite di N modelli, ognuno dei quali utilizza solo una delle N caratteristiche del nostro set di dati come input. Viene selezionata la caratteristica che produce il modello con le migliori prestazioni. Nell'iterazione successiva, viene creata un'altra suite di N-1 modelli con due caratteristiche in ingresso: quella selezionata nell'iterazione precedente e un'altra delle N-1 caratteristiche rimanenti. E così via.
- ▶ **Cost Sensitive Analysis** : la libreria weka permette anche l'utilizzo di un **Cost Sensitive Classifier (J48)**.
Per questo studio ha perfettamente senso valutare l'utilizzo di un **classificatore sensibile ai costi**!
 - ▶ Setting di una **matrice di costo** in modo da apportare una specifica **penalità** all'occorrenza di falsi negativi (FN), ogni **errore di tipo FN verrà pesato tanto quanto indicato nella matrice**!

Valutazione dei modelli

Ibk, RF, NB, J48

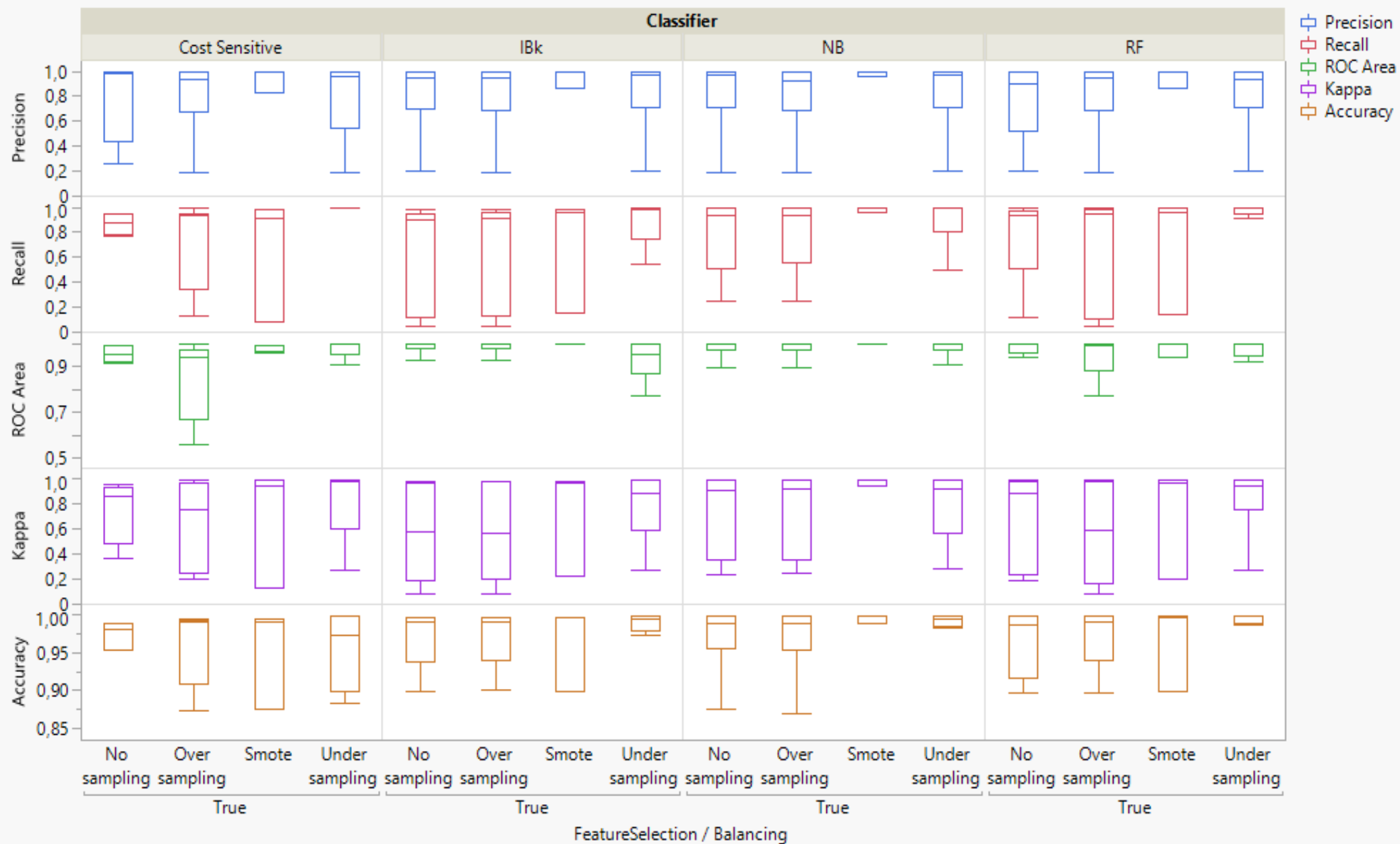
- ▶ I modelli di Machine Learning (classificatori) addestrati per lo studio sono:
 - ▶ **IBK, Random Forest, Naive Bayes, Cost Sensitive Classifier (J48).**
- ▶ La **valutazione dei predittori** tiene conto di alcune **metriche** che abbiamo affrontato a lezione, quali:
 - ▶ **Precision** : $\frac{TP}{TP+FP}$ (quanti dei punti che ho classificato come positivi lo erano davvero?)
 - ▶ **Recall** : $\frac{TP}{TP+FN}$ (quanti positivi sono riuscito a «beccare»?)
 - ▶ **Precision e Recall devono essere osservate insieme!** Se ho alta Precision e bassa Recall, significa che i punti classificati come positivi erano giusti... ma vuol dire che ho tanti falsi negativi! Viceversa, se la Recall è alta e la Precision bassa, significa che magari i positivi li ho beccati tutti... ma ho tanti falsi positivi!
 - ▶ **AUC** : è l'area sotto la curva ROC. Indica la probabilità che una istanza di positivo presa randomicamente sia classificata meglio di un'istanza di negativo presa randomicamente. In generale la vogliamo alta (>0.5), dato è molto più grave un falso negativo rispetto ad un falso positivo!
 - ▶ **Kappa** : indica quanto il classificatore si discosta da un classificatore Dummy (perfettamente randomico)
 - ▶ **Accuracy** : 1- error rate, indica la percentuale di punti classificati correttamente, indipendentemente dal fatto che questi siano classificati come positivi o negativi.



