

Metody przetwarzania dużych ilości danych

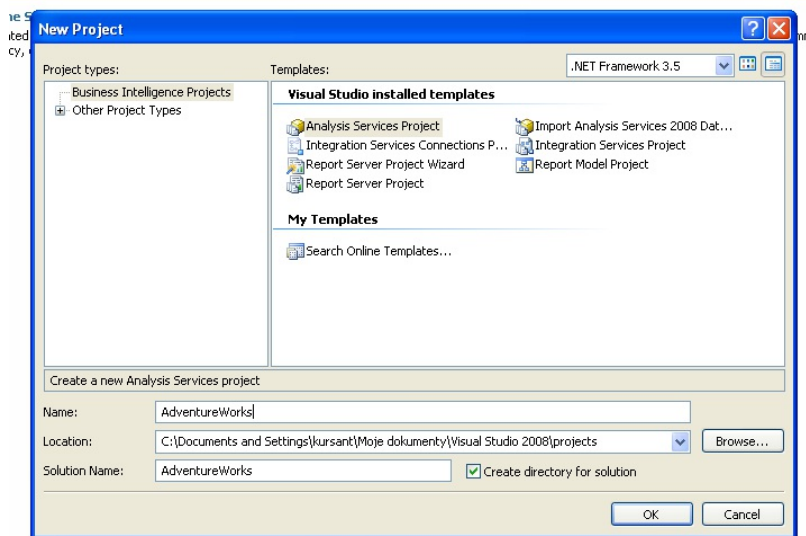
Sprawozdanie z laboratorium

Data	Tytuł zajęć	Uczestnicy
15.12.2020 18:55	Data Mining	Adam Hyjek (234987) Bartosz Rodziewicz (226105)

