

Имплементација софтвера за филтрирање нежељених порука употребом класификационих алгоритама

Проблем



Решења

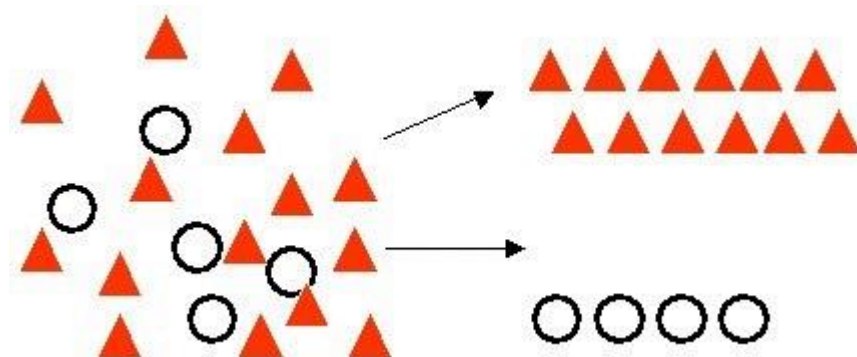
Нежељене мејл адресе



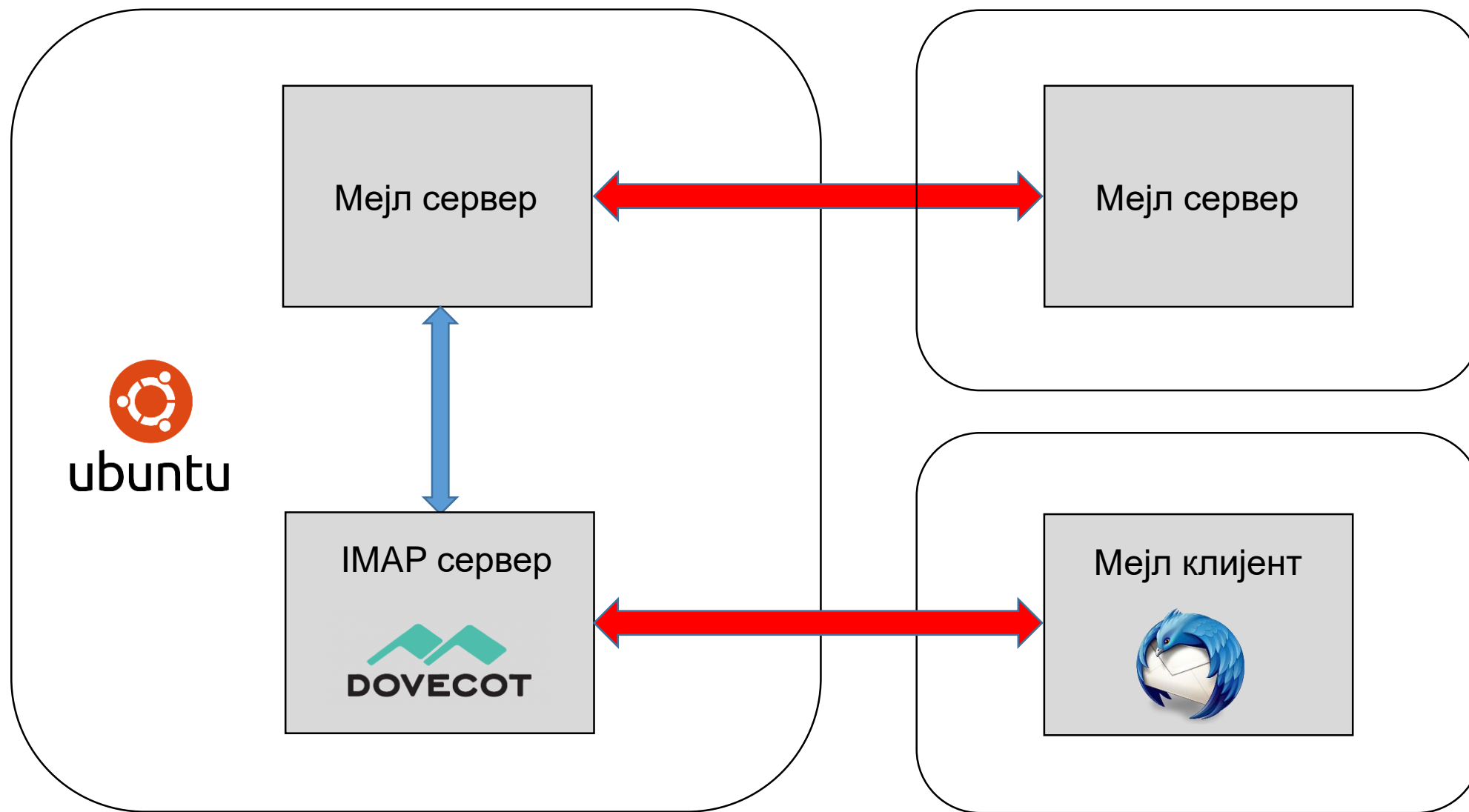
Нежељени мејл домени



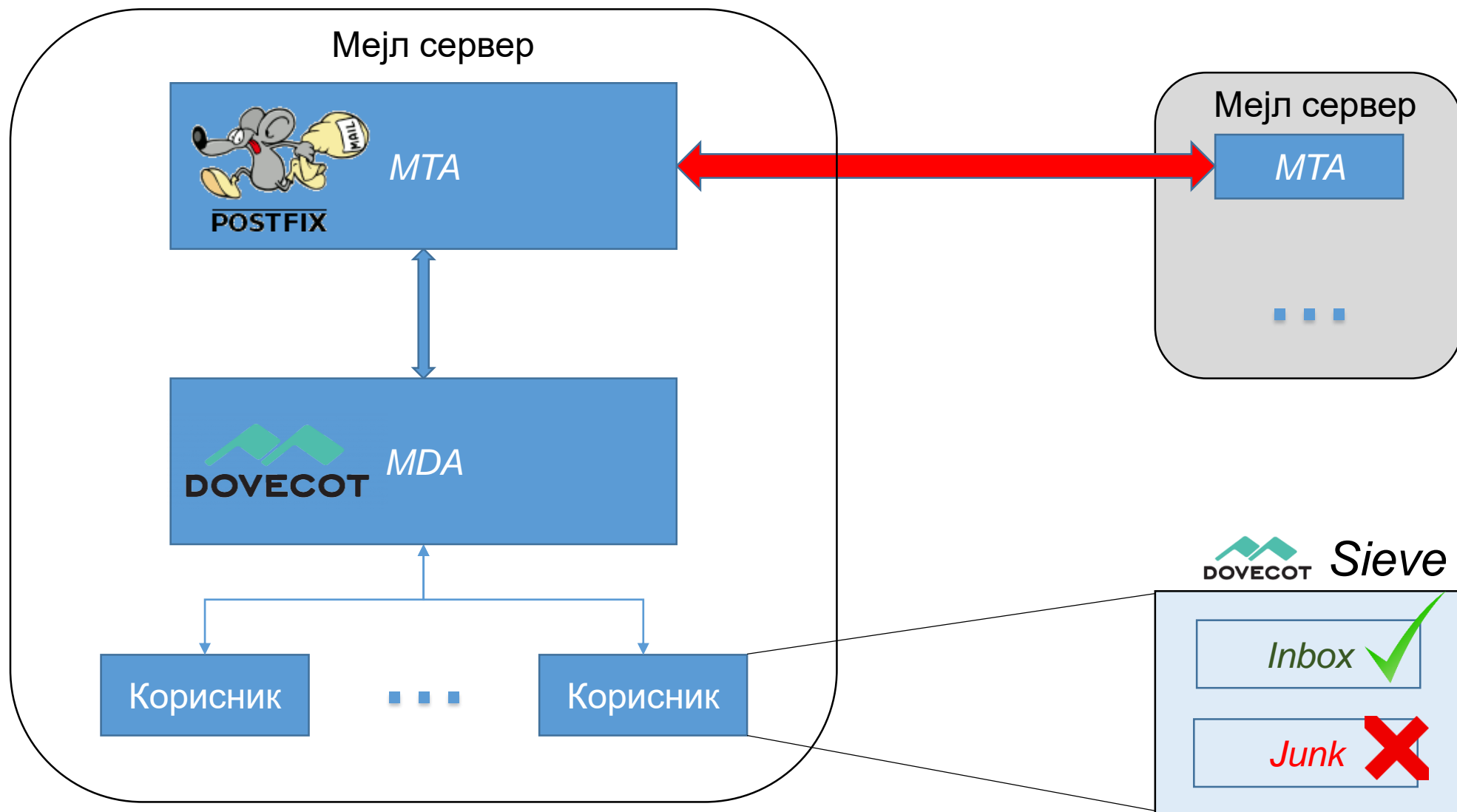
