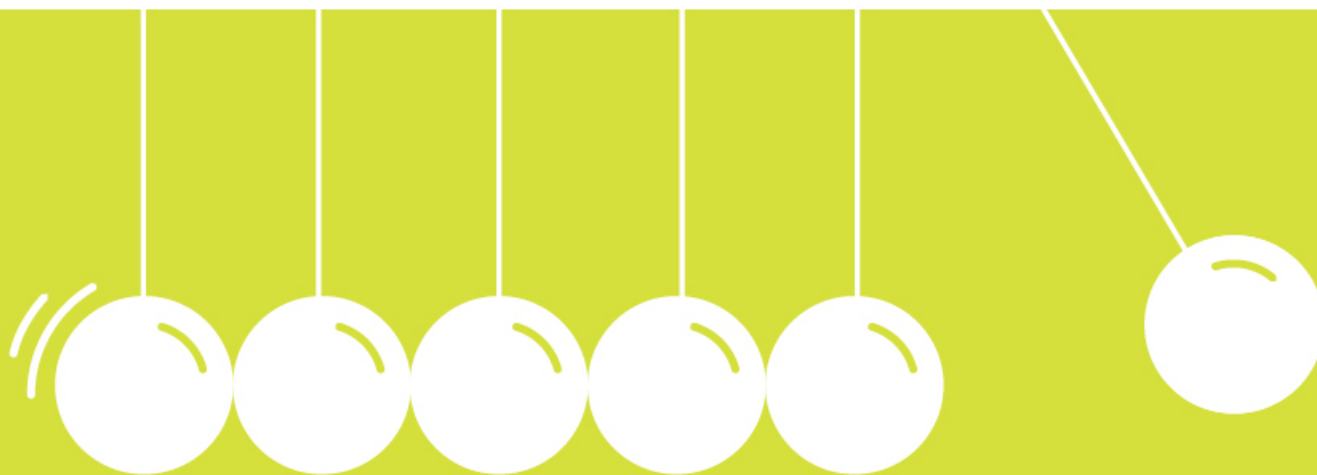




# Metoda *NoiseRank* za odkrivanje anomalij v podatkih



**Borut Sluban, *Institut „Jožef Stefan“***

# Anomalije?



# Anomalije?



# Anomalije so ...



- **Napake v podatkih - šum**

**živali bele barve:**



# Anomalije so ...



- **Napake v podatkih - šum**

**živali bele barve:**



- **Izjeme ali osamelci v podatkih**

**čreda ovc:**



# Namen dela



- **Razvoj metodologije za odkrivanje anomalij v podatkih**

# Namen dela



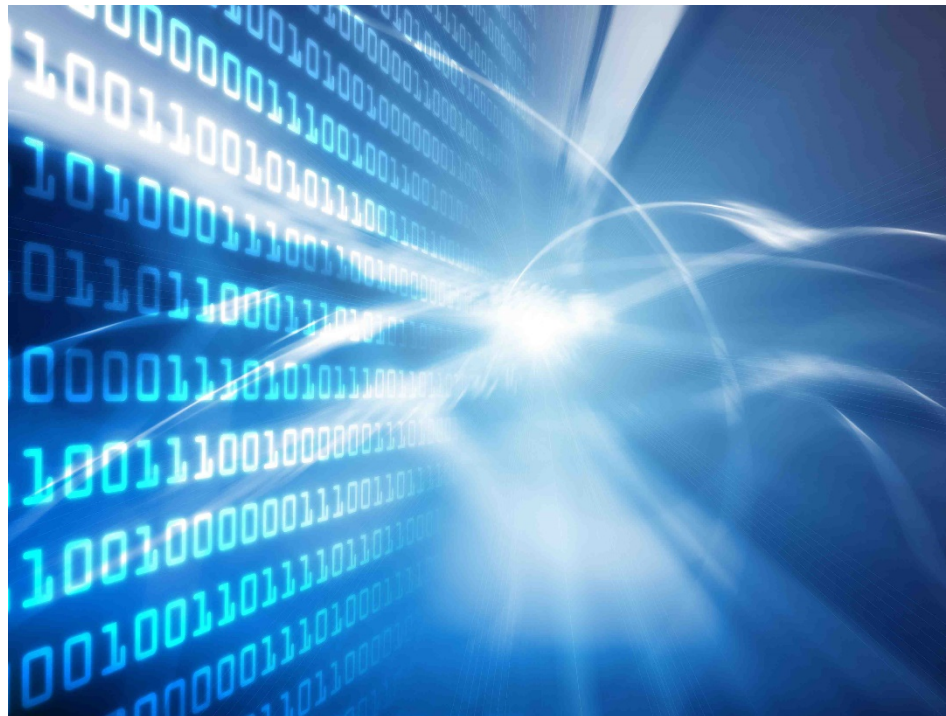
- **Razvoj metodologije za odkrivanje anomalij v podatkih**
- **za potrebe:**
  - **čiščenje podatkov**
  - **razumevanje podatkov / domene**
  - **odkrivanje novih posebnih primerov v podatkih**

