

Relazioni Quantitative Struttura-Attività (QSAR) Mediante Algoritmi di Machine Learning. Sviluppo di Modelli per la Predizione di Attività di inibitori di Janus Chinasi 2 (JAK2)

**Tesi sperimentale
in
Chimica Farmaceutica**



SAPIENZA
UNIVERSITÀ DI ROMA

**Facoltà di Farmacia e Medicina
Corso di Laurea in Chimica e Tecnologia Farmaceutiche
Tesi Sperimentale in Chimica Farmaceutica
a.a. 2020/2021**

**Laureanda: Giorgia Canini
Matricola: 1583394**

Relatore: prof. Rino Ragno



Obiettivo della ricerca



Dataset inibitori JAK 2

QSAR

Machine Learning

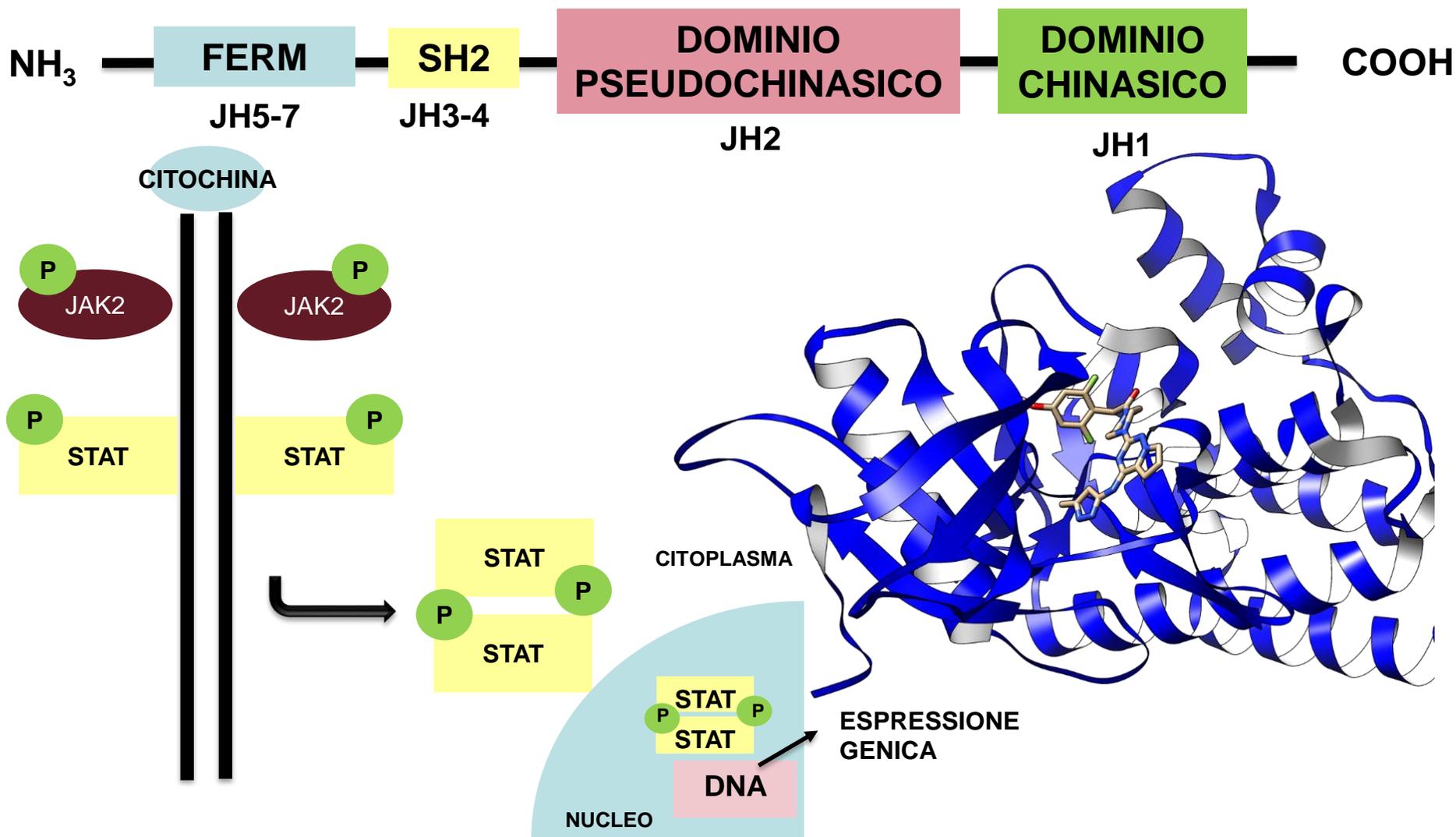
Virtual Screening

Nuovi potenziali
inibitori JAK2

Scopo
del
lavoro

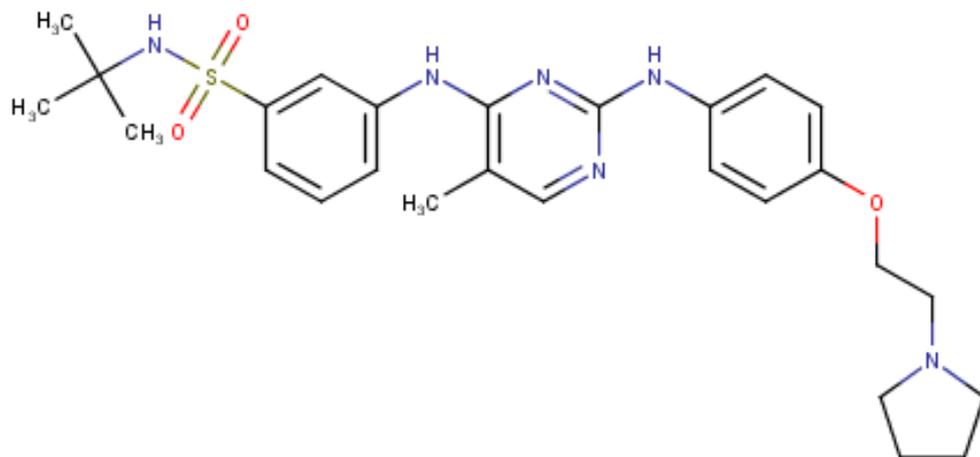


Target: JAK2



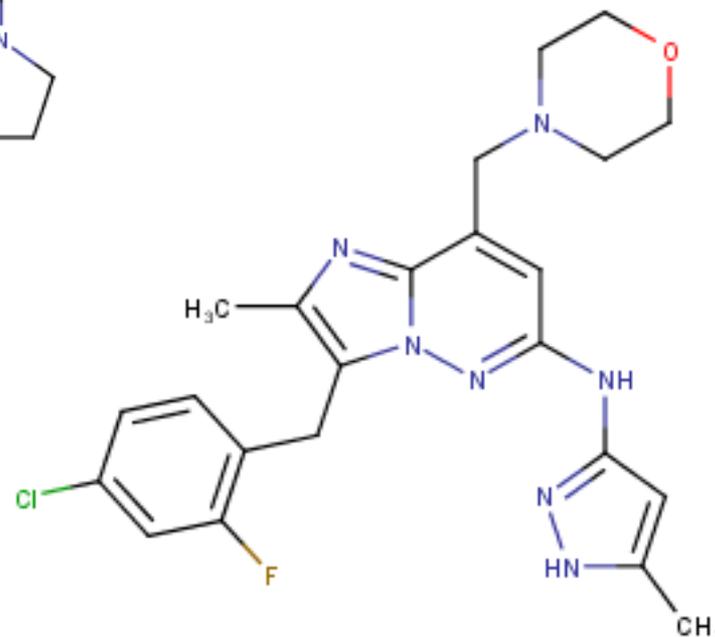


Inibitori JAK2



Fedratinib

- Inibitore selettivo di JAK2
- Approvato nel 2019 negli Stati Uniti per mielofibrosi primaria o secondaria di livello intermedio o ad alto rischio



