

# Detecting Packed Executable Files

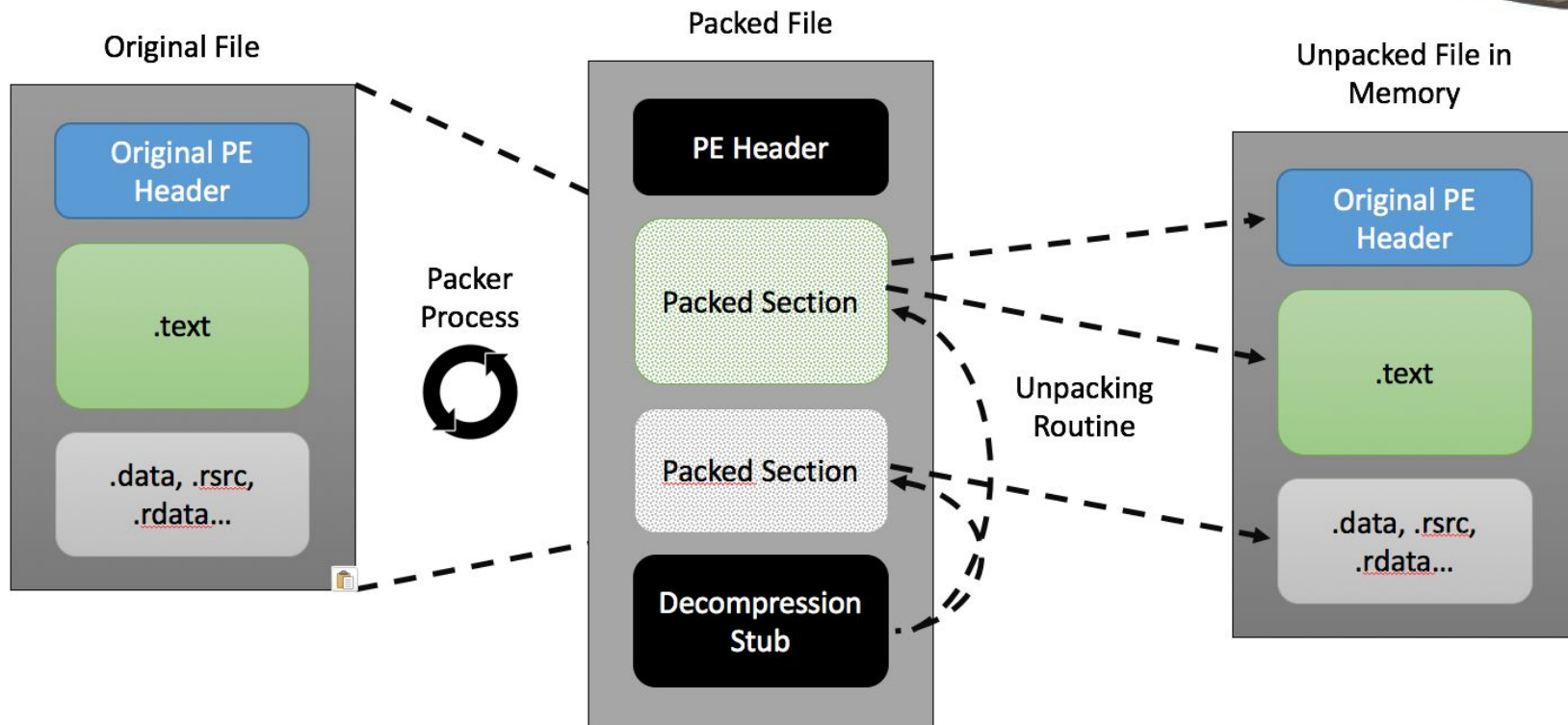
Projekt za kolegij Strojno učenje

Ivan Čeh

Sebastijan Horvat

PMF, Zagreb, 27. lipnja 2018.