# Podatki?

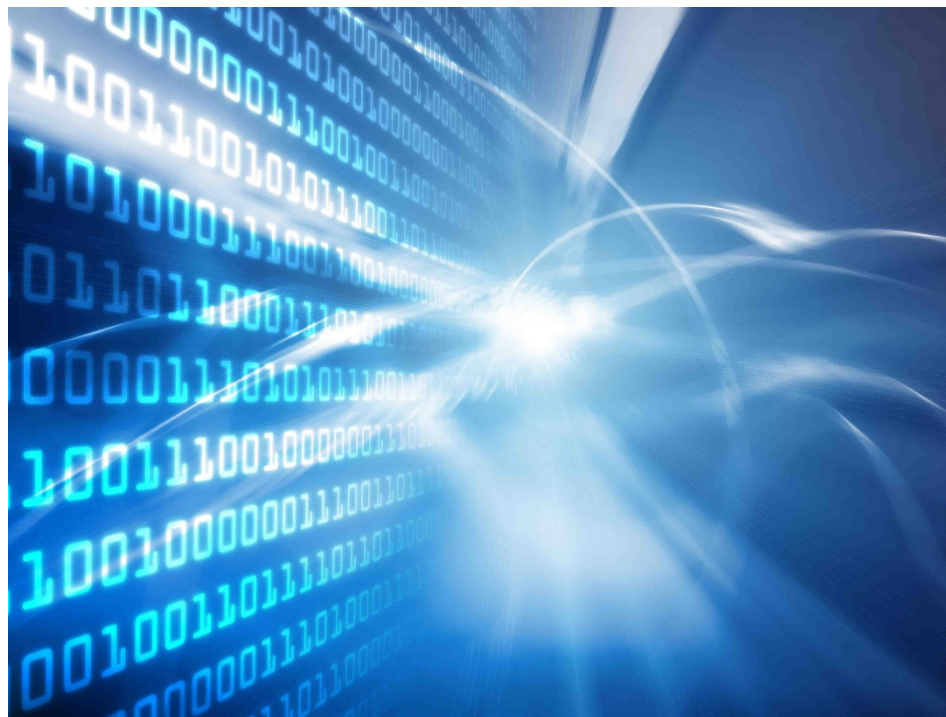




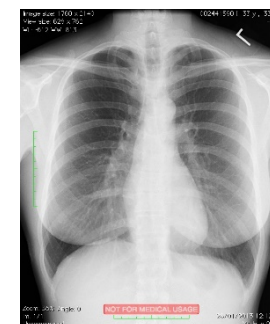
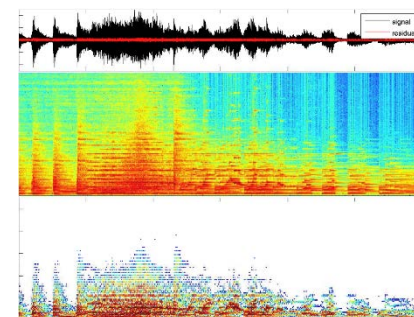
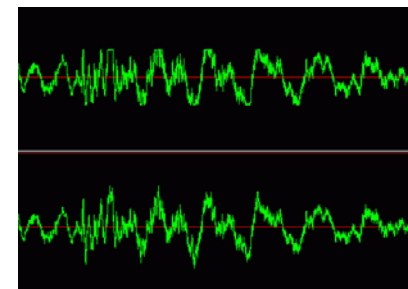
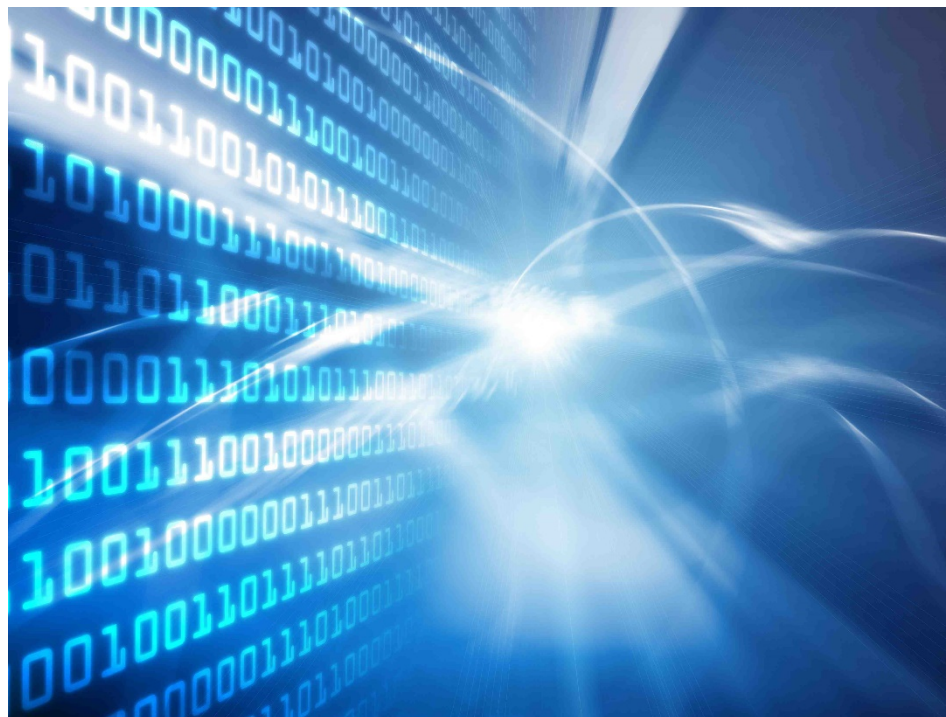
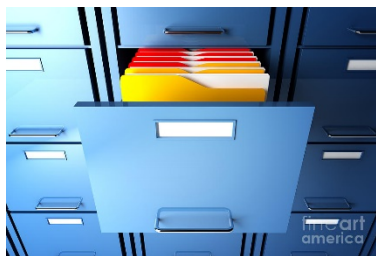
# Podatki?