Gandotinib

- LY2784544, inibitore selettivo JAK2
- Studi fase II nel 2018 per neoplasie mieloproliferative



Procedura sperimentale



1

Generazione
training set
senza attività
binarizzata

2

Modelli di
classificazione
e
regressione

3

Validazione
esterna
modelli

M
E
T
O
D
O
L
O
G
I
A

S
P
E
R
I
M
E
N
T
A
L
E

Virtual
Screening

6

Studio di
selettività
JAK2 / JAK1

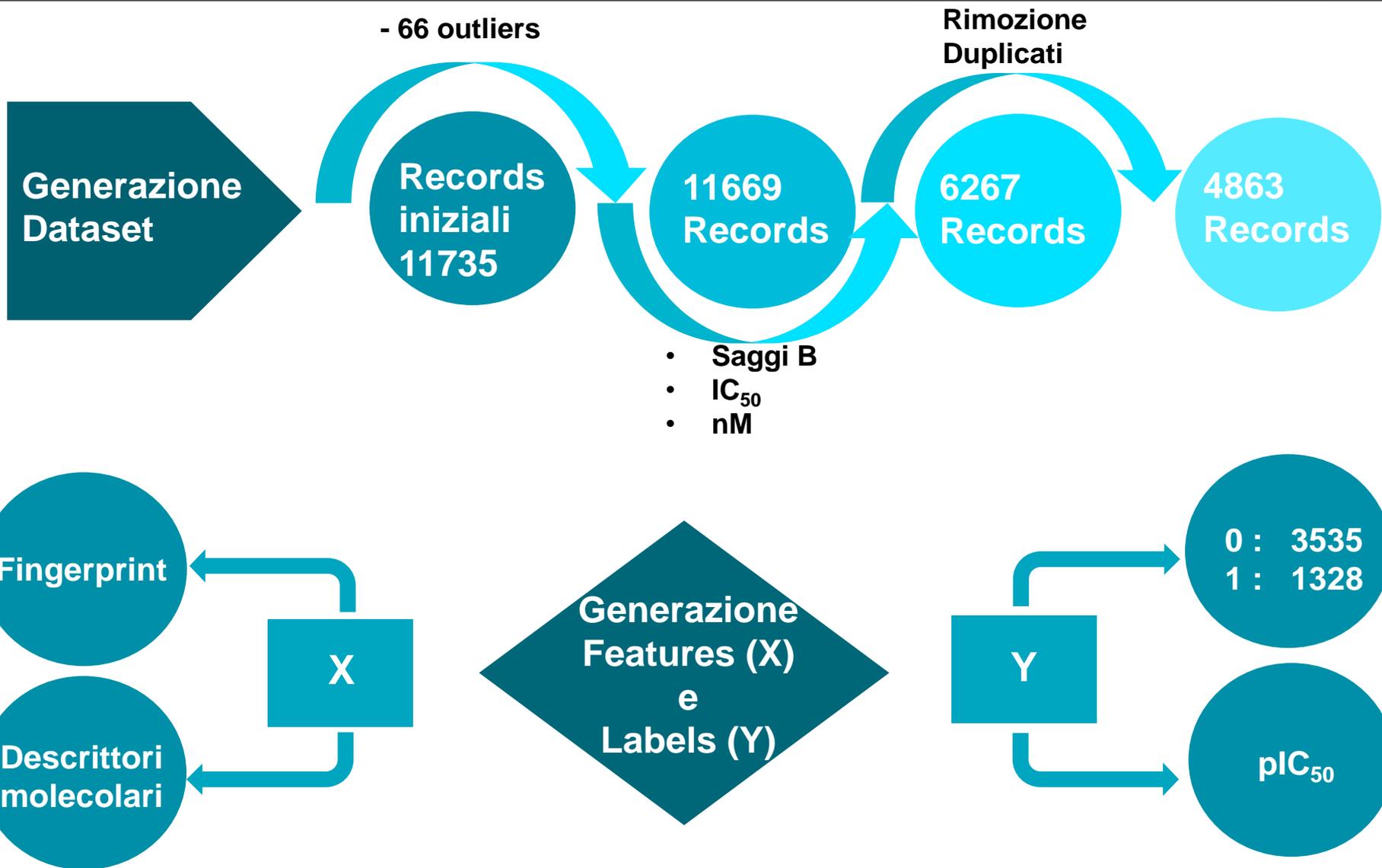
5

Analisi
similarità
molecolare

4

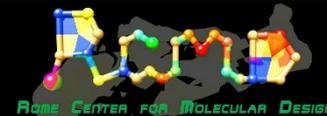


Generazione train set





Modelli di classificazione



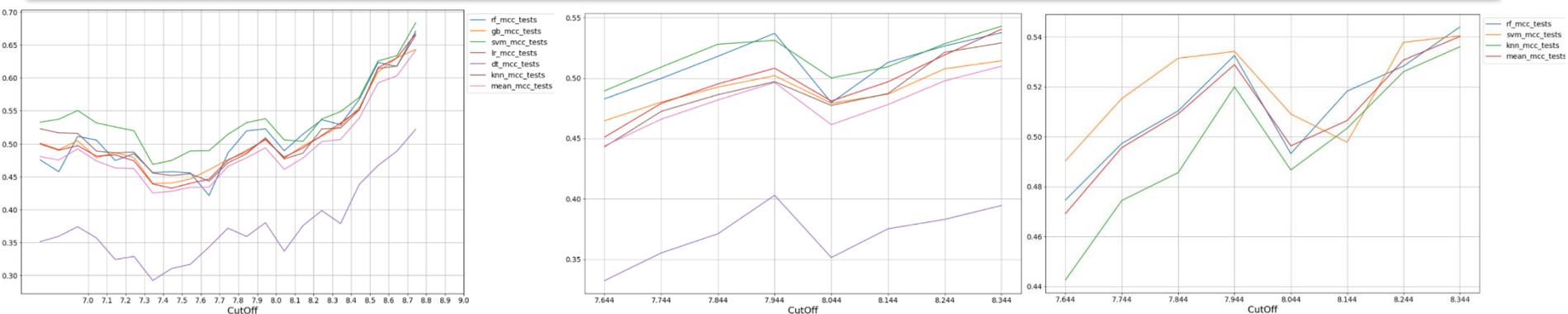
Stratified K-Fold
Train set 10%
Test set 90%

