Haskell na produkcji

czyli czy warto sięgać po niestandardowe narzędzia

Tomasz Tylec 11 listopada 2018

IF Research Polska (Gdańsk) / Mazars-Zettafox (Paris) | tomasz.tylec@ifresearch.pl

Outline

Reguły klasyfikujące

Implementacja w Haskellu

Dlaczego Haskell?

Kluczowe elementy

Dlaczego nie Python

Integracja z ekosystemem Pythona

Podsumowanie

Reguły klasyfikujące

Reguly

Definicja

Regułą nazywamy predykat, wyrażenie logiczne, które każdemu elementowi w zbiorze danych przypisuje wartość prawda lub fałsz.

Przykład

Dane:

index	А	В	C
0	low	high	low
1	low	low	low
2	high	low	high
3	medium	low	high

Reguly

Definicja

Regułą nazywamy predykat, wyrażenie logiczne, które każdemu elementowi w zbiorze danych przypisuje wartość *prawda* lub *fałsz*.

Przykład

Dane:

index	A	В	C
0	low	high	low
1	low	low	low
2	high	low	high
3	medium	low	high

Reguła: $A \in \{low, medium\} \land B \in \{low\}$

Rozważać będziemy reguły postaci:

$$c_1 \wedge c_2 \wedge \ldots \wedge c_n$$

gdzie c_i jest postaci $X \in \{...\}$, lub $X \in (a, b)$.

Reguly

Definicja

Regułą nazywamy predykat, wyrażenie logiczne, które każdemu elementowi w zbiorze danych przypisuje wartość *prawda* lub *fałsz*.

Przykład

Dane:

index	А	В	C
0	low	high	low
1	low	low	low
2	high	low	high
3	medium	low	high

Reguła: $A \in \{low, medium\} \land B \in \{low\}$

Rozważać będziemy reguły postaci:

$$c_1 \wedge c_2 \wedge \ldots \wedge c_n$$

gdzie c_i jest postaci $X \in \{...\}$, lub $X \in (a, b)$.

Algebra reguł

Reguły tworzą zbiór częściowo uporządkowany

 $r \leq q$ jeśli dla każdego wiersza dla którego qjest prawdziwe, rteż jest prawdziwe. Przykłady:

- \cdot $A \in \{low\} \land B \in \{low\} \le A \in \{low, medium\} \land B \in \{low\}$
- · ale $A \in \{low\}$ oraz $B \in \{high\}$ jest są w relacji porządku.

Algebra reguł

Reguły tworzą zbiór częściowo uporządkowany

 $r \leq q$ jeśli dla każdego wiersza dla którego q jest prawdziwe, r też jest prawdziwe. Przykłady:

- \cdot $A \in \{low\} \land B \in \{low\} \le A \in \{low, medium\} \land B \in \{low\}$
- ale $A \in \{low\}$ oraz $B \in \{high\}$ jest są w relacji porządku.

Reguły można ze sobą łączyć

Jeśli r, q są regułą to $r \land q$ też jest regułą. Wtedy $r \land q \leq r$ oraz $r \land q \leq q$.

Rozważać będziemy reguły postaci:

$$c_1 \wedge c_2 \wedge \ldots \wedge c_n$$

gdzie c_i jest postaci $X \in \{...\}$, lub $X \in (a, b)$.

Jeśli r,q są regułą to $r\lor q$ definiujemy jako najmniejszą (w sensie \leq) regułę spełniającą: $r\le r\lor q$ oraz $q\le r\lor q$.

Zadanie

Mając do dyspozycji zbiór sklasyfikowanych danych, znaleźć niewielki zbiór *najlepszych* reguł, względem pewnej funkcji oceniającej (*score function*).

Zaczynając od zbioru najogólniejszych nietrywialnych reguł a_1,a_2,\ldots,a_k , tworzymy reguły postaci:

- · $a_2 \wedge a_5$,
- $a_5 \wedge a_{10} \wedge a_{143}$,
- $a_{242} \wedge a_{743} \wedge a_{342} \wedge a_{1233}$,
- · itd.

wybierając tylko te najlepsze.