# Podatki?



# Podatki?



# Podatki



Ime	Starost	Poklic	Prosti čas	Najljubša jed
Janko	9	šolar/ka	Sprehodi po gozdu	sladkarije
Metka	7	šolar/ka	Sprehodi po gozdu	sladkarije
Sneguljčica	17	manekenka	Sprehodi po gozdu	rdeče jabolko
Volk	467	gozdar	Sprehodi po gozdu	Rdeča kapica
Špicparkelj	?	škrat	Sprehodi po gozdu	borovnice

# Odkrivanje anomalij



# Odkrivanje anomalij



- **Stvari v naravi:**
  - so urejene po določenih vzorcih



# Odkrivanje anomalij



- **Stvari v naravi:**
  - so urejene po določenih vzorcih
  - sledijo fizikalnim zakonom

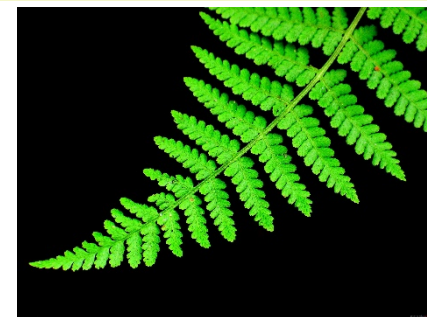


# Odkrivanje anomalij



- **Stvari v naravi:**

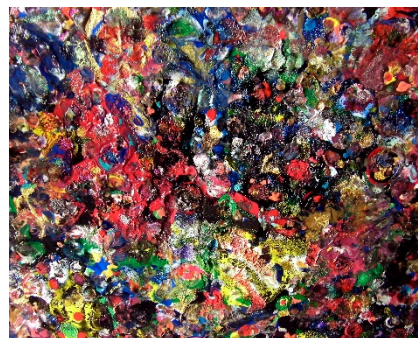
- so urejene po določenih vzorcih



- sledijo fizikalnim zakonom



- niso naključne





# Odkrivanje anomalij



- **Napake ali izjeme zaznamo kot:**

# Odkrivanje anomalij



- **Napake ali izjeme zaznamo kot:**
  - neskladja z ustaljenimi vzorci



# Odkrivanje anomalij

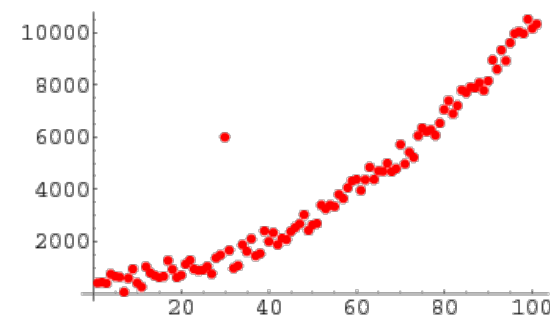


- **Napake ali izjeme zaznamo kot:**

- neskladja z ustaljenimi vzorci



- velika odstopanja od pričakovanih vrednosti



# Odkrivanje anomalij



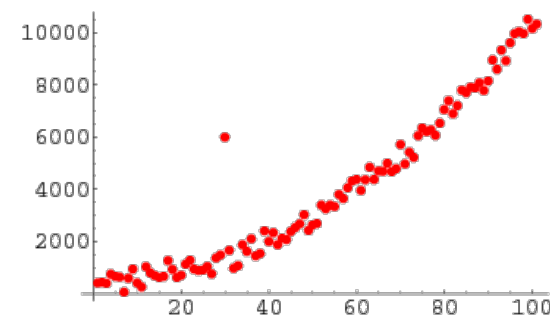
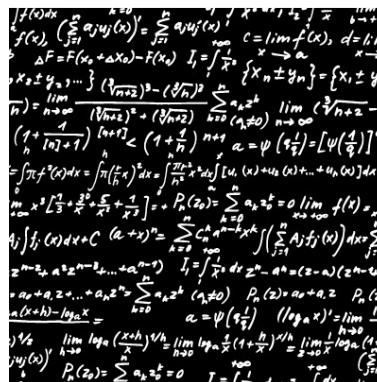
- **Napake ali izjeme zaznamo kot:**