Streszczenie laboratorium

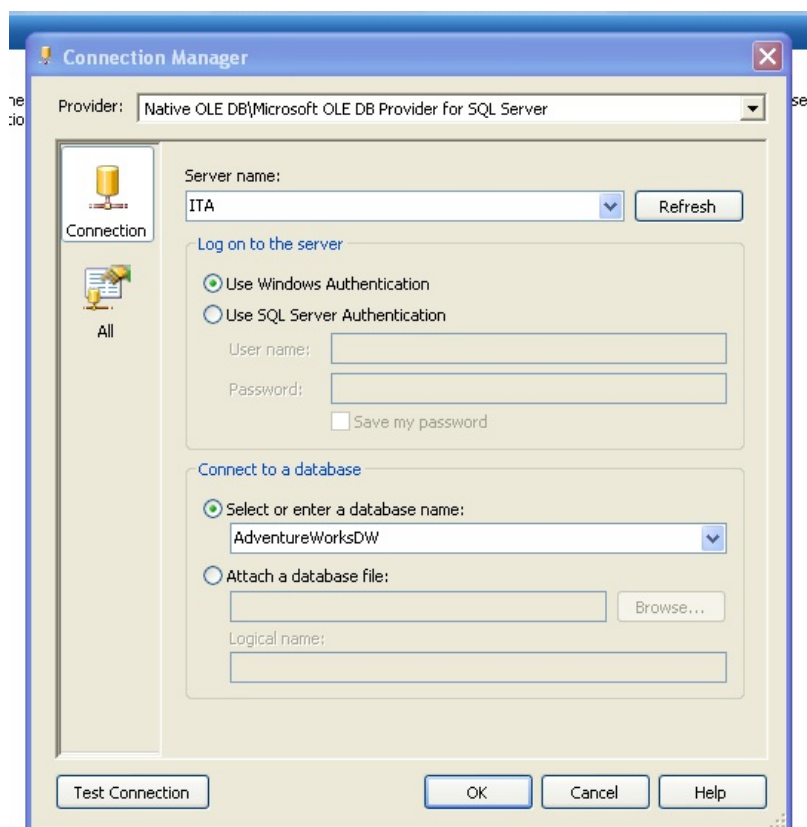
Problem 1

Stworzenie nowego projektu Analysis Services w SS BIDS

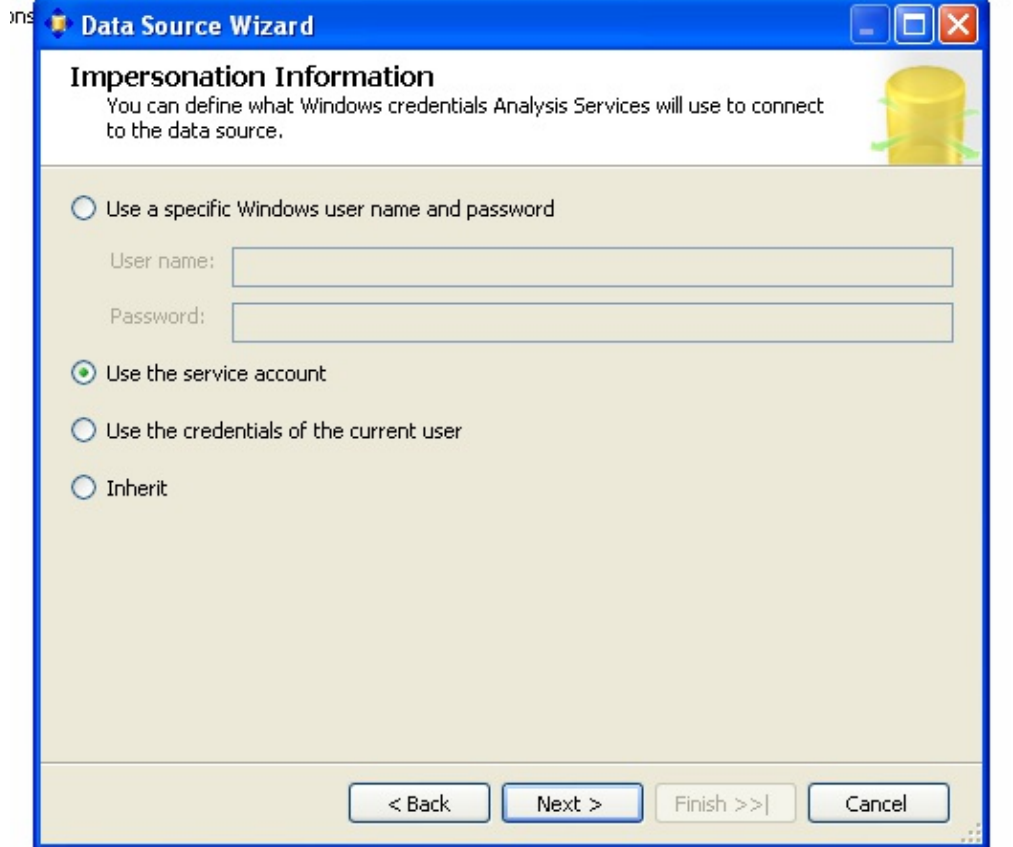


Stworzenie nowego projektu.