Zaczynając od zbioru najogólniejszych nietrywialnych reguł a_1,a_2,\ldots,a_k , tworzymy reguły postaci:

- · $a_2 \wedge a_5$,
- $a_5 \wedge a_{10} \wedge a_{143}$,
- $a_{242} \wedge a_{743} \wedge a_{342} \wedge a_{1233}$,
- · itd.

wybierając tylko te najlepsze.

index	A	В	С
0	low	hig	low
1	low	low	low
3	high	low	high
4	medium	low	high

$$a_1 = A \in \{\text{low, medium}\}, a_3 = B \in \{\text{low}\}\$$

 $r = a_1 \land a_3 = A \in \{\text{low, medium}\} \land B \in \{\text{low}\}\$

 $\mathsf{Dom}(A) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

 $Dom(B) = \{low, high\}$

 $\mathsf{Dom}(\mathit{C}) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

Zaczynając od zbioru najogólniejszych nietrywialnych reguł a_1,a_2,\ldots,a_k , tworzymy reguły postaci:

- · $a_2 \wedge a_5$,
- $a_5 \wedge a_{10} \wedge a_{143}$,
- $a_{242} \wedge a_{743} \wedge a_{342} \wedge a_{1233}$,
- · itd.

wybierając tylko te najlepsze.

index	A	В	C
0	low	hig	low
1	low	low	low
3	high	low	high
4	medium	low	high

$$a_1 = A \in \{\text{low, medium}\}, a_3 = B \in \{\text{low}\}\$$

 $r = a_1 \land a_3 = A \in \{\text{low, medium}\} \land B \in \{\text{low}\}\$

 $\mathsf{Dom}(A) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

 $Dom(B) = \{low, high\}$

 $\mathsf{Dom}(\mathit{C}) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

Zaczynając od zbioru najogólniejszych nietrywialnych reguł a_1,a_2,\ldots,a_k , tworzymy reguły postaci:

- · $a_2 \wedge a_5$,
- $a_5 \wedge a_{10} \wedge a_{143}$
- $a_{242} \wedge a_{743} \wedge a_{342} \wedge a_{1233}$,
- · itd.

wybierając tylko te najlepsze.

index	A	В	С
0	low	hig	low
1	low	low	low
3	high	low	high
4	medium	low	high

$$a_1 = A \in \{\text{low, medium}\}, a_3 = B \in \{\text{low}\}\$$

 $r = a_1 \land a_3 = A \in \{\text{low, medium}\} \land B \in \{\text{low}\}\$

 $\mathsf{Dom}(A) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

 $Dom(B) = \{low, high\}$

 $\mathsf{Dom}(\mathit{C}) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

Zaczynając od zbioru najogólniejszych nietrywialnych reguł a_1,a_2,\ldots,a_k , tworzymy reguły postaci:

- · $a_2 \wedge a_5$,
- $a_5 \wedge a_{10} \wedge a_{143}$
- $a_{242} \wedge a_{743} \wedge a_{342} \wedge a_{1233}$,
- · itd.

wybierając tylko te najlepsze.

index	A	В	С
0	low	hig	low
1	low	low	low
3	high	low	high
4	medium	low	high

$$a_1 = A \in \{\text{low, medium}\}, a_3 = B \in \{\text{low}\}\$$

 $r = a_1 \land a_3 = A \in \{\text{low, medium}\} \land B \in \{\text{low}\}\$

 $Dom(A) = \{low, medium, high\}$

 $Dom(B) = \{low, high\}$

 $\mathsf{Dom}(\mathit{C}) = \{\mathsf{low}, \mathsf{medium}, \mathsf{high}\}$

Każdy element zbioru danych x_i odpowiada regule $r(x_i)$ – takiej, która wybiera dokładnie ten element. Tworzymy reguły postaci:

- $\cdot r(x_{10}) \vee r(x_{31}),$
- $\cdot r(x_1) \vee r(x_{45}) \vee r(x_{321}) \vee r(x_{1421}) \vee r(x_{2352}) \vee r(x_{31002})$
- · itp.

wybierając najlepsze.