# 0 problemu



# Skup podataka

- 43784 izvršivih datoteka i pripadnih izvještaja (označenih)
- u sklopu Mozgala, još 6928 (neoznačenih)

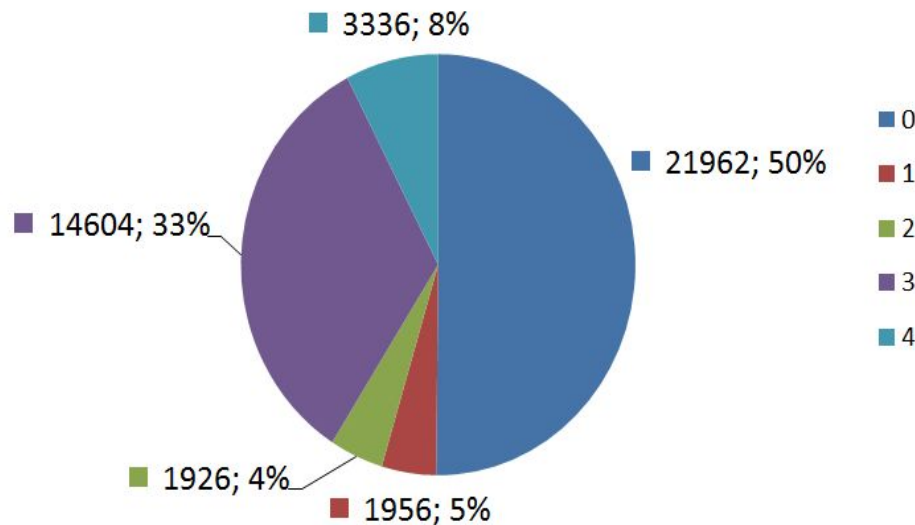


mozgalo

REVERSING™  
LABS

Oznaka	Nepakirano ili vrsta pakiranja
0	nije pakirana
1	pakirana - <i>Overlay</i>
2	pakirana - <i>Crypter</i>
3	pakirana - <i>Compressing packer</i>
4	pakirana - <i>Protector</i>

Broj datoteka



# Kako smo tražili značajke?

# Literatura

# Proučavanje izvještaja

## Classification of Packed Executables for Accurate Computer

Roberto Perdisci<sup>a,\*</sup>, Andrea I

\*Damballa, Inc., Atlanta.

<sup>b</sup>Georgia Tech Information Security Center

*<sup>c</sup>Dipartimento di Informatica e Comunicazione  
Milano, Ita*

# Using Entropy Analysis to Find Encrypted and Packed Malware

In statically analyzing large sample collections, packed and encrypted malware pose a significant challenge to automating the identification of malware attributes and functionality. Entropy analysis examines the statistical variation in malware executables, enabling analysts to quickly and efficiently identify packed and encrypted samples.



## Abstract

Executable packing is the most common technique to obfuscate malicious code and evade detection. Unpackers have been proposed that can detect packed executables, therefore potentially revealing code hidden by traditional signature-based anti-virus engines. However, these unpackers are computationally expensive and searching for virus infections may take several hours.

In this paper we apply pattern recognition to executables. The objective is to efficiently and compactly represent packed and non-packed executables, so that only one universal unpacker is needed to unpack to an universal unpacker, thus saving a significant amount of space. We will show that our system achieves very high detection rates with a low average processing time.

**Key words:** Computer Security, Pattern Recognition, Computer Virus Detection.

\* Corresponding author.

Email addresses: roberto.perdischi@gmail.com  
andrew@idea.sec.dico.unimi.it (Andrea L.  
(Wenke Lee).

