

UNIVERSITÀ DEGLI STUDI DI PADOVA

INGEGNERIA DEL SOFTWARE 1 - 20/21

Tribuo

Analisi di un progetto open source

Professori:

Mauro Migliardi

Alberto Sillitti

Gruppo:

Alberto Castagnaro, 1219641

Antonio Franzoso, 1219610

Christian Marchiori, 1218317

Francesco Visentin, 1216419

Indice

1	Introduzione	2
2	Analisi produzione	3
2.1	Team e organizzazione	3
2.2	Sviluppo software	4
2.2.1	Requirements elicitation	4
2.2.2	Code development	5
2.2.3	Software validation & delivery	5
3	Analisi metriche	6
3.1	Lines Of Code (LOC)	7
3.2	Chidamber and Kemerer's metrics suite (CK)	8
3.3	McCabe Cyclomatic Complexity (CC)	13
4	Conclusione	14

1 Introduzione

Il progetto che abbiamo scelto di analizzare è Tribuo.

Tribuo è una libreria open source scritta in Java da un team di ricercatori facenti parte dell'Oracle Labs Machine Learning Research Group. È stato utilizzato internamente ai laboratori Oracle per numerosi anni, per poi essere rilasciato in versione open source da Oracle nell'agosto 2020 su GitHub, sotto una licenza Apache 2.0.

Tribuo nasce con l'obiettivo di fornire una libreria di Machine Learning in ambiente Java che sia in linea con le funzionalità e i bisogni richiesti per lo sviluppo di grandi sistemi software.

La prima versione 1.0.0 di Tribuo viene sviluppata nel 2016 e rilasciata nell'autunno dello stesso anno. Le versioni open source, rilasciate a partire da agosto 2020, sono:

- v4.0.0 (13 agosto 2020 First OS release)
- v4.0.1 (1 settembre 2020 Patch release)
- v4.0.2 (5 novembre 2020 Patch release)
- v4.1.0 (26 maggio 2021 Minor release)

2 Analisi produzione

L'analisi del processo di produzione del software si è basata sulle informazioni ricavate dalla repository oracle/tribuo di GitHub, dal sito, dalla documentazione ufficiale di Tribuo e da alcuni articoli relativi al progetto.

Abbiamo inoltre contattato Adam Pocock, lead developer del progetto, che ha fornito ulteriori informazioni riguardo l'organizzazione del team e l'evoluzione di Tribuo.

2.1 Team e organizzazione

Il team di sviluppo è formato da un gruppo di ricercatori nel campo del Machine Learning che hanno come focus primario research & data science.

Tribuo nasce inizialmente come progetto di due persone e cresce nel tempo coinvolgendo altri membri vista la sua utilità all'interno degli altri progetti seguiti dal gruppo di ricerca.

L'ultima release di Tribuo ha coinvolto 6 sviluppatori interni e un paio di sviluppatori esterni. Molti sviluppatori di Tribuo lavorano contemporaneamente ad altri progetti interni ad Oracle, perciò il lavoro all'interno del gruppo è suddiviso sia sulla base delle aree di maggior competenza di ogni membro, sia considerando gli impegni legati ai diversi progetti di ricerca che questi stanno attualmente seguendo.

I meeting giornalieri sono previsti solo nei momenti in cui il team sta lavorando al design ed implementazione di features che interessano una grossa parte del progetto. Non sono previsti meetings ricorrenti ma solo in caso di necessità.

I membri della squadra si coordinano utilizzando Slack, mentre la scelta dell'IDE è facoltativa, sebbene la maggior parte degli sviluppatori utilizzi IntelliJ IDEA.

Vengono utilizzati inoltre dei plugin come *JaCoCo* per verificare la test coverage e per eseguire analisi statiche sul codice al fine di trovare eventuali bug prima della release del prodotto.

I contributors esterni possono collaborare al progetto aprendo issues su GitHub, segnalando bug, proponendo bug-fixes, features e scrivendo documentazione aggiuntiva.

2.2 Sviluppo software

Lo sviluppo di Tribuo è motivato dalla necessità di avere uno strumento a supporto dei lavori di data science all'interno di Oracle e il progredire del progetto non è quindi legato a regolari deadline. Per questo Tribuo non segue rigorosamente un processo standard di sviluppo software.