Класификација података



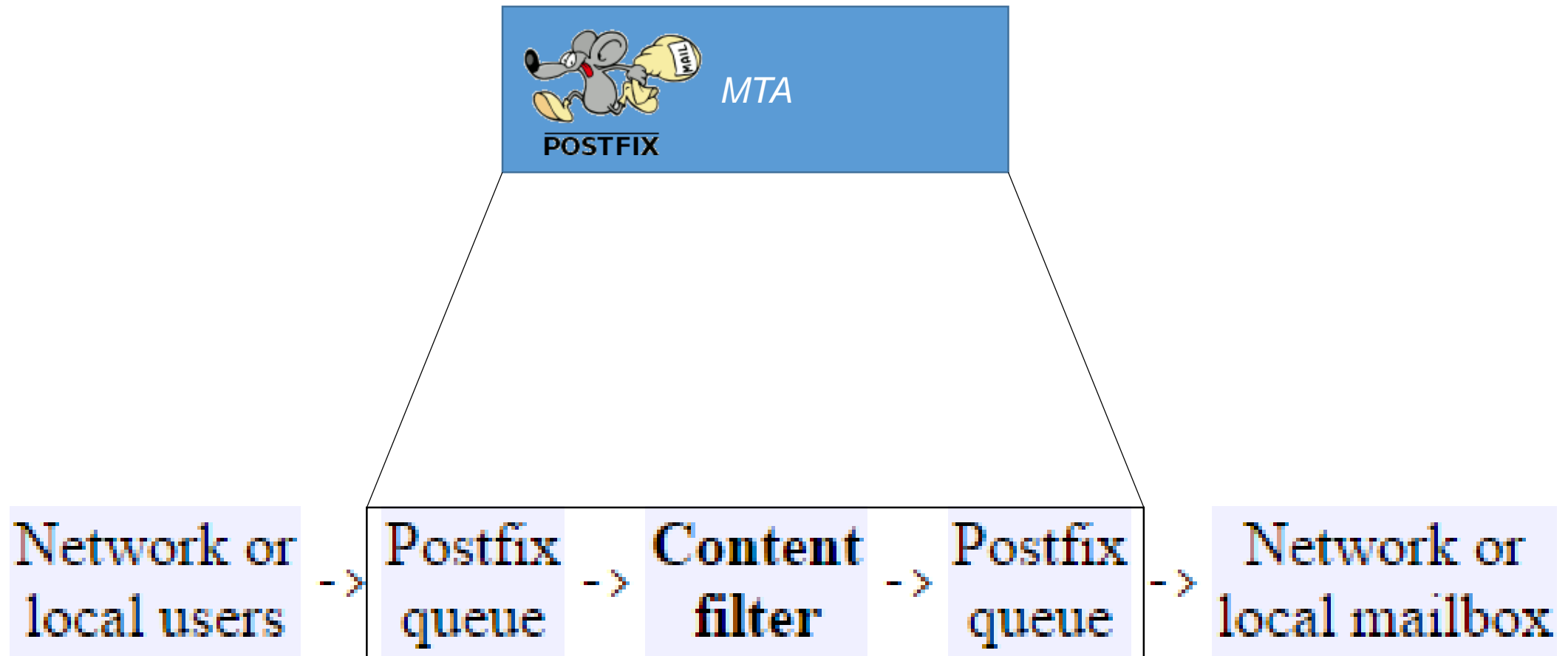
Платформа



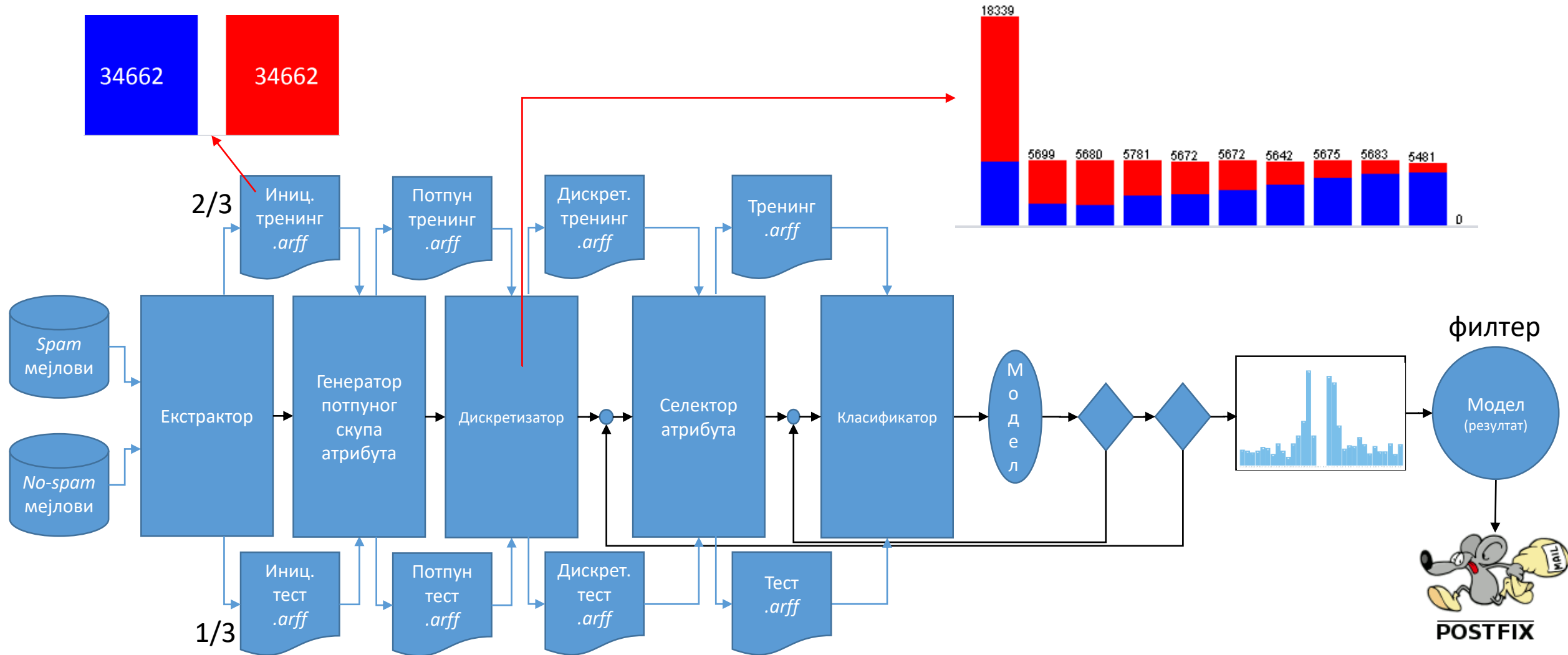
Мејл сервер