- neskladja z ustaljenimi vzorci



- velika odstopanja od pričakovanih vrednosti

- težko opisljive





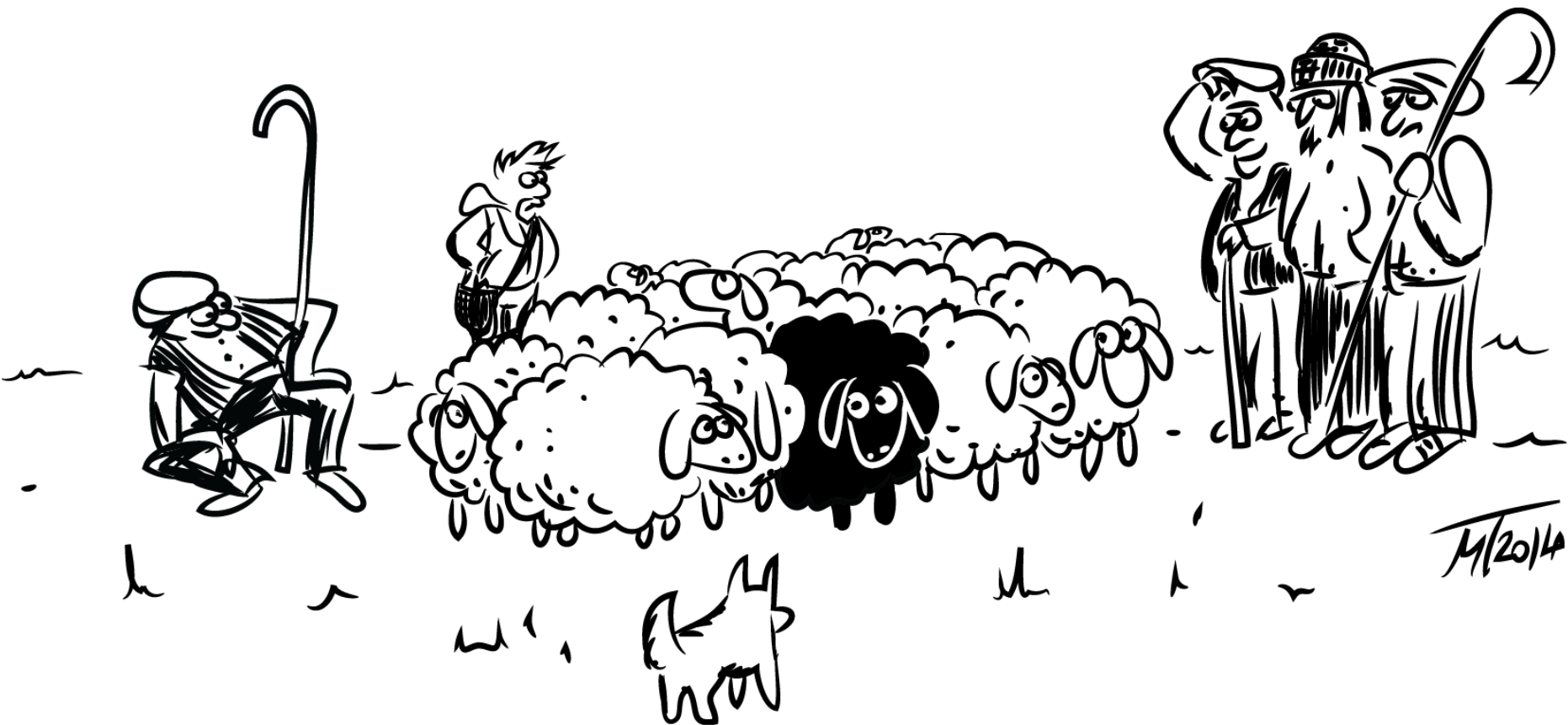








# Odkrivanje anomalij



# NoiseRank



- **Ansambelska metoda za odkrivanje šuma in osamelcev v podatkih**

# NoiseRank



- **Ansambelska metoda za odkrivanje šuma in osamelcev v podatkih**

- uporaba poljubnega števila modelov



# NoiseRank



- **Ansambelska metoda za odkrivanje šuma in osamelcev v podatkih**

- uporaba poljubnega števila modelov
- združevanje njihovih rezultatov

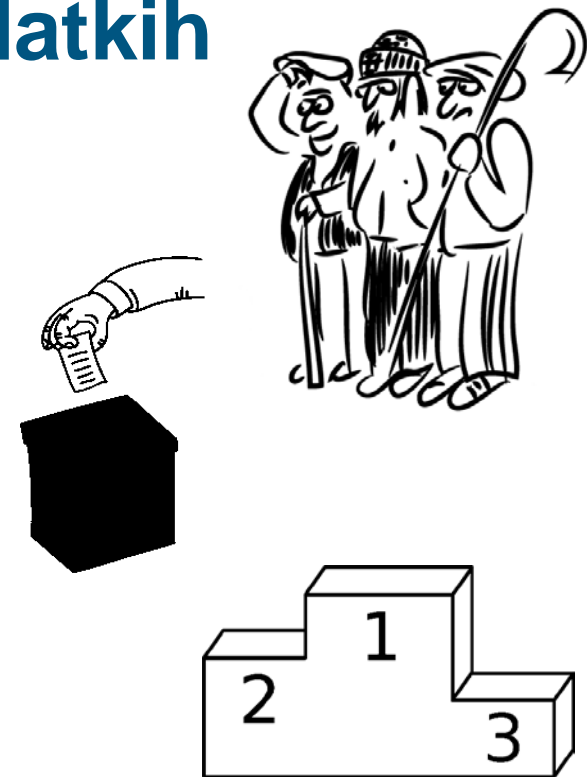


# NoiseRank



- **Ansambelska metoda za odkrivanje šuma in osamelcev v podatkih**

- uporaba poljubnega števila modelov
- združevanje njihovih rezultatov
- rangiranje primerov po „šumnosti“



# Uporaba



## ClowdFlows

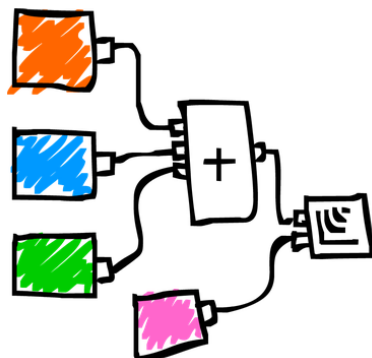
Data mining workflows on the cloud.