APACHE STORM

BEST FIRST

21

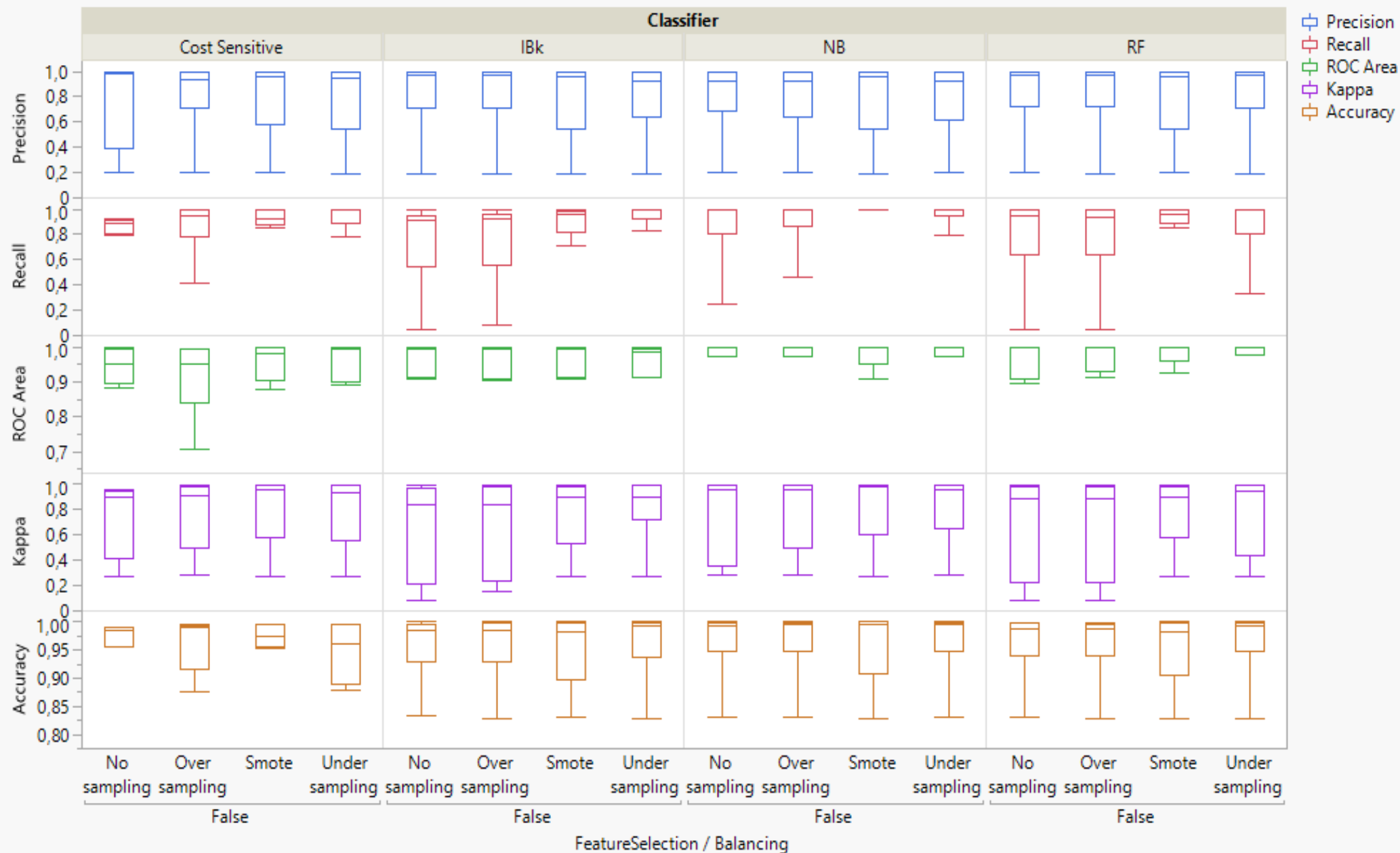




APACHE STORM 2

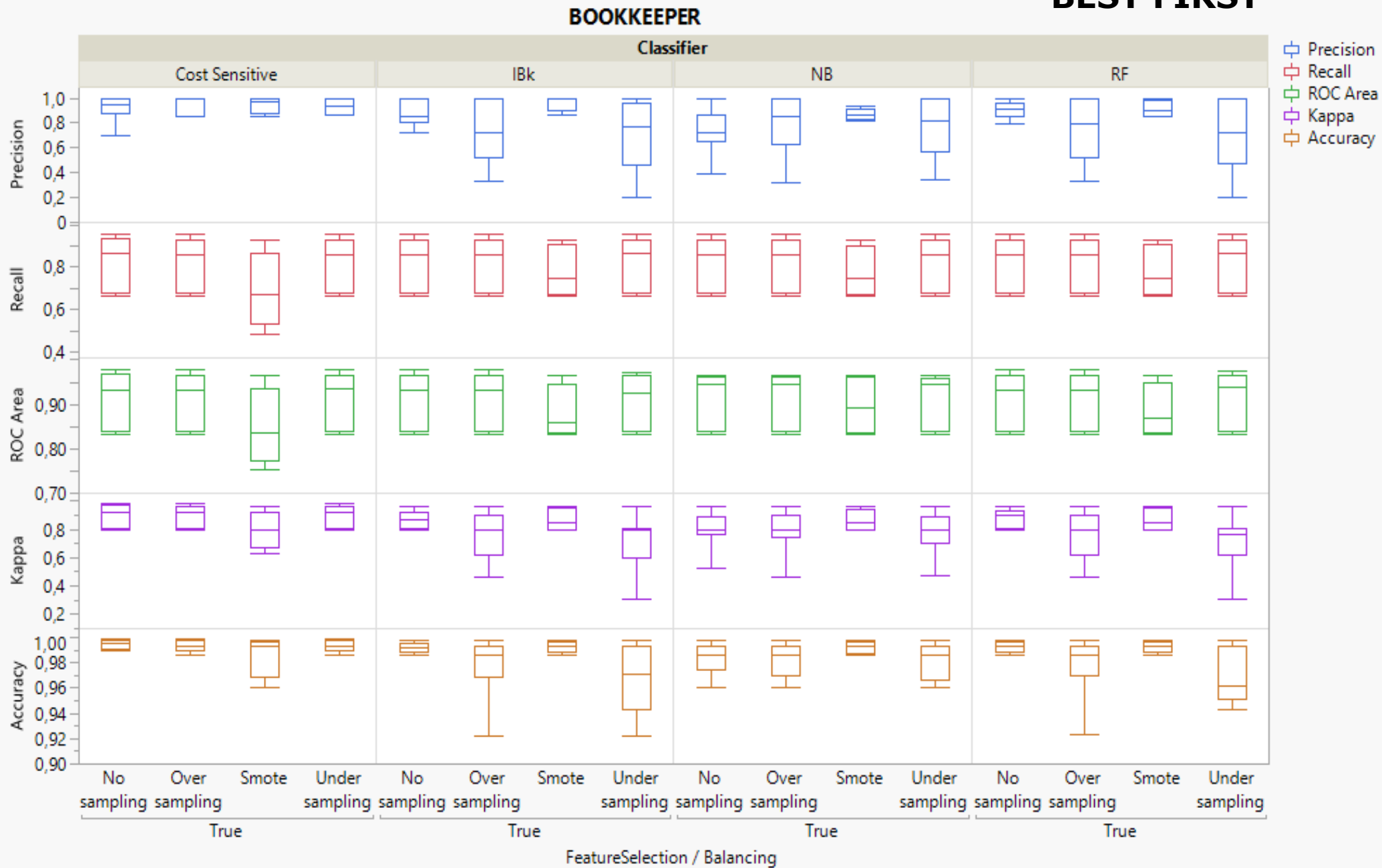
NO FEATURE SELECTION

22





- Analizziamo i risultati ottenuti per il **progetto Storm**, facendo **focus sui punti chiave**.
- i. **Problemi con Smote** : Molto spesso la **classe minoritaria** è risultata **insufficiente** per applicare smote correttamente, e i valori delle metriche sono stati settati a **NaN**, e perciò li ho dovuti **scartare nella valutazione finale**.
- ii. **Feature Selection** : La tecnica di feature selection, salvo pochi casi, si è rivelata **controproducente**, aumentando in particolare la **variabilità dei risultati sulla Recall**.
Interpretazione : le features nativamente utilizzate apportavano migliori risultati in fase di addestramento se utilizzate assieme. Si può dedurre che le features avessero una **bassa inter-correlazione** (tra loro), ma una **buona correlazione con l'output** del predittore.
- iii. **No Silver Bullet** : dai risultati ottenuti si evince chiaramente che **non esiste il proiettile d'argento**.
Alcuni classificatori, in determinate configurazioni, si comportano meglio di altri per alcune caratteristiche, e peggio per altre.