Stworzenie źródła danych

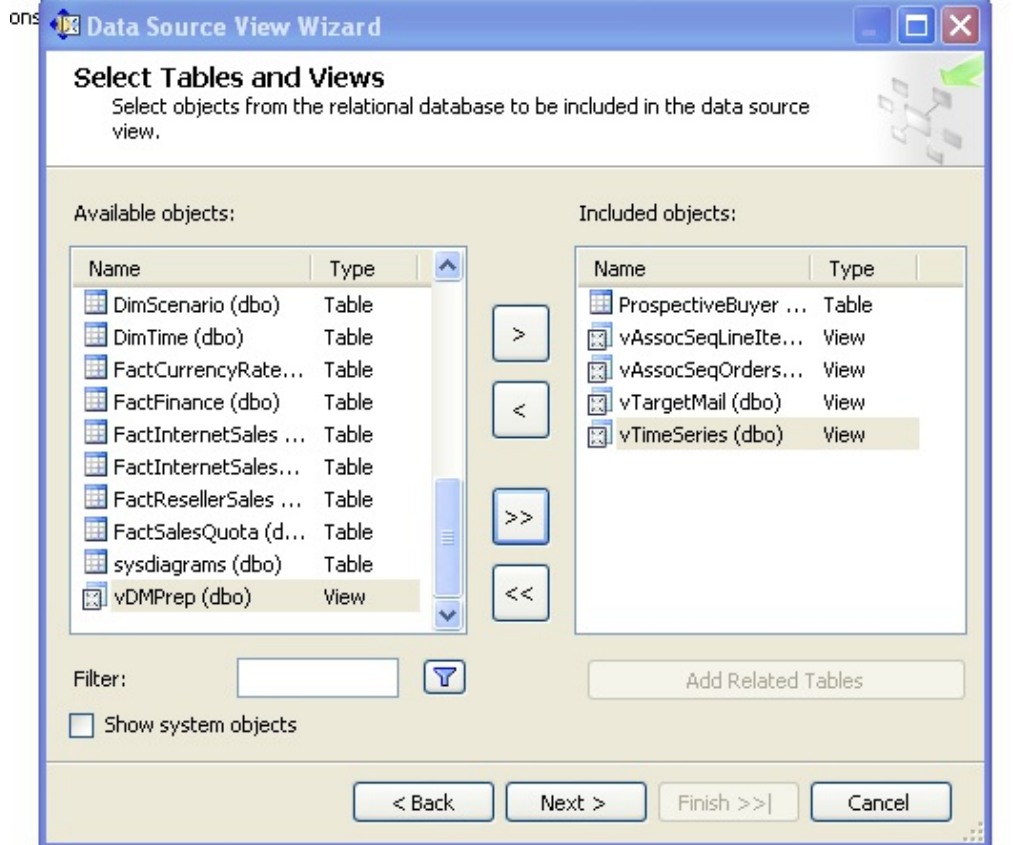


Stworzenie połączenia z bazą danych (różnica do instrukcji, serwer to nie localhost, a ITA).

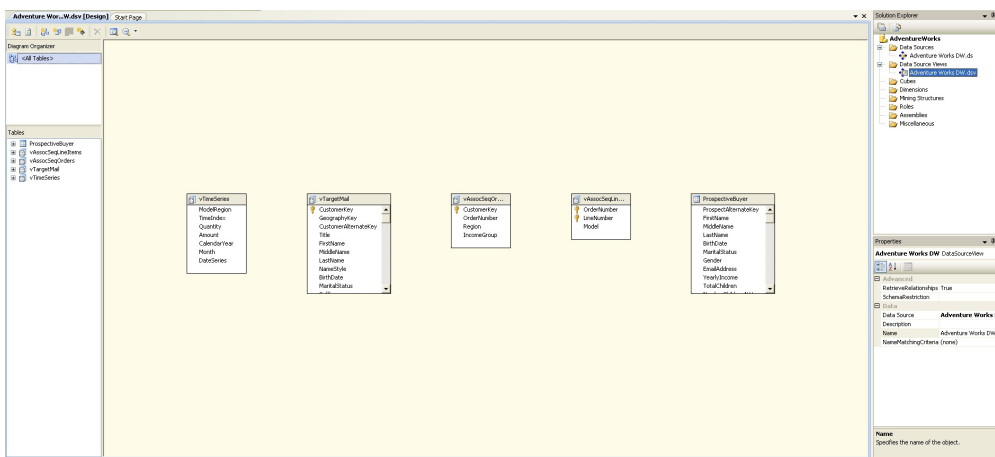


Wybór metody logowania.

Stworzenie widoku źródła danych

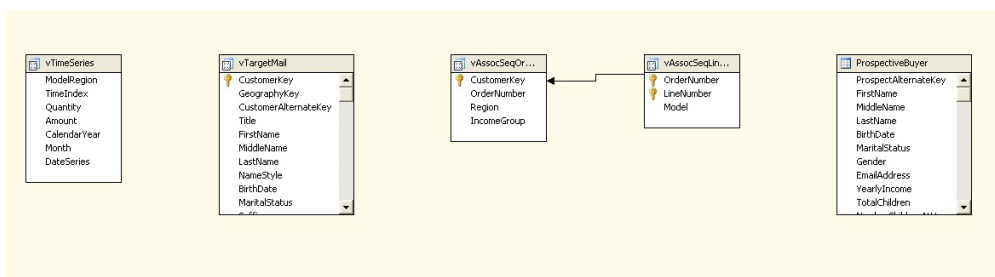


Wybór tabel do widoku źródła danych.