10 iterazioni
Stabilire l'intorno
del cutoff

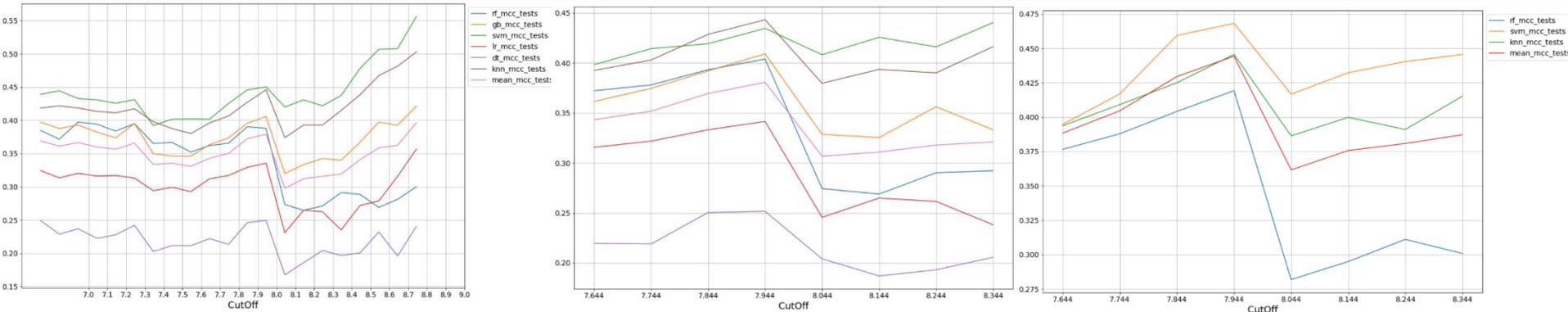
100 iterazioni
Selezione
algoritmi migliori

