

Strojno učenje in podatkovno rudarjenje

Nada Lavrač

Odsek za tehnologije znanja

Institut Jožef Stefan



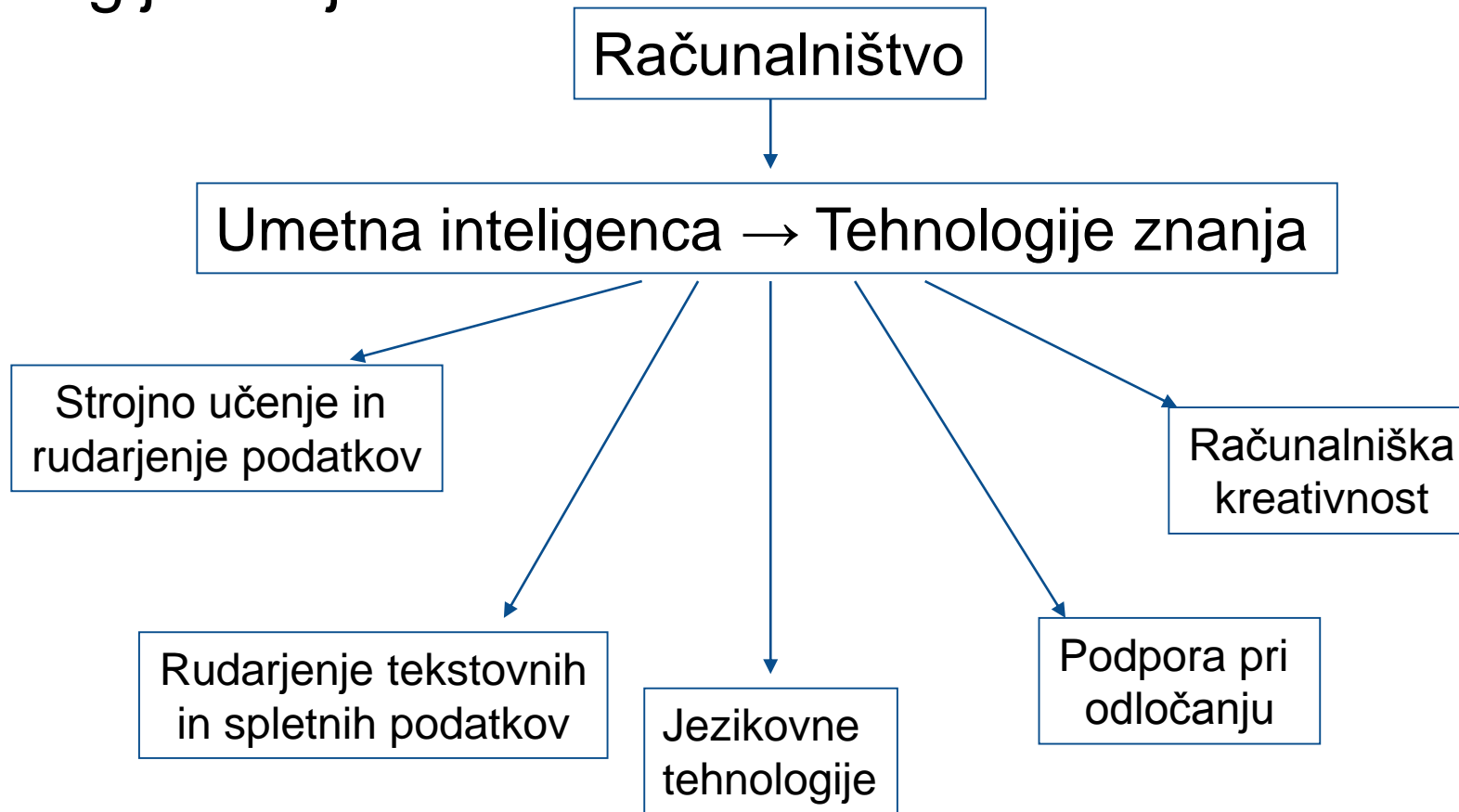
Odsek za tehnologije znanja IJS

- IJS (~ 700 raziskovalcev): naravoslovje in tehnike
- Odsek (~35 raziskovalcev): Raziskave in aplikacije tehnologij znanja



Odsek za tehnologije znanja IJS

- IJS (~ 700 raziskovalcev): naravoslovje in tehnike
- Odsek (~40 raziskovalcev): Raziskave in aplikacije tehnologij znanja



Kdo smo in kaj delamo




Tehnologi znanja IJS

Kdo smo in kaj delamo

- odkrivanje zakonitosti v podatkih
- rudarjenje tekstovnih in spletnih podatkov
- računalniška podpora pri odločanju
- jezikovne tehnologije in računalniško jezikoslovje
- računalniška kreativnost

Vsebina

- **Tehnologije znanja na IJS**
-  **Uvod: Klasične tehnike strojnega učenja**
- **Rudarjenje podatkov**
 - Izbrani sistemi in algoritmi druge generacije
 - Izbrane biomedicinske aplikacije
- **Napredne tehnike rudarjenja podatkov**
 - Relacijsko in semantično podatkovno rudarjenje
 - Izbrane biomedicinske aplikacije
- **Tekoče delo in zaključki**

Klasične tehnike strojnega učenja

Oseba	Starost	Dioptrija	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	MEHKE
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	TRDE
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	MEHKE
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE

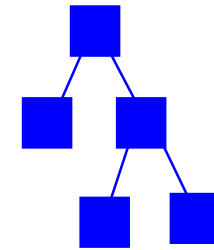
podatki

Vhod: tabela podatkov, relacijska podatkovna baza, dokumenti, spletne strani

Izhod: klasifikacijski model, posamezni vzorci

odkrivanje zakonitosti v podatkih

Strojno učenje
Rudarjenje podatkov



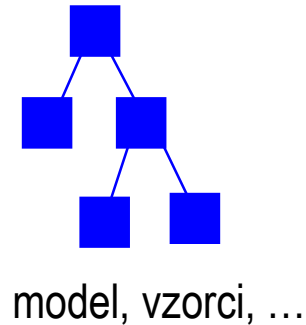
model, vzorci, ...

Klasične tehnike strojnega učenja

Oseba	Starost	Dioptrija	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	MEHKE
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	TRDE
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	MEHKE
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE

odkrivanje zakonitosti v podatkih

Strojno učenje
Rударjenje podatkov

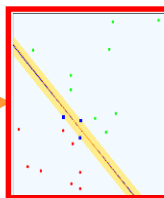


podatki

Vhod: tabela podatkov, relacijska podatkovna baza, dokumenti, spletne strani

Izhod: klasifikacijski model, posamezni vzorci

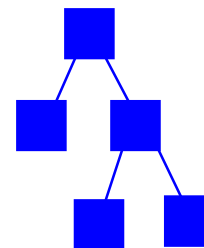
nov neklasificiran primer



klasificiran primer



črna škatla
ni razlage



simbolični model
simbolični vzorci



razlaga



Klasične tehnike strojnega učenja

1. PODATKI

starost	dioptriya	astigmatizem	solzenje	kontaktne leče
17	kratkovidnost	ne	zmanjšano	ne
23	kratkovidnost	da	normalno	trde
22	daljnovidnost	ne	normalno	mehke
27	kratkovidnost	ne	normalno	mehke
19	kratkovidnost	da	zmanjšano	ne
23	kratkovidnost	da	normalno	trde
20	daljnovidnost	ne	zmanjšano	ne
25	daljnovidnost	ne	normalno	mehke
28	daljnovidnost	da	zmanjšano	ne
21	daljnovidnost	da	normalno	trde
35	kratkovidnost	ne	zmanjšano	ne
36	kratkovidnost	ne	normalno	mehke
40	kratkovidnost	da	zmanjšano	ne
...	kratkovidnost	da	normalno	trde
...	daljnovidnost	ne	zmanjšano	ne
...	daljnovidnost	da	zmanjšano	ne
...	daljnovidnost	da	normalno	ne
...	kratkovidnost	ne	zmanjšano	ne
...	kratkovidnost	ne	normalno	ne
...	kratkovidnost	da	zmanjšano	ne
...	daljnovidnost	ne	zmanjšano	ne
...	daljnovidnost	ne	normalno	mehke
...	daljnovidnost	da	zmanjšano	ne
...	daljnovidnost	da	normalno	ne

Klasične tehnike strojnega učenja

starost	dioptriya	astigmatizem	solzenje	kontaktne leče
17	kratkovidnost	ne	zmanjšano	ne
23	kratkovidnost	da	normalno	trde
22	daljnovidnost	ne	normalno	mehke
27	kratkovidnost	ne	normalno	mehke
19	kratkovidnost	da	zmanjšano	ne
23	kratkovidnost	da	normalno	trde
20	daljnovidnost	ne	zmanjšano	ne
25	daljnovidnost	ne	normalno	mehke
28	daljnovidnost	da	zmanjšano	ne
21	daljnovidnost	da	normalno	trde
35	kratkovidnost	ne	zmanjšano	ne
36	kratkovidnost	ne	normalno	mehke
40	kratkovidnost	da	zmanjšano	ne
...	kratkovidnost	da	normalno	trde
...	daljnovidnost	ne	zmanjšano	ne
...	daljnovidnost	da	zmanjšano	ne
...	daljnovidnost	da	normalno	ne
...	kratkovidnost	ne	zmanjšano	ne
...	kratkovidnost	ne	normalno	ne
...	kratkovidnost	da	zmanjšano	ne
...	daljnovidnost	ne	zmanjšano	ne

2. VZOREC

PRAVILO

ČE

solzenje
zmanjšano

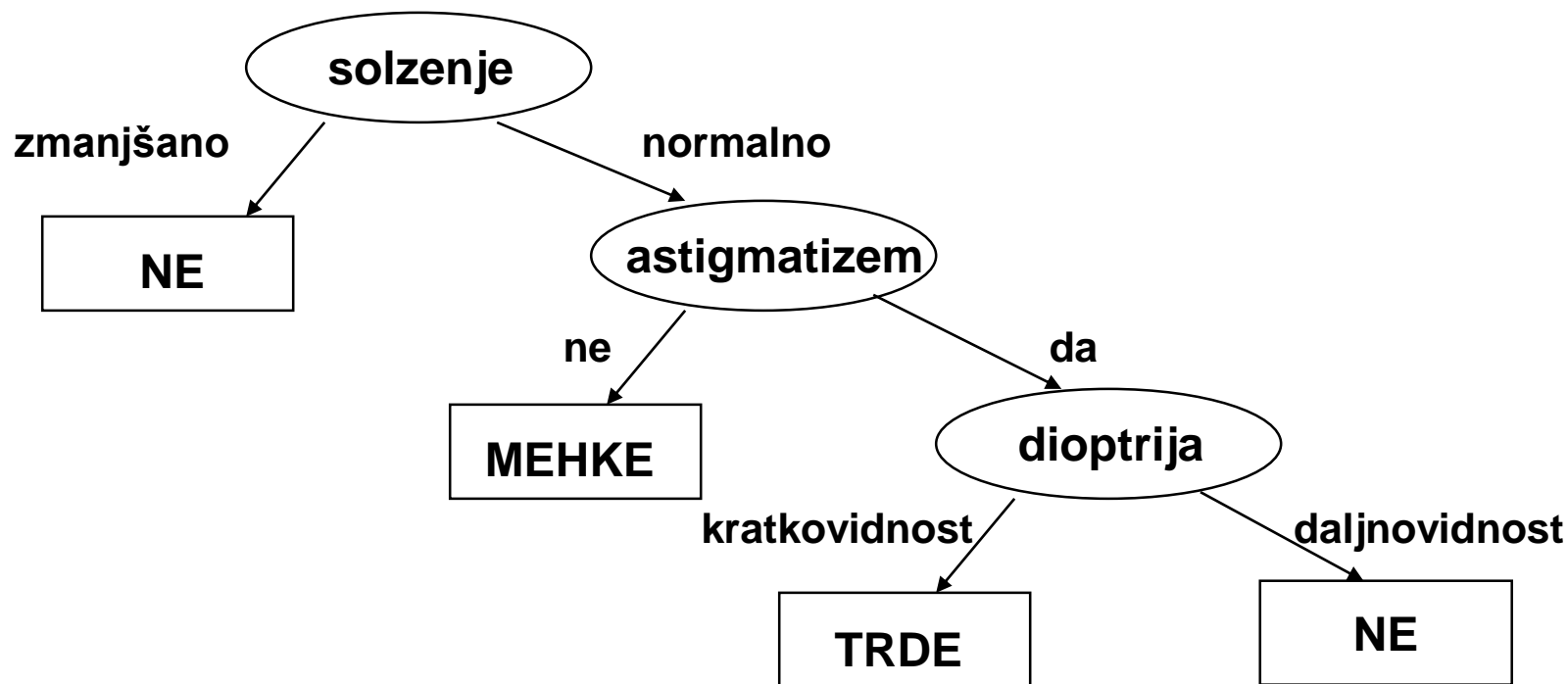
POTEM

kontaktne
leče=ne

Klasične tehnike strojnega učenja

... Odkrivanje znanja v podatkih

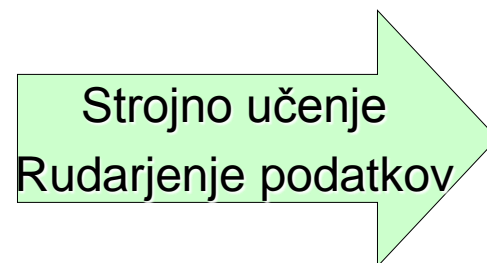
3. VZOREC / MODEL, ki ga potrdi človek = ZNANJE



4. UPORABA ZNANJA kot dodatno ekspertno mnenje

Klasične tehnike strojnega učenja

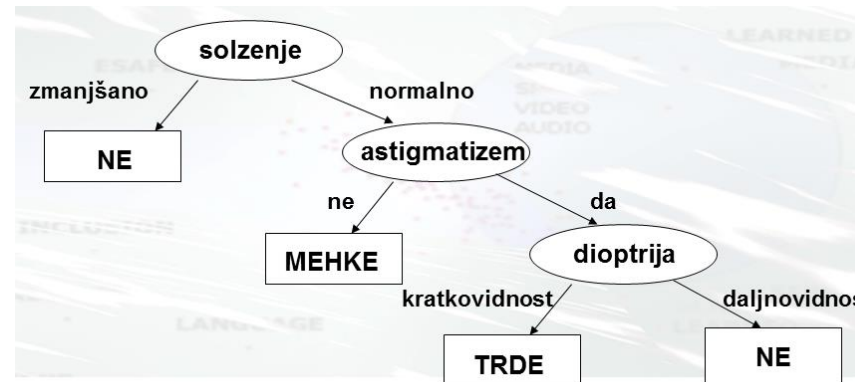
Oseba	Starost	Dioptrija	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	MEHKE
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	TRDE
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	MEHKE
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE



Klasične tehnike strojnega učenja

Oseba	Starost	Dioptrija	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	MEHKE
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	TRDE
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	MEHKE
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE

Strojno učenje
Rudarjenje podatkov



Klasične tehnike strojnega učenja

Oseba	Starost	Dioptrija	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	MEHKE
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	TRDE
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	MEHKE
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE

Strojno učenje
Rudarjenje podatkov



$$Gain(S, A) = E(S) - \sum_{v \in Values(A)} p_v \cdot E(S_v)$$

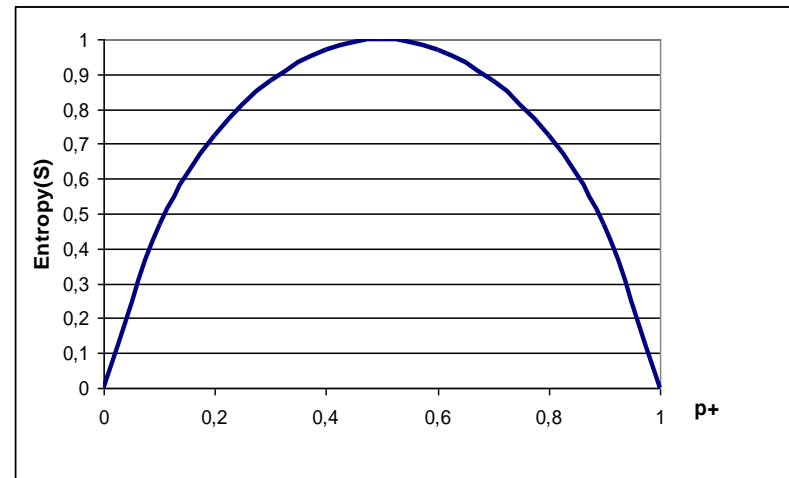
Hevristika za določanje najinformativnejšega atributa: $Gain(S,A)$
ocena zmanjšanja entropije množice S zaradi razbitja na podmnožice glede na vrednosti v atributa A

Klasične tehnike strojnega učenja

... Ocenjevanje informativnosti atributov

- **Glavna heuristika:** Kateri atribut izbrati kot test v danem vozlišču odločitvenega drevesa ? Atribut, ki je najkoristnejši za čimtočnejšo klasifikacijo primerov.
- Definiramo **statistično oceno informativnosti atributa**, ki meri kako dobro atribut ločuje med primeri, ki pripadajo različnim razredom
- **Informativnost** merimo kot **zmanjšanje entropije učne množice** primerov
- **Entropija** je mera “nečistosti” učne množice: $E(S) = - p_+ \log_2 p_+ - p_- \log_2 p_-$.
- **Najinformativnejši atribut:**
 - uporabi atribut v vozlišču drevesa, razbij S na S_1, S_2, \dots, S_v
 - izberi A , ki maksimizira informacijski prispevek $\text{Max Gain}(S, A)$

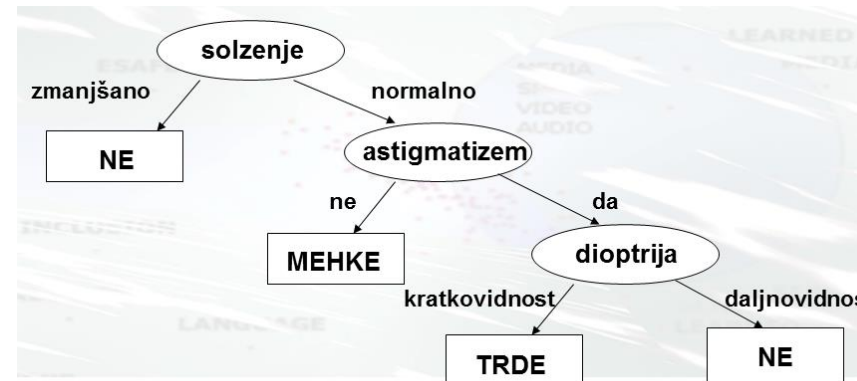
$$\text{Gain}(S, A) = E(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \cdot E(S_v)$$



Klasične tehnike strojnega učenja

Oseba	Starost	Dioptriya	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	MEHKE
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	TRDE
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	MEHKE
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE

Strojno učenje
Rudarjenje podatkov



leče=NE ← solzenje=zmanjšano

leče=NE ← solzenje=normalno & astigmatizem=da & dioptriya=daljnovidnost

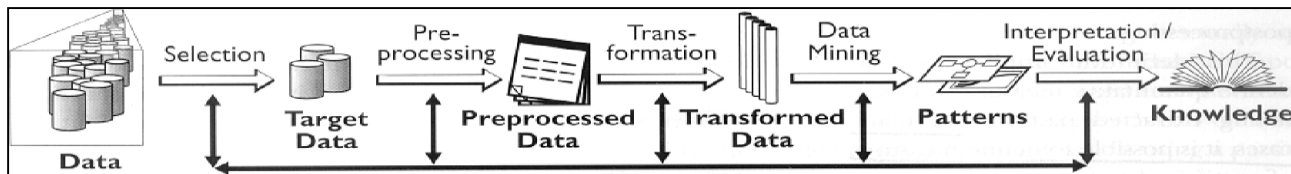
leče=MEHKE ← solzenje=normalno & astigmatizem=ne

leče=TRDE ← solzenje=normalno & astigmatizem=da & dioptriya=kratkovidnost

Druga generacija sistemov

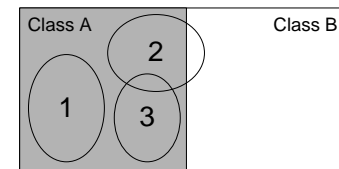
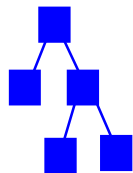
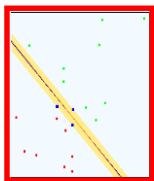
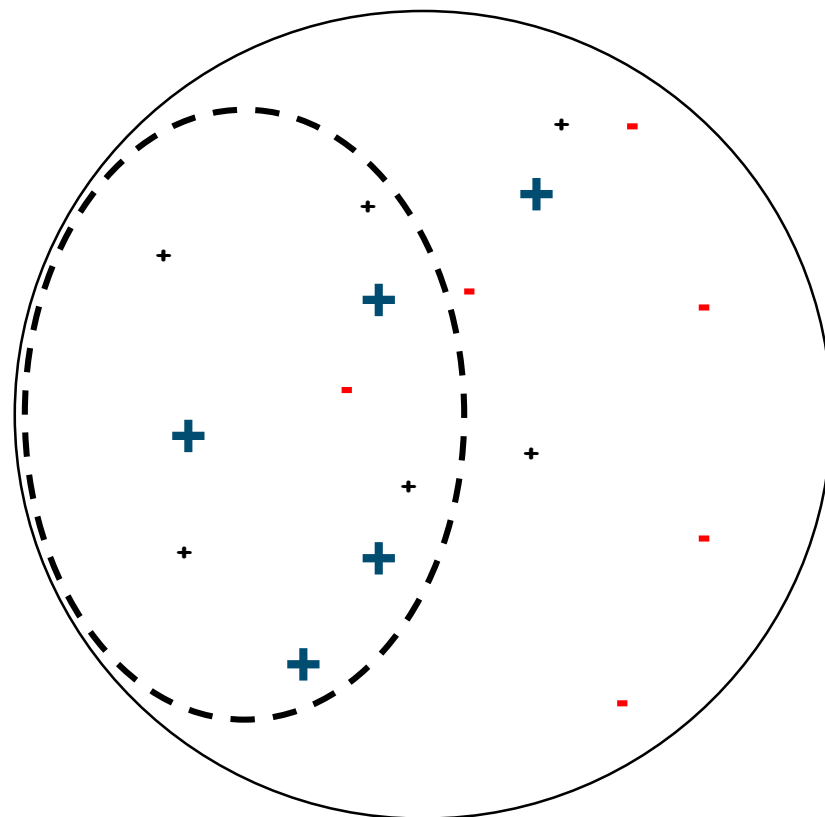
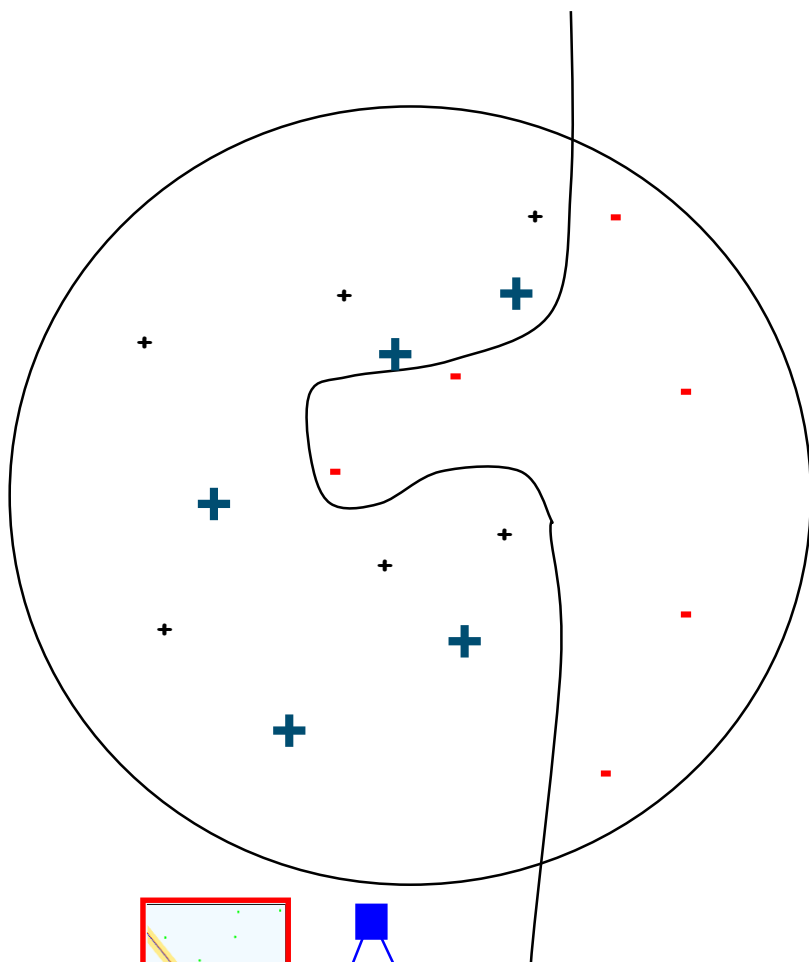
Po letu 1990 ...

- **Odkrivanje zakonitosti v podatkih** iz velikih podatkovnih baz, opisanih z velikim številom atributov



- Industrijski standard: CRISP-DM (1997)
- Nove konference o praktičnih aspektih rudarjenja podatkov in odkrivanja zakonitosti v podatkih: KDD, PKDD, ...
- **Nove naloge in učinkoviti algoritmi:**
 - **Učenje klasifikacijskim modelov:** Bayesovske mreže, večrelacijsko učenje, statistično relacijsko učenje, metode podpornih vektorjev (SVM), ...
 - **Odkrivanje vzorcev in učenje opisnih pravil:** učenje povezovalnih pravil, odkrivanje podskupin, ...

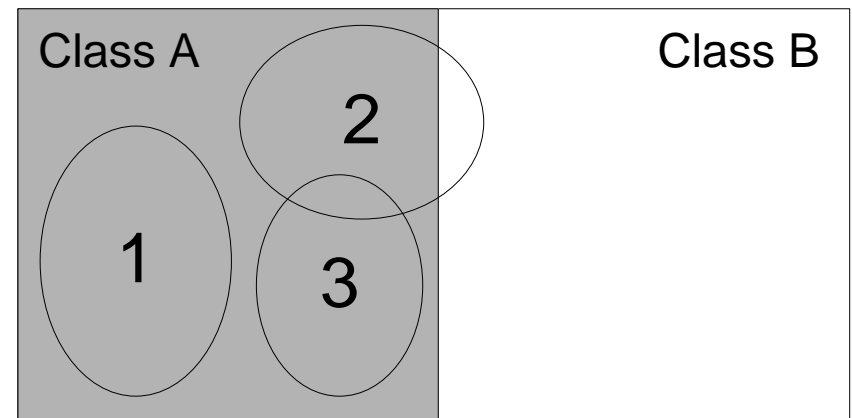
Klasifikacija vs. odkrivanje vzorcev



Odkrivanje podskupin

- Učenje opisnih pravil:
 - učenje posameznih pravil, ki opisujejo izbrano lastnost podatkov
- Algoritmi za odkrivanje podskupin se naučijo več neodvisnih pravil, vsako opiše podskupino primerov s podobnimi lastnostmi
 - odkrivanje in opis velikih signifikantnih podmnožic primerov

Oseba	Starost	Dioptrija	Astigmat.	Solzenje	Leče
O1	mlad	kratko	ne	zmanjšano	NE
O2	mlad	kratko	ne	normalno	DA
O3	mlad	kratko	da	zmanjšano	NE
O4	mlad	kratko	da	normalno	DA
O5	mlad	daleko	ne	zmanjšano	NE
O6-O13
O14	pr_st_dal	daleko	ne	normalno	DA
O15	pr_st_dal	daleko	da	zmanjšano	NE
O16	pr_st_dal	daleko	da	normalno	NE
O17	st_daleko	kratko	ne	zmanjšano	NE
O18	st_daleko	kratko	ne	normalno	NE
O19-O23
O24	st_daleko	daleko	da	normalno	NE



Primer: Odkrivanje rizičnih skupin pacientov

Vhod: Podatki o pacientih, opisanih z atributi **A** (anamneza), **B** (anamneza & laboratorij), in **C** (an., lab. & ECG)

Naloga: V dani populaciji poišči in opiši rizične skupine za CHD (dovolj velike podskupine z signifikantno različno distribucijo ciljnega razreda glede na distribucijo v celotni populaciji)

Izmed avtomatsko najdenih pravil je ekspert izbral tiste, ki imajo najboljši potencial za seznanjanje pacientov o rizičnosti:

A1: rizičnaSkupina ← moški & pozitivna družinska anamneza & starost > 46

A2: rizičnaSkupina ← ženska & indeks telesne teže > 25 & starost > 63

B1: rizičnaSkupina ← ...

B2: rizičnaSkupina ← ...

C1: rizičnaSkupina ← ...

Primer: Odkrivanje rizičnih skupin pacientov

Rizična skupina A2:

rizičnaSkupina ← ženska & indeks telesne teže > 25 & starost > 63

Podporni faktorji za dodatno karakterizacija skupin –
izračunani z uporabo χ^2 statističnega testa signifikance

Pozitivna družinska anamnreza in visok krvni pritisk. Ženske v tej skupini imajo tipično tudi povečano LDL vrednost holesterola in normalno a znižano vrednost HDL vrednost holesterola.

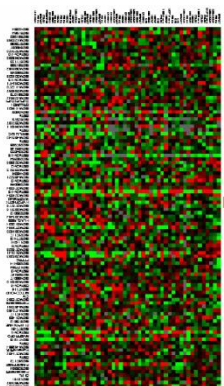
(Gamberger & Lavrač, JAIR 2002)

Primer: Analiza DNA mikromrež z odkrivanjem podskupin

- **Funkcijska genomika** je tipična domena za znanstveno odkrivanje zakonitosti, katere cilj je študij genov in njihovih funkcij. Za to domeno je značilno:
 - zelo veliko število atributov (genov) v primerjavi s številom primerov
 - tipična velikost: 7000-16000 atributov, 50-150 primerov
- **Primer problema:** Ločevanje med Akutno limfoblastično levkemijo (ALL, 27 primerov) in akutno mieloidno levkemijo (AML, 11 primerov), z dodatnimi 34 primeri v ločeni testni množici. Vsak primer je opisan z izraženostjo 7129 genov

<http://www-genome.wi.mit.edu/cgi-bin/cancer/datasets.cgi>

Podatki DNA mikromrež: format za rudarjenje pdoatkov



#1



#2



...