Филтер компонента



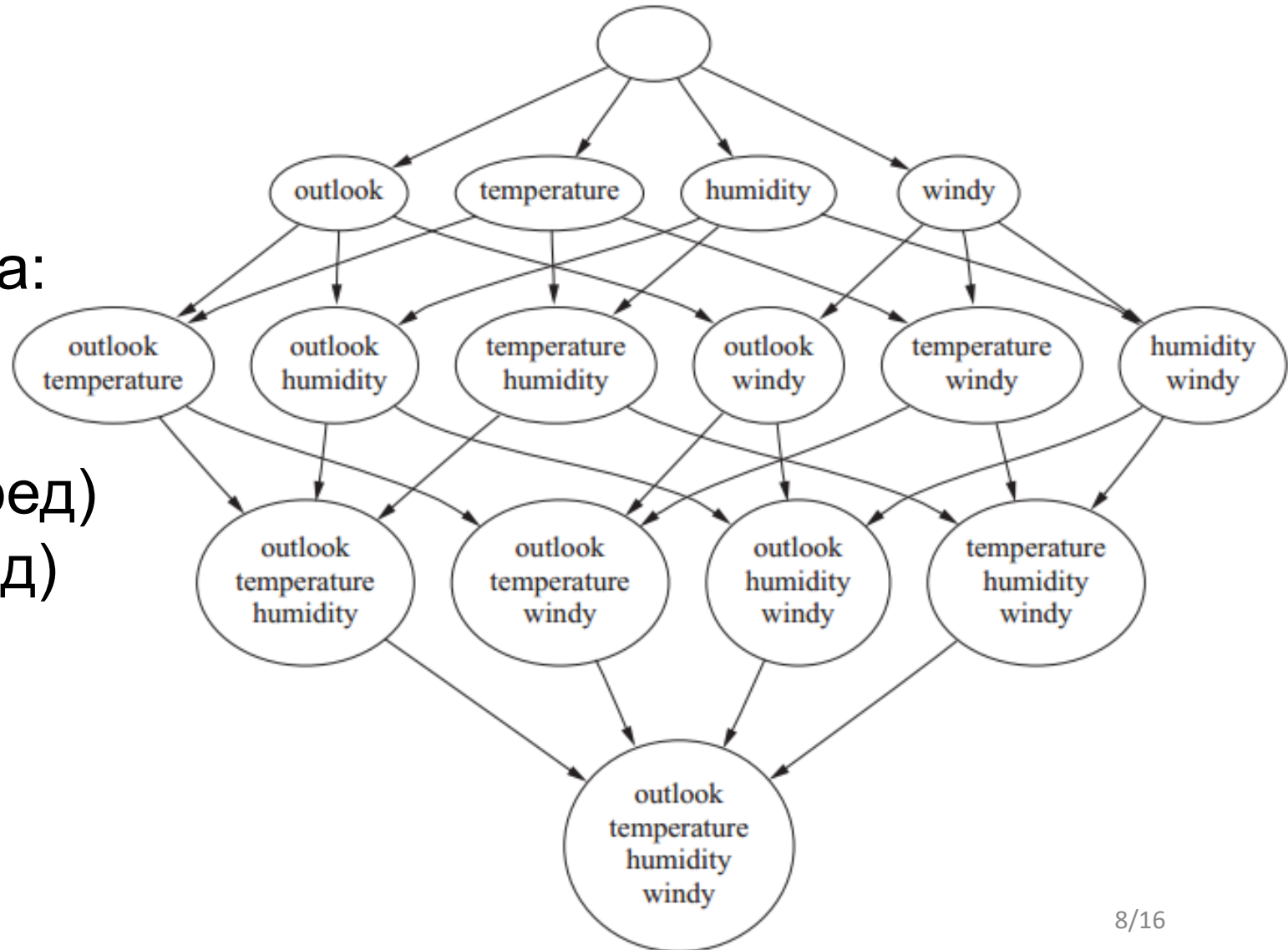
Филтер компонента



Алгоритми за претпроцесирање

Претрага скупа атрибута:

- Селекција унапред
- Елиминација уназад
- Прво најбољи (унапред)
- Прво најбољи (уназад)



Алгоритми за претпроцесирање

Евалуација скупа атрибута:

- Селекција својстава заснована на корелацији

$$\sum_j U(A_j, C) / \sqrt{\sum_i \sum_j U(A_i, A_j)}$$

Симетрична несигурност

$$U(A, B) = 2 \frac{H(A) + H(B) - H(A, B)}{H(A) + H(B)}$$

Алгоритми за претпроцесирање

Претрага скупа атрибута:

- Селекција унапред
- Елиминација уназад
- Прво најбољи (унапред)
- Прво најбољи (уназад)

Евалуација скупа атрибута:

- Селекција својстава заснована на корелацији

5

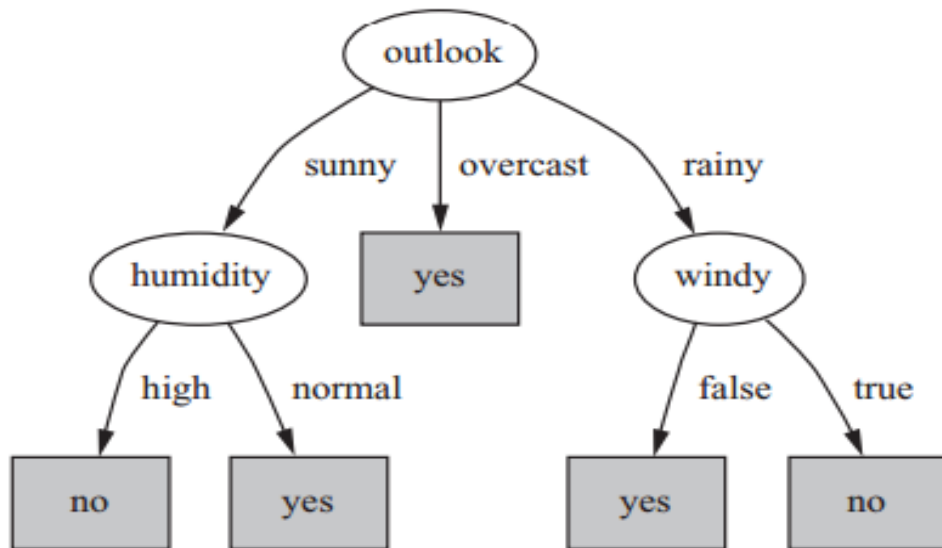
Рангирање атрибута

Евалуација атрибута:

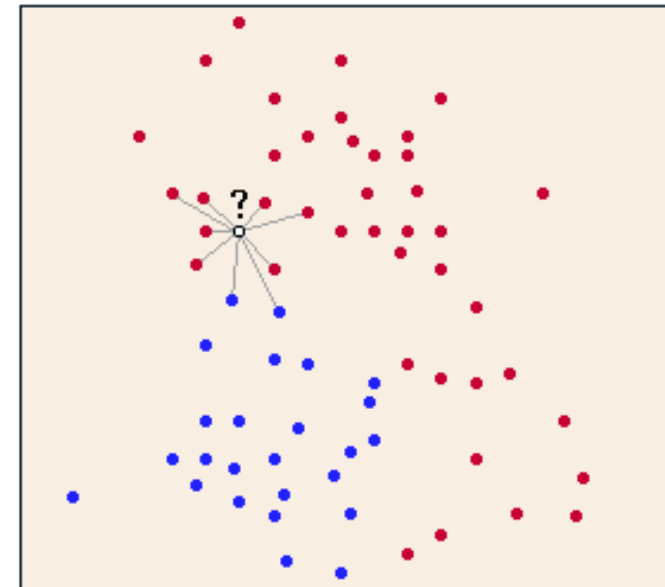
- Информациона добит

Алгоритми за класификацију

Стабло одлучивања (C4.5) **x4**



Учење базирано на инстанцама
(K најближих суседа) **x3**

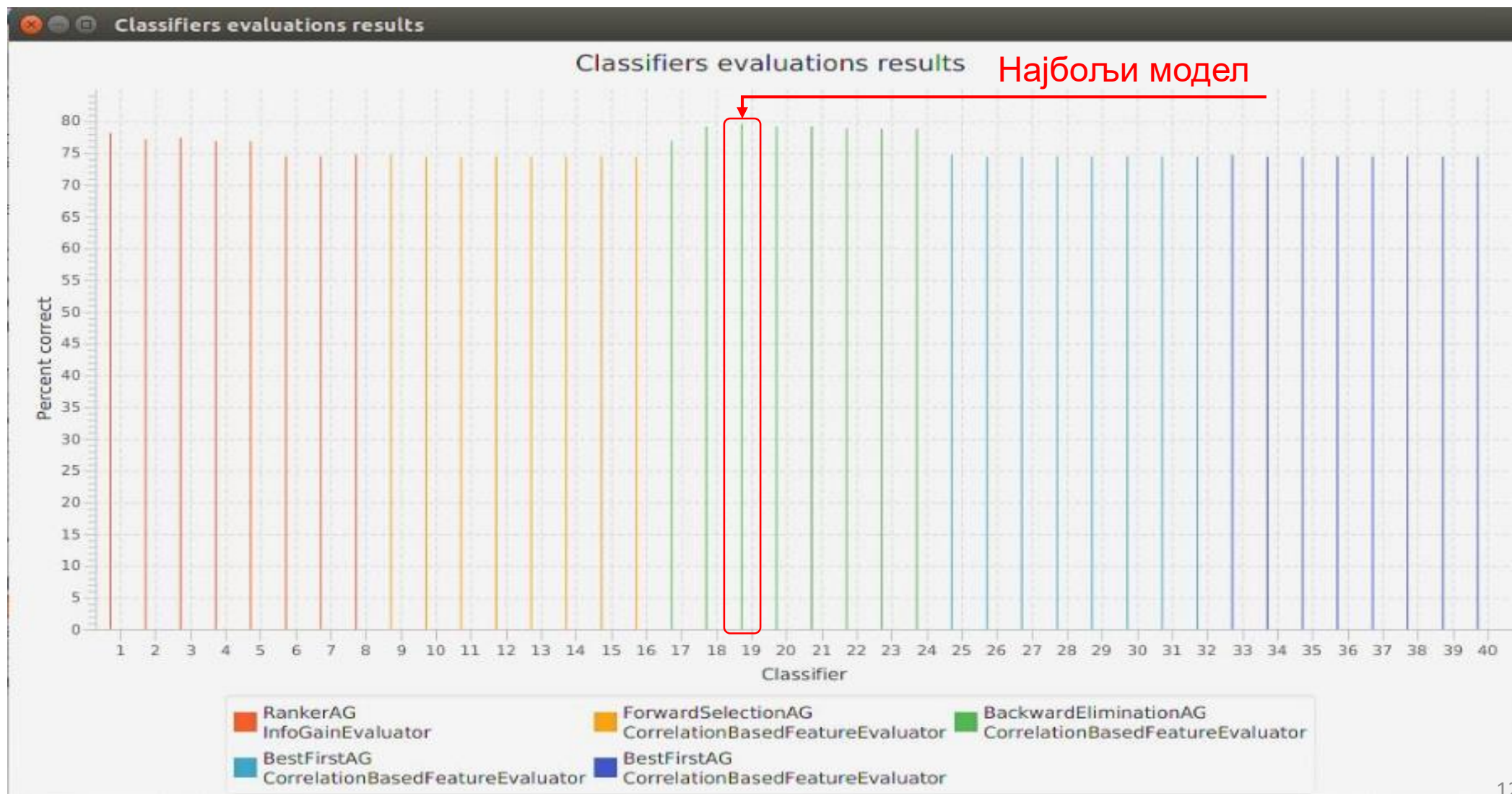


8

Наиван Bayes **x1**

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Резултат



Резултат

Најбољи модел:

- Елиминација уназад
- Стабло одлучивања (неорезано + уздизање подстабала)

The best attribute generator - classifier combination is:

Attribute generator: BackwardEliminationAG

Evaluator Name: CorrelationBasedFeatureEvaluator

Classifier Name: J48

=====

Correctly Classified Instances	21711	79.615	%
--------------------------------	-------	--------	---