1000 iterazioni
Scelta cutoff e
algoritmi ottimali

Fingerprint



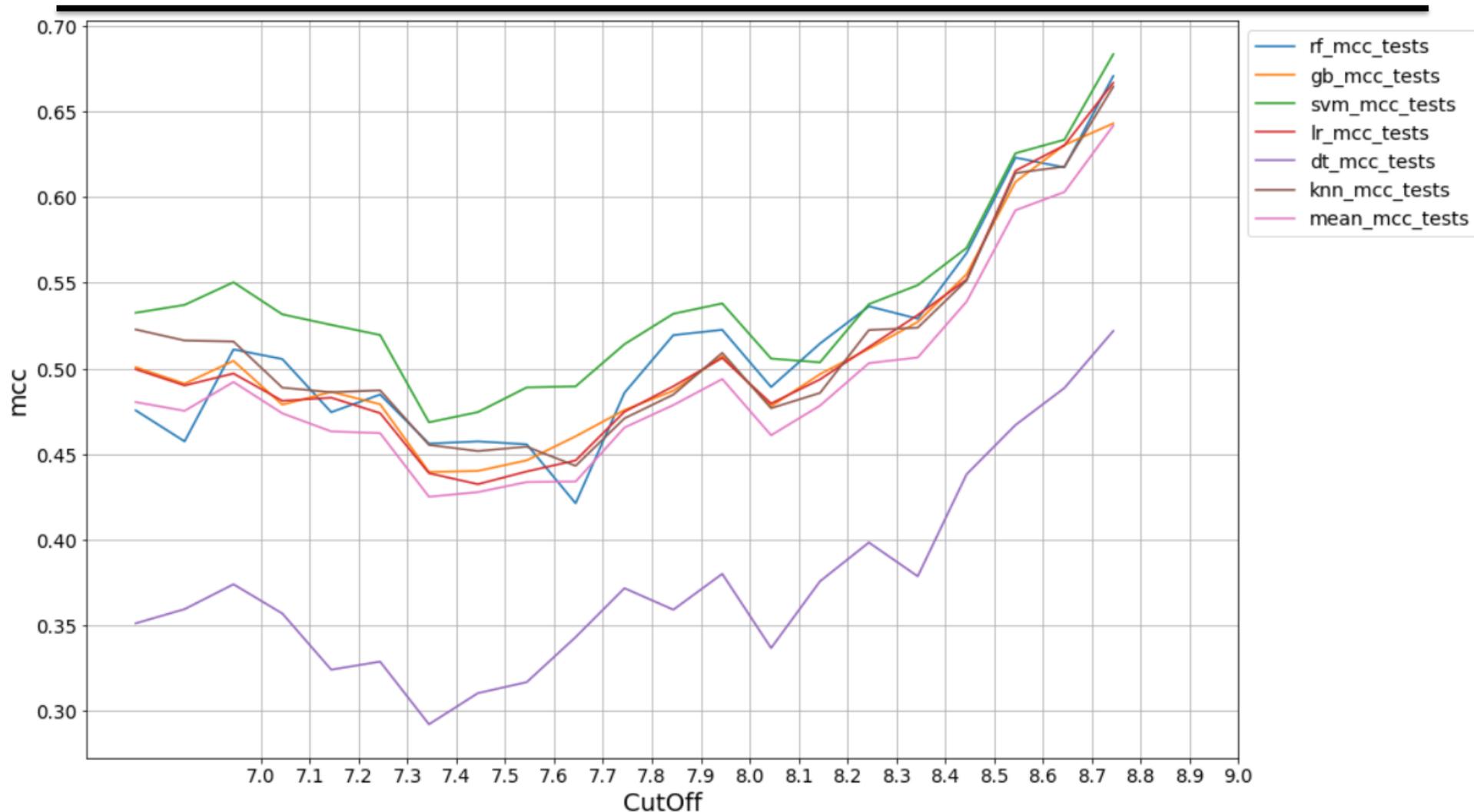
Descrittori





Modelli di classificazione

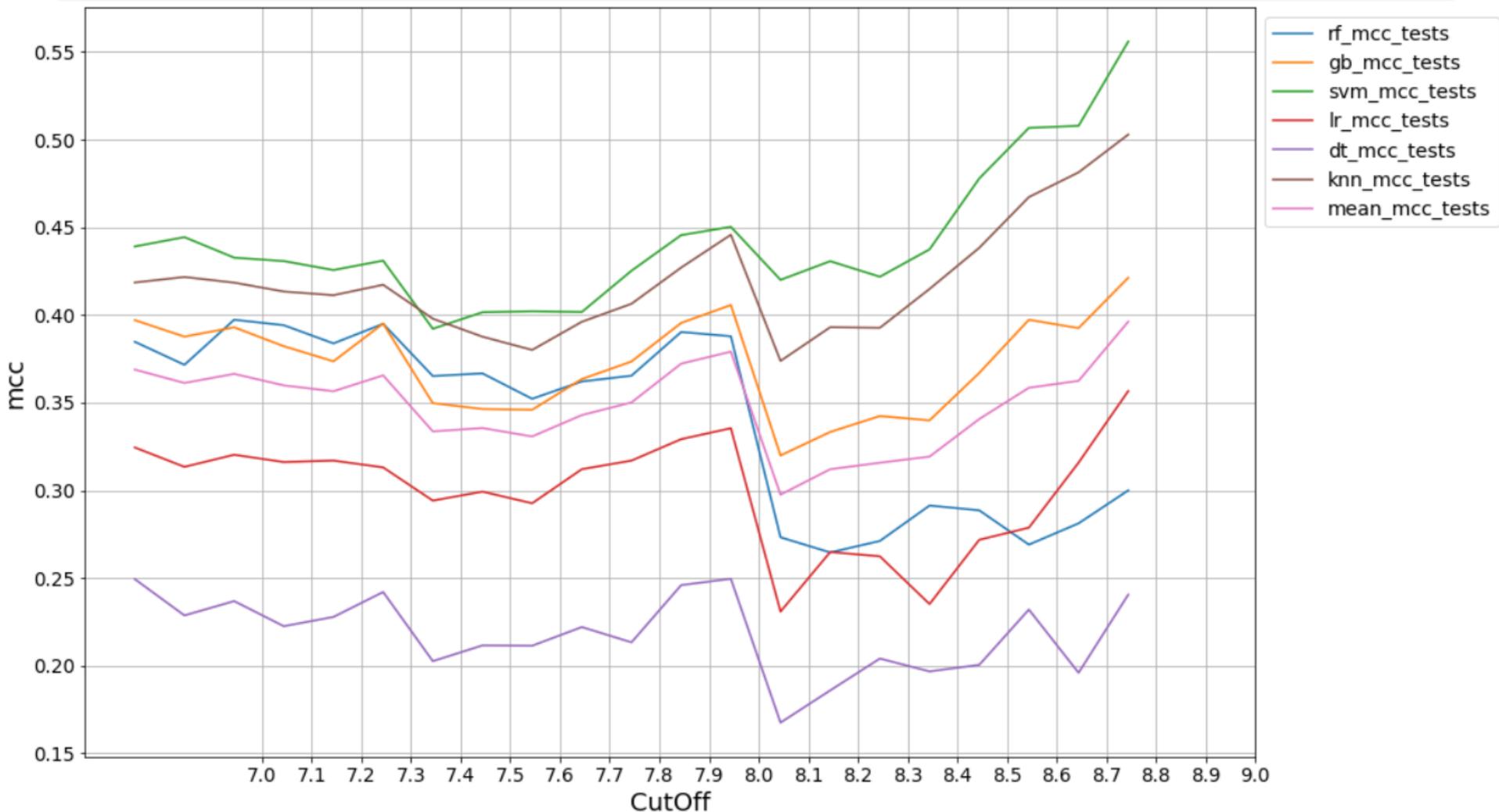
Fingerprint





Modelli di classificazione

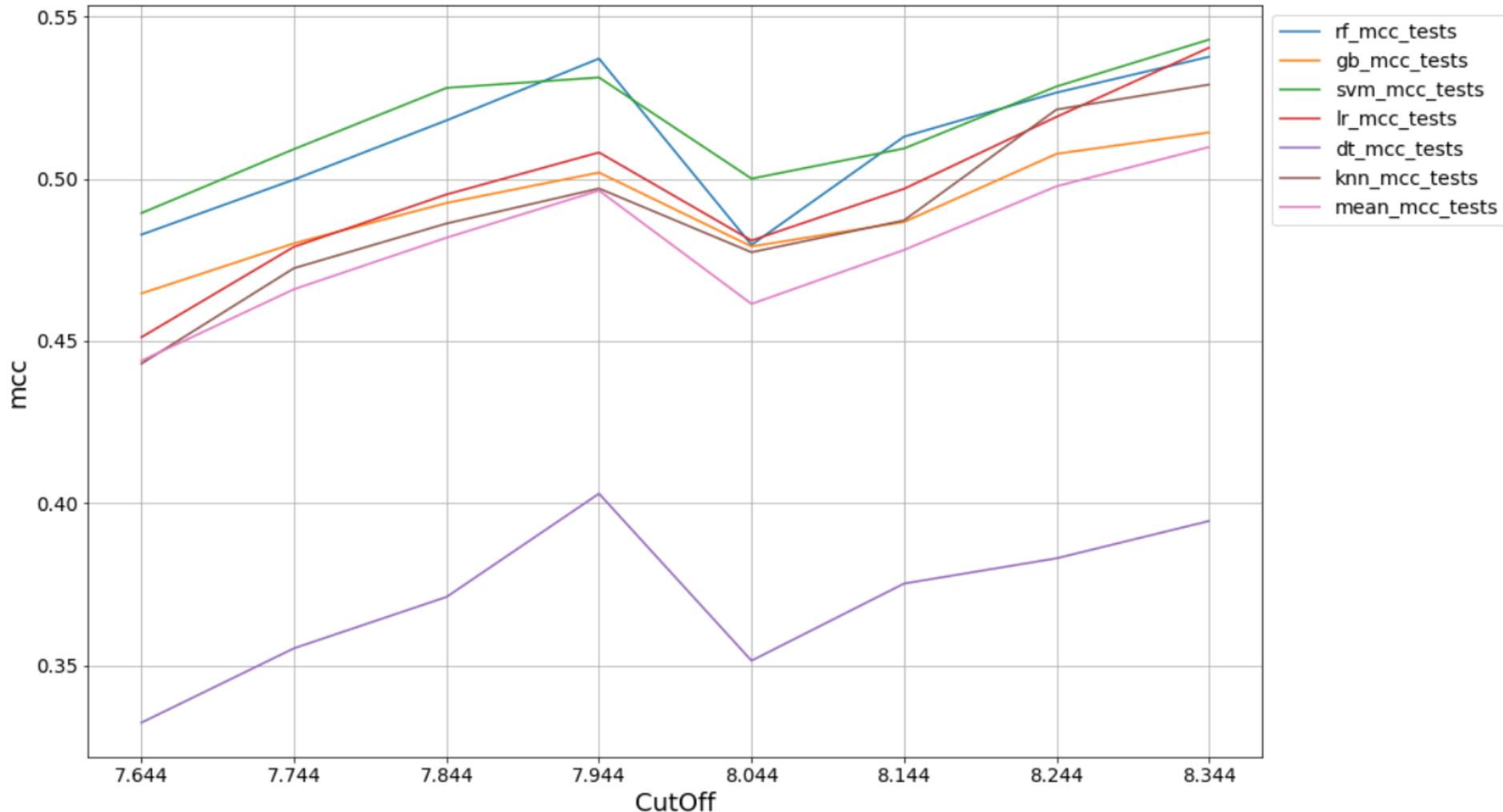
Descrittori





Modelli di classificazione

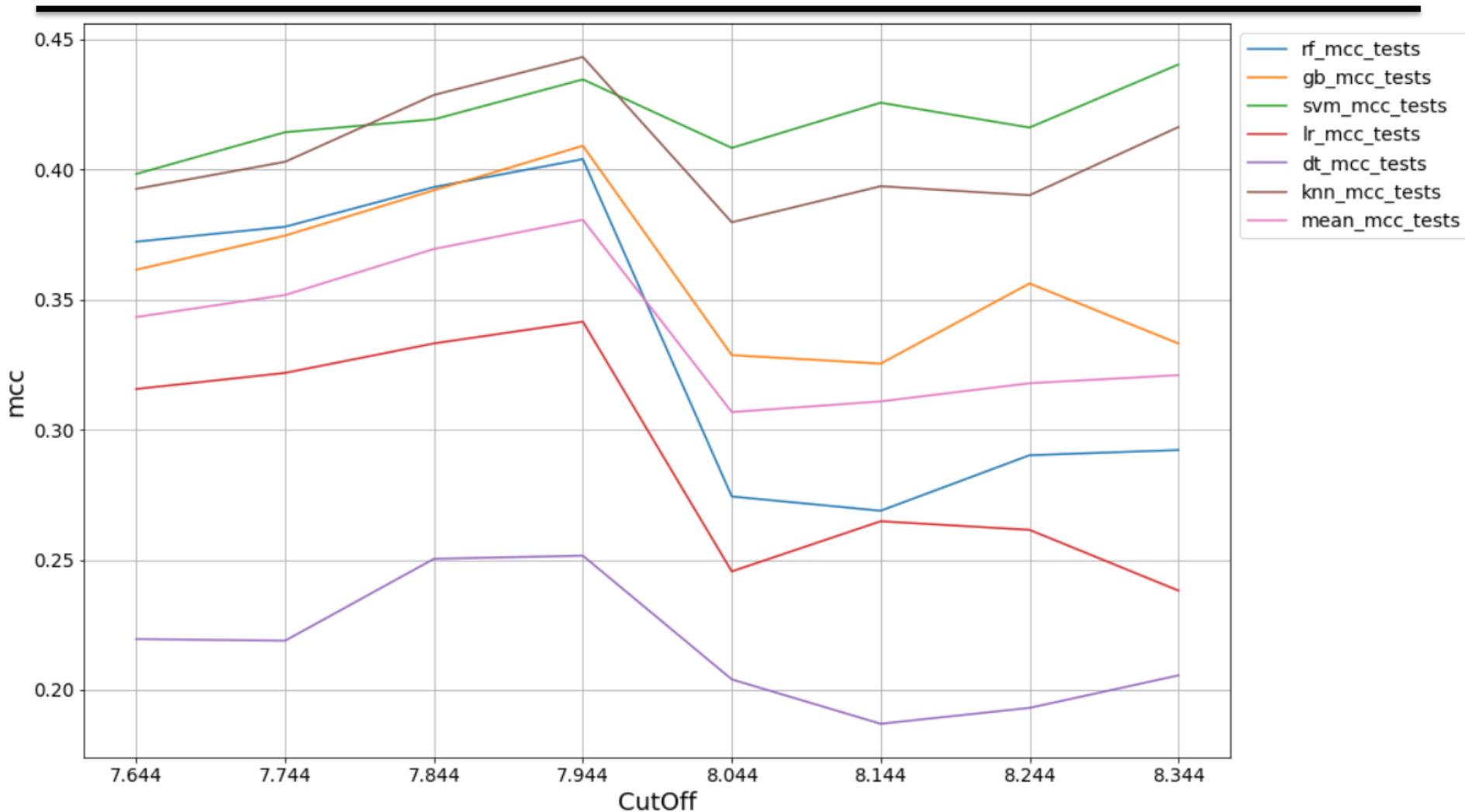
Fingerprint





Modelli di classificazione

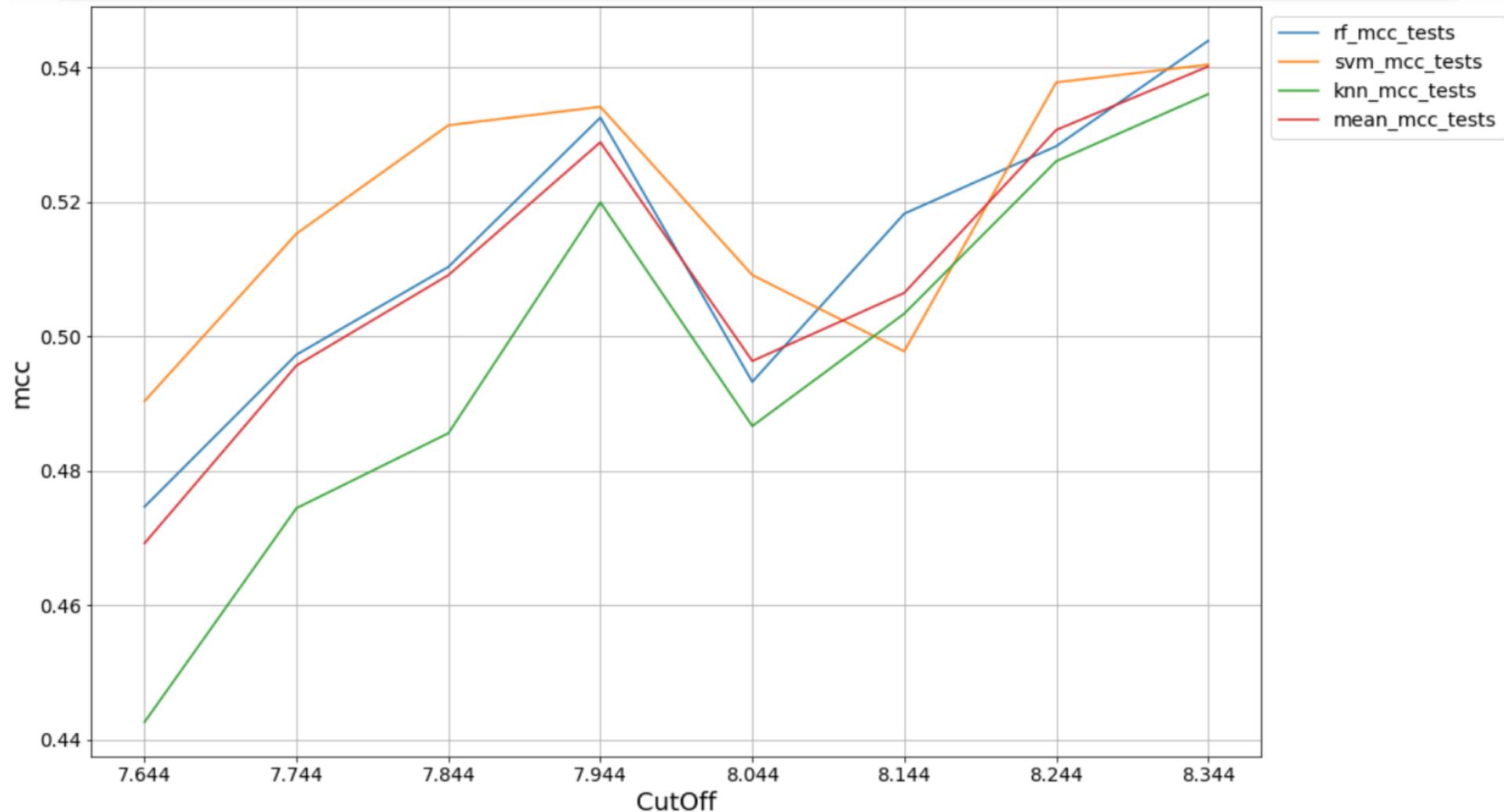
Descrittori





Modelli di classificazione

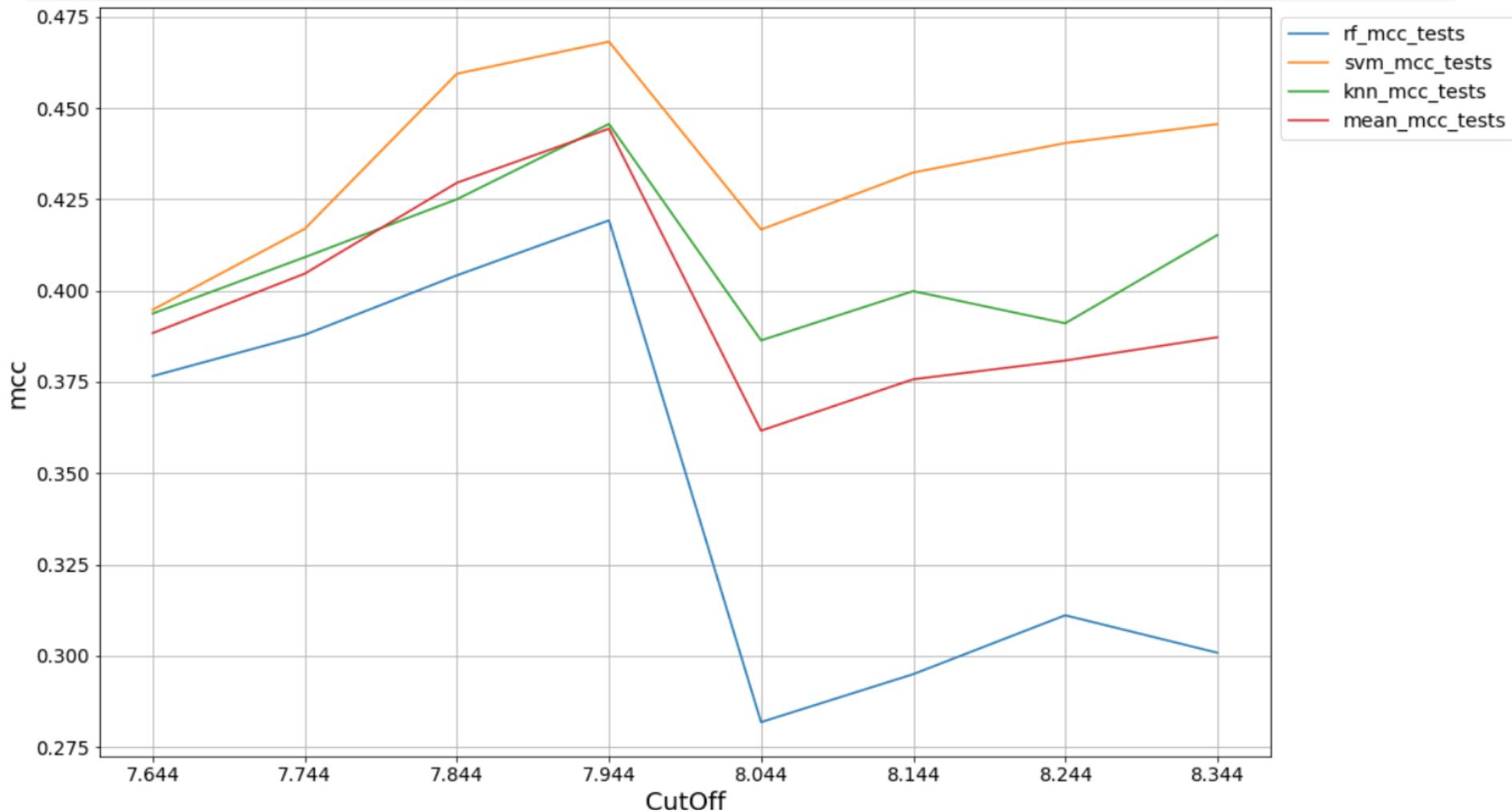
Fingerprint