#100

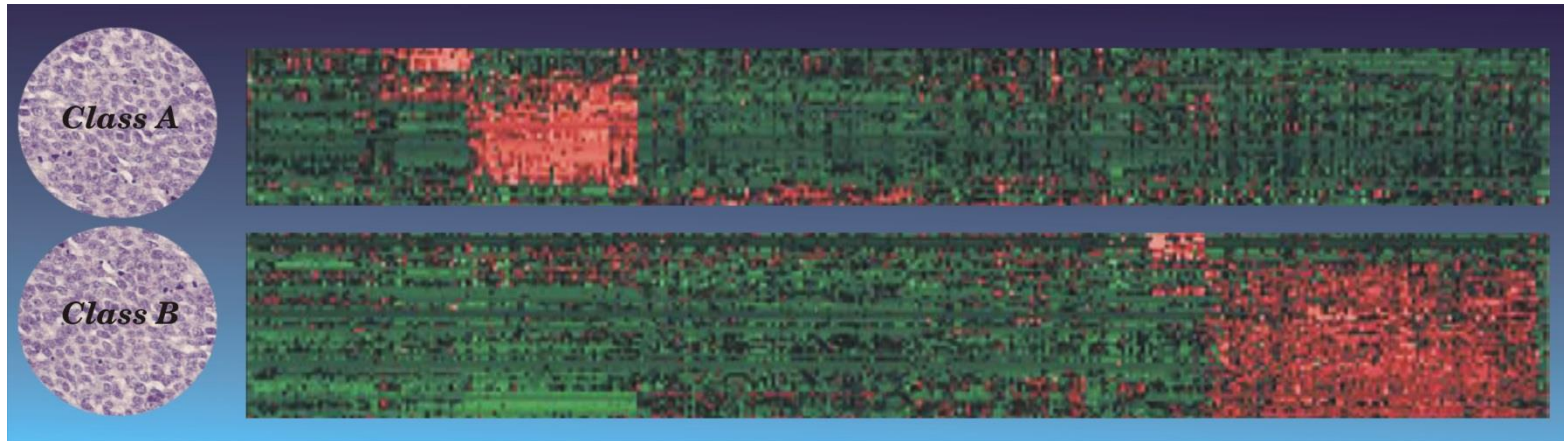


malo
primerov

Patient #	Tumor Type	Gene #1	Gene #2	Gene #3	...	Gene #10,000
1	A	5.00	1.33	3.45	...	4.22
2	A	0.98	0.87	1.04	...	?
3	B	0.33	1.40	0.42	...	0.24
...
100	B	0.89	0.90	1.00	...	0.66

veliko atributov

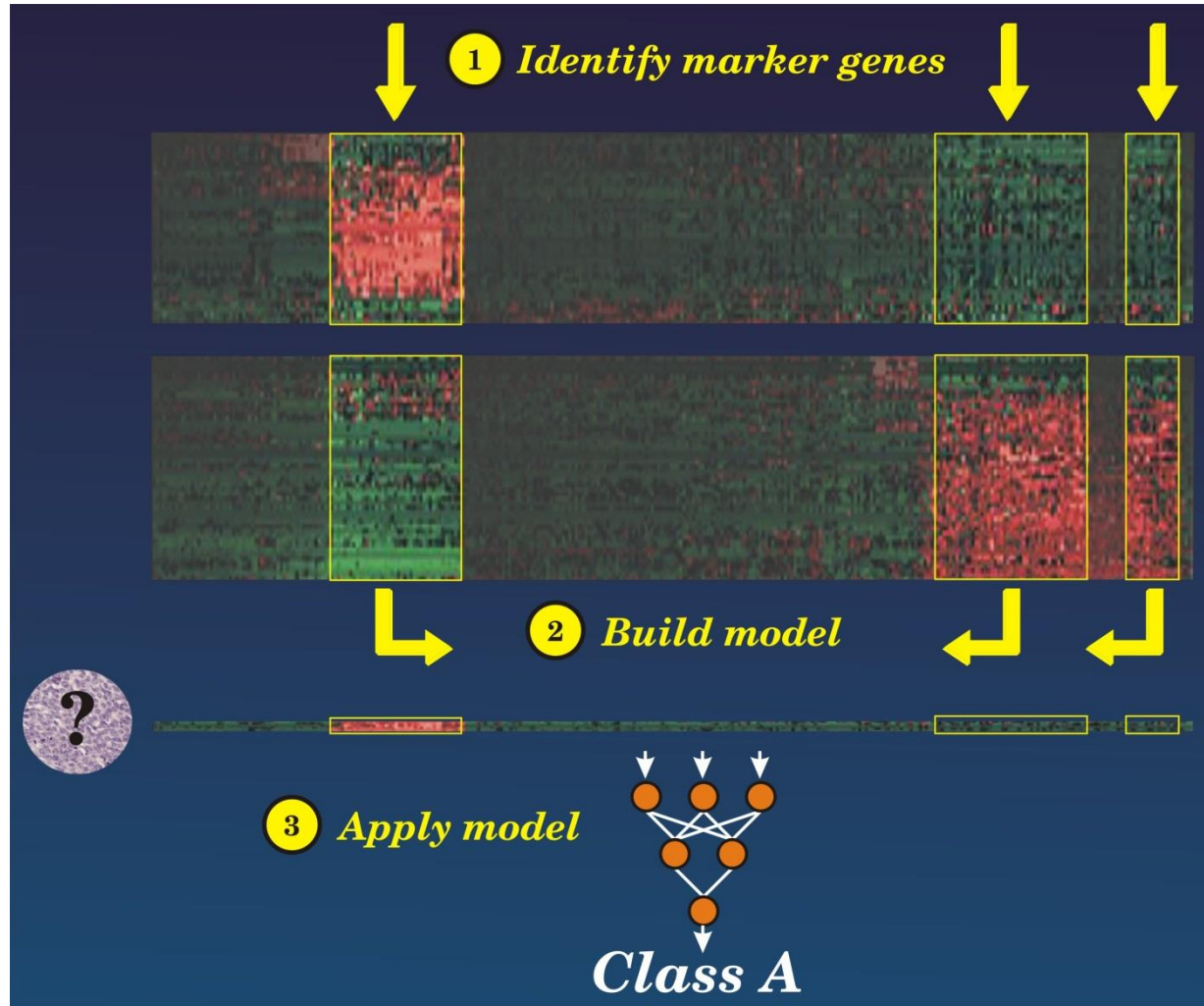
Podatki DNA mikromrež



- Dvo-razredni diagnostični problem ALL vs. AML
- Večrazredni diagnostični problem: 14 tipov raka, skupaj le 144 primerov v učni množici in 54 primerov v testni množici. Vsak primer opisan z izraženostjo 16063 genov.

<http://www-genome.wi.mit.edu/cgi-bin/cancer/datasets.cgi>.

Klasična naloga: Učenje klasifikacijskih modelov



Analiza DNA mikromrež z metodami odkrivanja podskupin

- Vzorci, odkriti z algoritmom SD za odkrivanje podskupin

AML ← gen_20056 = DIFF. IZRAŽEN &
gen_23984 = NI_DIFF. IZRAŽEN

Levkemija ← KIAA0128 = DIFF. IZRAŽEN &
prostoglandin_d2_synthase = NI_DIFF. IZRAŽEN

- Možna je biološka **interpretacija** pravil, ki opisujejo podskupine

D. Gamberger, N. Lavrač, F. Železný, J. Tolar,
Journal of Biomedical Informatics, 2004

Analiza DNA mikromrež z metodami odkrivanja podskupin

- Vzorci, odkriti z algoritmom SD za odkrivanje podskupin

AML ← gen_20056 = DIFF. IZRAŽEN &
gen_23984 = NI_DIFF. IZRAŽEN

Levkemija ← KIAA0128 = DIFF. IZRAŽEN &
prostoglandin_d2_synthase = NI_DIFF. IZRAŽEN

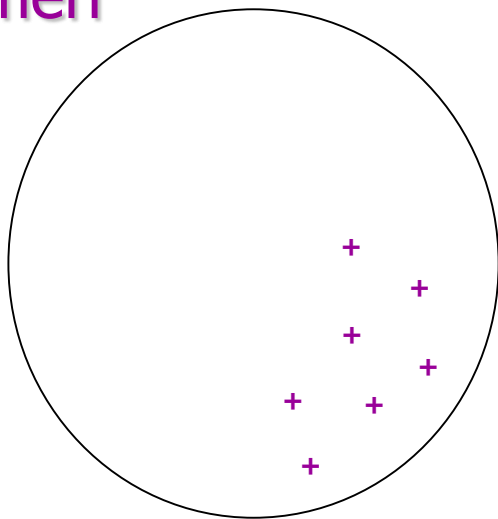
- Možna je biološka **interpretacija** pravil, ki opisujejo podskupine

D. Gamberger, N. Lavrač, F. Železný, J. Tolar,
Journal of Biomedical Informatics, 2004

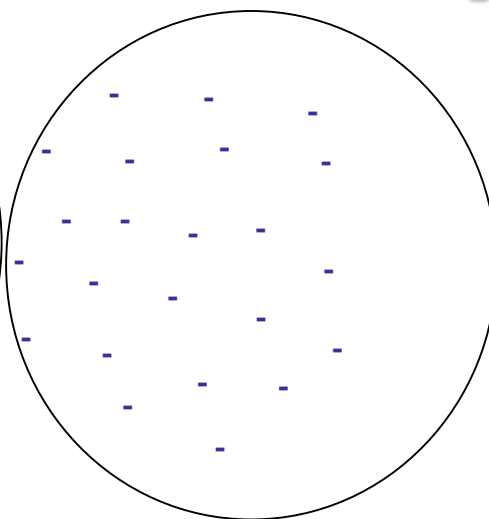
- nove hevristike za učenje opisnih pravil
- opisna pravila imajo manjšo klasifikacijsko točnost
- **trade-off med točnostjo in interpretabilnostjo**

Klasični pokrivni algoritem za učenje klasifikacijskih pravil

Pozitivni primeri



Negativni primeri

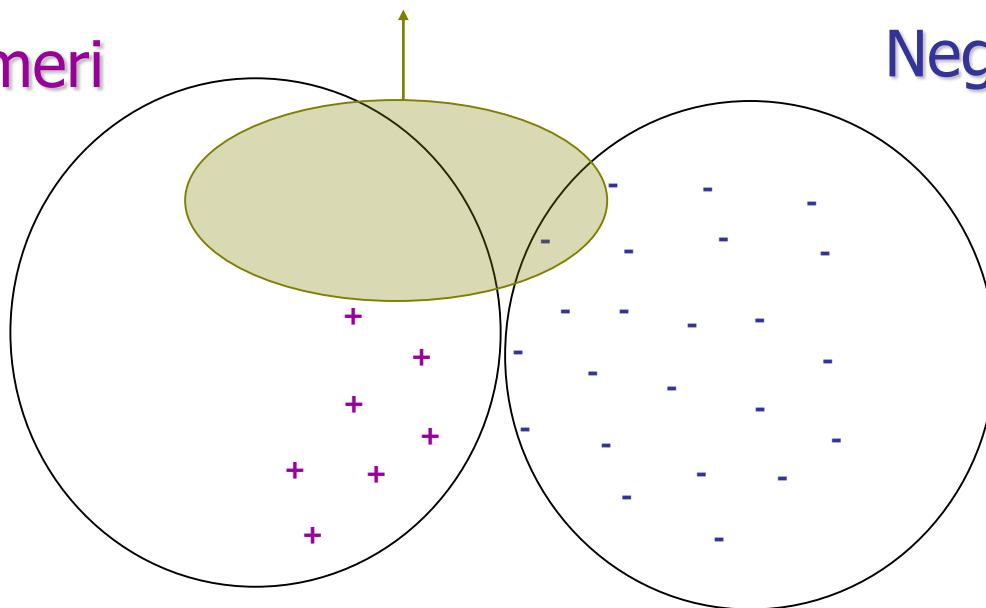


Klasični pokrivni algoritem za učenje klasifikacijskih pravil

Pravilo 1: Razred=Poz. ← Pogoju2 & Pogoju3

Pozitivni primeri

Negativni primeri

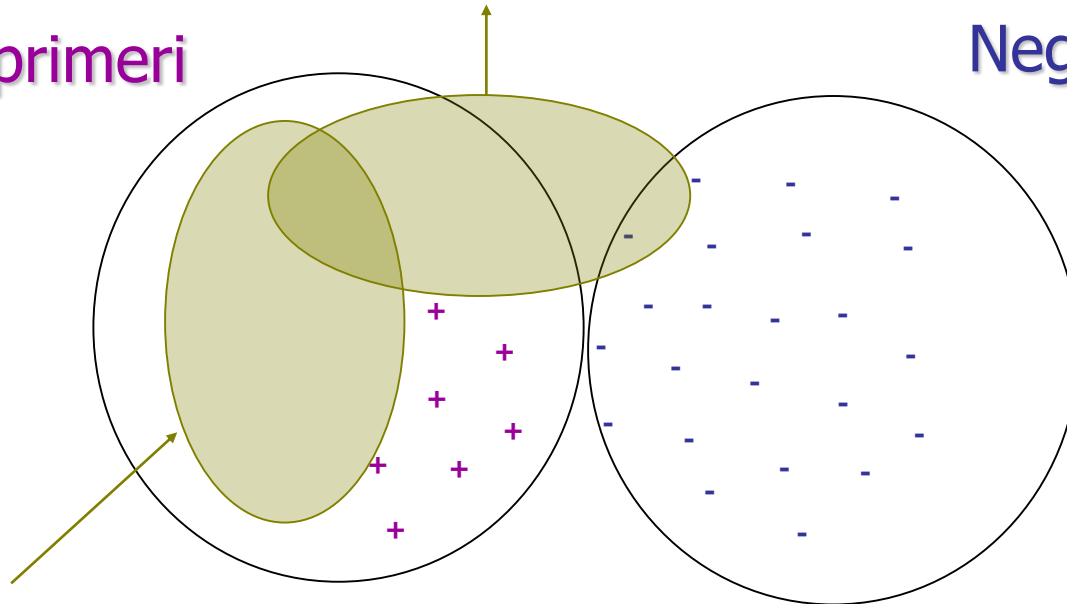


Klasični pokrivni algoritem za učenje klasifikacijskih pravil

Pravilo 1: Razred=Poz. ← Pogoju2 & Pogoju3

Pozitivni primeri

Negativni primeri



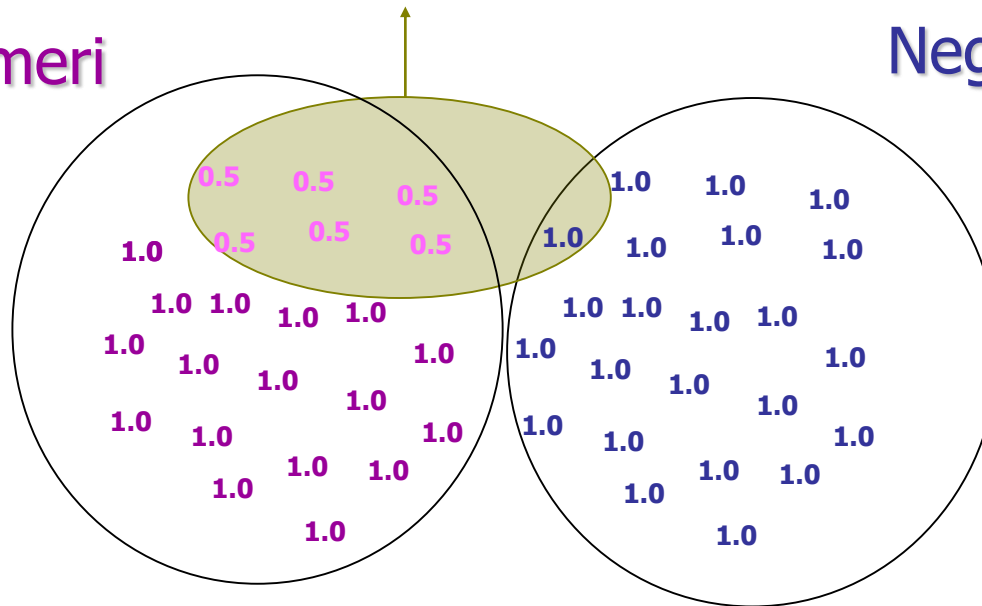
Pravilo 2: Razred=Poz. ← Pogoju8 AND Pogoju6