Try the tutorial

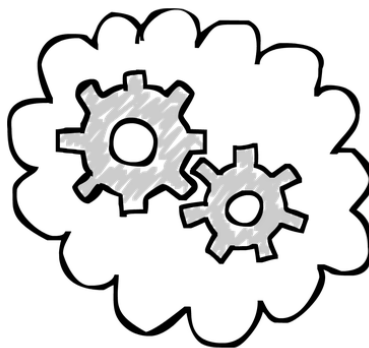
Start working

Explore existing workflows

Construct a workflow in the  
browser



Execute the workflow in the  
cloud

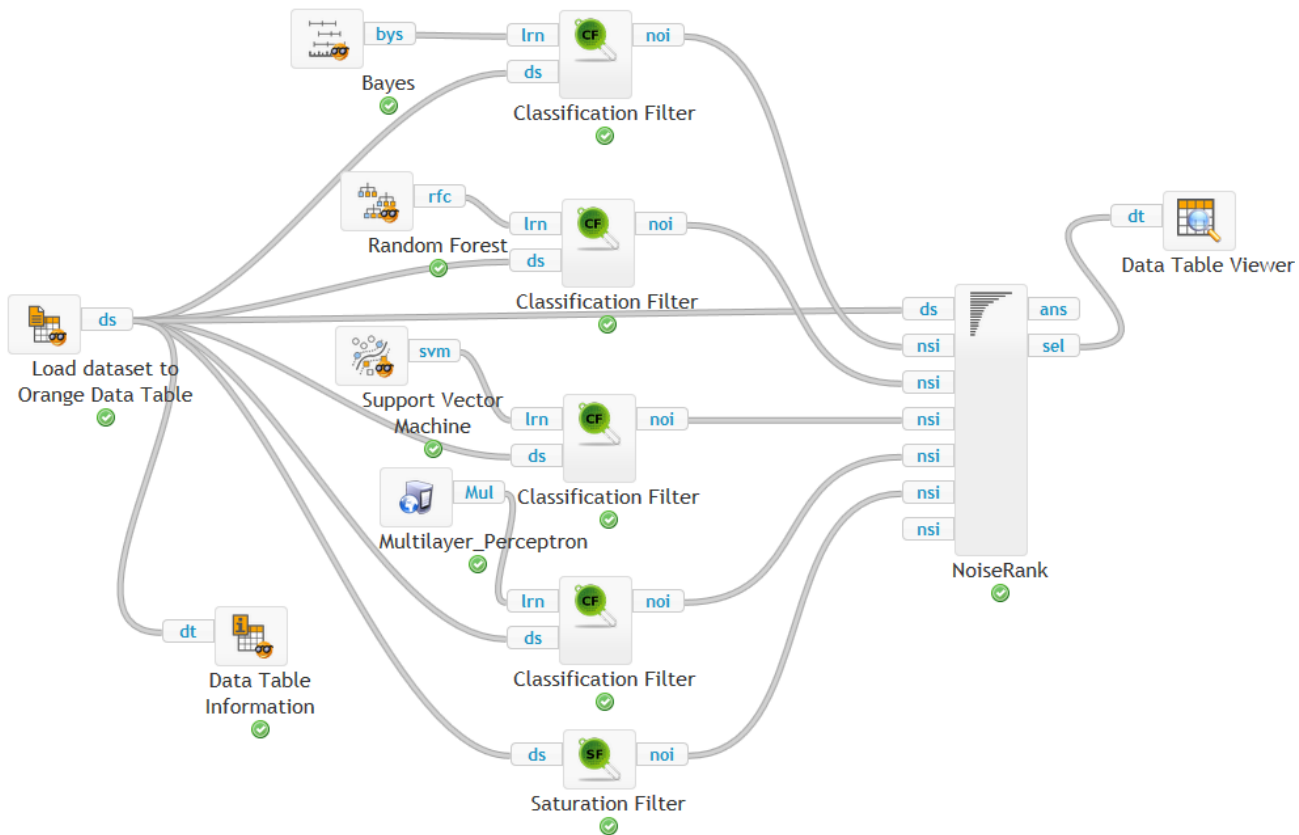


Share your workflows and  
results



<http://www.clowdflows.org>

# Uporaba



# Uporaba

NoiseRank wants your input!

Select the data instances that you want to examine in more detail.

Selected	Rank	Class	ID	Detected by:				
<input checked="" type="checkbox"/>	1.	non-CHD	51	Naive Bayes (Orange)	RF500 (Orange)	SVM (Orange)	Multilayer Perceptron	SF
<input checked="" type="checkbox"/>	2.	CHD	229	RF500 (Orange)	SVM (Orange)	Multilayer Perceptron	SF	
<input checked="" type="checkbox"/>	3.	CHD	0	SVM (Orange)	Multilayer Perceptron	SF		
<input checked="" type="checkbox"/>	4.	non-CHD	27	RF500 (Orange)	Multilayer Perceptron	SF		
<input checked="" type="checkbox"/>	5.	non-CHD	39	Naive Bayes (Orange)	SVM (Orange)	Multilayer Perceptron		
<input checked="" type="checkbox"/>	6.	CHD	176	Naive Bayes (Orange)	SVM (Orange)	Multilayer Perceptron		
<input checked="" type="checkbox"/>	7.	CHD	194	Naive Bayes (Orange)	SVM (Orange)	Multilayer Perceptron		
<input checked="" type="checkbox"/>	8.	CHD	213	RF500 (Orange)	SVM (Orange)	Multilayer Perceptron		
<input type="checkbox"/>	9.	CHD	42	SVM (Orange)	Multilayer Perceptron			
<input type="checkbox"/>	10.	non-CHD	120	Naive Bayes (Orange)	SVM (Orange)			
<input type="checkbox"/>	11.	non-CHD	164	Naive Bayes (Orange)	RF500 (Orange)			