Preprint submitted to Elsevier

James  
McDonald  
Bradley

**M**alware analysts often use encryption or packing (‘compression’) methods to conceal their malicious ‘cousins’<sup>1</sup>’s erasing data and code. These methods—which transform some or all of the original bytes into a series of random-looking data bytes—appear in 80 to 90 percent of malware samples.<sup>2</sup> This fact creates special challenges for analysts who use static methods to analyze large malware collections, as they must quickly and efficiently identify the samples and unpack or decrypt them before analysis can begin. Many tools, including the packing tools themselves, are generally successful at automatically unpacking or decrypting files, but they are often sensitive to file size. Often, analysts use tools to reorganize and remove the transformation subroutines or functions from the

areas. Specifically, it sums the frequency of each observed fixed-length data, and uses a formula to generate a trend to correlate the data with the season. Furthermore, it packs and encodes the transformed data using a binary tree structure. Our methodology for training BinaryTrend uses the training data to derive statistical parameters for each season, and then uses these parameters to generate a trend for each season.

s. Higher entropy scores

```

    "address": "0x00000000",
    "size": "0x00000000"
  },
  {
    "address": "0x00000000",
    "size": "0x00000000"
  }
},
"sections": {
  "section_list": [
    {
      "name": ".text",
      "flags": {
        "flag_list": [
          "IMAGE_SCN_CNT_CODE",
          "IMAGE_SCN_MEM_EXECUTE",
          "IMAGE_SCN_MEM_READ"
        ],
        "size": "0x00001800",
        "address": "0x00001000",
        "offset": "0x00000400",
        "entropy": "5.887627060638998"
      },
    },
    {
      "name": ".rdata",
      "flags": {
        "flag_list": [
          "IMAGE_SCN_CNT_INITIALIZED_DATA",
          "IMAGE_SCN_MEM_READ"
        ],
        "size": "0x00000140",
        "address": "0x00000000",
        "offset": "0x00000000",
        "entropy": "5.148888888888889"
      },
    },
    {
      "name": ".data",
      "flags": {
        "flag_list": [
          "IMAGE_SCN_CNT_INITIALIZED_DATA",
          "IMAGE_SCN_MEM_READ",
          "IMAGE_SCN_MEM_WRITE"
        ],
        "size": "0x00000000",
        "address": "0x00000000",
        "offset": "0x00000000",
        "entropy": "0.000000000000000"
      },
    }
  ]
}

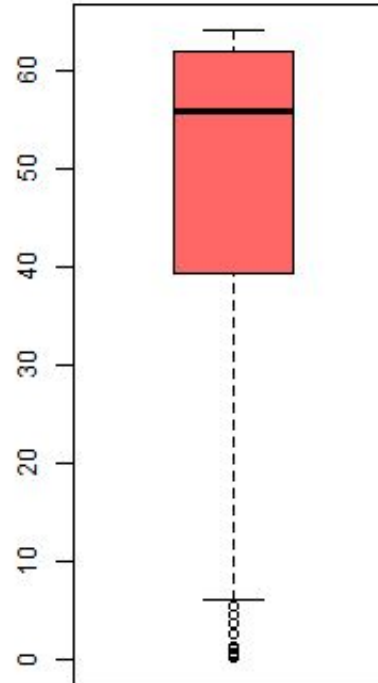
```

# Težinska sredina kvadrata entropija sekcija

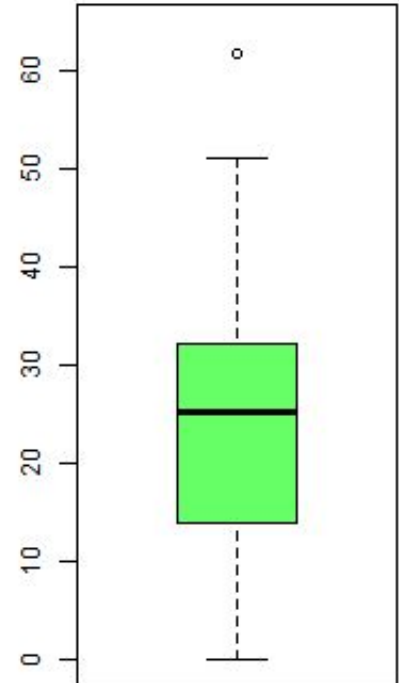
		veličina	entropija
Section 1	→	0x1800	5.8876
Section 2	→	0x1400	5.1488
Section 3	→	0x200	0.4524
Section 4	→	0x200	3.7829
Section 5	→	0x200	0.7613