Incorrectly Classified Instances	5559	20.385	%
----------------------------------	------	--------	---

Kappa statistic	0.5725
-----------------	--------

Mean absolute error	0.3023
---------------------	--------

Root mean squared error	0.3988
-------------------------	--------

Relative absolute error	60.4622	%
-------------------------	---------	---

Root relative squared error	79.7655	%
-----------------------------	---------	---

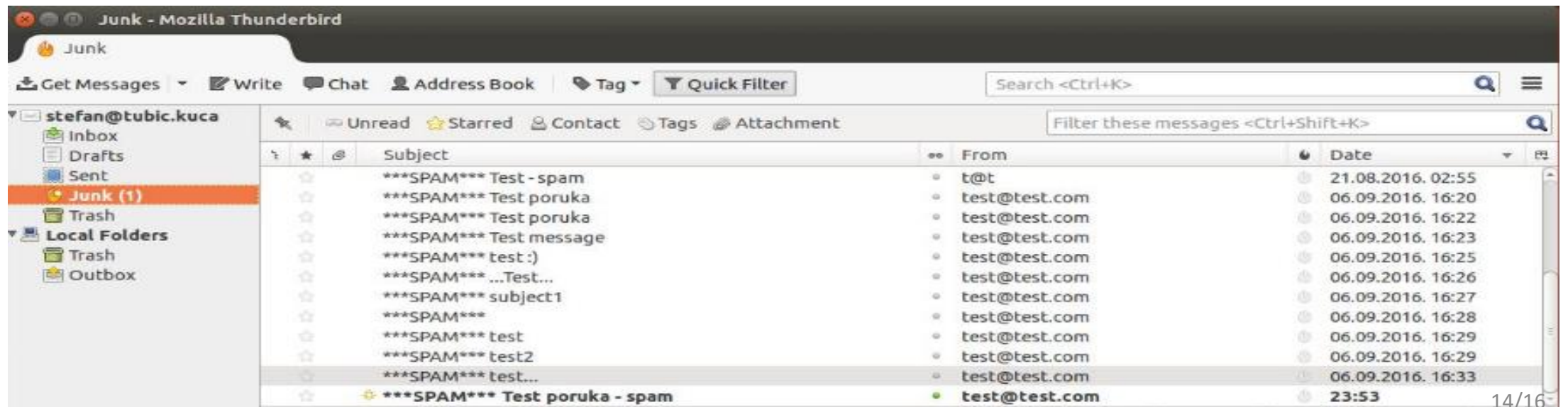
Total Number of Instances	27270
---------------------------	-------

=== Confusion Matrix ===

	a	b	<-- classified as
7887	2310		a = 0
3249	13824		b = 1

Резултат

```
root@stefan-LIFEB00K-E751:/home/stefan# telnet localhost 25
Trying 127.0.0.1...
Connected to localhost.
Escape character is '^]'.
220 tubic.kuca ESMTP Postfix (Ubuntu)
mail from: test@test.com
250 2.1.0 Ok
rcpt to: stefan@tubic.kuca
250 2.1.5 Ok
data
354 End data with <CR><LF>.<CR><LF>
subject: Test poruka - spam
BAIT, YOU HAVE BEEN APPROVED. CASH GRANT AMOUNT: $10,000-$5,000,000 DID YOU KNOW? -EACH YEAR THE U.S. GOVER
MENT GIVES AWAY BILLIONS IN CASH GRANTS? -THERE ARE NO SPECIAL REQUIREMENTS TO OBTAIN THESE GRANTS. -THESE
ARE FREE CASH GRANTS THAT YOU NEVER HAVE TO REPAY! BAIT,YOU QUALIFY! CLICK HERE LIMITED TIME OFFER
.
250 2.0.0 Ok: queued as 1EE5628032F
```



Закључак

- Додатак (*plugging*) за постојећи мејл сервер
- Примена алгоритама машинског учења за класификацију порука

Унапређења:

- Повећање фокуса на заглавље мејл порука
- Повећање спектра алгоритама за претпроцесирање и класификацију

Хвала на пажњи!