2.2.1 Requirements elicitation

Le priorità e gli obiettivi di Tribuo sono quindi definiti sulla base di:

- necessità interne dei ricercatori di Oracle
- ciò che la community richiede
- ciò che gli sviluppatori ritengono sia utile
- elementi presenti nel backlog del progetto (gestito tramite JIRA e GitHub)

Tribuo aderisce al concetto di *semantic versioning*. Si distinguono pertanto:

- **Major release:** sono le release più sostanziose, nelle quali le modifiche apportate riguardano solitamente i pacchetti che stanno alla base della libreria e possono quindi compromettere la retro compatibilità del codice con le versioni precedenti. La struttura di Tribuo è stata volutamente resa il più modulare possibile così da poter offrire la possibilità ai fruitori della libreria di usare unicamente i pacchetti di cui effettivamente hanno bisogno. Questa scelta di design ha portato però con sé ulteriori difficoltà nello sviluppo delle funzionalità riguardanti il core della libreria, rendendo questo tipo di release il meno frequente.
- **Minor release:** garantiscono la retro compatibilità del codice e si concentrano sull'implementazione o sul rework di una singola funzionalità principale (scelta effettuata per mantenere ridotta la complessità) accompagnata da una serie di piccoli bug fixes e features minori presenti nel backlog del progetto.
- **Patch release:** retro compatibili, risolvono eventuali bug e problemi di sicurezza quando questi vengono scoperti.

I requisiti da implementare durante lo sviluppo della prossima release, vengono selezionati a seguito di meeting coordinativi iniziali e seguendo una priorità decisa dal team. In questi meeting viene steso un documento (solitamente in linguaggio informale) contenente le linee guida per la progettazione della prossima release.

L'implementazione delle features avviene su branch separati del progetto, uniti previa approvazione del lead developer, al branch principale quando lo sviluppo e la validazione di queste si è concluso.

2.2.2 Code development

Non seguendo in maniera precisa una metodologia di sviluppo software, né Agile né Plan-based, lo sviluppo non è scandito in fasi ben delineate, ma segue il documento (accessibile a tutto il gruppo) sulla base del quale viene effettuato lo sviluppo software. Esso viene gestito ed aggiornato tramite Slack, software di collaborazione aziendale utilizzato per scambiare messaggi e documenti in modo istantaneo all'interno del team, in modo da fornire una panoramica riguardante lo stato corrente del progetto.

Il compito di ogni membro del team, a differenza di un modello Plan-based, può includere mansioni diverse all'interno del progetto, potenzialmente anche variabili nel tempo. Queste vengono definite in relazione, oltre che alle competenze individuali, anche agli impegni legati ad altri progetti di ricerca. L'assegnamento e la coordinazione generale di questi compiti viene supervisionata dal lead developer del progetto, il quale detenendone la responsabilità si assicura, tramite un costante rapporto diretto tra i membri del team, che il lavoro venga svolto correttamente e efficientemente.

Nel caso uno dei membri si trovi in difficoltà riguardo una sezione particolarmente problematica del codice, egli ha modo di contattare immediatamente gli altri e, se necessario, viene organizzato un meeting al fine di trovare una soluzione adeguata. Durante la progettazione di alcuni degli elementi più complessi all'interno del progetto, gli sviluppatori hanno organizzato sessioni di **pair programming**, durante la quale due membri del team lavorano contemporaneamente all'implementazione di una stessa parte di codice scambiandosi regolarmente il ruolo di “*driver*” (colui che scrive il codice e espone una possibile implementazione funzionante) e “*navigator*” (colui che lo controlla e propone soluzioni alternative).

Nonostante si faccia largo uso di test all'interno del progetto, Tribuo non aderisce rigorosamente al modello del **test driven development**. Infatti, a seconda della natura della feature da implementare, in alcuni casi i test, effettuati attraverso *JUnit 5*, vengono scritti *prima* dell'effettiva implementazione di una specifica funzionalità (come previsto dal modello TDD), mentre in altri casi i test vengono ideati solamente *dopo* che gran parte del codice relativo alla feature è stato già scritto, al fine di validare il suo comportamento. In altri casi ancora può capitare di scrivere test per una libreria già esistente (e già testata in precedenza) al fine di estendere la *test coverage* di quella porzione di codice, sebbene ciò avvenga solitamente in preparazione ad un imminente *refactoring* dello stesso.

2.2.3 Software validation & delivery

Il codice appartenente alla prossima release deve quindi superare una fase di testing (**validazione**) durante la quale viene verificata la conformità ai requisiti ed il corretto funzionamento delle singole feature in differenti modalità d'uso.

Terminato il processo di validazione la release è pronta per essere rilasciata. Viene quindi effettuato il **delivery** del prodotto.

3 Analisi metriche

In questa sezione si analizza l'evoluzione del codice di Tribuo, prendendo in considerazione le metriche LOC, CK, CC.

I dati analizzati vanno a confrontare tra loro le quattro releases open source di Tribuo partendo dalla versione 4.0.0 fino alla 4.1.0.

Le metriche delle versioni 4.0.1 e 4.0.2, essendo patch release, evidenzieranno una variazione minima tra le due versioni mentre quelle delle minor release 4.0.0 e 4.0.1 saranno invece caratterizzate da differenze maggiori.