Każdy element zbioru danych x_i odpowiada regule $r(x_i)$ – takiej, która wybiera dokładnie ten element. Tworzymy reguły postaci:

 $r(x_{10}) \lor r(x_{31}),$ $r(x_1) \lor r(x_{45}) \lor r(x_{321}) \lor r(x_{1421}) \lor r(x_{2352}) \lor r(x_{31002})$ \cdot itp.

wybierając najlepsze.

index	A	В	C				
0	low	high	low	$Dom(A) = \{low, medium, high\}$			
1	low	low	low	$Dom(B) = \{low, high\}$			
3	low	low	high	$Dom(C) = \{low, medium, high\}$			
4	medium	low	high	$Dom(c) = \{low, mediani, mgn\}$			
$a_1 = A \in \{low\} \land B \in \{low\} \land C \in \{low\}$ $a_2 = A \in \{medium\} \land B \in \{low\} \land C \in \{high\},$							
r = c	$r = a_1 \land a_2 = A \in \{low, medium\} \land B \in \{low\} \land C \in \{low, high\}$						

Każdy element zbioru danych x_i odpowiada regule $r(x_i)$ – takiej, która wybiera dokładnie ten element. Tworzymy reguły postaci:

 $r(x_{10}) \lor r(x_{31}),$ $r(x_1) \lor r(x_{45}) \lor r(x_{321}) \lor r(x_{1421}) \lor r(x_{2352}) \lor r(x_{31002})$ \cdot itp.

wybierając najlepsze.

Α	В	C					
low	high	low	$Dom(A) = \{low, medium, high\}$				
low	low	low	$Dom(B) = \{low, high\}$				
low	low	high	$Dom(C) = \{low, medium, high\}$				
medium	low	high	Dom(c) = {tow, medium, mgm}				
$a_1 = A \in \{low\} \land B \in \{low\} \land C \in \{low\}$ $a_2 = A \in \{medium\} \land B \in \{low\} \land C \in \{high\},$							
$r = a_1 \land a_2 = A \in \{\text{low, medium}\} \land B \in \{\text{low}\} \land C \in \{\text{low, high}\}$							
	$\begin{array}{c} \text{low} \\ \text{low} \\ \text{low} \\ \text{medium} \end{array}$ $A \in \{\text{low}\}$ $A \in \{\text{medium}\}$	$\begin{array}{c c} \text{low} & \text{high} \\ \text{low} & \text{low} \\ \text{low} & \text{low} \\ \text{low} & \text{low} \\ \text{medium} & \text{low} \\ A \in \{\text{low}\} \land B \\ A \in \{\text{medium}\} \end{array}$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$				

Każdy element zbioru danych x_i odpowiada regule $r(x_i)$ – takiej, która wybiera dokładnie ten element. Tworzymy reguły postaci:

- $r(x_{10}) \lor r(x_{31}),$ $r(x_1) \lor r(x_{45}) \lor r(x_{321}) \lor r(x_{1421}) \lor r(x_{2352}) \lor r(x_{31002})$ \cdot itp.
- wybierając najlepsze.

	index	A	В	C			
	0	low	high	low	$Dom(A) = \{low, medium, high\}$		
	1	low	low	low	$Dom(B) = \{low, high\}$		
	3	low	low	high	$Dom(C) = \{low, medium, high\}$		
	4	medium	low	high	Dom(c) = \text{tow, mediani, mgn}		
$a_1 = A \in \{low\} \land B \in \{low\} \land C \in \{low\}$							
$a_2 = A \in \{\text{medium}\} \land B \in \{\text{low}\} \land C \in \{\text{high}\},$							
	$r = a_1 \land a_2 = A \in \{low, medium\} \land B \in \{low\} \land C \in \{low, high\}$						

Każdy element zbioru danych x_i odpowiada regule $r(x_i)$ – takiej, która wybiera dokładnie ten element. Tworzymy reguły postaci:

- $r(x_{10}) \lor r(x_{31}),$ $r(x_1) \lor r(x_{45}) \lor r(x_{321}) \lor r(x_{1421}) \lor r(x_{2352}) \lor r(x_{31002})$ \cdot itp.
- wybierając najlepsze.

index	А	В	С			
0	low	high	low	$Dom(A) = \{low, medium, high\}$		
1	low	low	low	$Dom(B) = \{low, high\}$		
3	low	low	high	$Dom(C) = \{low, medium, high\}$		
4	medium	low	high	$Bom(e) = \{low, mediam, mgn\}$		
$a_1 = A \in \{low\} \land B \in \{low\} \land C \in \{low\}$ $a_2 = A \in \{medium\} \land B \in \{low\} \land C \in \{high\},$						
$r = a_1 \land a_2 = A \in \{\text{low, medium}\} \land B \in \{\text{low}\} \land C \in \{\text{low, high}\}$						

Problemy

Obie metody wymagają:

- · szybkiego liczenia funkcji score (wiele reguł do sprawdzenia);
- · efektywnej implementacji algebry reguł (wiele reguł do skonstruowania);
- · stosowania heurystycznych optymalizacji.