- iv. **Precision e Recall** : a livello di Precision e Recall, la configurazione migliore è stata ottenuta dal **Cost Sensitive Classifier in assenza di feature selection**, e in particolare con **over-sampling**. Mi aspettavo un risultato di questo tipo, sebbene sia importante sottolineare che sto valutando la bontà del risultato **non solamente sui valori medi** assunti, bensì **anche** sul fatto che è importante avere una **bassa varianza sui risultati**, ai fini della validazione del modello.
È interessante notare anche che la **classificazione orientata ai costi produce alti valori di Recall**, ma inevitabilmente **impatta negativamente sulla Precision**, perché il modello è portato a classificare più positivi che negativi.
Volendo guardare alla Precision da sola, infatti, i risultati migliori sono stati ottenuti da Random Forest/NoFS.
- v. **AUC** : i risultati migliori a livello di AUC sono quelli ottenuti per **Naive Bayes/NoFS** e **Random Forest/NoFS**, in entrambi i casi in modo analogo con l'applicazione dell'under-sampling o dell'over-sampling.
- vi. **Kappa** : è interessante notare che, per la metrica di valutazione kappa, i risultati migliori sono stati ottenuti dal classificatore **Naive Bayes/NoFS**, pur essendo questo tra i peggiori nel contesto di Precision e di Recall.