	Number of classes	Number of packages	Number of external packages	Number of external classes
Tribuo 4.0.0	782	101	23	90
Tribuo 4.0.1	787	101	23	90
Tribuo 4.0.2	788	101	23	90
Tribuo 4.1.0	873	109	31	116

Tabella 1: Numero di classi e pacchetti per ogni release di Tribuo.

A causa dell'elevato numero di dati ottenuti, il confronto tra le metriche delle varie versioni è stato effettuato principalmente per via grafica in modo da poterlo rendere più visibile e immediato.

3.1 Lines Of Code (LOC)

Le linee di codice sono state misurate utilizzando l'estensione *VS Code Counter* dell'editor Visual Studio Code. Vengono riportate le linee di codice di tutti i linguaggi presenti all'interno del progetto, differenziate tra **LOC** (*Lines Of Code*), **CLOC** (*Comment Lines Of Code*), **ELOC** (*Empty Lines Of Code*) e **TLOC** (*Total Lines Of Code*).

language	files	code	comment	blank	total
Java	683	52,956	25,543	12,162	90,661
XML	70	4,148	1,109	324	5,581
Markdown	14	1,359	0	302	1,661
Python	5	240	111	58	409
JSON	1	91	0	1	92
Shell Script	3	22	3	12	37
YAML	1	21	3	4	28

(a) v4.0.0

language	files	code	comment	blank	total
Java	685	53,522	25,691	12,252	91,465
XML	70	4,150	1,109	324	5,583
Markdown	14	1,360	0	302	1,662
Python	5	240	111	58	409
JSON	1	93	0	1	94
Shell Script	3	22	3	12	37
YAML	1	21	3	4	28

(b) v4.0.1

language	files	code	comment	blank	total
Java	687	54,135	27,191	12,360	93,686
XML	71	4,349	1,152	330	5,831
Markdown	16	1,498	1	337	1,836
JSON	4	558	0	6	564
Python	5	240	111	58	409
YAML	3	63	9	12	84
Shell Script	3	22	3	12	37

(c) v4.0.2

language	files	code	comment	blank	total
Java	752	61,188	32,016	13,836	107,040
JSON	6	31,445	0	8	31,453
XML	77	4,918	1,257	364	6,539
Markdown	18	1,613	1	364	1,978
Python	4	178	61	37	276
YAML	3	63	9	12	84
Shell Script	3	22	3	12	37
HTML	1	16	0	1	17

(d) v4.1.0

Figura 1: *Lines Of Code* relative alle varie releases di Tribuo.

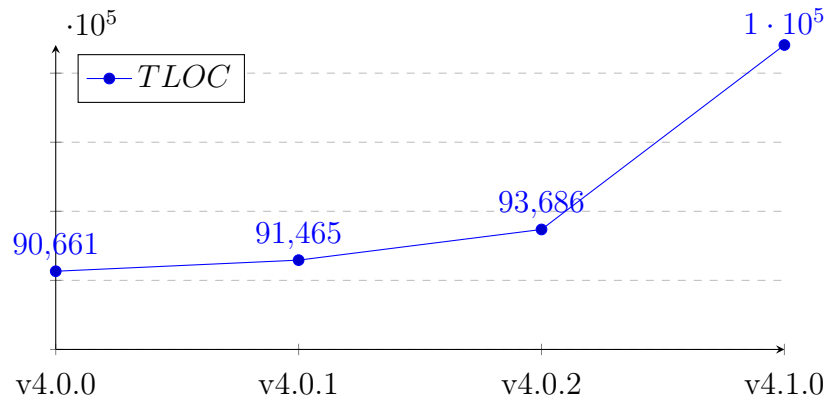


Figura 2: Andamento della metrica *TLOC* lungo lo sviluppo di Tribuo.

3.2 Chidamber and Kemerer's metrics suite (CK)

Una della suite di metriche più conosciute ed utilizzate in ambito Object Oriented è la **CK Suite** (*Chidamber and Kemerer's metrics suite*). Essa fornisce delle misure numeriche per ogni classe del progetto, facili da ottenere e utili per valutare la qualità del codice. Queste metriche, ottenute attraverso il plugin *CodeMR* dell'IDE IntelliJ IDEA, comprendono:

- *Weighted Methods per Class* (**WMC**)
- *Depth of Inheritance Tree* (**DIT**)
- *Number Of Children* (**NOC**)
- *Coupling Between Objects* (**CBO**)
- *Response For a Class* (**RFC**)
- *Lack of Cohesion in Methods* (**LCOM**)¹

Essendo questa una misura numerica ottenuta a livello di classe e non di progetto, si è deciso di riportare, per ogni release di Tribuo, i valori di tali metriche in tre formati differenti.