27.8467

**Pakirane**



**Nepakirane**



# Nestandardna imena sekcija

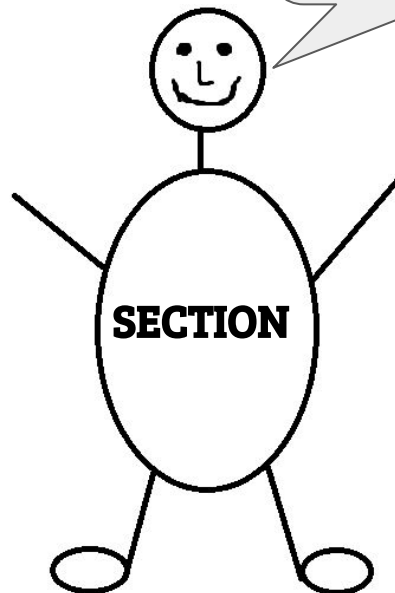
<b>.text</b>	<b>.reloc</b>
<b>.data</b>	<b>.sdata</b>
<b>.rdata</b>	<b>.srdata</b>
<b>.idata</b>	<b>.pdata</b>
<b>.edata</b>	<b>.debug\$\$</b>
<b>.rsrc</b>	<b>.debug\$T</b>
<b>.bss</b>	<b>.debug\$P</b>
<b>.crt</b>	<b>.drectve</b>
<b>.tls</b>	<b>.didat</b>

+

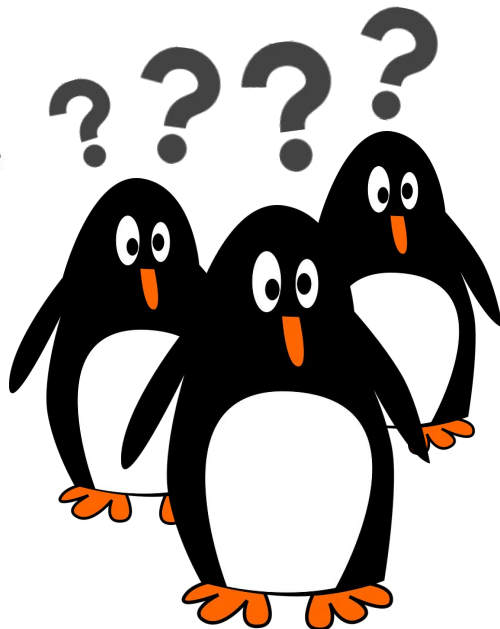
**.imports**



**VS**



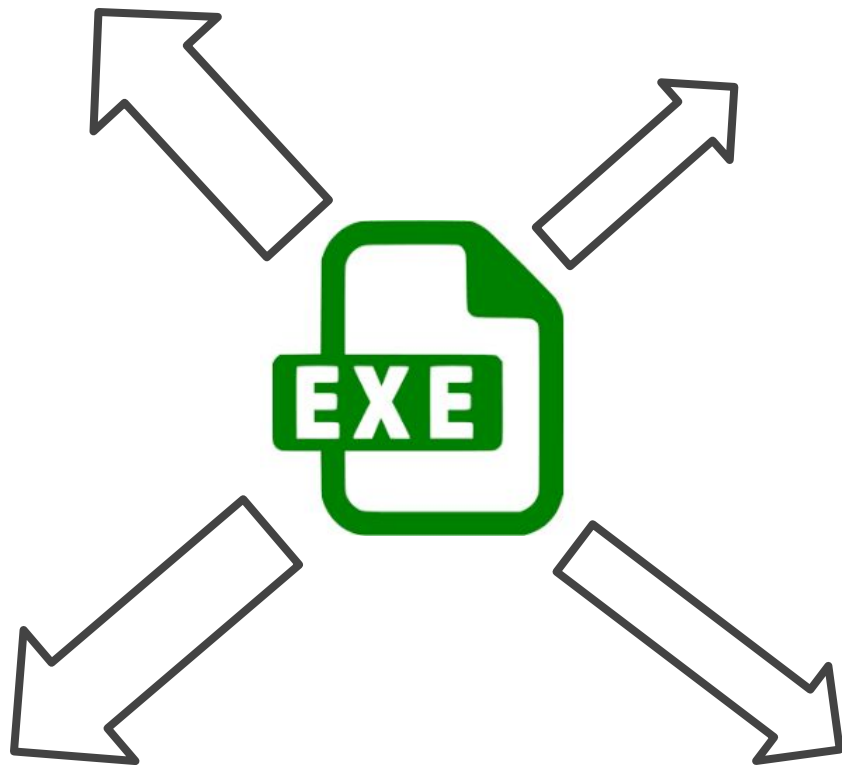
Ja sam  
Pingvin.