Modelli di classificazione

Descrittori



Training set 80%
Test set 20%

Ottimizzazione
iperparametri

Ricerca modello
migliore

Dataset 100%

Validazione
interna

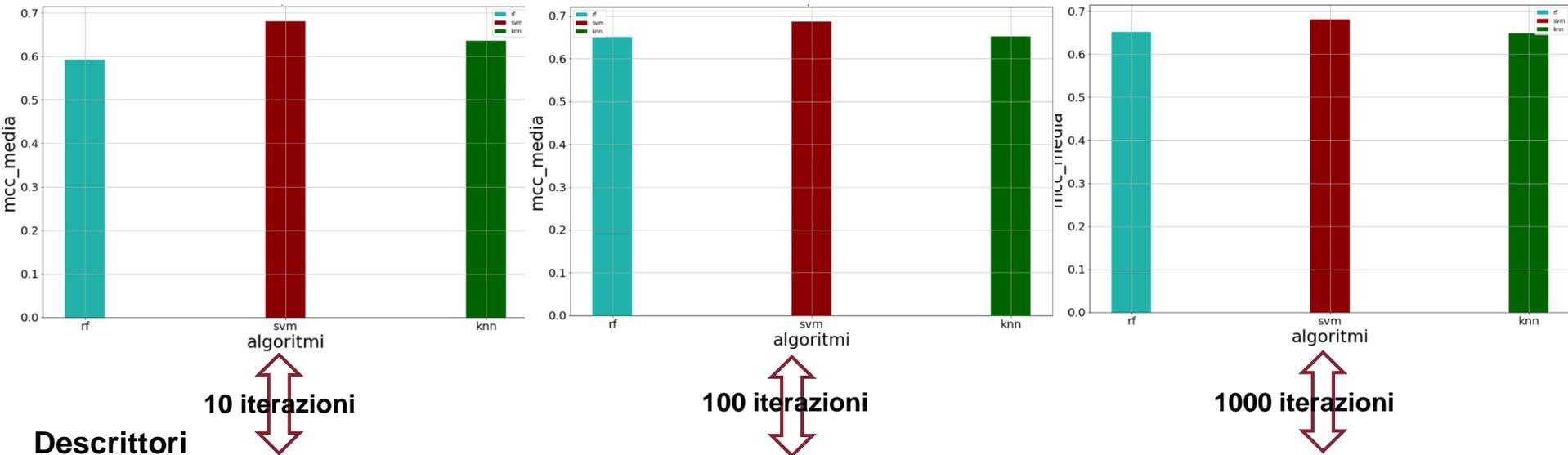
$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

- $-1 < MCC < 1$;
- Misura della qualità delle classificazioni binarie;
- Metrica per validare modelli di classificazione.
- veri positivi (TP): composti attivi correttamente predetti come attivi;
- falsi negativi (FN): composti attivi erroneamente classificati come inattivi;