Il primo (figura 3) rappresenta ogni metrica di CK sotto forma di grafico a torta. Il colore e la grandezza di ogni fetta è proporzionale alla quantità di classi del progetto aventi quella determinata metrica a quel determinato livello (si veda la legenda di tabella 2).

Il secondo (figura 4) e il terzo (figura 5) confrontano, attraverso un istogramma e dei grafici, le metriche della CK Suite per l'intero progetto in modo da poterne tracciare la progressione durante lo sviluppo delle differenti releases.

	Low	Low-Medium	Medium-High	High	Very High
WMC	0-20	20-50	50-101	101-120	>120
DIT	0-1	1-3	3-10	10-20	>20
NOC	0-1	1-5	5-10	10-15	>15
CBO	0-5	5-10	10-20	20-30	>30
RFC	0-50	50-100	100-150	150-200	>200
LCOM	0-0.5	0.5-0.7	0.7-0.8	0.8-1	>1

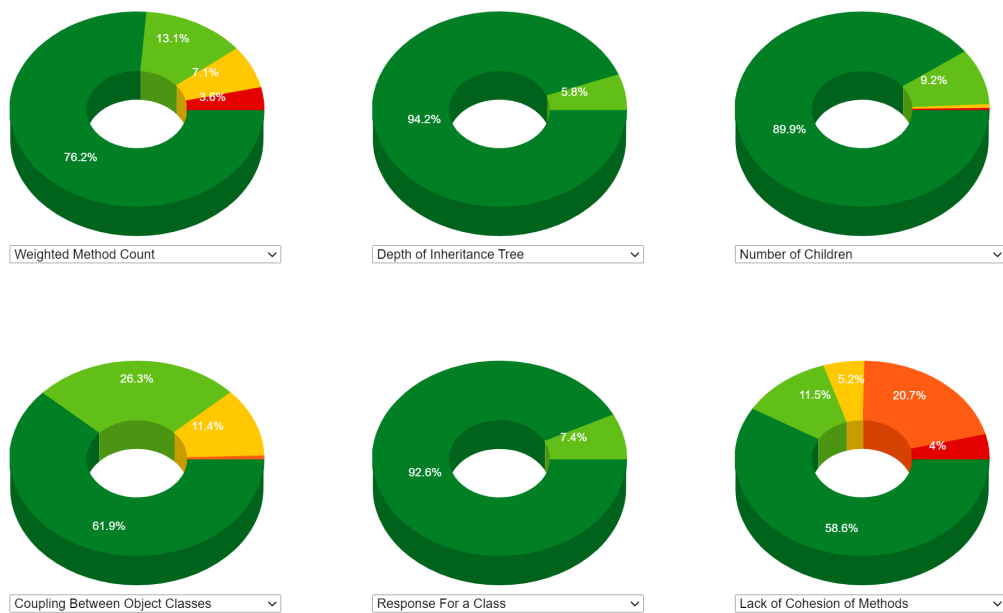
Tabella 2: Legenda di grafici a torta in figura 3.

¹LCOM è compreso tra 0 e 2, dove valori superiori a 1 sono considerati allarmanti poiché indicano una mancanza di coesione critica. Il suo valore è definito come:

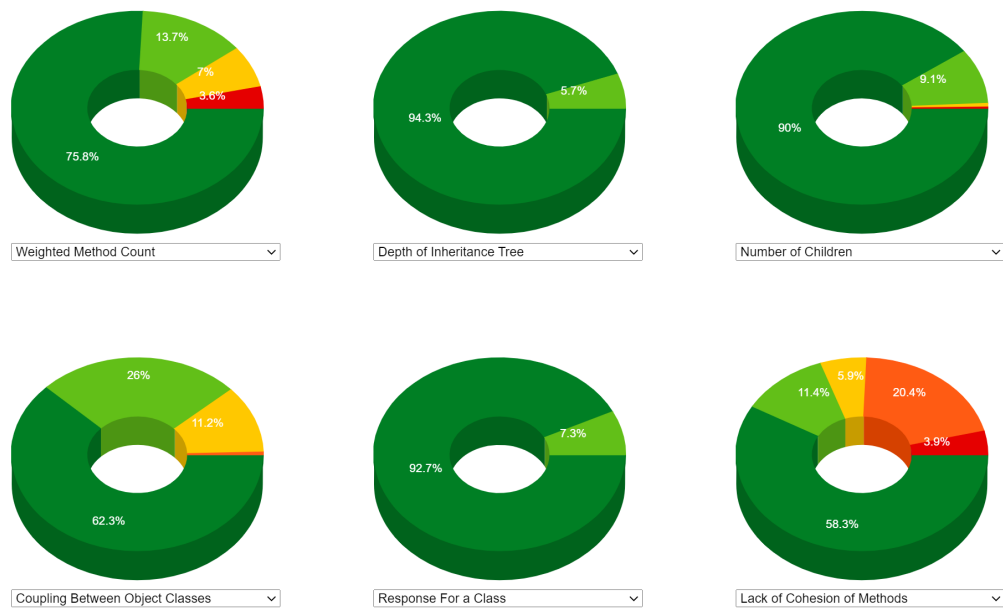
$$LCOM = \frac{m - \frac{\sum mA}{a}}{m - 1} \quad (1)$$

dove:

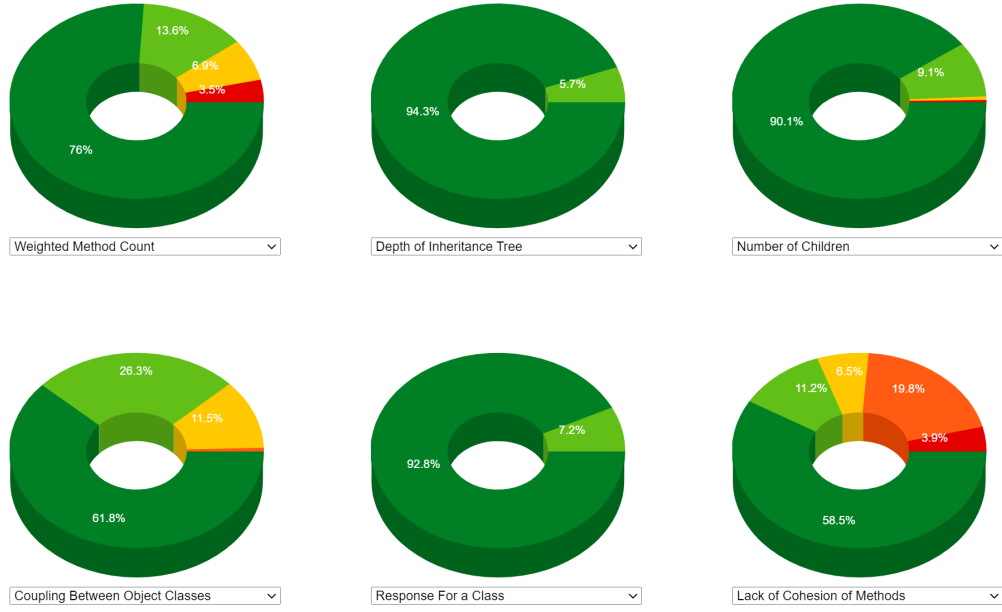
- m , numero di metodi della classe
- a , numero di variabili della classe
- mA , numero di metodi che accedono a una variabile
- $\sum mA$, sommatoria degli mA per ogni variabile della classe



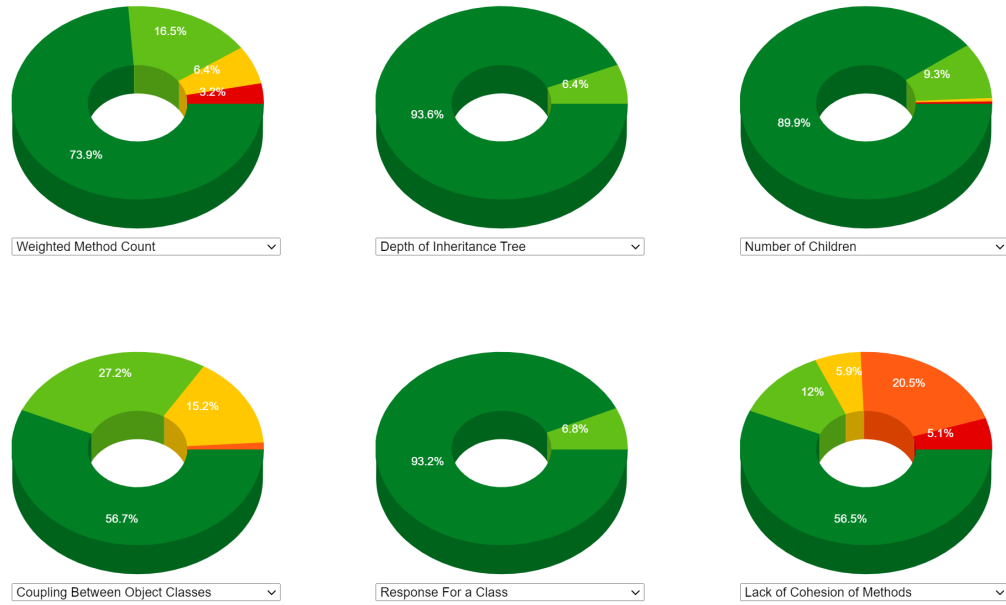
(a) v4.0.0



(b) v4.0.1



(c) v4.0.2



(d) v4.1.0

Figura 3: Metriche *CK* relative alle varie releases di Tribuo.