Nowy widok źródła danych.

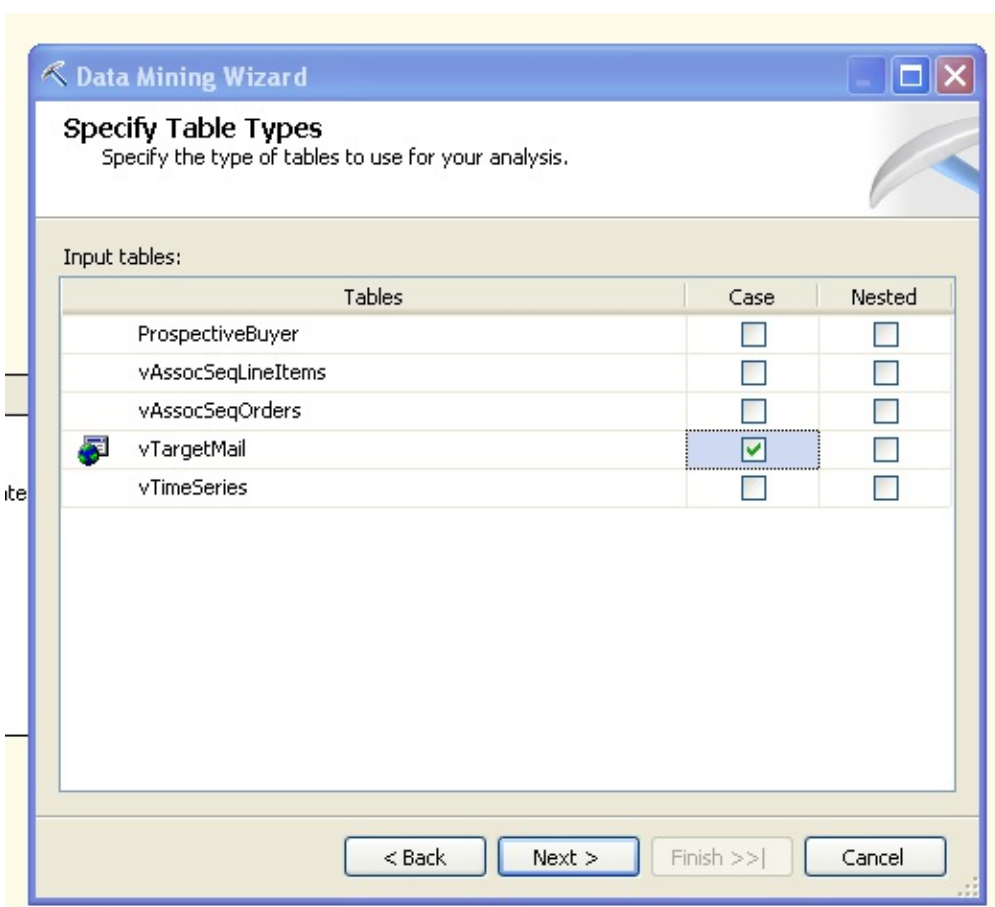
Modyfikacja widoku źródła danych



Utworzenie nowego związku między tabelami.

Problem 2

Stworzenie struktury eksploracji danych dla modelu Targeted Mailing



Wybór tabeli do eksploracji.