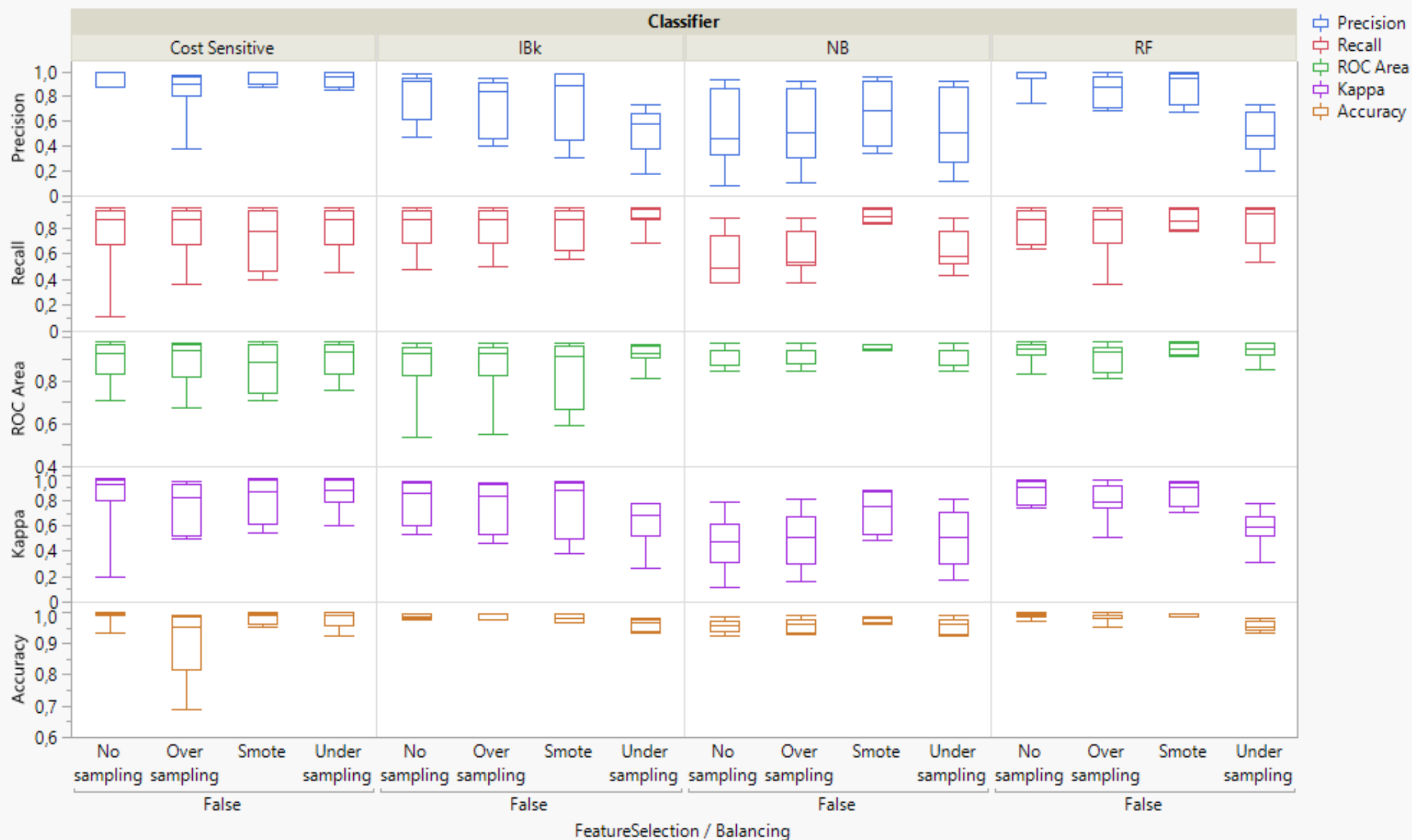




NO FEATURE SELECTION

25

BOOKKEEPER 2





- Analizziamo i risultati ottenuti per il **progetto Bookkeeper**, facendo **focus sui punti chiave**.
- i. **Problemi con Smote** : in questo caso i problemi con la tecnica smote sono stati minori, poiché solamente poche volte la classe minoritaria è risultata insufficiente per applicare la tecnica. Ad ogni modo, è stato comunque necessario eliminare qualche riga dal file di *evaluation* a causa dei valori NaN, come nel caso di Storm, e la variabilità sui valori con smote è comunque bassa. Ad ogni modo, **nel dataset di Bookkeeper la classe buggy era molto più rappresentata** rispetto a quanto non lo fosse in Storm.
- i. **Feature Selection** : la tecnica di feature selection ha in generale avuto un impatto **di poca rilevanza**, per quanto riguarda le metriche Precision, Recall e AUC. Un netto miglioramento dato dall'utilizzo della feature selection è tuttavia osservabile sui valori medi e sulla variabilità della metrica Kappa, e questo vale per tutti i predittori.