Nei grafici seguenti si utilizzano le metriche della CK Suite (definite $\forall classe \in progetto$) le quali sono state estese per ogni release e di cui vengono riportati i valori medi.

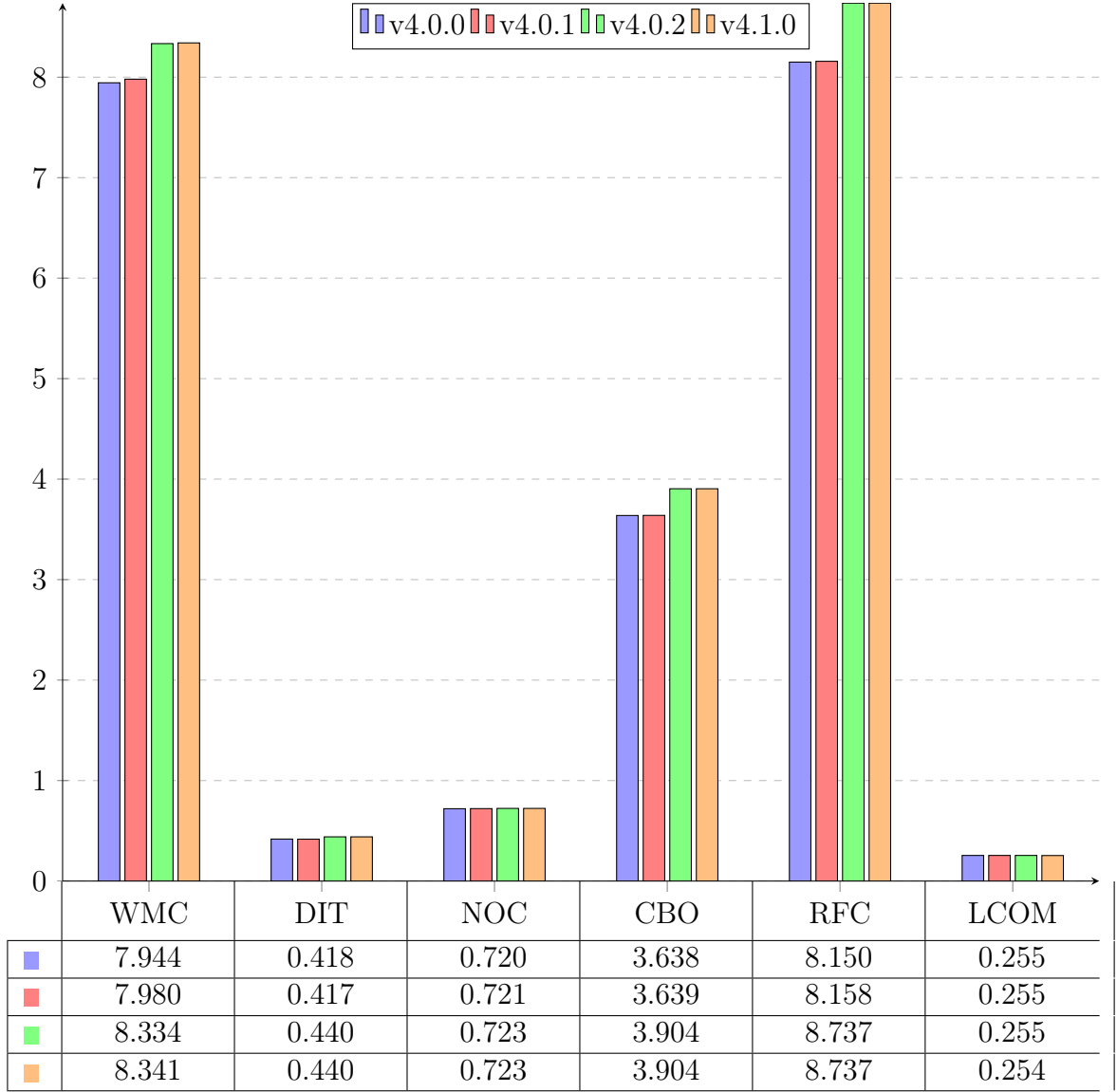
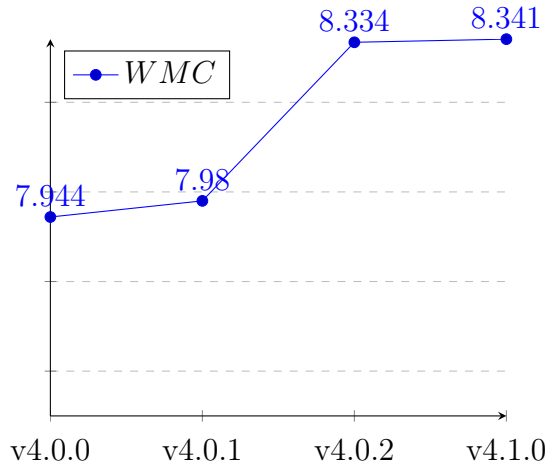
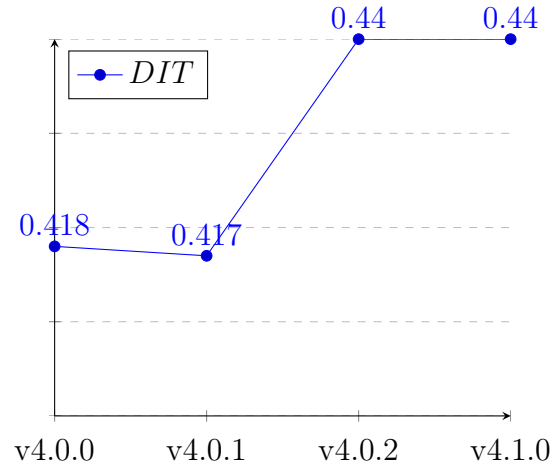