# Flagovi



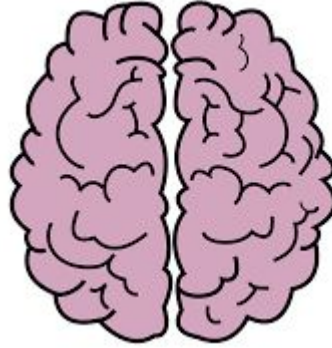
# API-ji



# Složenije značajke

**sizeofInitializedData**  
(suma po svim Entrijima)

**sizeofUninitializedData**  
(suma po svim Entrijima)



$$\log \left( \frac{1 + \textit{initalized}}{1 + \textit{uninitialized}} \right)$$

entropy (section#1)

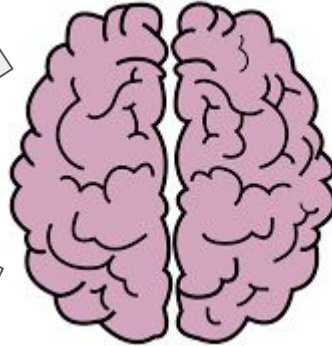
⋮

entropy (section#n)

size (section#1)

⋮

size (section#n)

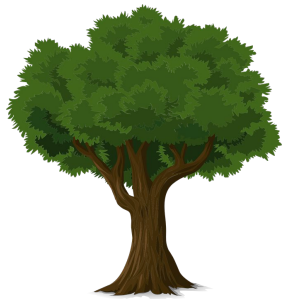


$$\max_{s \text{ je sekcija}} \left( \frac{\text{size}(s)}{\text{entropy}^2(s)} \right)$$



# SLUČAJNA ŠUMA





× 800



**RANDOM  
STATE**

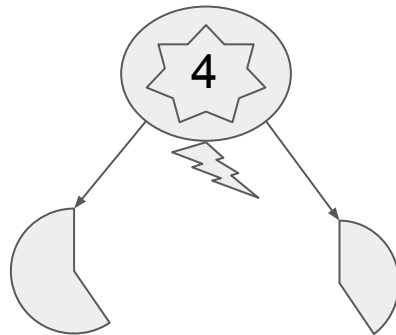
# PARAMETRI MODELA

**BOOTSTRAP = False**



Ponavljáš  
se!

