

# Metody zbiorów przybliżonych w uczeniu się podobieństwa z wielowymiarowych zbiorów danych

Andrzej Janusz



WMIM, Uniwersytet Warszawski  
ul. Banacha 2, 02-097 Warszawa, Polska

*andrzejjanusz@gmail.com*

13.06.2013

# Dlaczego właśnie podobieństwo?

Myślenie...



... i formowanie  
pojęć

Podjęmowanie  
decyzji



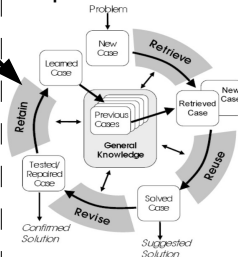
sampleID	AFFX_3_at	3322_l_at	4969_s_at	...	22095_s_at	22379_l_at	Diagnosis
GSM1	4.010	12.434	32.443	...	1.665	12.44	3
GSM2	5.314	43.765	5.763	...	3.567	7.645	2
GSM3	3.275	17.567	23.842	...	0.657	12.46	2
GSM4	2.112	8.432	54.849	...	87.656	45.32	1
...	...	...	...	...	...	...	...
GSM14	8.453	10.087	8.678	...	2.986	9.656	3

**Podobieństwo**

Uczenie się



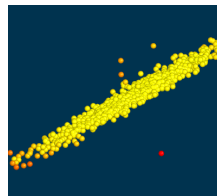
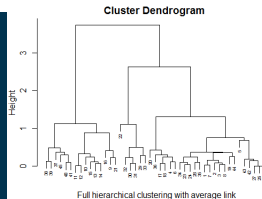
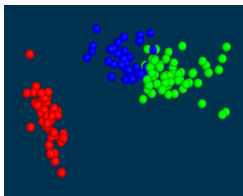
Rozwiązywanie  
problemów



# Zastosowania modeli podobieństwa

## Przykłady:

- klasyfikacja i regresja,
- segmentacja danych,
- planowanie, rozwiązywanie problemów,
- wykrywanie nietypowych obiektów,
- wizualizacja i streszczenie danych.



## Podstawowa zasada:

Podobne obiekty powinny być traktowane podobnie (np. należeć do tej samej klasy decyzyjnej, czy grupy).

# Czym tak naprawdę jest podobieństwo?

Trudności ze ścisłą definicją podobieństwa:

- relacja, czy funkcja?
- obiektywne, czy subiektywne?
- bezkontekstowe, czy kontekstowe?
- globalne, czy lokalne?

Czynniki, które wpływają na kontekst to:

- cel lub zadanie, któremu służy ewaluacja podobieństwa,
- wiedza o innych znanych obiektach.



# Czym tak naprawdę jest podobieństwo?

Trudności ze ścisłą definicją podobieństwa:

- relacja, czy funkcja?
- obiektywne, czy subiektywne?
- bezkontekstowe, czy kontekstowe?
- globalne, czy lokalne?

Czynniki, które wpływają na kontekst to:

- cel lub zadanie, któremu służy ewaluacja podobieństwa,
- wiedza o innych znanych obiektach.



## Model kontrastu cech:

- obiekty postrzegane są jako zbiory cech jakościowych,
- cechy są zazwyczaj na wyższym poziomie abstrakcji niż dane “sensoryczne”, np.  
*dwa samochody są podobne ponieważ są małe i szybkie,*
- ważne są zarówno wspólne jak i wyróżniające cechy obiektów,
- $S(a, b) = \theta f(A \cap B) - \alpha f(A \setminus B) - \beta f(B \setminus A)$ , gdzie  $\theta, \alpha, \beta \geq 0$

## Model kontrastu cech:

- obiekty postrzegane są jako zbiory cech jakościowych,
- cechy są zazwyczaj na wyższym poziomie abstrakcji niż dane “sensoryczne”, np.  
*dwa samochody są podobne ponieważ są małe i szybkie,*
- ważne są zarówno wspólne jak i wyróżniające cechy obiektów,
- $S(a, b) = \theta f(A \cap B) - \alpha f(A \setminus B) - \beta f(B \setminus A)$ , gdzie  $\theta, \alpha, \beta \geq 0$

Model Tversky-ego trudno jest zaaplikować do rzeczywistych danych:

- jak definiować wysokopoziomowe cechy?
- jak wybrać te istotne w danym kontekście?