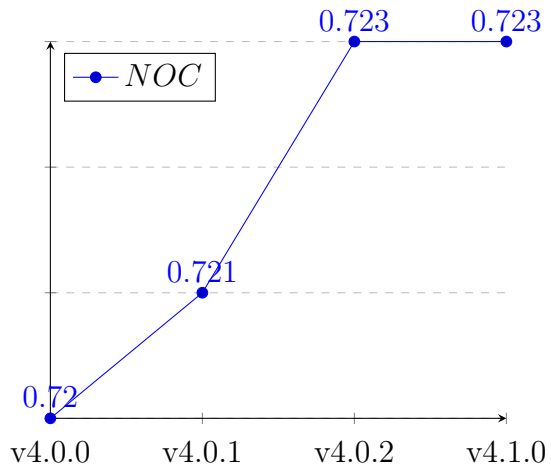
Figura 4: Andamento delle metriche *CK* lungo lo sviluppo di Tribuo.



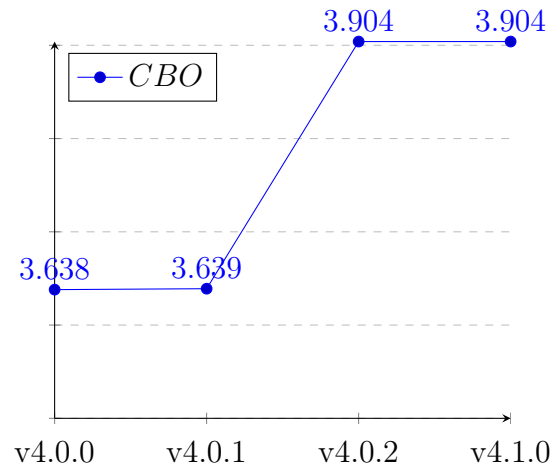
(a) WMC



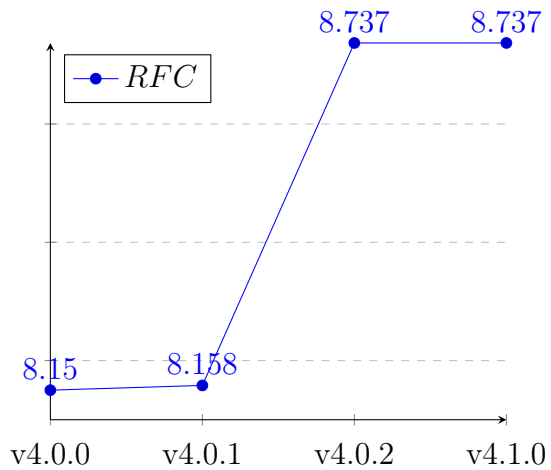
(b) DIT



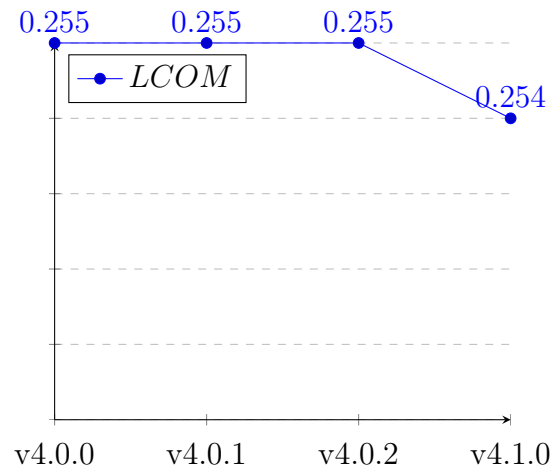
(c) NOC



(d) CBO



(e) RFC



(f) LCOM

Figura 5: Andamento delle metriche *CK* lungo lo sviluppo di Tribuo.