Novi uteženi pokrivni algoritem in nove heuristike za učenje opisnih pravil

Pravilo 1: Razred=Poz. ← Pogoje2 & Pogoje3

Pozitivni primeri

Negativni primeri



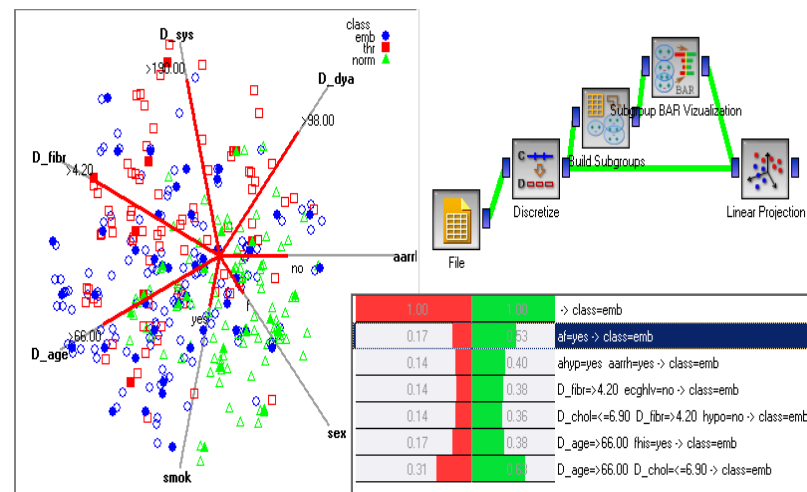
Heuristika - mera kvalitete pravila v SD: $q(\text{Razred} \leftarrow \text{Pogoji}) = TP / (FP + g)$

Heuristika v CN2-SD: $q(\text{Razred} \leftarrow \text{Pogoji}) = \sim p(\text{Pogoji}) \times (p(\text{Razred} | \text{Pogoji}) - p(\text{Razred}))$
 $\sim \text{Pokritost} \times (\text{Preciznost} - \text{Default})$

Pokritost = delež pokritih primerov, **Preciznost** = delež pravilno pokritih pozitivnih primerov

Napredna orodja za rudarjenje podatkov v platformah Orange (FRI) in Orange4WS (IJS)

- **Orange**
 - klasifikacijski algoritmi
 - algoritmi za odkrivanje podskupin
 - vizualizacija podatkov
 - delotoki podatkovnega rudarjenja



- **Algoritmi za odkrivanje podskupin v Orange in Orange4WS**

SD (Gamberger & Lavrač, JAIR 2002)

CN2-SD (Lavrač et al., JMLR 2004)

APRIORI-SD (Kavšek & Lavrač, AAI 2006)

Metodologija SegMine implementirana v platformi Orange4WS (BMC Bioinformatics 2011)

SegMine overview

Segment1-P1	Segment2-P2	Segment3-P3	Segment4-P4	Segment5-P5	Segment6-P6
25.73	41.29	33.13	49.53	54.89	38.59
8.15	11.84	12.05	8.7	7.81	9.82
7.69	108.73	291.82	9.71	105.98	84.38
65.48	86.82	130.19	119.27	82.59	118.26
1.33	1.11	15.98	1.91	1.25	5.03
50.84	53.27	38.18	43.25	75.51	32.19
2.89	0.68	4.24	1.43	4.91	4.41
184.58	150.82	119.35	141.87	155.45	157.76
5.45	1.51	30.72	0.34	2.89	0.65
292.55	359.93	485.48	289.12	344.46	291.91
9.34	12.14	9.67	7.82	9.39	8.37
7.84	52.98	47.63	89.49	55.46	40.43
4.41	39.8	17.72	26.42	19.17	12.15
0.35	0.65	2.5	0.34	0.41	1.95
31.09	43.62	151.49	25.51	101.89	26.77

raw data from a microarray experiment (expression of genes)

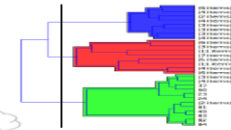


knowledge from ontologies

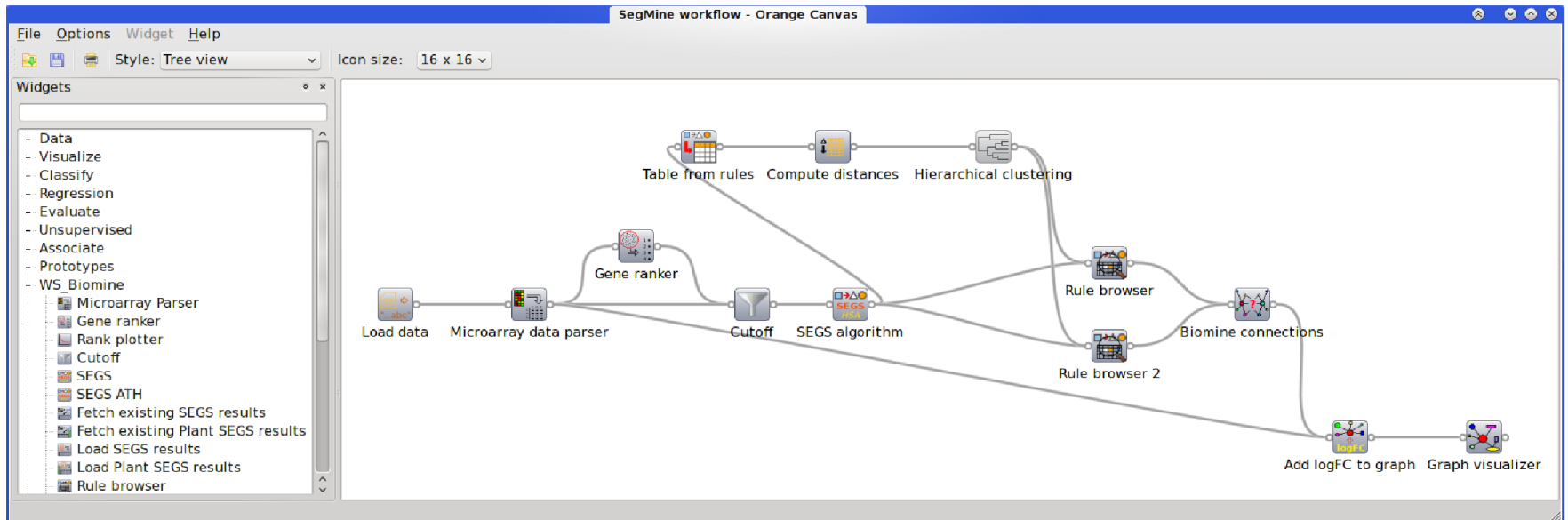
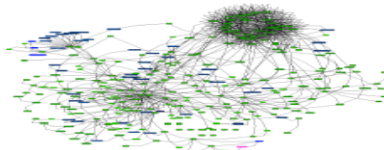
RULE 1 := organelle organization
AND INTERACT: transcription coactivator activity
AND nuclear part

RULE 2 := cellular macromolecule metabolic process
AND INTERACT: chromatin binding

RULE 3 := cellular response to stimulus
AND macromolecule organelle part
AND INTERACT: RNA binding



interpretation of gene expression data: rules, clusters, genesets



Primer kompleksne metodologije podatkovnega rudarjenja: SegMine

SegMine overview

1donor1-P2	1donor2-P2	1donor3-P2	2donor1-P11	2donor2-P7	2donor3-P8
25.71	41.29	33.11	49.53	54.89	36.59
8.15	11.84	12.85	6.7	7.61	9.82
7.69	108.73	291.82	9.71	105.98	84.38
95.46	86.82	110.13	118.57	92.53	118.26
1.53	1.11	15.98	1.41	1.25	5.03
50.04	53.07	36.16	43.25	73.51	32.19
2.89	0.64	4.24	1.63	6.91	4.41
184.58	150.62	119.35	141.87	155.45	157.76
5.45	1.51	0.72	0.34	2.83	0.65
292.55	359.03	465.48	289.12	344.66	291.91
9.34	12.14	9.67	7.82	5.39	8.37
7.04	52.98	47.63	89.49	55.46	40.43
4.41	39.9	17.72	26.42	19.17	12.15
0.35	0.65	2.2	0.34	0.41	1.95
31.09	43.62	151.49	25.51	101.89	26.77

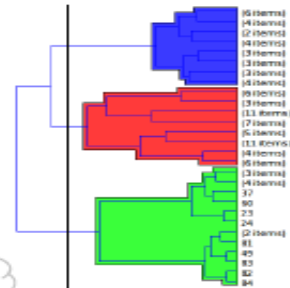
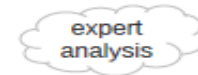
raw data from a microarray experiment (expression of genes)

SEGS

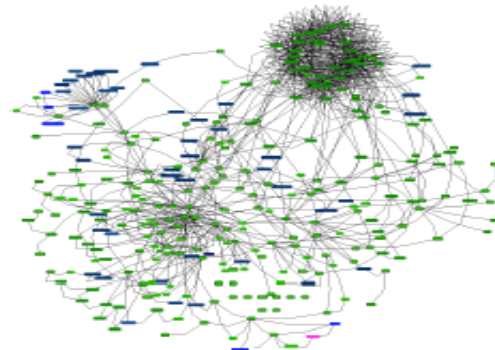


knowledge from ontologies

- RULE 1 := organelle organization
AND intracellular non-membrane-bounded organelle
AND INTERACT: transcription coactivator activity
- RULE 2 := cellular macromolecule metabolic process
AND nuclear part
AND INTERACT: chromatin binding
- RULE 3 := cellular response to stimulus
AND intracellular organelle part
AND INTERACT: RNA binding



interpretation of gene expression data:
rules, clusters, genesets

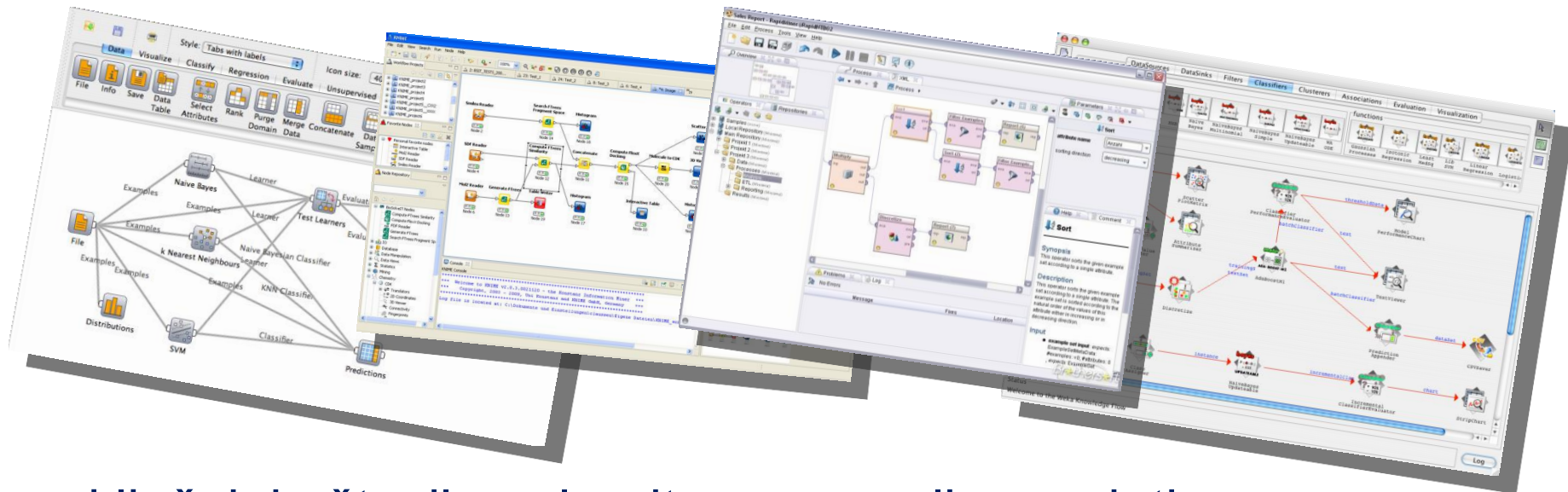


Biomine



public databases

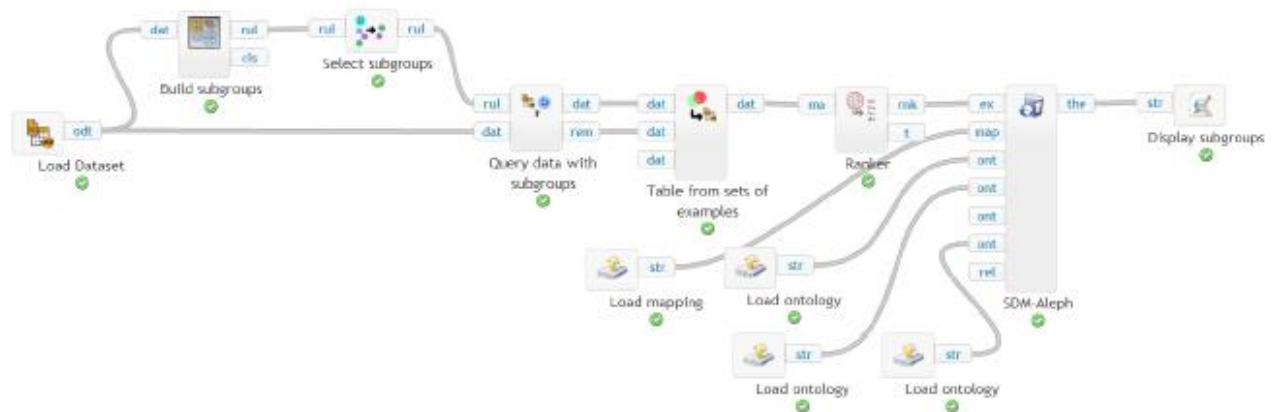
Platforme WEKA, KNIME, RapidMiner, Orange (FRI) in Orange4WS (IJS)



- vključujejo številne algoritme za analizo podatkov
- omogočajo analizo in vizualizacijo podatkov
- omogočajo gradnjo delotokov
- **Orange4WS omogoča tudi vključevanje spletnih servisov**