- i. **Precision e Recall** : anche in questo caso, a livello di Precision e Recall, la configurazione migliore è stata ottenuta dal **Cost Sensitive Classifier** a cui **non è stata applicata feature selection e nemmeno il sampling**. Si può spiegare il risultato col fatto che, data la presenza di una migliore rappresentanza della classe buggy, in questo caso il sampling non è stato particolarmente necessario.
- i. **AUC e Kappa** : i valori migliori per AUC e Kappa sono stati ottenuti con il predittore **Naive Bayes**, con feature selection **best first** e **no sampling**.

Info e Links

Ambiente di sviluppo e frameworks utilizzati :

- ▶ **VSCode** come IDE per lo sviluppo dell'applicazione
- ▶ **Maven** come framework per la gestione delle dipendenze e per la build del progetto
- ▶ **Java 11** come linguaggio di programmazione (principale)
- ▶ **Excel** e **JMP (GUI)** per lavorare dati e grafici



Maven™



Istruzioni per l'utilizzo del codice :

- ▶ Scaricare il codice sorgente dal link alla github repository https://github.com/CLOCKWORK95/ML_4_SE
- ▶ Posizionarsi da terminale nella directory : `ML_4_SE/my-app`
- ▶ lanciare il comando bash : `sh buildAndLaunch.sh` e seguire le indicazioni da CLI

Link a **Sonarcloud** per la repository del progetto **ML_4_SE** :

- ▶ https://sonarcloud.io/summary/overall?id=CLOCKWORK95_ML_4_SE

sonarcloud 