3.3 McCabe Cyclomatic Complexity (CC)

L'analisi della complessità ciclomatica viene effettuata attraverso il software JArchitect al fine di valutare l'andamento della complessità del codice di Tribuo all'interno delle varie releases.

La complessità ciclomatica è definita per i singoli metodi come:

$$CC_{method} = 1 + \{decision\ points\ found\ in\ the\ body\ of\ the\ method\} \quad (2)$$

È possibile estendere la complessità ciclomatica di una singola classe con la seguente formula:

$$CC_{class} = \sum_{\forall method \in class} CC_{method} \quad (3)$$

Perciò, al fine di ottenere una metrica di progetto, si considera la metrica ottenuta dalla somma della complessità ciclomatica delle singole classi facenti parti del progetto:

$$CC_{project} = \sum_{\forall class \in project} CC_{class} \quad (4)$$

Viene di seguito riportato un grafico raffigurante l'andamento delle CC risalenti alle quattro differenti releases di Tribuo.

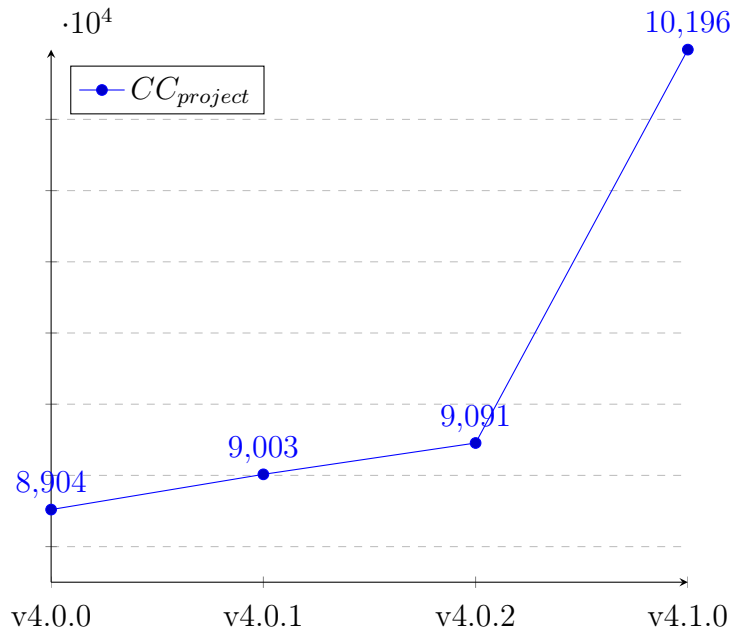


Figura 6: Andamento della CC lungo lo sviluppo di Tribuo.

4 Conclusione

L'analisi del codice ha evidenziato lievi cambiamenti tra le metriche descrittive del progetto. Ciò è dovuto al fatto che Tribuo, essendo un software “giovane”, non ha ancora ricevuto major release, ma solo patch e minor release, dopo la prima OS rilasciata nel Settembre 2020.

Versione	LOC	CK						CC
		WMC	DIT	NOC	CBO	RFC	LCOM	
v4.0.0 → v4.0.1	0,89%	0,45%	-0,34%	0,13%	0,03%	0,10%	0,18%	1,11%
v4.0.1 → v4.0.2	2,43%	4,45%	5,51%	0,25%	7,27%	7,09%	-0,01%	0,98%
V4.0.2 → v4.1.0	14,25%	0,08%	0,00%	0,00%	0,00%	0,00%	-0,24%	12,15%
v4.0.0 → v4.1.0	18,07%	5,00%	5,15%	0,39%	7,31%	7,20%	-0,07%	14,51%

Tabella 3: Incremento percentuale delle varie metriche tra tutte le versioni.

Nel dettaglio, le principali differenze riscontrate tra le varie versioni sono:

- **LOC**: Incremento complessivo tra le quattro versioni di circa il 18% con un picco legato a v4.1.0. Questo dato trova riscontro nel significativo aumento di classi e pacchetti che ha caratterizzato quest'ultima release.
- **CK**: Incremento medio complessivo delle metriche *WMC*, *DIT*, *CBO*, *RFC* con valori compresi tra il 5% e il 7% con picco legato a v4.1.0. L'incremento di queste metriche aumenta le difficoltà legate allo sviluppo, testing e manutenzione del codice. Le metriche *NOC* e *LCOM*, contrariamente alle precedenti, sono invece rimaste costanti durante le release analizzate.
- **CC**: Incremento complessivo tra le quattro versioni del 14.5% con un picco legato a v4.1.0. L'incremento del valore della *CC* denota una crescente complessità assoluta ma, considerato il parallelo aumento della metrica *LOC*, la complessità relativa del progetto (CC/LOC) è rimasta pressapoco costante.