- falsi positivi (FP): composti inattivi erroneamente classificati come attivi;
- veri negativi (TN): categoria che include i composti inattivi correttamente predetti inattivi

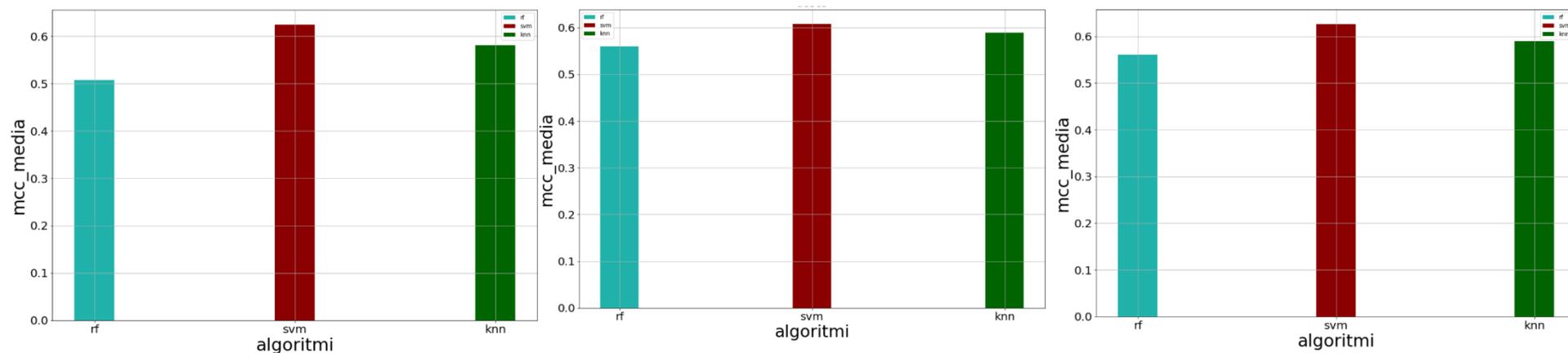


Modelli di classificazione

Fingerprint



Descrittori



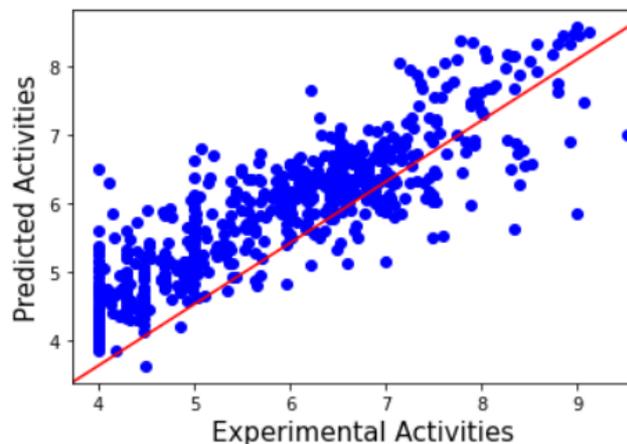


Training regressione

| R^2 | Fingerprint | Descrittori | Descrittori scalati |
|-----------------------------|-------------|-------------|---------------------|
| SupportVectorMachine | 0.77 | 0.64 | 0.63 |
| RandomForest | 0.69 | 0.54 | 0.65 |
| Ridge | 0.67 | 0.23 | 0.34 |