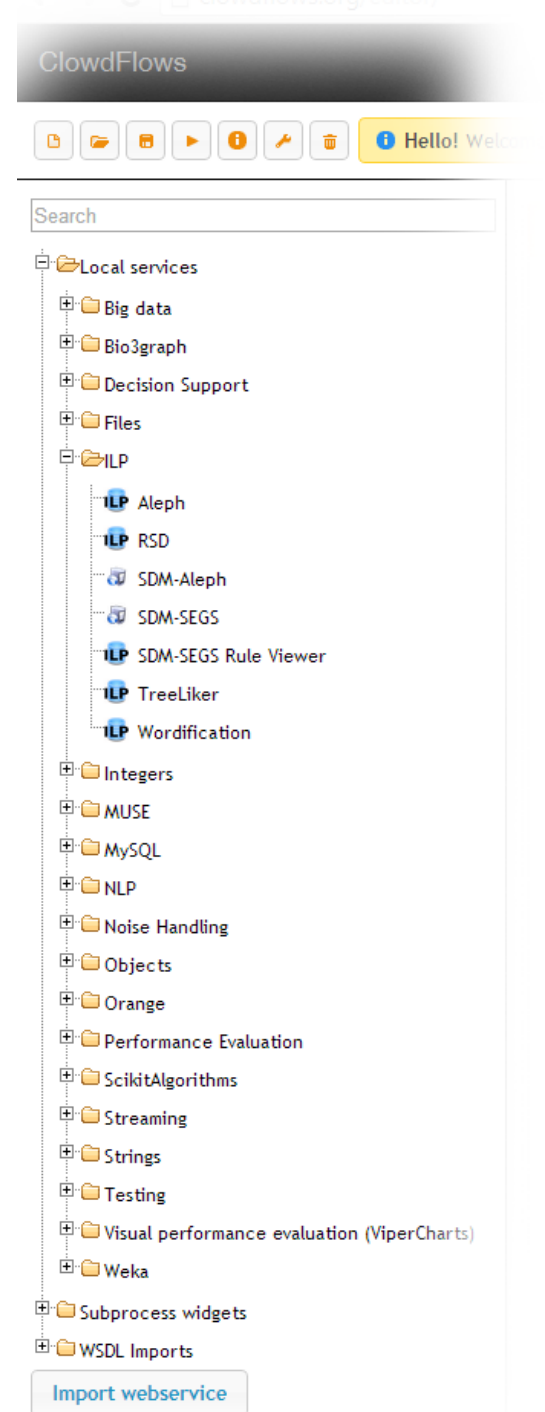
IJS platforma za rudarjenje podatkov CloudFlows

- Platforma tretje generacije za kreiranje in izvajanje kompleksnih postopkov rudarjenja podatkov
 - algoritmi so spletni servisi (v “oblaku”)
 - ni potrebna instalacija platforme
 - delotoki so dostopni vsakomur kar iz brskalnika s preprostim klikom na spletni naslov:
npr. <http://clowdflows.org/workflow/910/>



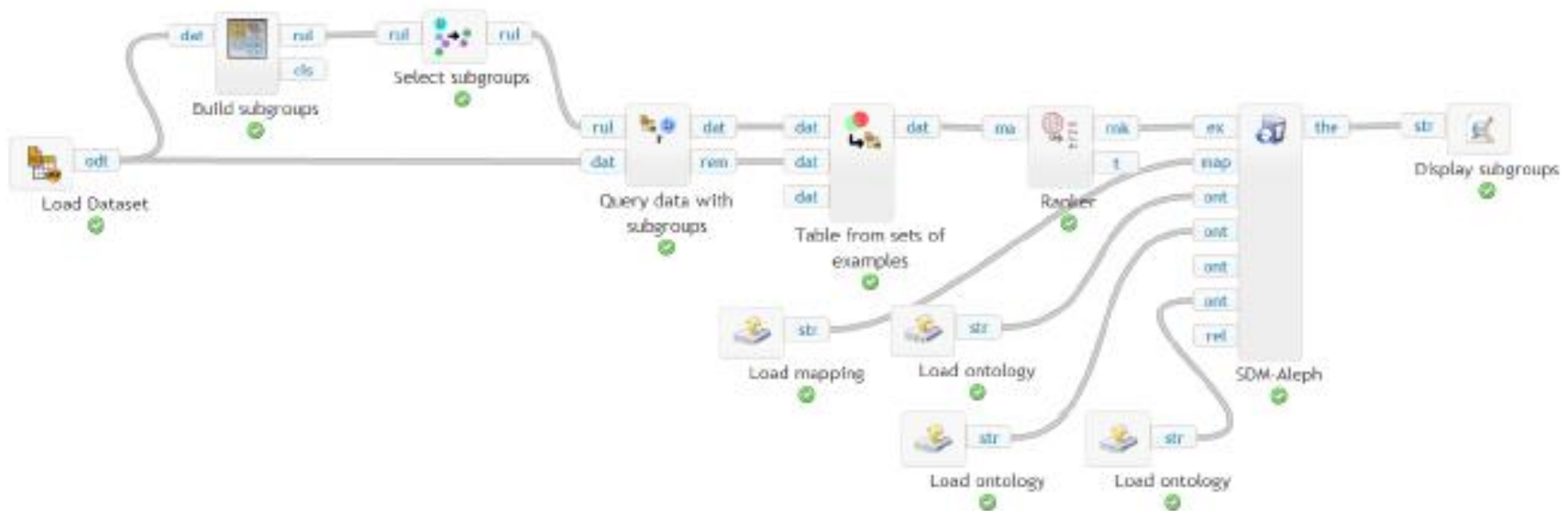
Platforma CloudFlows

- **Velik repozitorij algoritmov**
 - analiza relacijskih podatkovnih baz
 - vsu algoritmi iz platforme Orange
 - WEKA algoritmi kot spletni servisi
 - vizualizacija podatkov in rezultatov analiz
 - analiza tekstovnih podatkov
 - analiza socialnih omrežij
 - analiza velikih tokov podatkov
- **Velik repozitorij delotokov**
 - omogča dostop do naše tehnološke dediščine



Primer: Semantično rudarjenje podatkov

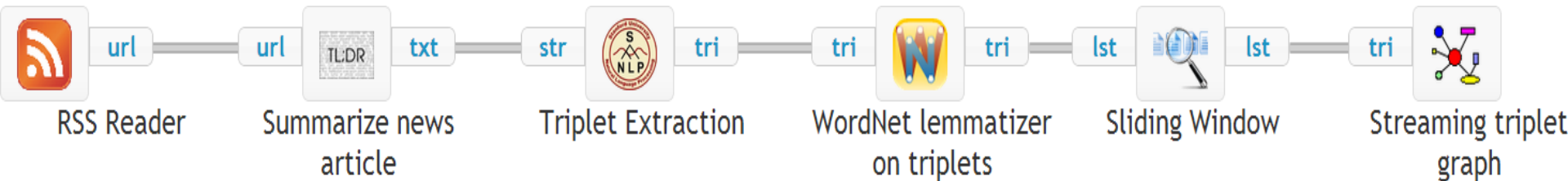
- Odkrivanje zanimivih podskupin v podatkih ter njihova razlaga s pomočjo biološkega predznanja - ontologij



- <http://clowdflows.org/workflow/910/>

Primer: “Big Data”

- Analiza velikih količin podatkov v realnem času.
Primer: izgradnja semantičnega grafa iz tokov spletnih novic. <http://clowdfwos.org/workflow/1729/>.



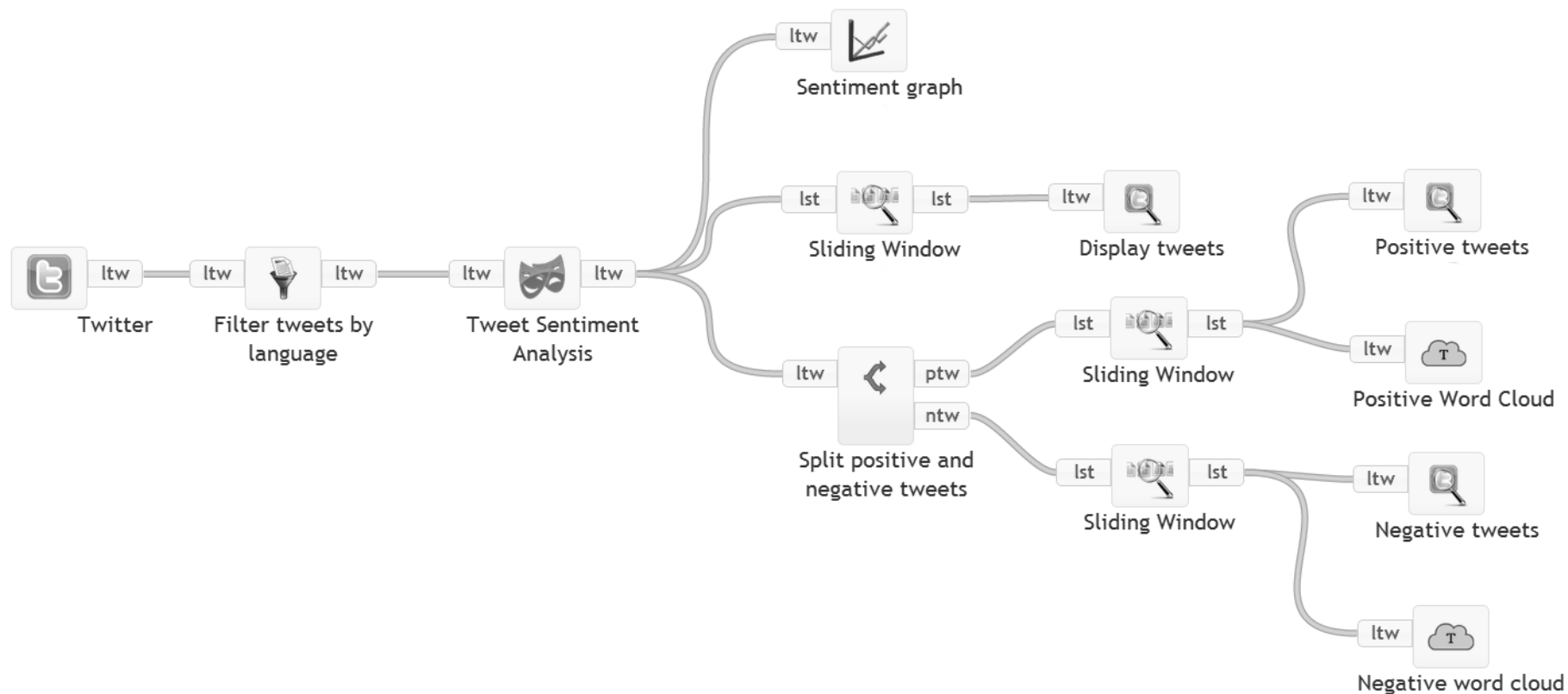
- Primer: monitoring novic z vizualizacijo grafa zgrajenega iz športnih novic (vir: CNN RSS feeds)



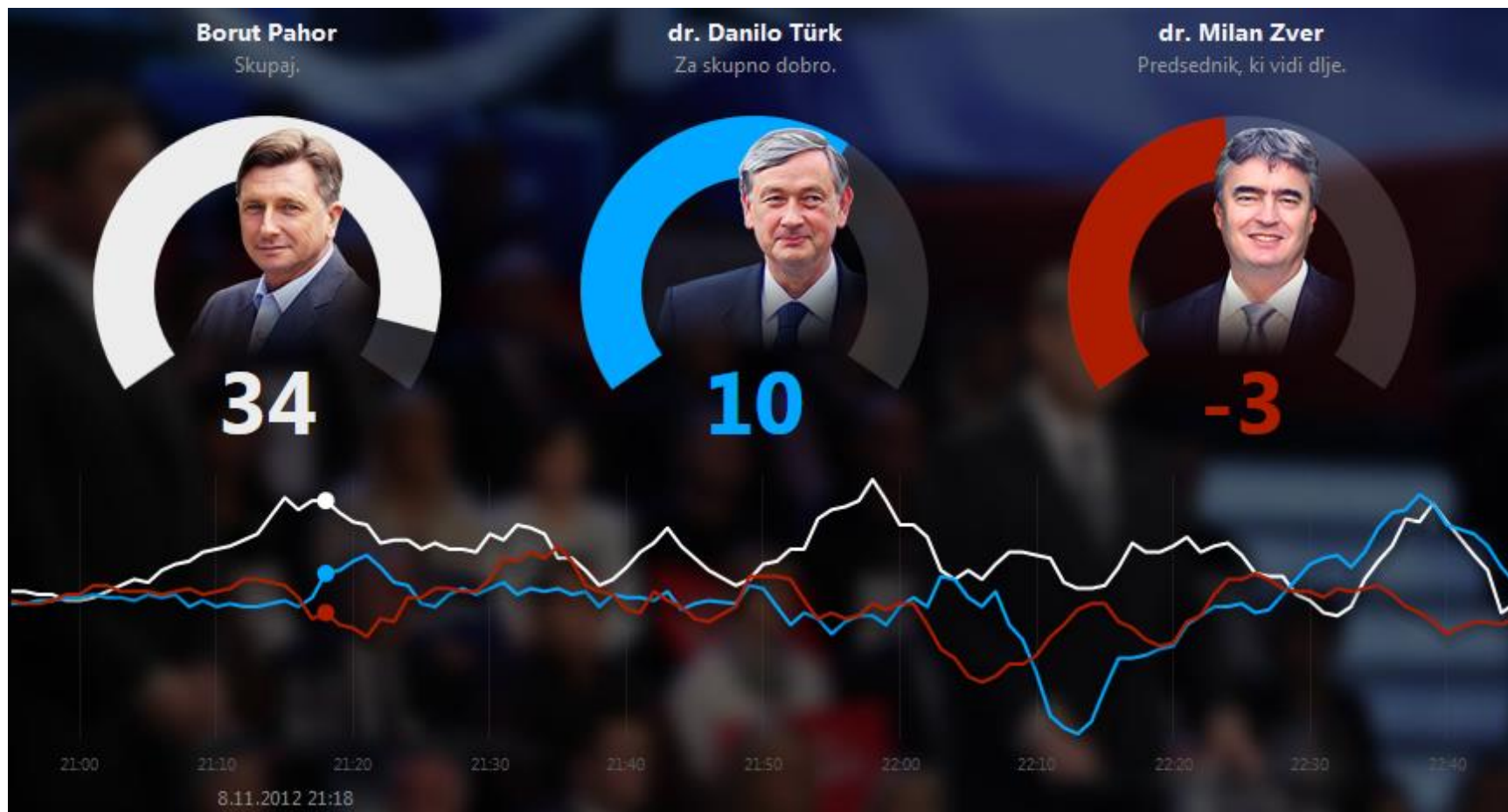
<http://clowdfwos.org/streams/data/31/15524/>.