Problemy

Obie metody wymagają:

- · szybkiego liczenia funkcji score (wiele reguł do sprawdzenia);
- · efektywnej implementacji algebry reguł (wiele reguł do skonstruowania);
- · stosowania heurystycznych optymalizacji.

W konsekwencji:

- · kod będzie się często zmieniał na etapie tworzenia prototypu
- · nie wiadomo z góry, która reprezentacja reguł będzie optymalna
- · może istnieć więcej niż jeden sposób heurystycznej optymalizacji przeszukania
- · może nie istnieć globalnie najefektywniejsze rozwiązanie

Z drugiej strony, opis algorytmu jest bardzo prosty.

Implementacja w Haskellu

General-purpose, purely functional programming language with non-strict semantics and strong static typing [wikipedia]

General-purpose, purely functional programming language with non-strict semantics and strong static typing [wikipedia]

 strongly typed: wszystko ma określony typ; kompilator statycznie sprawdza zgodność typów; brak domyślnych konwersji między typami.

```
pi :: Double sin :: Double -> Double podobnie jak w matematyce: \pi \in \mathbb{R}, sin \colon \mathbb{R} \to \mathbb{R}
```

General-purpose, purely functional programming language with non-strict semantics and strong static typing [wikipedia]

 strongly typed: wszystko ma określony typ; kompilator statycznie sprawdza zgodność typów; brak domyślnych konwersji między typami.

```
pi :: Double sin :: Double -> Double podobnie jak w matematyce: \pi \in \mathbb{R}, sin \colon \mathbb{R} \to \mathbb{R}
```

· pure: brak zmiennych i efektów ubocznych

General-purpose, purely functional programming language with non-strict semantics and strong static typing [wikipedia]

 strongly typed: wszystko ma określony typ; kompilator statycznie sprawdza zgodność typów; brak domyślnych konwersji między typami.

```
pi :: Double sin :: Double -> Double podobnie jak w matematyce: \pi \in \mathbb{R}, sin \colon \mathbb{R} \to \mathbb{R}
```

· pure: brak zmiennych i efektów ubocznych

· lazy evaluation: wartości są liczone dopiero wtedy, gdy są potrzebne:

```
-- lista kwadratów wszystkich liczb naturalnych
squares = [x*x | x <- [0..]]
take 5 squares -- pierwszych 5 wartości
-- [0, 1, 4, 9, 16]</pre>
```

Typy i klasy typów

type Rule a = Map Text a

- · ułatwia prototypowanie: kompilator sam podpowiada co i gdzie jeszcze zmienić;
- · mniej błędów w wersji finalnej;
- · porządkują proces myślowy.

Typy i klasy typów

```
type Rule a = Map Text a
```

- · ułatwia prototypowanie: kompilator sam podpowiada co i gdzie jeszcze zmienić;
- · mniej błędów w wersji finalnej;
- · porządkują proces myślowy.

Łatwo możemy tworzyć nietrywialne typy:

```
data Tree a = Node a [Tree a]
```

Więcej niż jeden typ może być regułą:

class RuleType a where

```
(.<=.) :: a -> a -> Bool -- poset
(/\) :: a -> a -> a -- and for rules
(\/) :: a -> a -> a -- sup for rules
```

Więcej niż jeden typ może być regułą:

class RuleType a where

```
(.<=.) :: a -> a -> Bool -- poset
(/\) :: a -> a -> a -- and for rules
(\/) :: a -> a -> a -- sup for rules
```

Zbiory jako reguły:

```
instance RuleType (Set a) where
  a .<=. b = a `isSubsetOf` b
  a /\ b = a `intersect` b
  a \/ b = a `union` b</pre>
```

Więcej niż jeden typ może być regułą:

```
class RuleType a where
  (.<=.) :: a -> a -> Bool -- poset
  (/\) :: a -> a -> a -- and for rules
  (\/\) :: a -> a -> a -- sup for rules
Zbiory jako reguly:
instance RuleType (Set a) where
  a .<=. b = a 'isSubsetOf' b
  a /\ b = a 'intersect' b
  a /\ b = a 'union' b</pre>
```

instance Rule a => RuleType (Rule a) where
a .<=, b = ...</pre>

a /\ b = ... a \/ b = ...