**MIN. SAMPLES SPLIT**



**CRITERION**

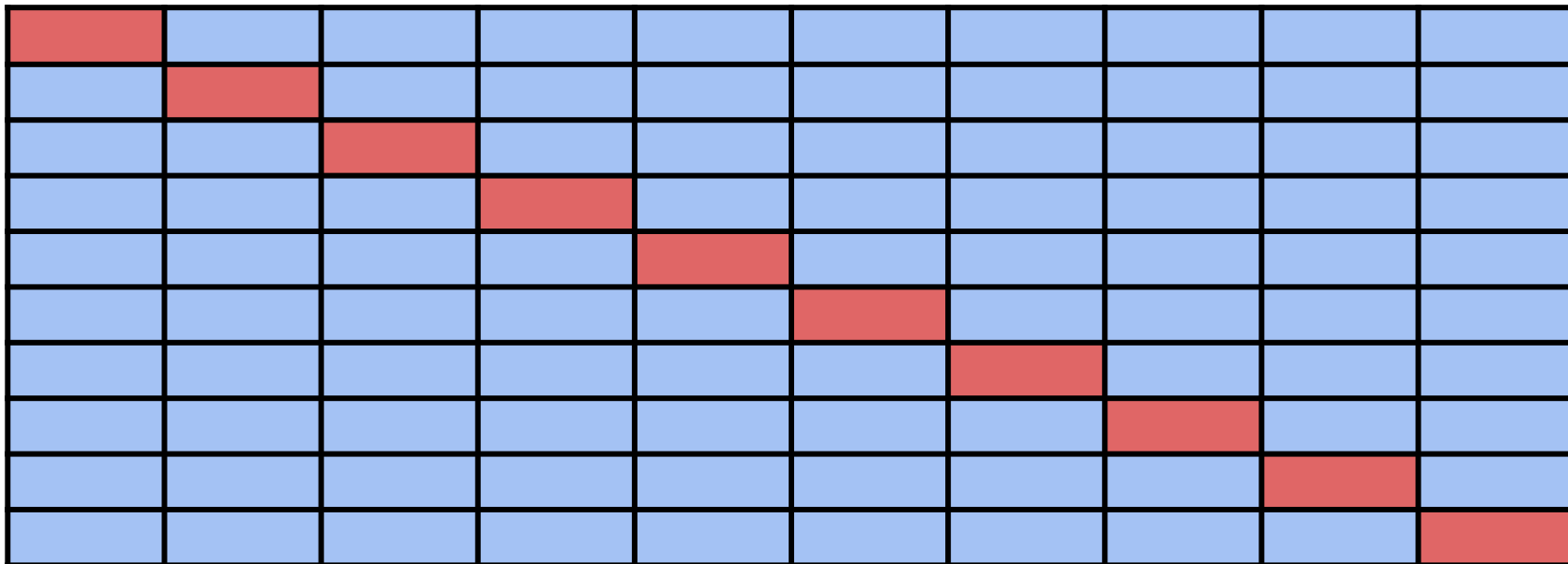
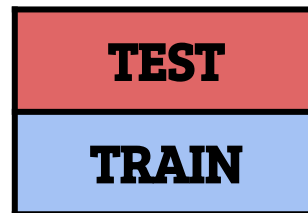
*gini*