Data Mining Wizard

Specify the Training Data

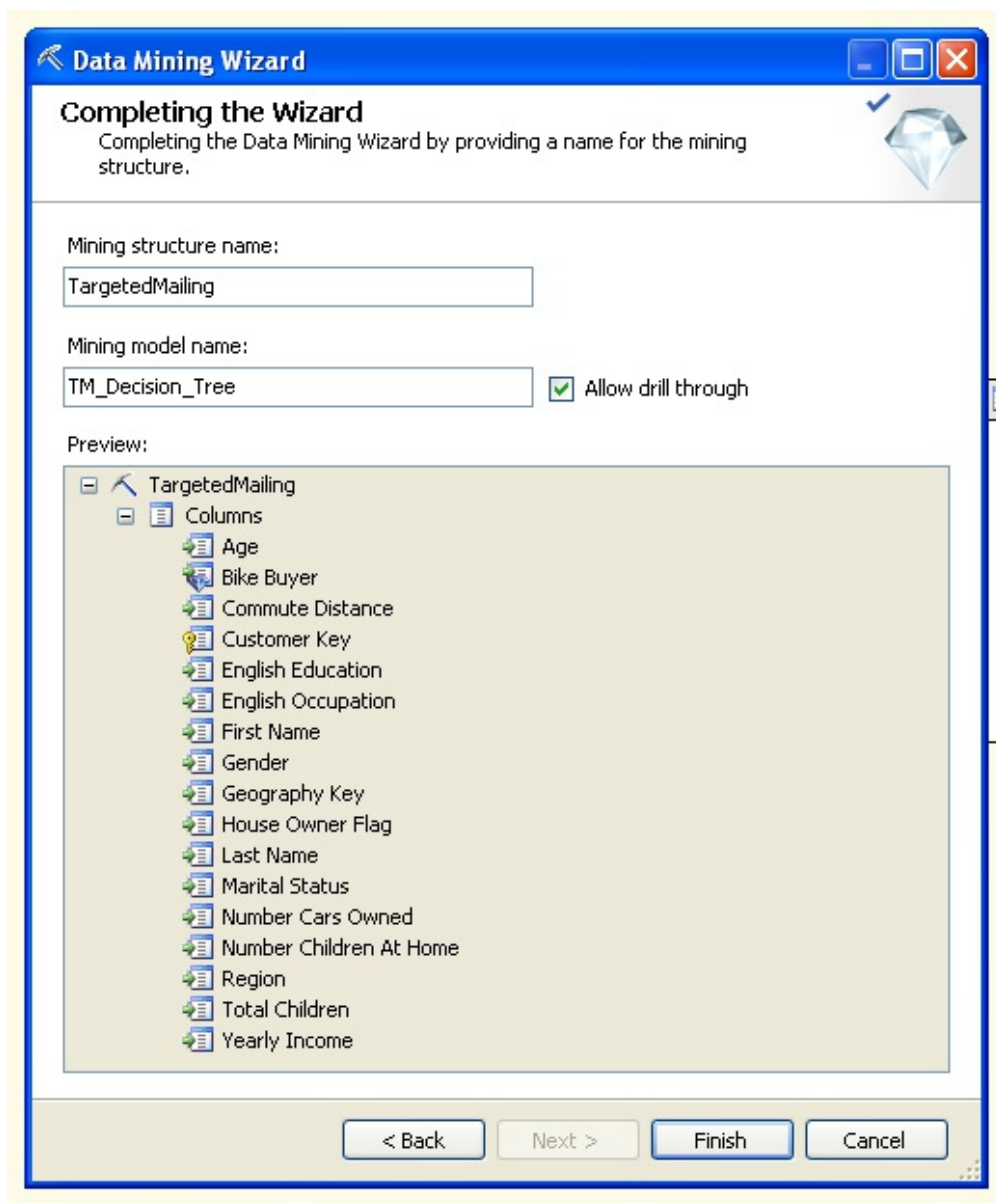
Specify the columns used in your analysis.

Mining model structure:

<input type="checkbox"/>	Tables/Columns	Key	<input type="checkbox"/> Input	<input type="checkbox"/> Predic...
	vTargetMail			
<input type="checkbox"/>	AddressLine1	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	AddressLine2	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	Age	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	BikeBuyer	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
<input type="checkbox"/>	BirthDate	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	CommuteDistance	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	CustomerAlternateKey	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	CustomerKey	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	DateFirstPurchase	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	EmailAddress	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	EnglishEducation	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	EnglishOccupation	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	FirstName	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	FrenchEducation	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	FrenchOccupation	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	Gender	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	GeographyKey	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	HouseOwnerFlag	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	LastName	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	MaritalStatus	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	MiddleName	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	NameStyle	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	NumberCarsOwned	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	NumberChildrenAtHome	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	Phone	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	Region	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	SpanishEducation	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	SpanishOccupation	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	Suffix	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	Title	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	TotalChildren	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
<input checked="" type="checkbox"/>	YearlyIncome	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>