Primer: “Big Data”

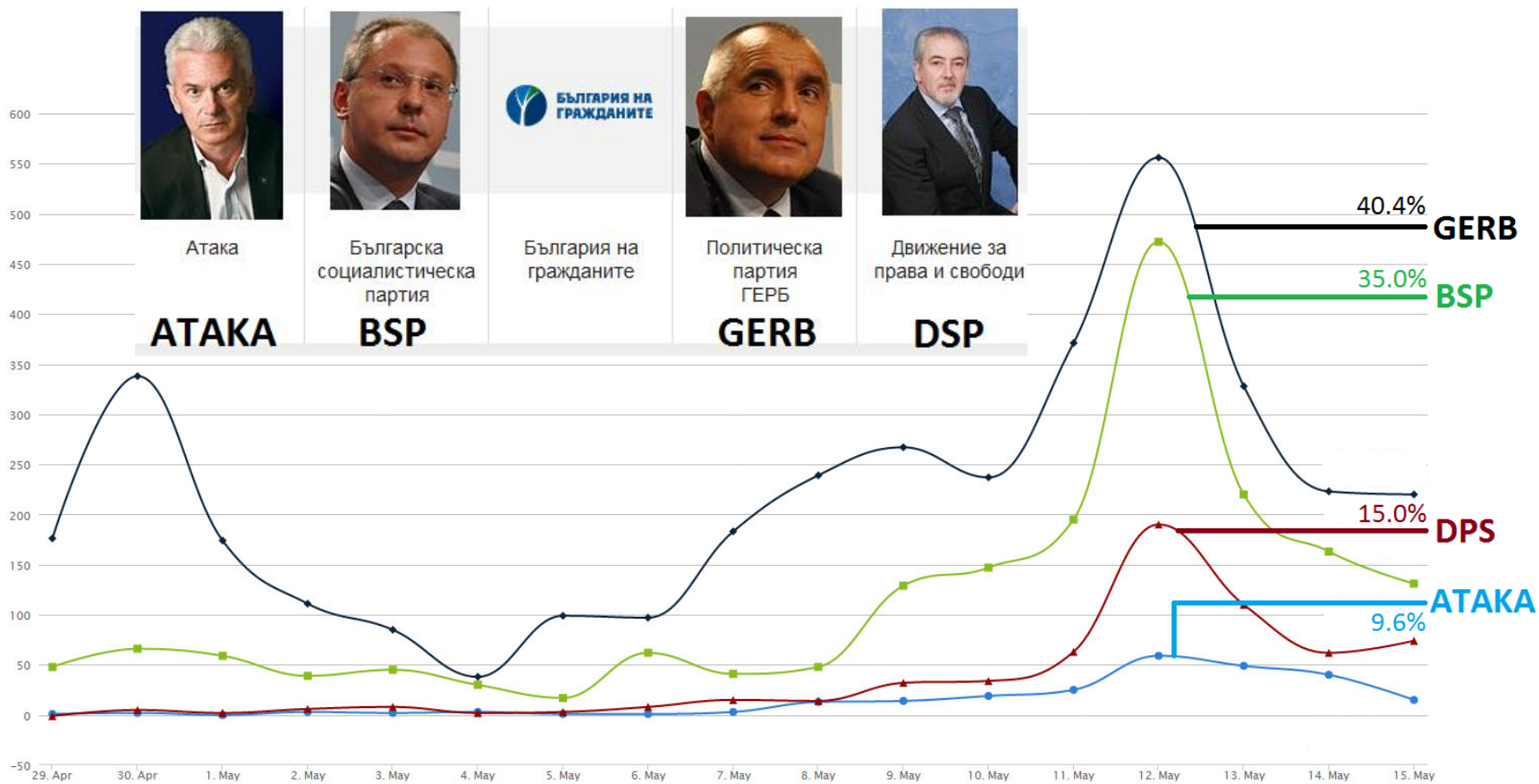
- Analiza pozitivnega/negativnega sentimenta v tvitih v realnem času: <http://clowdflows.org/workflow/1041/>.



Analiza sentimenta v tvitih: Predsedniške volitve 2012

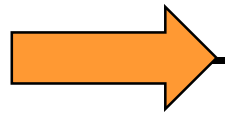


Parlamentarne volitve v Bolgariji



Vsebina predavanja

- odkrivanje zakonitosti in anomalij v podatkih



- **rudarjenje tekstovnih in spletnih podatkov**

- računalniška podpora pri odločanju

- jezikovne tehnologije in računalniško jezikoslovje

- **računalniška kreativnost**

Tekstovno rudarjenje :

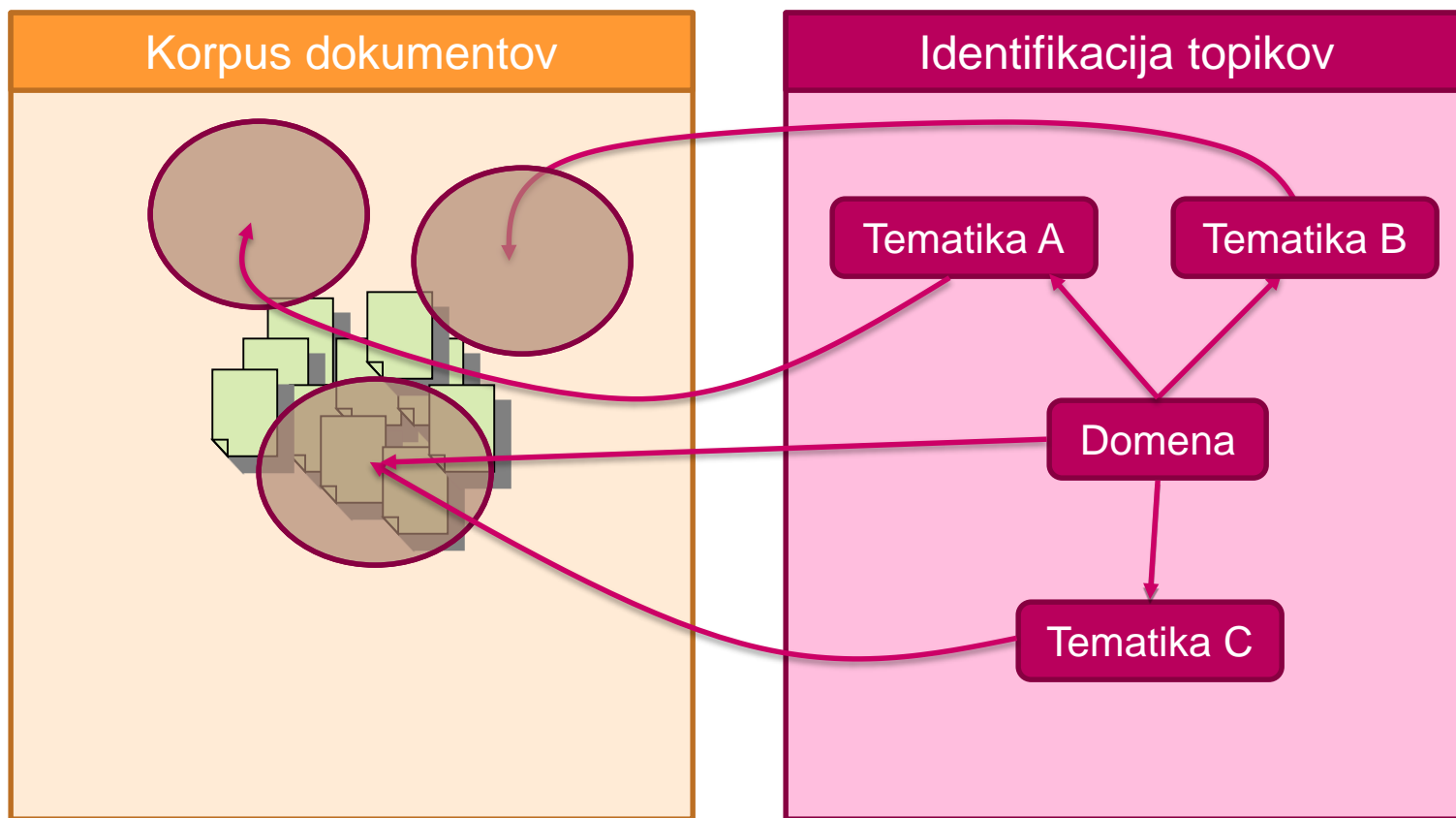
Primeri aplikacij

- Analiza sentimenta v tvitih
- Klasifikacija dokumentov
- Ugotavljanje avtorstva dokumentov
- Razvrščanje dokumentov v skupine glede na njihovo podobnost
- Pomoč pri preiskovanju svetovnega spleta
- Analiza profilov uporabnikov spleta
- Iskanje anomalij in osamelcev
- Iskanje skritih povezav med znanstvenimi domenami
- ...

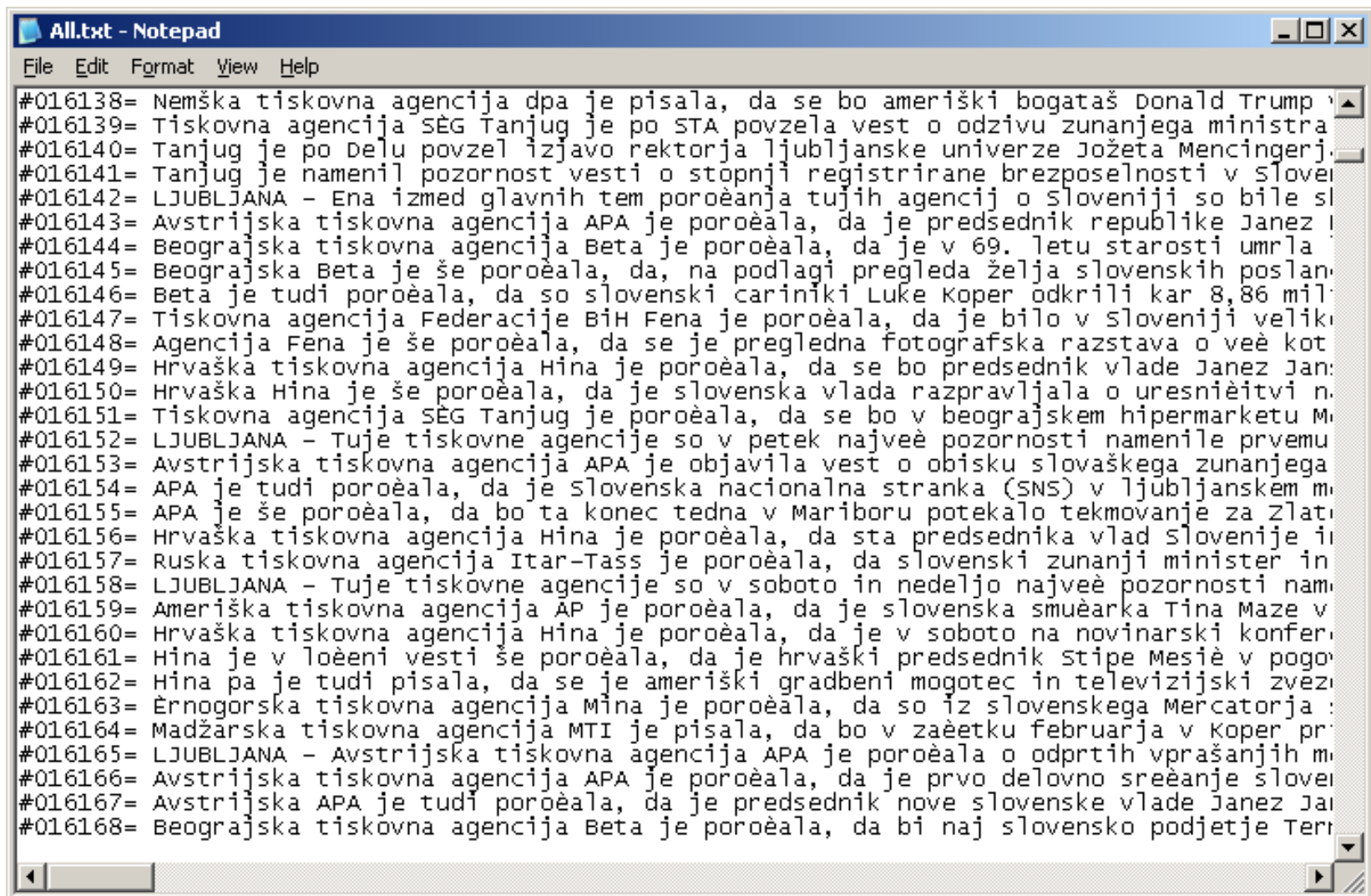
Primer: Analiza javnega mnenja z razvrščanjem novic STA v skupine

- Podatki: STA agencijske novice o Sloveniji povzete v tujih medijih
- Uporaba sistema za pol-avtomatsko razvrščanje v skupine
- V vsakem koraku postopka razcepimo množico dokumentov na k podskupin
 - Podatki znotraj skupine naj bodo čimbolj podobni
 - Podatki med skupinami naj bodo čim bolj različni
- Pri razvrščanju uporabimo mero podobnosti ali razdalje med dokumenti