<input type="checkbox"/>	12.	non-CHD	173	RF500 (Orange)	SF			
<input type="checkbox"/>	13.	CHD	196	Naive Bayes (Orange)	SVM (Orange)			
<input type="checkbox"/>	14.	non-CHD	226	RF500 (Orange)	SF			
<input type="checkbox"/>	15.	non-CHD	30	SVM (Orange)				
<input type="checkbox"/>	16.	CHD	45	Multilayer Perceptron				





# Rezultati



# Rezultati



- **Medicinska domena\***

\* Institute of Cardiovascular Prevention and Rehabilitation, Zagreb, Hrvaška

# Rezultati

- **Medicinska domena\***
- **Bolezni srca (CHD)**
  - 1. Napačna diagnoza
  - 2. Zapleten primer
  - 3. Športnik

\* Institute of Cardiovascular Prevention and Rehabilitation, Zagreb, Hrvaška

Rank	Class	ID	Detected by:
1.	non-CHD	51	__Bayes__ __RF__ __SVM__ __NN__ __SF__
2.	CHD	229	__RF__ __SVM__ __NN__ __SF__
3.	CHD	0	__SVM__ __NN__ __SF__
3.	non-CHD	27	__RF__ __NN__ __SF__
3.	non-CHD	39	__Bayes__ __SVM__ __NN__
3.	CHD	176	__Bayes__ __SVM__ __NN__
3.	CHD	194	__Bayes__ __SVM__ __NN__
3.	CHD	213	__RF__ __SVM__ __NN__
4.	CHD	42	__SVM__ __NN__
4.	non-CHD	120	__Bayes__ __SVM__
4.	non-CHD	164	__Bayes__ __RF__
4.	non-CHD	173	__RF__ __SF__
4.	CHD	196	__Bayes__ __SVM__
4.	non-CHD	226	__RF__ __SF__
5.	non-CHD	30	__SVM__
5.	CHD	45	__NN__
5.	non-CHD	59	__RF__
5.	non-CHD	62	__Bayes__
5.	non-CHD	63	__SF__
5.	CHD	67	__SVM__
5.	non-CHD	72	__SVM__
5.	CHD	97	__SF__
5.	non-CHD	135	__SVM__
5.	non-CHD	177	__NN__
5.	CHD	181	__Bayes__
5.	non-CHD	189	__SVM__
5.	CHD	195	__SVM__
5.	CHD	200	__NN__
5.	CHD	205	__SVM__
5.	CHD	231	__Bayes__