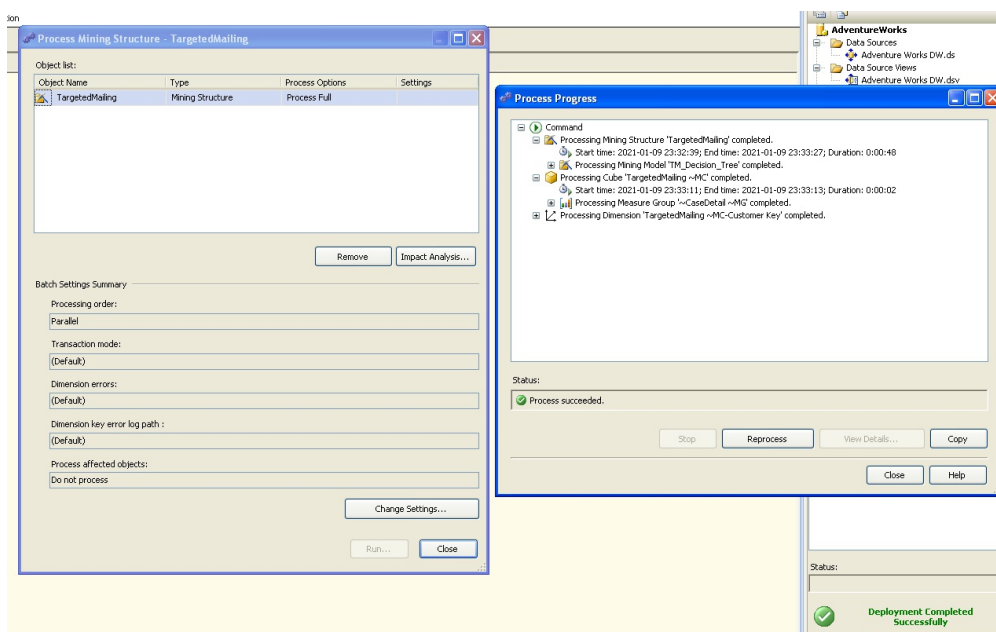
Recommend inputs for currently selected predictable:

Wybór danych treningowych.



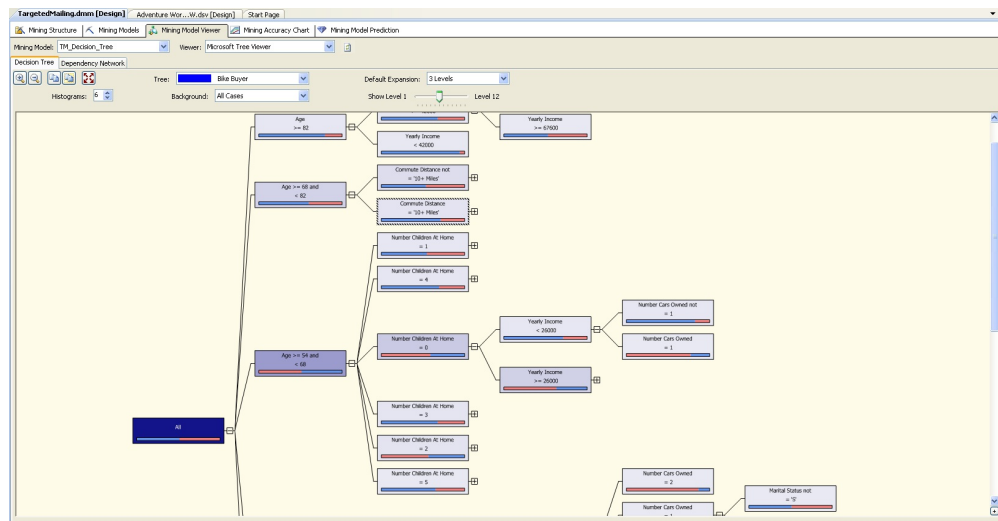
Podsumowanie tworzenia struktury eksploracji danych.