Razvrščanje člankov v skupine

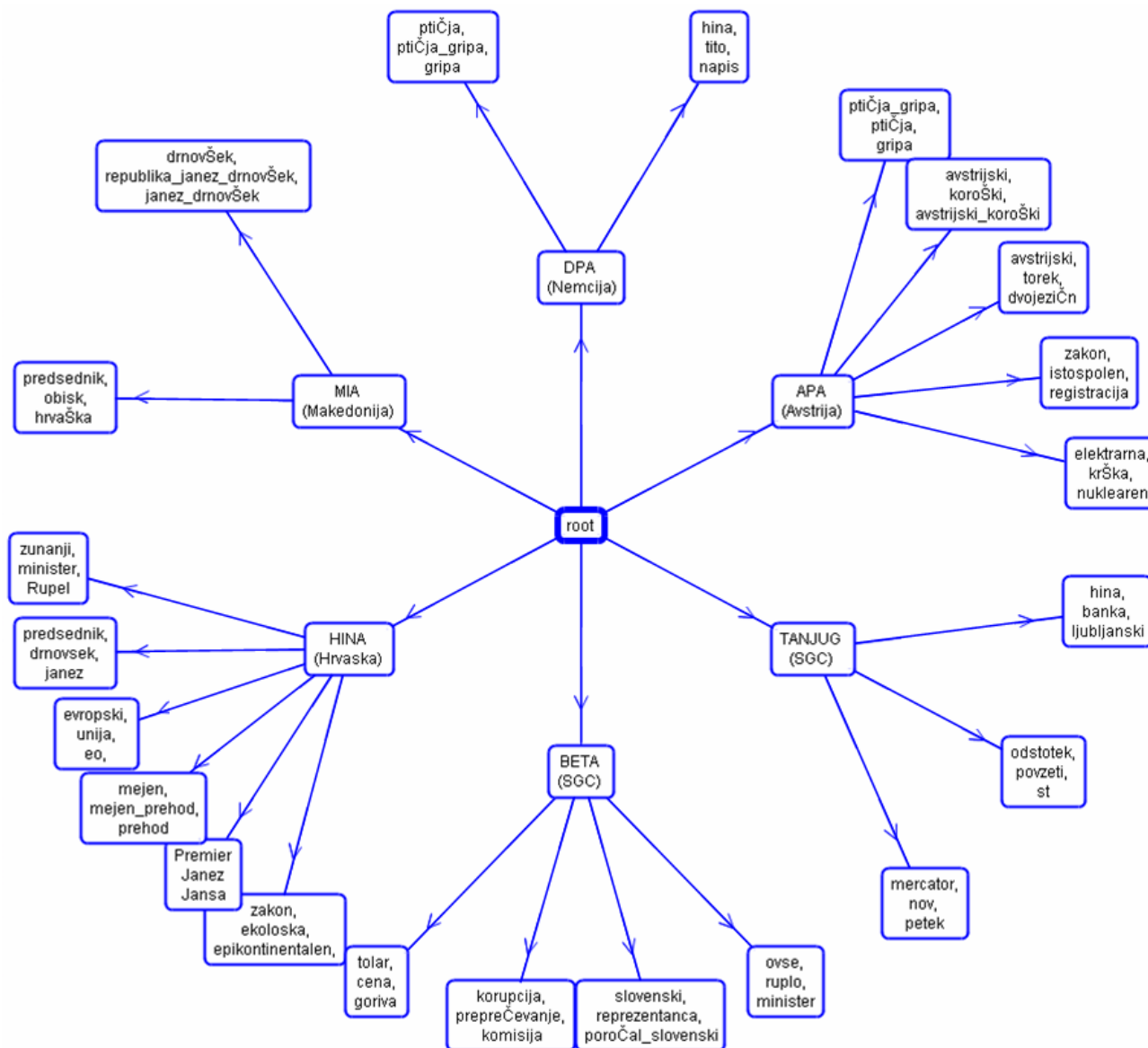


Vhod v program



```
All.txt - Notepad
File Edit Format View Help
#016138= Nemška tiskovna agencija dpa je pisala, da se bo ameriški bogataš Donald Trump
#016139= Tiskovna agencija S&E Tanjug je po STA povzela vest o odzivu zunanjega ministra
#016140= Tanjug je po Delu povzel izjavo rektorja ljubljanske univerze Jožeta Mencingerj.
#016141= Tanjug je namenil pozornost vesti o stopnji registrirane brezposelnosti v Sloveni
#016142= LJUBLJANA - Ena izmed glavnih tem poro&anja tujih agencij o Sloveniji so bile sl
#016143= Avstrijska tiskovna agencija APA je poro&ala, da je predsednik republike Janez I
#016144= Beograjska tiskovna agencija Beta je poro&ala, da je v 69. letu starosti umrla
#016145= Beograjska Beta je še poro&ala, da, na podlagi pregleda želja slovenskih poslan
#016146= Beta je tudi poro&ala, da so slovenski cariniki Luke Koper odkrili kar 8,86 mil
#016147= Tiskovna agencija Federacije BiH Fena je poro&ala, da je bilo v Sloveniji veliki
#016148= Agencija Fena je še poro&ala, da se je pregledna fotografska razstava o ve& kot
#016149= Hrvaška tiskovna agencija Hina je poro&ala, da se bo predsednik vlade Janez Jan
#016150= Hrvaška Hina je še poro&ala, da je slovenska vlada razpravljala o uresni&itvi n.
#016151= Tiskovna agencija S&E Tanjug je poro&ala, da se bo v beograjskem hipermarketu M
#016152= LJUBLJANA - Tuje tiskovne agencije so v petek najve& pozornosti namenile prvemu
#016153= Avstrijska tiskovna agencija APA je objavila vest o obisku slovaškega zunanjega
#016154= APA je tudi poro&ala, da je slovenska nacionalna stranka (SNS) v ljubljanskem m
#016155= APA je še poro&ala, da bo ta konec tedna v Mariboru potekalo tekmovanje za Zlati
#016156= Hrvaška tiskovna agencija Hina je poro&ala, da sta predsednika vlad Slovenije in
#016157= Ruska tiskovna agencija Itar-Tass je poro&ala, da slovenski zunanji minister in
#016158= LJUBLJANA - Tuje tiskovne agencije so v soboto in nedeljo najve& pozornosti nam
#016159= Ameriška tiskovna agencija AP je poro&ala, da je slovenska smu&arka Tina Maze v
#016160= Hrvaška tiskovna agencija Hina je poro&ala, da je v soboto na novinarski konfer
#016161= Hina je v lo&eni vesti še poro&ala, da je hrvaški predsednik stipe Mesi& v pogo
#016162= Hina pa je tudi pisala, da se je ameriški gradbeni mogotec in televizijski zvez
#016163= &rnogorska tiskovna agencija Mina je poro&ala, da so iz slovenskega Mercatorja :
#016164= Madžarska tiskovna agencija MTI je pisala, da bo v za&etku februarja v Koper pr
#016165= LJUBLJANA - Avstrijska tiskovna agencija APA je poro&ala o odprtih vpra&anjih m
#016166= Avstrijska tiskovna agencija APA je poro&ala, da je prvo delovno sre&anje slove
#016167= Avstrijska APA je tudi poro&ala, da je predsednik nove slovenske vlade Janez Jan
#016168= Beograjska tiskovna agencija Beta je poro&ala, da bi naj slovensko podjetje Ter
```

Analiza STA novic o Sloveniji v tujini



Znanstvena literatura kot vir znanja

Primer:

- Biomedicinska bibliografska podatkovna baza PubMed
- US National Library of Medicine
- Več kot 21M citatov
- Več kot 5.600 revij
- 2.000–4.000 referenc dodanih na vsak delovni dan!

The screenshot displays the PubMed website interface. At the top, the NCBI logo and 'PubMed' branding are visible, along with the text 'A service of the National Library of Medicine and the National Institutes of Health' and the URL 'www.pubmed.gov'. A search bar contains the query 'autism', with 'Go', 'Clear', and 'Save Search' buttons. Below the search bar, there are tabs for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. The search results are displayed in a list format, showing the number of items (11008) and the number of reviews (1632). The first three results are highlighted:

- 1:** [Fazzi E, Rossi M, Signorini S, Rossi G, Bianchi PE, Lanzi G](#) Related Articles
Leber's congenital amaurosis: is there an autistic component?
Dev Med Child Neurol. 2007 Jul;49(7):503-7.
AbstractID: 17592121 [PubMed - in process]
- 2:** [Paya B, Fuentes N](#) Related Articles
Neurobiology of autism: neuropathology and neuroimaging studies.
Actas Esp Psiquiatr. 2007 Jul-Aug;35(4):271-6.
PMID: 17592791 [PubMed - in process]
- 3:** [Hayashi ML, Rao BS, Seo JS, Choi HS, Dolan BM, Choi SY, Chattarji S, Tonegawa S](#) Related Articles
Inhibition of p21-activated kinase rescues symptoms of fragile X syndrome in mice.
Proc Natl Acad Sci U S A. 2007 Jun 25; [Epub ahead of print]
PMID: 17592139 [PubMed - as supplied by publisher]

The interface also includes a sidebar with navigation options such as 'About Entrez', 'Text Version', 'Entrez PubMed', 'PubMed Services', and 'Related Resources'. The bottom of the page shows the page number '1' of 23 Next.

Iskanje bisociativnih povezav med članki iz različnih podpodročij

Argument 1 (magnesium literature)

- Mg is a natural calcium channel blocker.
- Stress and Type A behavior can lead to body loss of Mg.
- Magnesium has anti-inflammatory properties.
- . . .

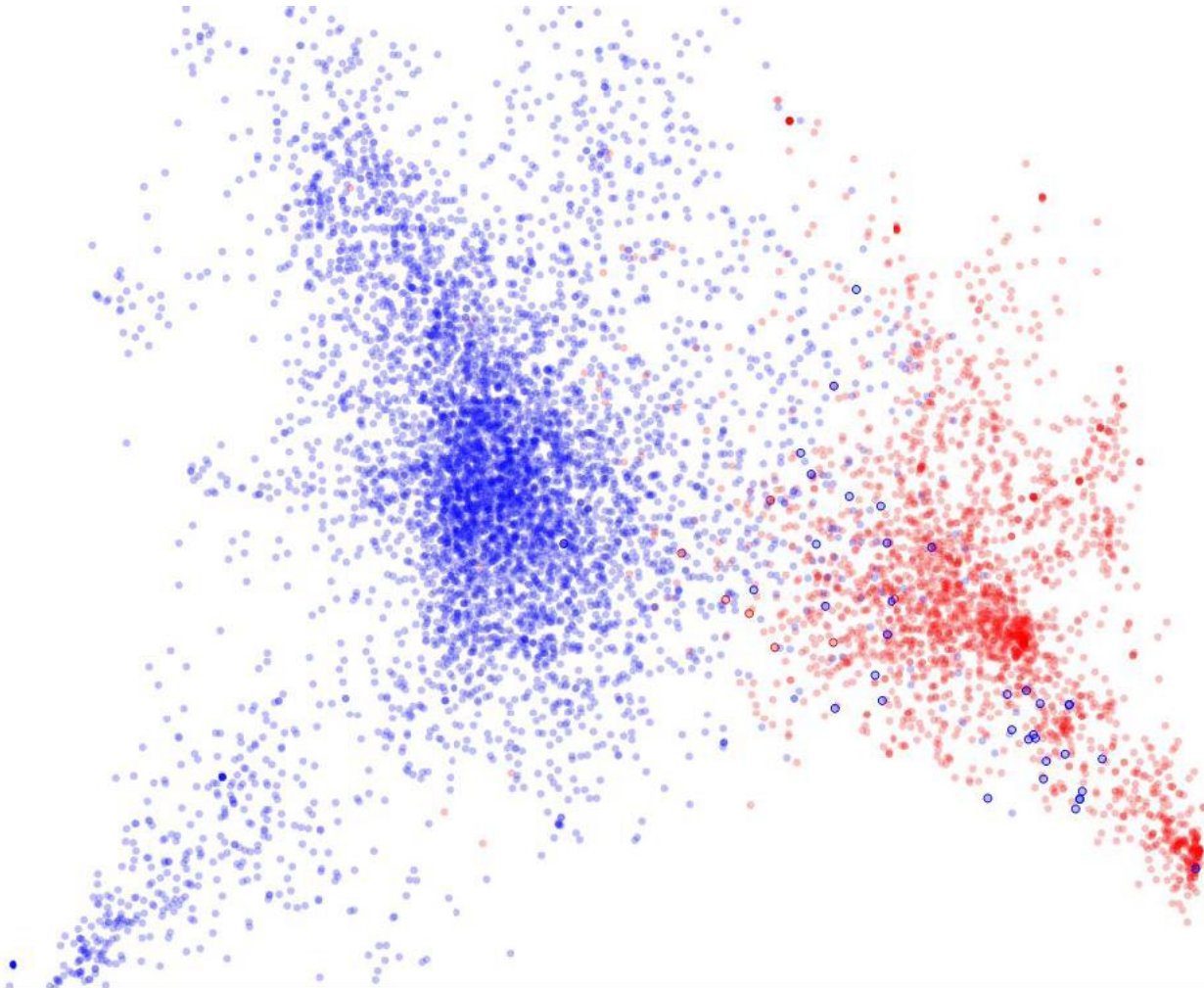
Argument 2 (migraine literature)

- Calcium channel blockers can prevent migraine attacks.
- Stress and Type A behavior are associated with migraine.
- Migraine may involve sterile inflammation of the cerebral blood vessels.
- . . .

Iskanje izjem in anomalij v podatkih



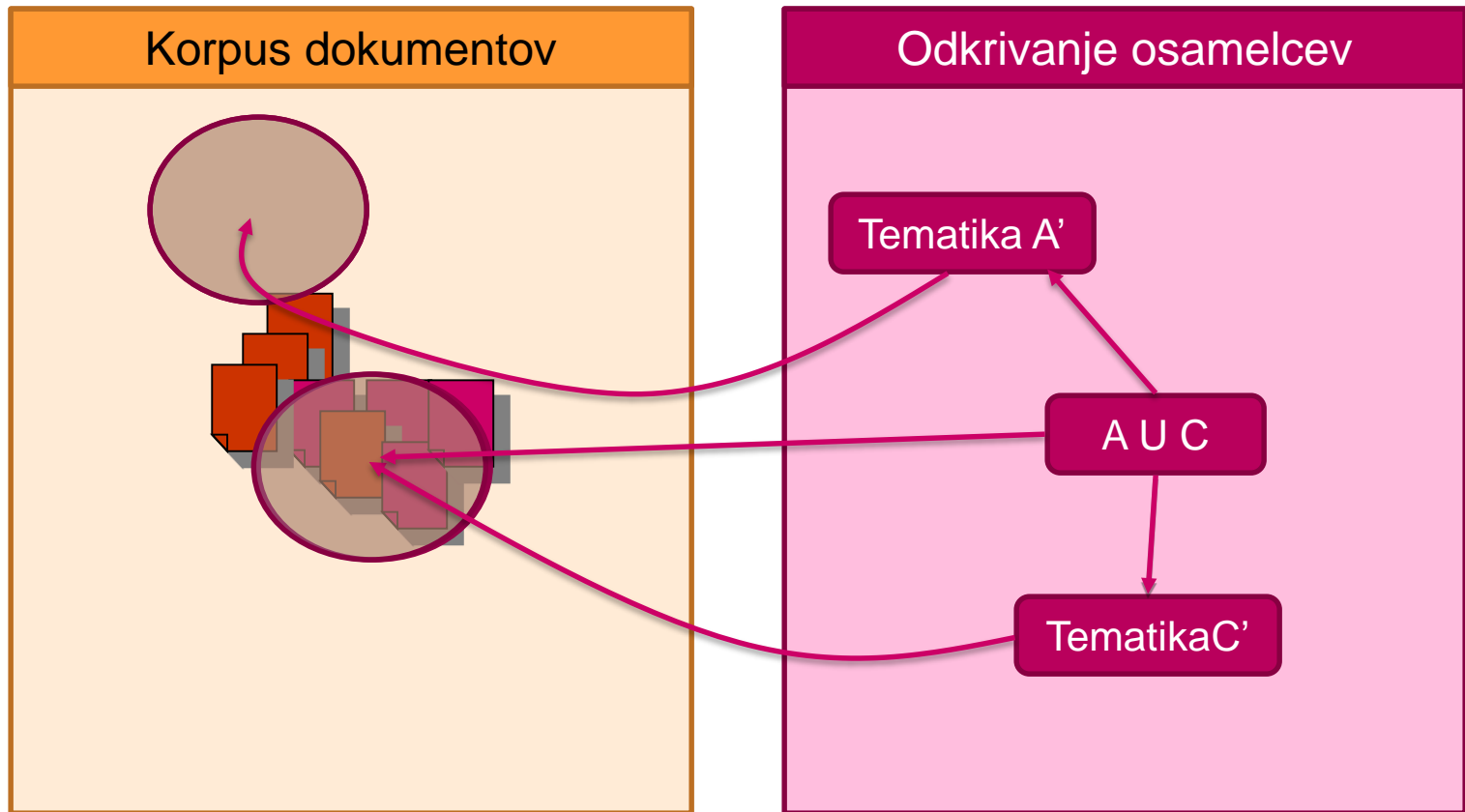
Iskanje izjem v korpusih besedil



2-dimenzionalna
projekcija
dokumentov z
dveh različnih
domen

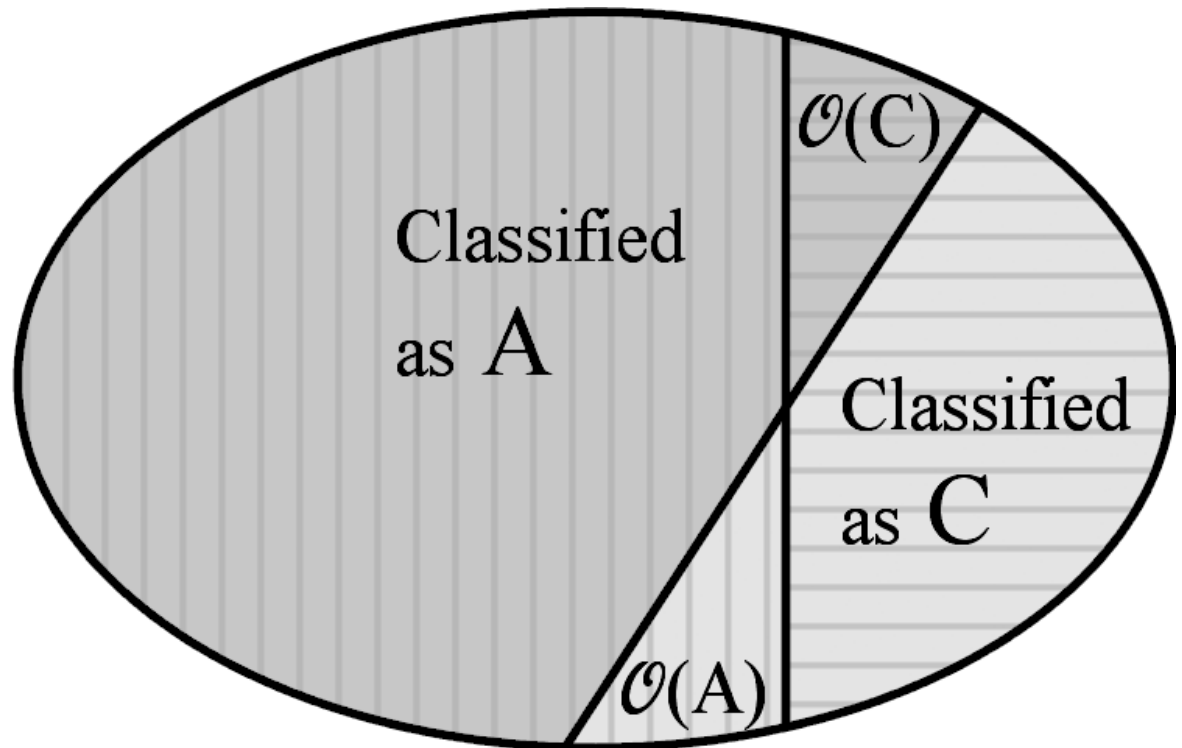
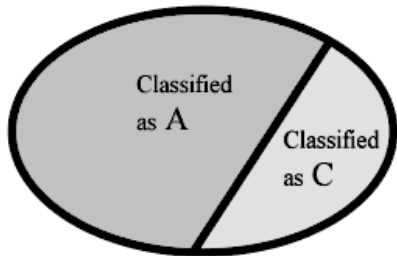
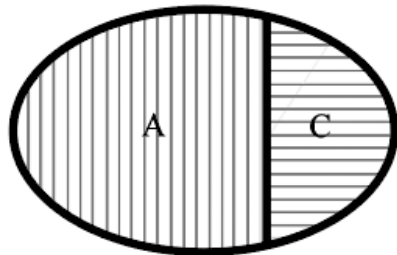
**Iskanje
nenavadnih
člankov v
korpusih**