*entropy*

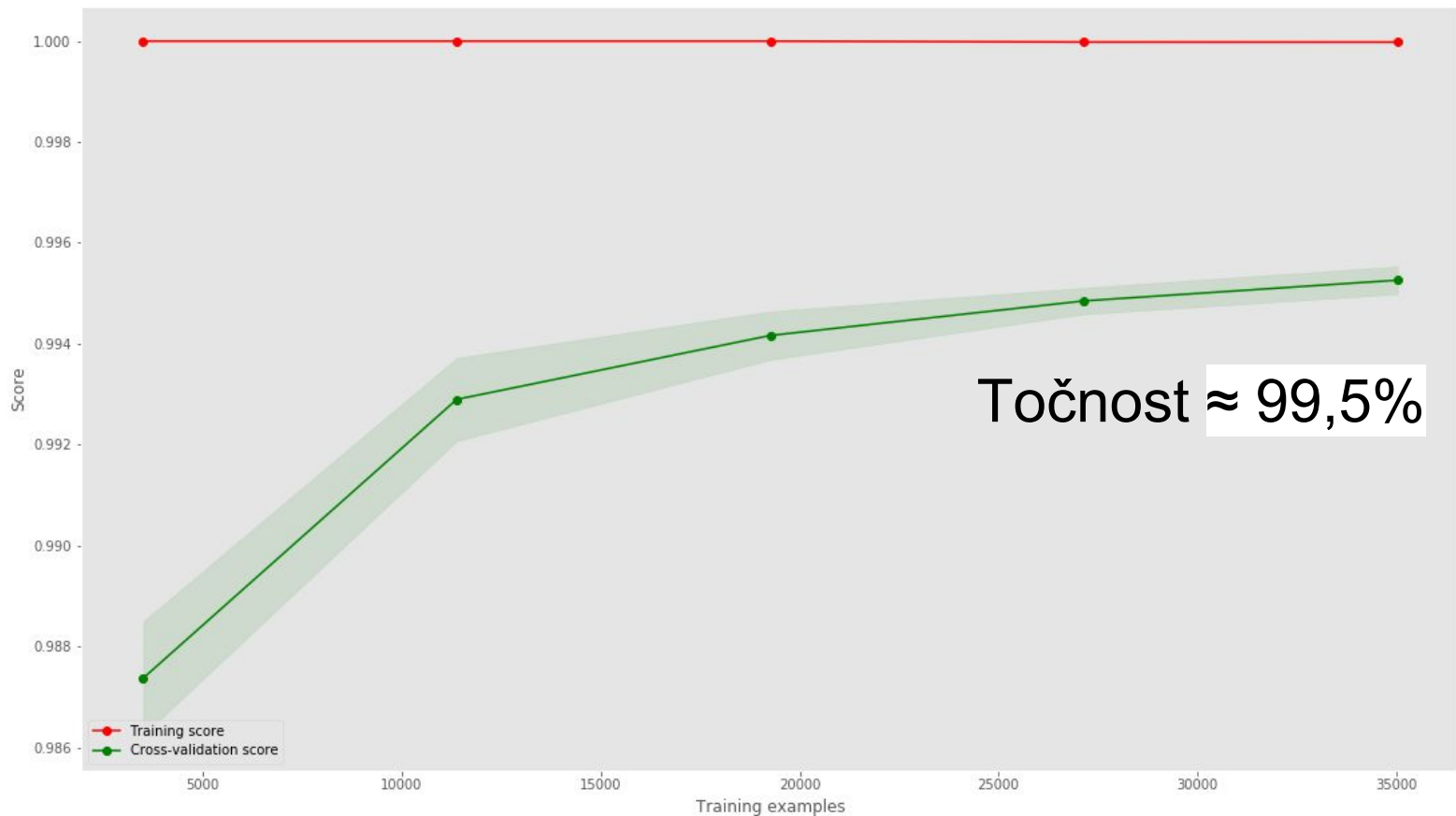
# 10-fold



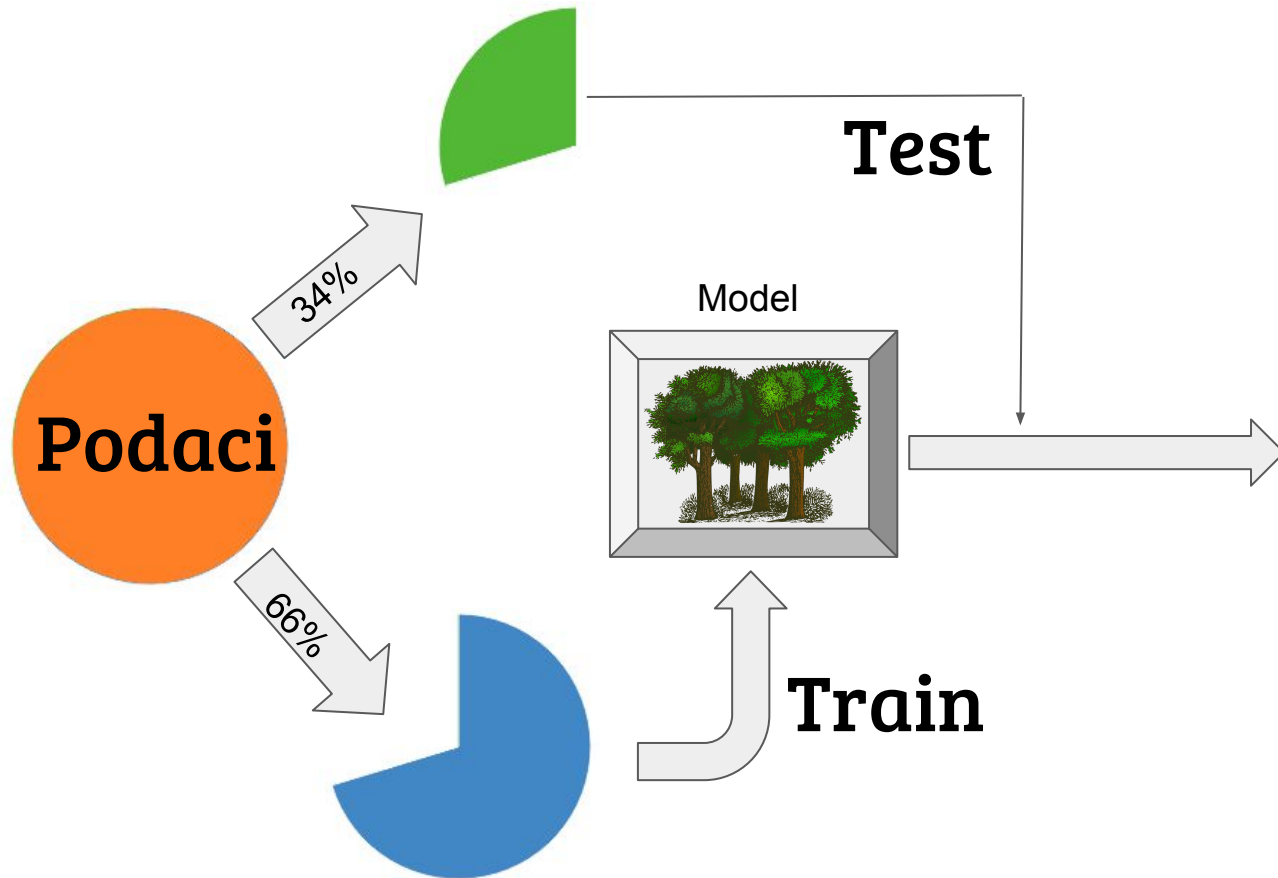
43784



# KRIVULJA UČENJA

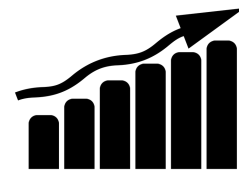


# Jedna matrica konfuzije



prediction	False	True
original		
t0	7371	18
t1	33	640
t2	10	688
t3	5	5013
t4	2	1107

# VRIJEME IZVRŠAVANJA



Učitavanje podataka:  
12 minuta

Trening: 1 minuta

Predikcija (14 000 datoteka):  
5 sekundi



Windows 8.1, Pentium Dual-Core  
2.00 GHz, 32b, 4GB RAM



# Rezultati natjecanja Mozgalo 2018.



6928 testnih primjeraka

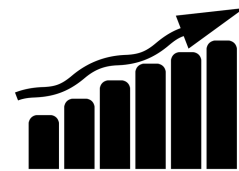
točnost: 97.97921%

točnost na neviđenim primjerima:

$536/556 = 96.40288\%$

prediction		False	True
original			
t0	3399	63	
t1	6	479	
t2	46	625	
t3	3	1147	
t4	22	1138	

# Rezultati natjecanja Mozgalo 2018.



Naziv tima	Točnost (%)	FP+FN
Laganica ③	98.239	122
NewTeam (mi) ②	97.979	140
Analysis Paralysis	97.546	170
GMO Lazo ①	97.373	182
505	95.901	284



# Pitanja?