Więcej niż jeden typ może być regułą:

```
class RuleType a where
  (.<=.) :: a -> a -> Bool -- poset
  (/\) :: a -> a -- and for rules
  (\/) :: a -> a -- sup for rules
Zbiory jako reguly:
instance RuleType (Set a) where
 a .<=. b = a `isSubsetOf` b
 a /\ b = a `intersect` b
 a \  \  \  \  \  b = a `union` b
instance Rule a => RuleType (Rule a) where
 a . <= . b = ...
 a / b = \dots
 a \/ b = ...
instance Rule a => RuleType [a] where
 a .<=. b = all $ zipWith (.<=.) a b
 a /\ b = all $ zipWith (/\) a b
 a \  \  ) b = all \  \   zipWith (\/) a b
```

Więcej niż jeden typ może być regułą:

```
class RuleType a where
  (.<=.) :: a -> a -> Bool -- poset
  (/\) :: a -> a -- and for rules
  (\/) :: a -> a -> a -- sup for rules
Zbiory jako reguly:
instance RuleType (Set a) where
  a \le b = a isSubsetOf b
  a /\ b = a `intersect` b
  a \  \  \  \  \  b = a `union` b
instance Rule a => RuleType (Rule a) where
  a . <= . b = ...
  a / b = \dots
  a \/ b = ...
instance Rule a => RuleType [a] where
  a <=. b = all $ zipWith (.<=.) a b
  a / b = all \pm zipWith (/ a b
  a \  \  ) b = all \  \   zipWith (\/) a b
Reprezentacja bitowa:
instance RuleType Bool where
  a <= . b = b implies a
 a /\ b = a && b
```

Więcej niż jeden typ może być regułą:

```
class RuleType a where
  (.<=.) :: a -> a -> Bool -- poset
  (/\) :: a -> a -- and for rules
  (\/\) :: a -> a -> a -- sup for rules
Zbiory jako reguly:
instance RuleType (Set a) where
 a \le b = a isSubsetOf b
 a /\ b = a `intersect` b
 a \ \ b = a \ union b
instance Rule a => RuleType (Rule a) where
 a . <= . b = ...
 a / b = \dots
 a \/ b = ...
instance Rule a => RuleType [a] where
 a <=. b = all $ zipWith (.<=.) a b
 a / b = all \pm zipWith (/ a b
 a \  \  ) b = all \  \   zipWith (\/) a b
Reprezentacia bitowa:
instance RuleType Bool where
 a <= . b = b implies a
 a /\ b = a && b
instance RuleType BitArray where
 a <= .b = not b \mid a
 a / b = a & b
```

Parametryczny polimorfizm + purity

Implementacja algorytmu wykorzystuje tylko ogólne cechy argumentów, nie konkretne reprezentacje, tu:

- algebra reguł (\leq , \wedge);
- · dodawanie nowego elementu do zbioru najlepszych reguł.

Parametryczny polimorfizm + purity

```
searchRules :: RuleType a, WithScore a =>
-> (a -> Bool) -- f. przycinająca
-> Best a -- najlepsze reguły
-> Tree a -- przestrzeń do przeszukania
-- codomain
-> Best a -- nowe najlepsze reguły
```

Implementacja algorytmu wykorzystuje tylko ogólne cechy argumentów, nie konkretne reprezentacje, tu:

- algebra reguł (\leq , \wedge);
- · dodawanie nowego elementu do zbioru najlepszych reguł.

Dzięki temu:

- łatwo eksperymentować z różnymi reprezentacji obiektu (korzystamy z dwóch, dwie czy trzy inne były przetestowane trakcie implementacji);
- łatwiej testować poprawność algorytmu (można np. wybrać najprostszą reprezentację reguł);
- bardziej przejrzysty kod.

Parametryczny polimorfizm + purity

Implementacja algorytmu wykorzystuje tylko ogólne cechy argumentów, nie konkretne reprezentacje, tu:

- algebra reguł (\leq , \wedge);
- · dodawanie nowego elementu do zbioru najlepszych reguł.