Propozycja: można wykorzystać teorię zbiorów przybliżonych!

# Założenia proponowanego modelu podobieństwa:

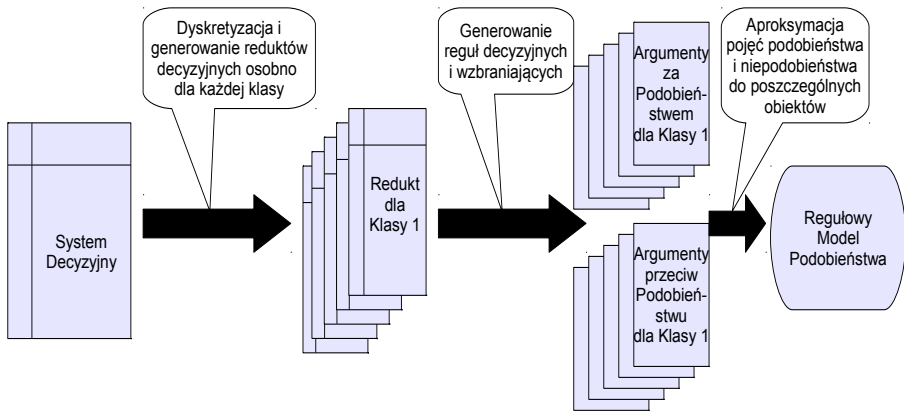
## Uczenie się podobieństwa w języku zbiorów przybliżonych:

wybór istotnych aspektów podobieństwa	↔	wybór przestrzeni aproksymacji
wysokopoziomowe cechy	↔	lewe strony reguł
agregacja argumentów za i przeciw podobieństwu	↔	aproksymacja pojęć bycia podobnym i niepodobnym do obiektu
funkcja podobieństwa	↔	funkcja przynależności do aproksymacji pojęcia

- Wysokopoziomowe cechy można traktować jak argumenty za lub przeciw podobieństwu obiektów!
- Aproksymacja podobieństwa do obiektu to zbiór obiektów, do który pasują argumenty za podobieństwem a nie pasują argumenty przeciwko.



# Konstrukcja proponowanego modelu podobieństwa



## Aproksymacja podobieństwa i niepodobieństwa:

$F_{(i)}^+$  oraz  $F_{(i)}^-$  – zbiory cech dla  $i$ -tej klasy decyzyjnej, wyznaczone przez reguły decyzyjne i wzbraniające;

$$F_{(i)}^+ = \{f : (f \rightarrow (d = i)) \in \text{RuleSet}(DR_i)\};$$

$$F_{(i)}^- = \{f : (f \rightarrow (d \neq i)) \in \text{RuleSet}(DR_i)\};$$

$$SIM_{(i)}(u) = \bigcup_{f \in F_{(i)}^+ \wedge f(u)=1} [u]_f \quad \left| \quad DIS_{(i)}^0(u) = \bigcup_{f \in F_{(i)}^- \wedge f(u)=0} U \setminus [u]_f \quad \right| \quad DIS_{(i)}^1(u) = \bigcup_{f \in F_{(i)}^- \wedge f(u)=1} [u]_f$$

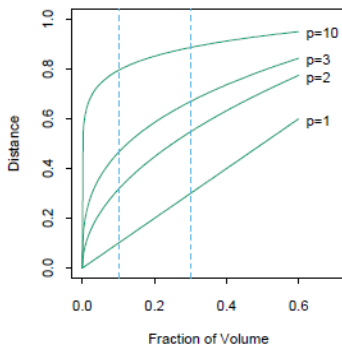
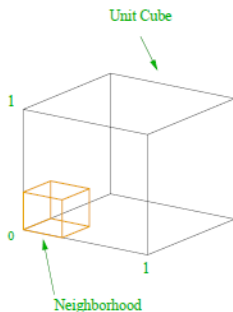
Przynależność do  $SIM_{d(u_1)}(u_1)$ :

$$\mu(u_2, SIM_i(u_1)) = \frac{|SIM_i(u_1) \cap SIM_i(u_2)|}{|SIM_i(u_1)|}$$

Przynależność do  $DIS_{d(u_1)}^0(u_1)$ :

$$\psi(u_2, DIS_i^0(u_1)) = \frac{|DIS_i^0(u_1) \cap DIS_i^1(u_2)|}{|DIS_i^0(u_1)|}$$

# Dlaczego dane wielowymiarowe?



**Rysunek:** Ilustracja “przekleństwa wielu wymiarów” (z książki *Elements of Statistical Learning: Data Mining, Inference and Prediction*).