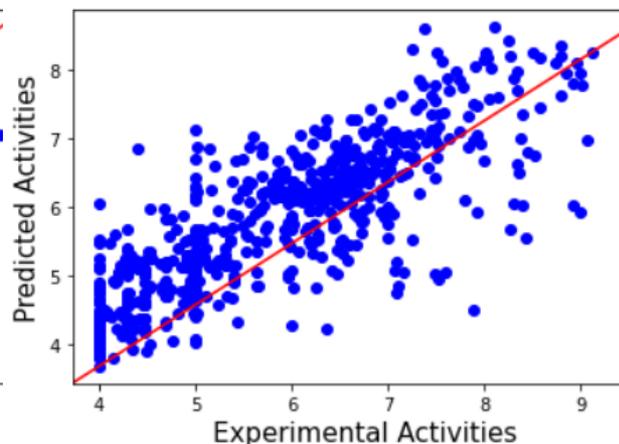
SVR Fingerprint

C=5, Kernel = rbf



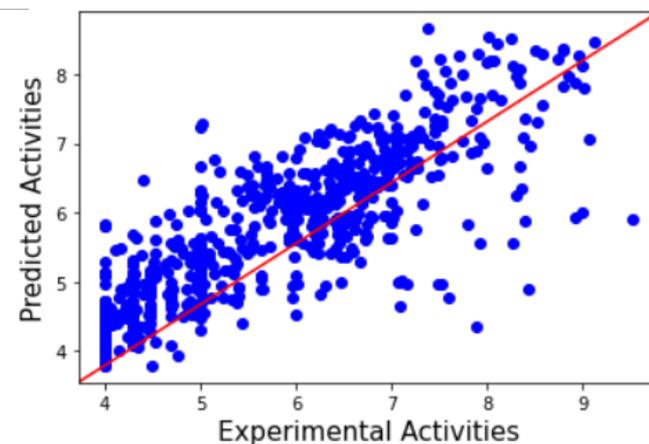
SVR Descrittori

C=9, Kernel = rbf



RF Descrittori scalati

Max_depth=20, n_estimators=104



CLASSIFICAZIONE

| | MCC_FIT | MCC_CV | MCC_EXT |
|------------------------|---------|--------|---------|
| Fingerprint | 0.95 | 0.67 | 0.72 |
| Descrittori molecolari | 0.80 | 0.59 | 0.65 |

REGRESSIONE

| | Q ² | SDEP |
|------------------------|----------------|------|
| Fingerprint | 0.77 | 0.58 |
| Descrittori molecolari | 0.66 | 0.71 |
| Descrittori scalati | 0.54 | 0.82 |

- Quanto un modello è predittivo
- SDEP (Standard Deviation Error in Prediction)
- Confrontabile con SDEP_CV



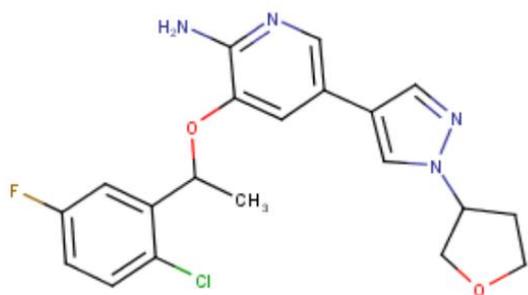
Analisi della similarità

Ricerca
test set

Confronto
train e
test set

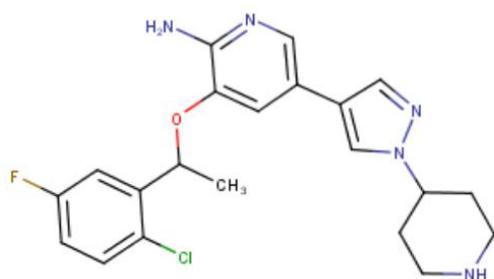
Generazione
features e
labels

62 molecole:
17 attive
45 inattive



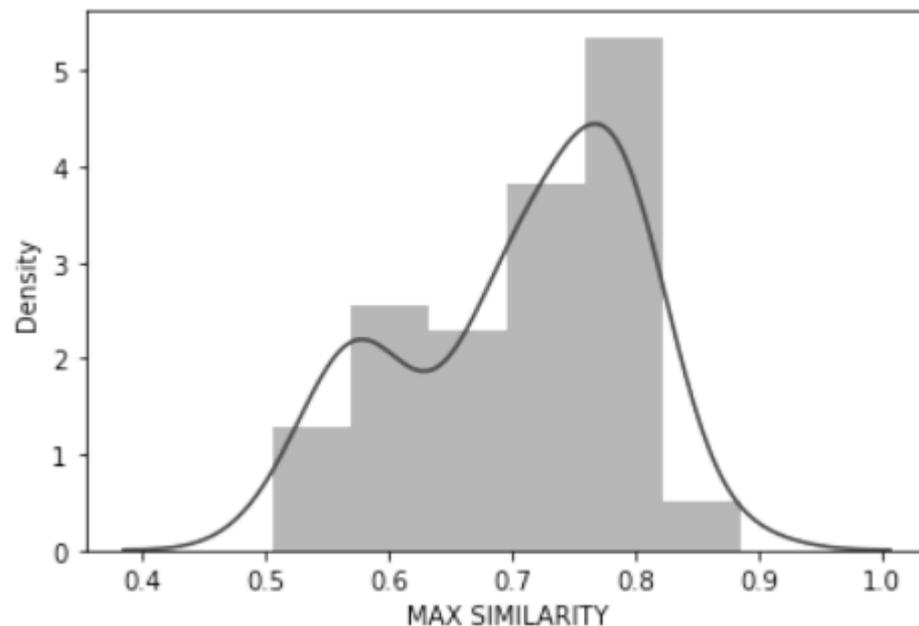
Test

Train



Max similarity: 0.73

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$



Classificazione

Fingerprint

Descrittori molecolari

0.56

Accuracy

0.55

0.53

F1 Score

0.48

[[20,25] [2,15]]

Confusion Matrix

[[21,24] [4,13]]