# Rezultati



- **Novice**
- **Izgredi po Kenijskih volitvah 2009\***

\* IPrA Research Center, University of Antwerp, Netherlands

# Rezultati

- Novice
- Izgredi po Kenijskih volitvah 2009
- Primerjava lokalnih in „zahodnih“ novic

DAILY  NATION

 THE STANDARD

THE  TIMES

The  INDEPENDENT

The New York Times

The Washington Post

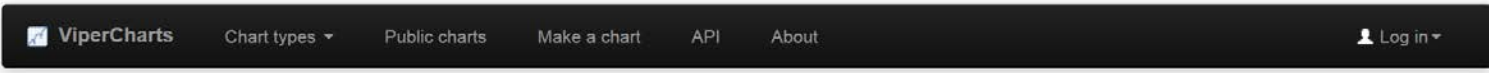
# Rezultati

- 1. Napake pri zajemu
- 2. Gostujoči novinarji, druge tematike, drugi žanr

Rank	Class	ID	Detected by:				
1.	WE	351	_Bayes_	_RF_	_SVM_	_NN_	_SF_
1.	WE	356	_Bayes_	_RF_	_SVM_	_PruneSF_	_SF_
1.	WE	357	_Bayes_	_RF_	_SVM_	_PruneSF_	_SF_
2.	LO	3	_Bayes_	_RF_	_SVM_	_PruneSF_	
2.	LO	24	_Bayes_	_RF_	_SVM_	_NN_	
2.	LO	100	_Bayes_	_RF_	_SVM_	_NN_	
2.	LO	161	_Bayes_	_RF_	_SVM_	_NN_	
2.	LO	172	_Bayes_	_RF_	_SVM_	_NN_	
2.	WE	325	_Bayes_	_RF_	_SVM_	_NN_	
2.	WE	409	_Bayes_	_RF_	_SVM_	_PruneSF_	
3.	LO	60	_Bayes_	_PruneSF_		_SF_	
3.	LO	152	_Bayes_	_RF_	_SVM_		
3.	LO	347	_Bayes_	_RF_	_SVM_		
3.	WE	363	_Bayes_	_RF_	_SVM_		
3.	WE	369	_Bayes_	_RF_	_SVM_		
3.	WE	463	_RF_	_SVM_	_NN_		
4.	LO	20	_RF_	_SVM_			
4.	LO	67	_RF_	_SVM_			
4.	LO	119	_Bayes_	_SF_			
4.	LO	157	_Bayes_	_PruneSF_			
4.	LO	158	_PruneSF_	_SF_			
4.	LO	184	_Bayes_	_PruneSF_			
4.	LO	200	_PruneSF_	_SF_			
4.	LO	214	_Bayes_	_NN_			
4.	LO	220	_PruneSF_	_SF_			
4.	LO	221	_Bayes_	_PruneSF_			
4.	WE	237	_RF_	_SVM_			
4.	WE	246	_Bayes_	_SVM_			
4.	WE	366	_Bayes_	_PruneSF_			
4.	WE	374	_Bayes_	_PruneSF_			
4.	WE	384	_Bayes_	_PruneSF_			
4.	WE	386	_SVM_	_NN_			
4.	WE	407	_Bayes_	_PruneSF_			
4.	WE	412	_RF_	_NN_			
4.	WE	415	_PruneSF_	_SF_			
4.	WE	420	_Bayes_	_SF_			
5.	LO	7	_PruneSF_				
5.	LO	17	_Bayes_				
5.	LO	70	_PruneSF_				
5.	LO	83	_SVM_				
5.	LO	129	_NN_				
5.	LO	131	_PruneSF_				
5.	LO	160	_Bayes_				
5.	LO	180	_SVM_				
5.	LO	186	_Bayes_				
5.	LO	213	_NN_				
5.	WE	245	_PruneSF_				
5.	WE	311	_RF_				
5.	WE	321	_PruneSF_				
5.	WE	327	_RF_				
...							