Uporaba razvrščanja v skupine za iskanje izjem v korpusih besedil



Klasifikacijski pristop za iskanje izjem v korpusih

- Iz podatkov se naučimo klasifikacijski model in ga uporabimo nad dokumenti iz dveh ločenih korpusov



Rangiranje izjem v člankih o Kenijskih volitvah, objavljenih v lokalnih in zahodnih medijih

- **Članek 352**

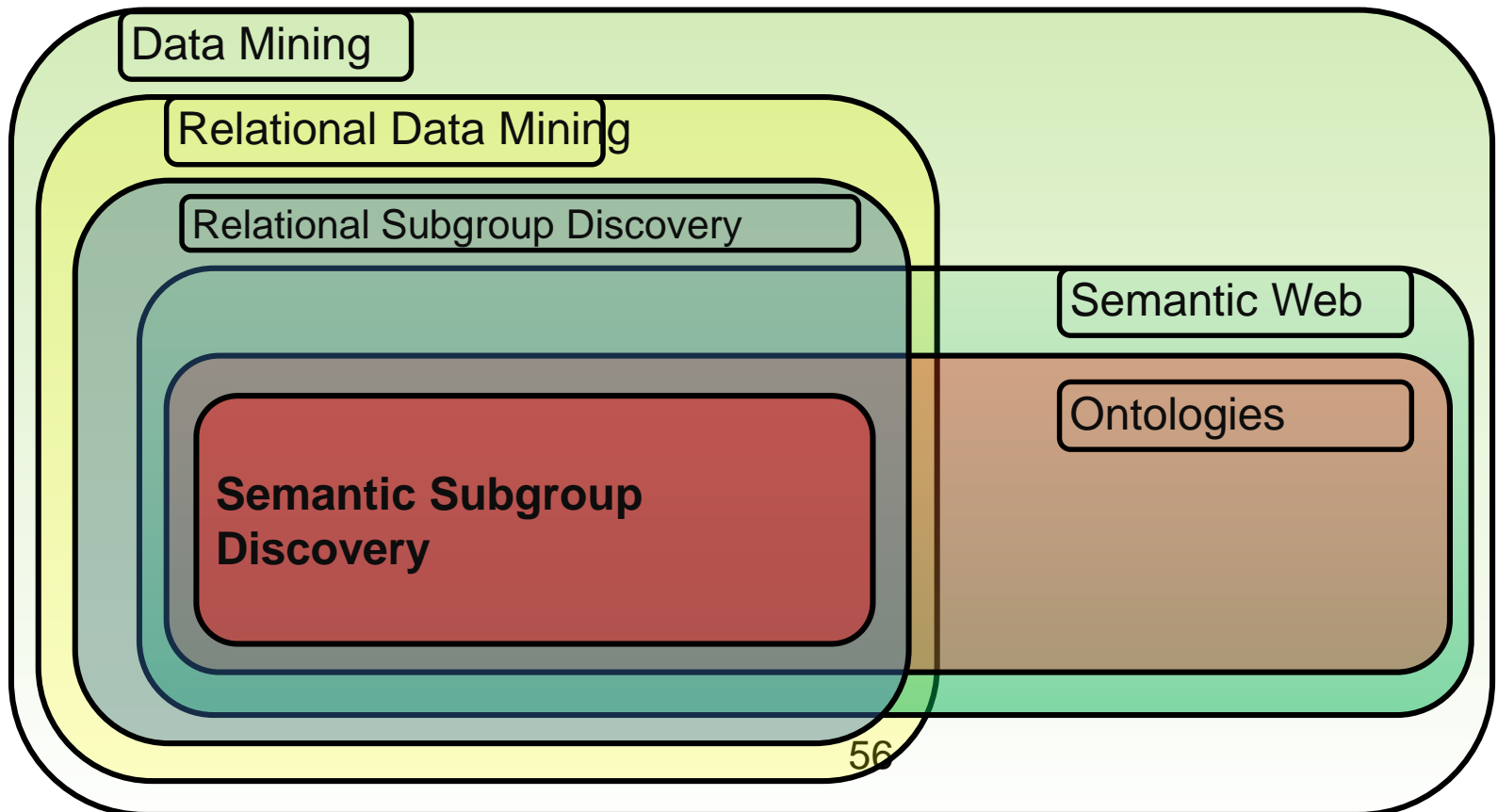
Ta članek so naknadno odstranili iz korpusa, ker ni govoril o politični in socialni problematiki, temveč o oropnem britanskem turistu..

- **Članek 173**

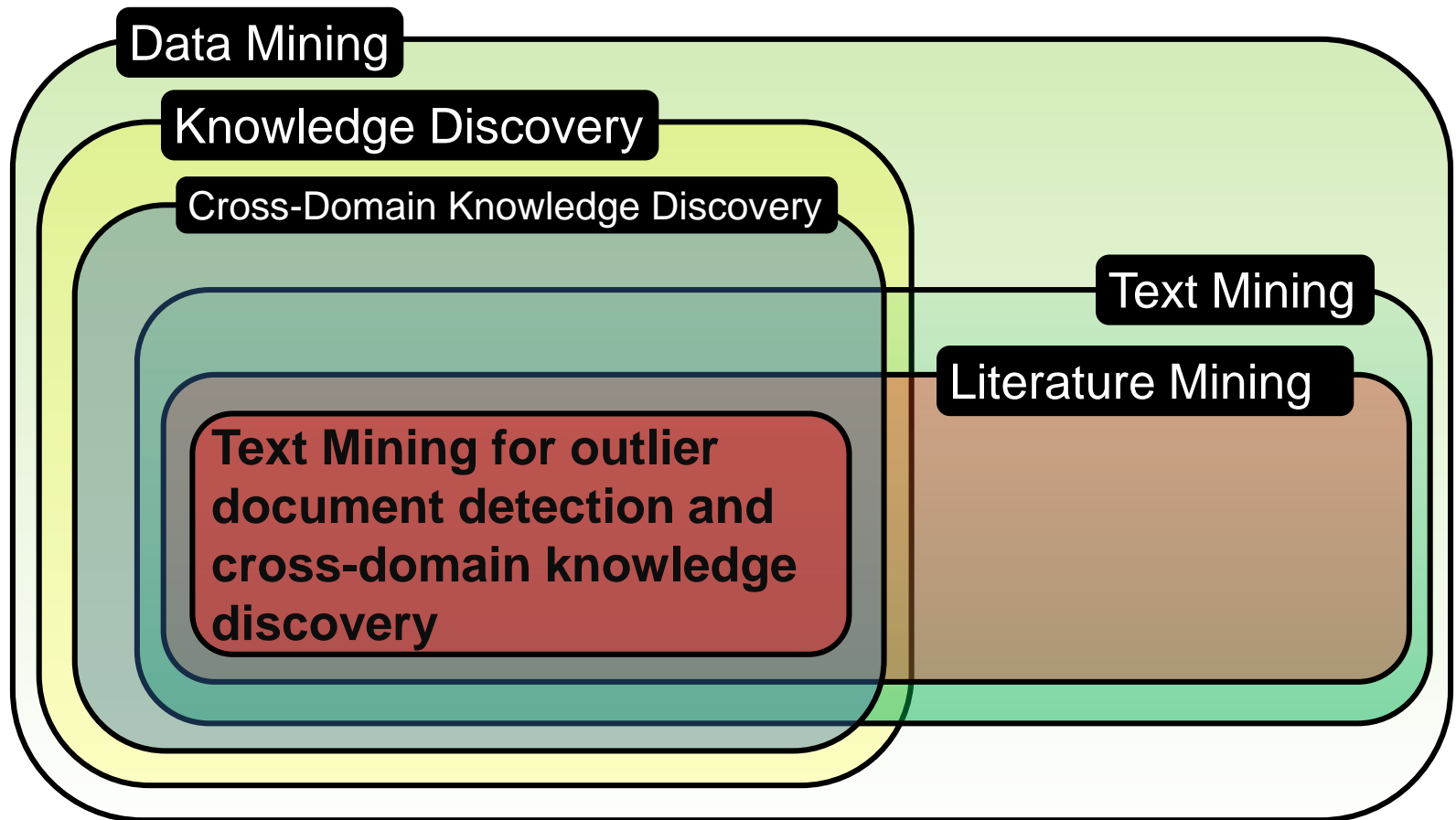
Ta članek je bil objavljen v lokalnem mediju a je bil avtomatsko klasificiran kot zahodni članek. Članek je dejansko napisal zahodni novinar, ki nima ustrezne senzitivnosti za pisanje (npr. omenja besede “pleme”, “plemenska pripadnost” itd. v negativnem kontekstu) – avtor je torej uporabljal zahodni stil pisanja.

Povzetek in tekoče delo

Semantično odkrivanje podskupin v podatkih (Vavpetič et al., 2012)

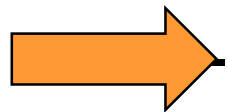


Povzetek in tekoče delo



Vsebina predavanja

- odkrivanje zakonitosti in anomalij v podatkih
- rudarjenje tekstovnih in spletnih podatkov
- računalniška podpora pri odločanju
- jezikovne tehnologije in računalniško jezikoslovje



računalniška kreativnost

Kaj je računalniška kreativnost

- Računalniška kreativnost (umetna kreativnost, mehanska kreativnost ali kreativno računstvo) je multidisciplinarno področje znanosti na preseku umetne inteligence, kognitivne psihologije, filozofije in umetnosti.
- Cilji področja so modeliranje in simulacija človeške ustvarjalnosti kot tudi računalniško kreativno ustvarjanje in/ali izboljšanje človeške lastne kreativnosti s pomočjo robota/računalnika



Naši tekoči projekti s področja računalniške kreativnosti

- **MUSE – Machine Understanding For Interactive Storytelling (STREP, 2012-2015)**
 - Koordinatorica: Sien Moens, KU Leuven
- **PROSECCO – Promoting The Scientific Exploration Of Computational Creativity (CSA, 2013-2016)**
 - koordinator Tony Veale, UCD, Dublin
- **Concrete – Concept Creation Technology (STREP, 2013-2016)**
 - Koordinator Geraint Wiggins, QMUL, London
- **WHIM – The What-if Machine (STREP, 2013-2016)**
 - Koordinator Simon Colton, Goldsmith Uni., London
- **Skupaj ~ 1.1 Mio EUR evropskih sredstev** pridobljenih za naše delo v teh evropskih projektih

Primer projekta: MUSE

Patient Guidelines

On completion of the preparatory phase, the multidisciplinary team makes one of three decisions:

- 1 The operation can go ahead. The team will then give you more information on the operative technique chosen, if you have decided to have the operation, you will be given an operation date and a request for your health insurance fund to agree to help with the operation costs (to find out more www.ameli.fr).
- 2 Your preparation for the operation is not sufficient. You will have to undertake additional preparations. On completion of these, the multidisciplinary team will re-examine your request and make a new decision.
- 3 Surgery is not suitable in your case. The multidisciplinary team will explain the reasons why and offer you another treatment (non-surgical).

NLP Analysis

Connexor Analysis



Feature Structures Extraction

```
(:MAIN "re-examine"
  (:OBJ "decision"
    (:ATTR "new")
    (:OBJ "request"
      (:ATTR "you")
      (:SUBJ "team"
        (:ATTR
          "multidisciplinary"))
      ("completion")))))
```

Plan-based Narrative Engine

Narrative Actions



World State

```
... (...) (can-provide-info drnicholson anaesthesia) (...)
(can-provide-info drdowell digestive-surgery) (...) (can-provide-info drrichards nutrition)
(...) (at-location mrroberts reception) (at-location drsmith or-room) (...)
(at-location msjones patient-room) (...) ...
```

update

staging



3D Visualization (Unreal® Engine - UDK)

Primer projekta: MUSE

- **MUSE – Machine Understanding for Interactive Storytelling (2012-2015)**
 - 3D animacija teksta
 - z uporabo metod jezikovnih tehnologij omogočimo razumevanje teksta
 - rendering za 3D vizualizacijo zgodbe
 - planiranje za iskanje možnih razpletov zgodbe (what-if)
 - mehanizmi za igranje “resnih iger” v 3D svetu omogočajo uporabniku vodeno sledenje igri in interakcijo z virtualnim svetom
- **Uporabnost:** izobraževanje, medicina, letališka varnost, ..., industrija iger

Področja računalniške kreativnosti

- Kreiranje pripovedi, zgodb, metafor, analogij, šal, neologizmov, poezije, glasbe, slik in drugih vizualnih artefaktov, marketinških oglasov, ...
- Primer: robot lahko napiše zgodbo, **Kaj-če stroj (What-If Machine)** pa lahko spremeni njen potek
- Uporabnost v industriji in znanosti

Naši raziskovalni cilji

- **Naš cilj:** Kreativno reševanje problemov in kreativno odkrivanje povezav, vzorcev in znanja iz podatkov.
 - Bisociativno preiskovanje biomedicinske literature za odkrivanje novih meddomenskih povezav
 - Kreativno kombiniranje algoritmov in avtomatsko generiranje novih
- **Naš pogled:**
 - Ni nujno, da je robot tudi sam kreativen, lahko je le inteligen ten pomočnik
- **Večinski pogled:**
 - Računalnik/robot naj bo avtonomen kreator



Povzetek in zahvala sodelavcem odseka za soustvarjanje raziskovalnih poti v novo in še neraziskano