88%

Recall

76%

Regressione

Fingerprint

Descrittori molecolari

Descrittori scalati

SDEP

1.18

1.20

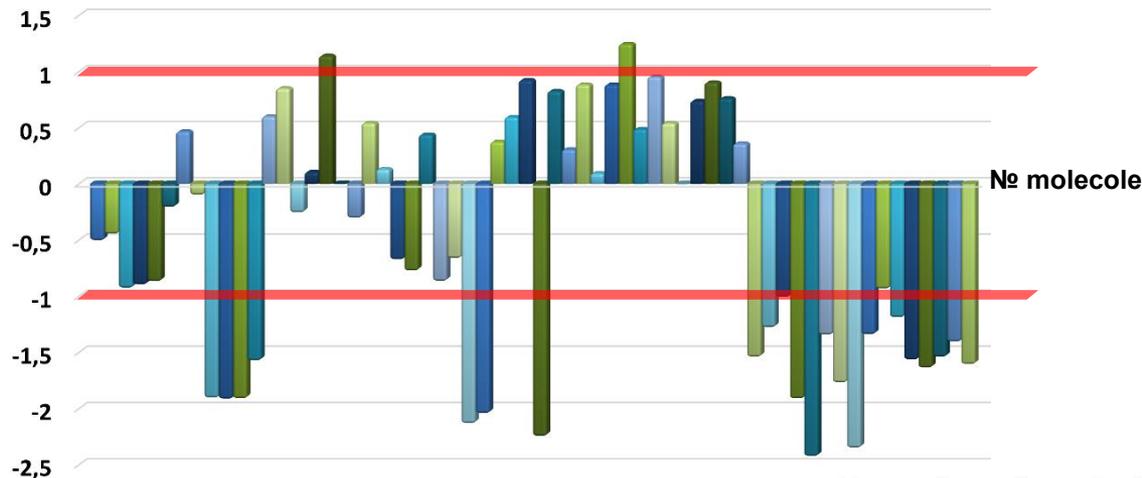
1.06



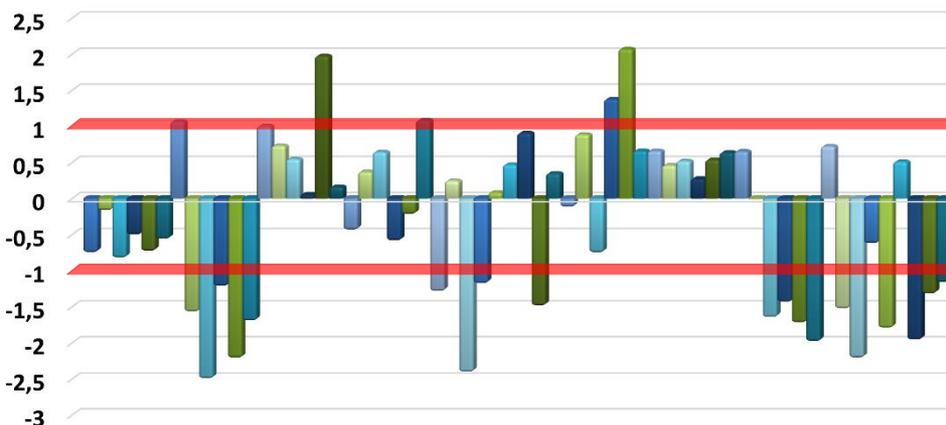
Validazione esterna regressione

Errori di Predizione

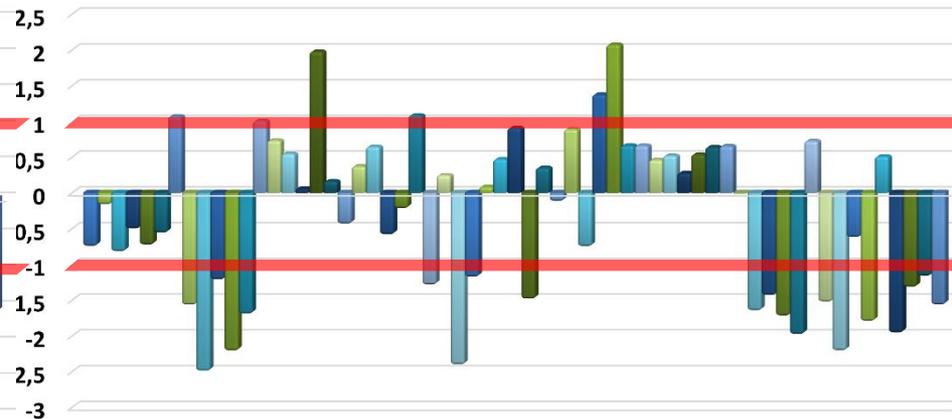
Fingerprint



Descrittori molecolari

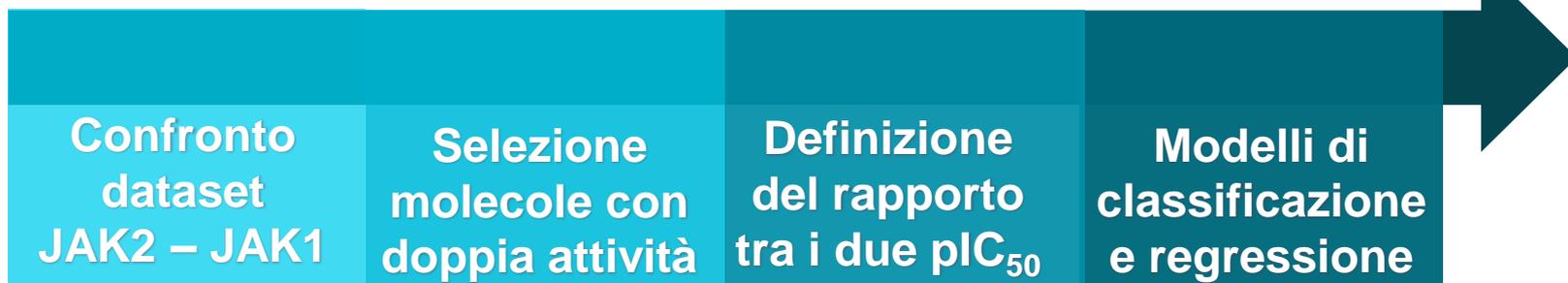


Descrittori scalati





Selettività JAK2 – JAK1



182 Molecole selettive JAK2
1756 Molecole selettive JAK1
944 Molecole non selettive

Modelli di classificazione e regressione

MCC

Q²

Fingerprint

0.90

0.80

Descrittori molecolari

0.90

0.79

Descrittori molecolari scalati

/

0.78

Virtual screening e sviluppi futuri

National Cancer Institute
252128 records

Features selection

Predizione attività

Consensus scoring

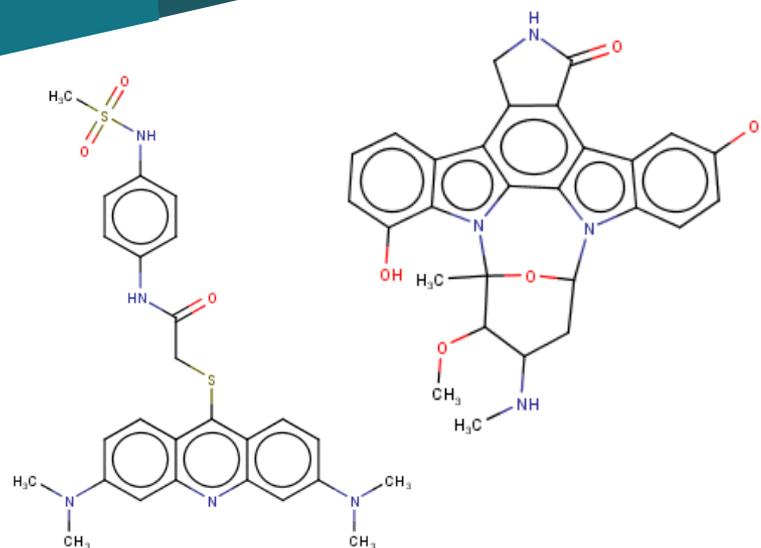
Selezione delle top 40 molecole più attive

Sviluppi futuri

Saggi biologici sulle top 40 molecole

Applicazione modelli elaborati a tutte le JAK

Individuare molecole che indicano la selettività verso tutte le altre JAK





RINGRAZIAMENTI

- Professor Rino Ragno
- Dottoressa Manuela Sabatino
- Dottore Lorenzo Antonini
- Dottoressa Giulia Fantera
- Dottoressa Paola Caprioli
- Beatrice Foti
- Maria Stella Delle Chiaie