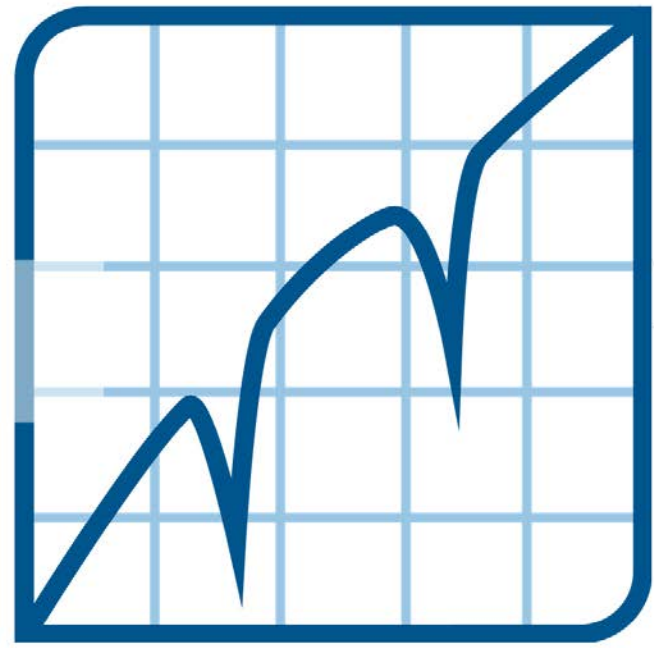


# Kvantitativna analiza



**Welcome to ViperCharts.**  
Visual performance evaluation made easy and intuitive.

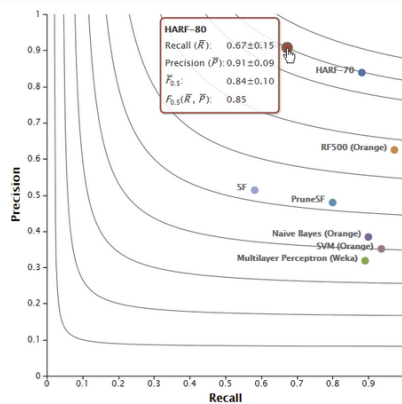
[See an example](#)



© 2014 Borut Sluban, Department of Knowledge Technologies, Jožef Stefan Institute, Slovenia.

<http://viper.ijs.si>

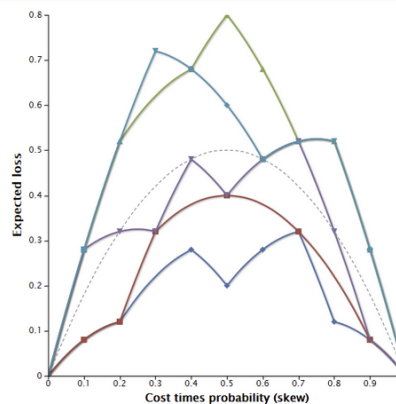
# Kvantitativna analiza



## Scatter Chart

Visual performance evaluation: compare the performance of your information retrieval, entity recognition or detection algorithms in the Precision-Recall or ROC space. [More »](#)

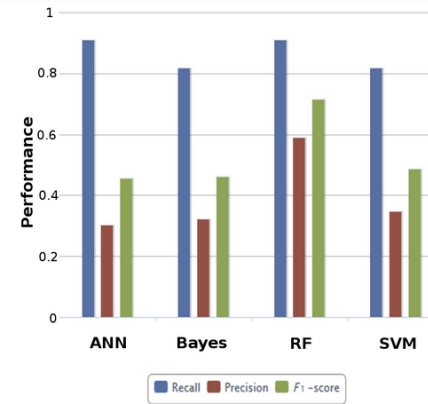
[Make a Scatter Chart »](#)



## Curve Chart

Easily chart the performance of your binary classifier systems. Included performance visualizations for PR curves, Lift curves, ROC curves, Cost curves, Rate-driven curves and Kendall curves. [More »](#)

[Make a Curve Chart »](#)



## Column Chart

Standard graphical presentation of algorithm performance. Visualizes the values of one or more performance measures on the evaluated algorithms. [More »](#)

[Make a Column Chart »](#)



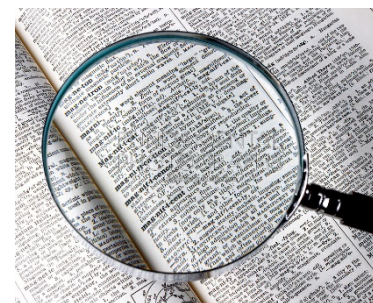
# Povzetek

- **Ansambelska metoda za odkrivanje anomalij v podatkih**



# Povzetek

- **Ansambelska metoda za odkrivanje anomalij v podatkih**
- **Podpora strokovnjakom pri analizi domenskih podatkov**



# Povzetek

- **Ansambelska metoda za odkrivanje anomalij v podatkih**



- **Podpora strokovnjakom pri analizi domenskih podatkov**



- **Javno dostopno, razširljivo, nadgradljivo**





**Hvala  
za Vašo pozornost**