- typowe metody nie radzą sobie z problemem *niewielu obiektów o dużej liczbie cech*,
- duża złożoność obliczeniowa algorytmów uczenia się podobieństwa z danych wielowymiarowych.

# Rozszerzenia modelu dla danych wielowymiarowych

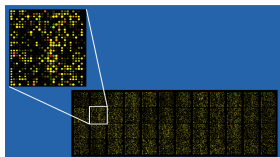
## Główna idea:

W przypadku danych wielowymiarowych konieczne jest rozpatrywanie wielu lokalnych modeli podobieństwa, które można interpretować jako autonomicznych agentów z własnymi preferencjami i doświadczeniem.

## Dwa typy wielowymiarowych danych

Dane mikromacierzowe:	Dane tekstowe:
uczenie z nadzorem	uczenie bez nadzoru
redukty dynamiczne	biredukty informacyjne
reguły decyzyjne i wzbraniające	pojęcia z ontologii dziedzinowej

# Opis eksperymentów na danych mikromacierzowych



**Microarray data:**  
*few-objects-many-attributes* problem

≈40k genes (attributes)

sampleID	AFFX-3_at	3322_i_at	4969_s_at	...	22095_s_at	22379_at	Diagnosis
GSM1.CEL	4.010	12.434	32.443	...	1.665	12.434	3
GSM2.CEL	5.314	43.765	5.763	...	3.567	7.645	2
GSM3.CEL	3.275	17.567	23.842	...	0.657	12.446	2
GSM4.CEL	2.112	8.432	54.849	...	87.656	45.324	1
...	...	...	...	...	...	...	...
GSM149.CEL	8.453	10.087	8.678	...	2.986	9.656	3

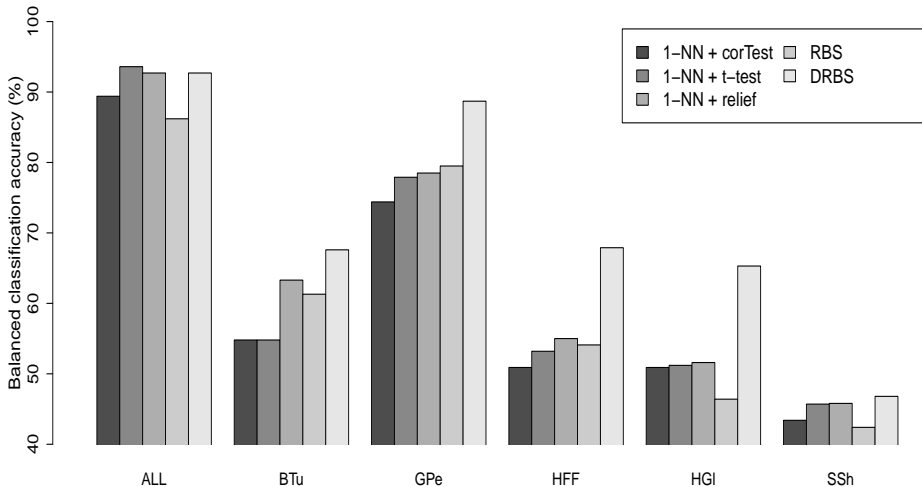
## Opis danych

- 11 zbiorów mikromacierzy,
- liczba obiektów: 124 – 284,
- liczba atrybutów: 22k – 61k,
- zbiory pochodzą z repozytorium ArrayExpress.