Dzięki temu:

- łatwo eksperymentować z różnymi reprezentacji obiektu (korzystamy z dwóch, dwie czy trzy inne były przetestowane trakcie implementacji);
- łatwiej testować poprawność algorytmu (można np. wybrać najprostszą reprezentację reguł);
- · bardziej przejrzysty kod.

```
searchRules :: RuleType a, WithScore a => Int -> (a -> Bool) -> Best a -> Tree a -> Best a
searchRules 0 _ _ bestK (Node r _) = update r bestK
searchRules d prune bestK (Node r cs) = foldl' (searchRules (d-1) prune) updated children
where
    children = filter interesting cs
    updated = update r bestK
    interesting (Node r _) = not $ prune r
```

Lazy evaluation

Możemy oddzielić tworzenie przestrzeni reguł do przeszukania od samego algorytmu przeszukującego:

```
hSpace :: RuleType r => [r] -> r -> Tree r
hSpace [] r = Node r []
hSpace pool r = Node r $ filter notNull . map subForest . init . tails $ pool
   where
       subForest (u:us) = hSpace us (r /\ u)
       notNull (Node d _) = not . isNull $ d
```

tworzy w praktyce nieskończone drzewo reguł do przeszukania. Możemy je transformować:

```
space = hSpace pool trivial
fmap toScored space -- infinite tree with scores!
```

Wersja Pythonowa?

- liczenie funkcji score, z wykorzystaniem pandas/numpy/numexpr było stosunkowo szybkie;
- · ale implementacja algebry reguł była okrutnie powolna...
- · ...lub bardzo nieczytelna, jeśli chcieć ją zoptymalizować (pythran, etc.);
- · mało wydajne zrównoleglanie kodu;
- przez większą złożoność kodu, znacznie trudniejsze eksperymentowanie z różnymi wariantami algorytmu.

Wersja Pythonowa?

- liczenie funkcji score, z wykorzystaniem pandas/numpy/numexpr było stosunkowo szybkie;
- · ale implementacja algebry reguł była okrutnie powolna...
- · ...lub bardzo nieczytelna, jeśli chcieć ją zoptymalizować (pythran, etc.);
- · mało wydajne zrównoleglanie kodu;
- przez większą złożoność kodu, znacznie trudniejsze eksperymentowanie z różnymi wariantami algorytmu.

Statystyki

- · liczba linii kodu: ok 1500 (Haskell) vs 3800 (Python)
- w tym kod bezpośrednio związany z implementacją algorytmu: 540 (Haskell) vs 3000 (Python)
- wersja w Haskellu średnio ok. 10 razy szybsza od pythonowej (single core) i znacznie lepiej zrównoleglająca się (do 32 rdzeni, w pythonie: kilka do kilkunastu).

Integracja z ekosystemem Pythona

Wymiana danych

Największy problem

Haskell ma ubogi ekosystem bibliotek odczytujących dane z różnych formatów do których łatwo można eksportować z pandas + brak standardowej biblioteki typu pandas

Wymiana danych

Największy problem

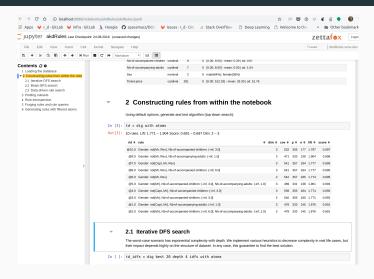
Haskell ma ubogi ekosystem bibliotek odczytujących dane z różnych formatów do których łatwo można eksportować z pandas + brak standardowej biblioteki typu pandas

ARFF

Własna implementacja formatu ARFF (od Weka):

- prosty format tekstowy niezależny od specyficznych sposób binarnego kodowania (pandas miewa z tym problemy);
- · w przeciwieństwie do csv kolumny mają ściśle określony typ i dziedzinę;
- efekt końcowy: czytanie plików ARFF niewiele wolniejsze od czytania np. HDF5 w pandas;
- · niestety konieczna była implementacja ARFF również w Pythonie...

Dedykowany jupyter kernel



- · znajomy interfejs
- pozwala wykorzystać algorytmy wyszukujące reguły do edycji istniejących reguł przydatne gdy dodaje się wiedzę domenową.

Podsumowanie

Podsumowanie

Dzięki dopasowaniu specyfiki problemu i cech narzędzia, udało się uzyskać:

- · wydajniejszy proces tworzenia narzędzia
- · kod łatwiejszy w utrzymaniu
- · łatwość tworzenia eksperymentalnych wariantów
- · większą szybkość działania
- · łatwiejsze testowanie
- · większa satysfakcja z pracy

Choć bywały trudniejsze momenty:

- · brak wsparcia w ekosystemie
- · nie jest to droga dla początkujących
- · trzeba uważać na przesadną abstrakcję