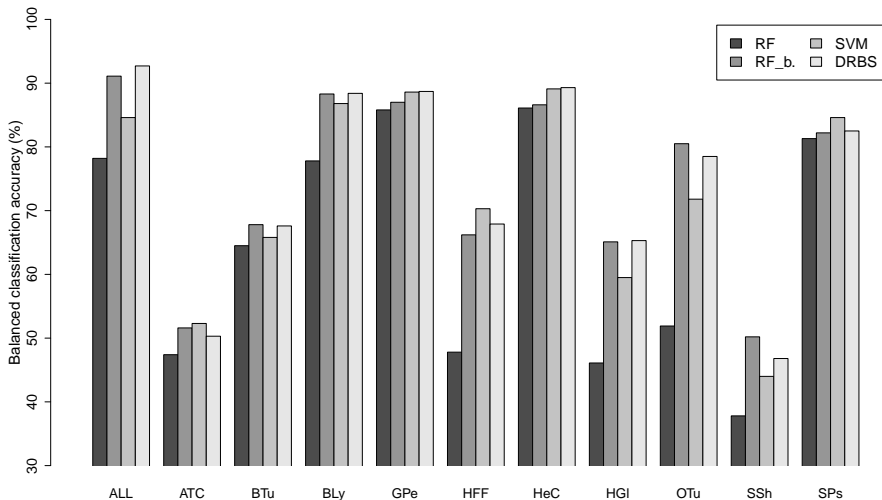
## Opis eksperymentu

- wielokrotnie powtarzana weryfikacja krzyżowa,
- miary jakości: ACC i BAC,
- porównywane klasyfikatory:  $k$ -NN\*, RF, SVM.

# Wyniki porównania z wybranymi modelami podobieństwa



# Wyniki porównania z wybranymi metodami klasyfikacji



# Opis eksperymentów na danych tekstowych

## Opis danych

- zbior 1000 artykułów naukowych z repozytorium PubMed Central,
- ontologia dziedzinowa MeSH ( $\approx 26k$  pojęć),
- metoda automatycznego etykietowania: ESA,
- zbiory etykiet nadanych przez ekspertów.

## Opis eksperymentu

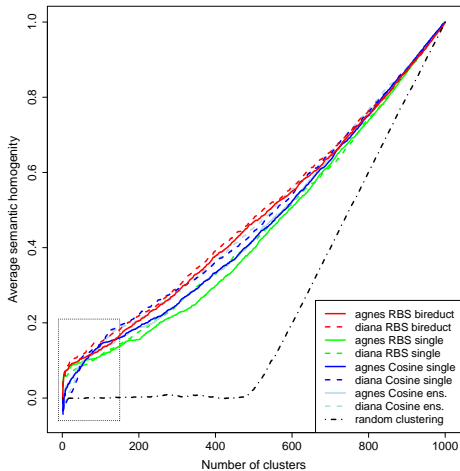
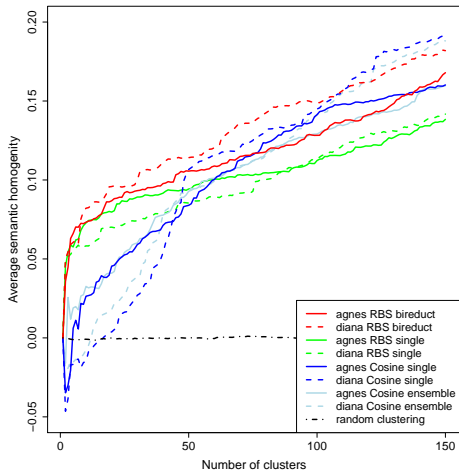
- grupowanie hierarchiczne artykułów,
- stosowane algorytmy: *agnes* i *diana*,
- porównywane modele: dwa oparte o miarę kosinusową,
- zewnętrzna miara oceny jakości grupowania.

## Ewaluacja wyników

Miara zgodności etykiet nadanych przez ekspertów wewnątrz grup.



# Wyniki ewaluacji modelu



## Co się udało?

- dokonano interpretacji problemu uczenia się podobieństwa z punktu widzenia teorii zbiorów przybliżonych,
- zaproponowano intuicyjny i elastyczny model uczenia się podobieństwa z danych,
- opracowano efektywne algorytmy działające dla wielowymiarowych zbiorów danych,
- przeprowadzono dokładną ewaluację zaproponowanego podejścia.

## Kierunki na przyszłość:

- lepsze wykorzystanie wiedzy dziedzinowej,
- optymalizacja wydajności obliczeniowej dla dużych zbiorów danych,
- stworzenie wysokopoziomowego środowiska do eksperymentów.

Dziękuję za uwagę!