Przetwarzanie modelu ekstrakcji danych



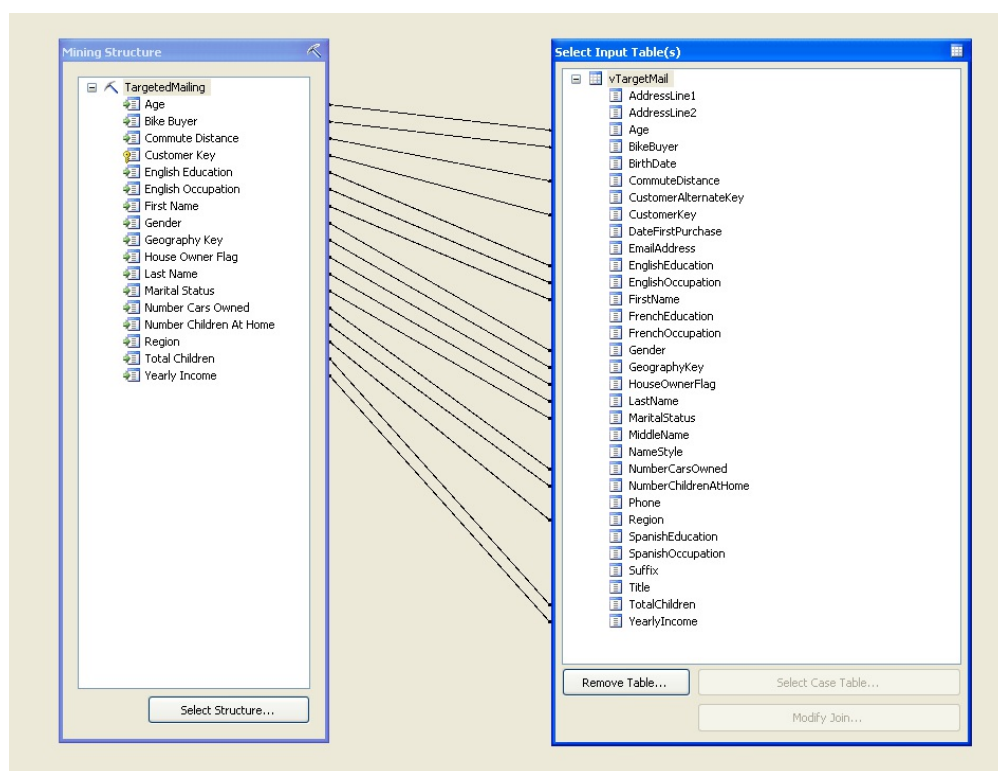
Poprawny wynik przetwarzania.

Analiza modelu opartego o drzewa decyzyjne



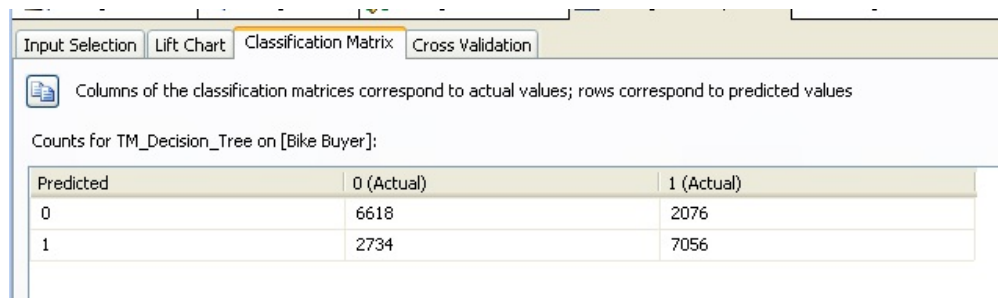
Wygląd widoku analizy modelu.

Mapowanie kolumn



Mapowanie kolumn.

Macierz klasyfikacji



Wyniki macierzy klasyfikacji.

Jaki jest całkowity błąd predykcji drzewa decyzyjnego, a także jaka część osób odrzuci ofertę z kampanii?

Wskaźnik	Działanie	Wynik
Całkowity błąd predykcji	$(2734 + 2076) / (6618 + 2734 + 2076 + 7056) = 4810 / 18484 = \sim 0.26$	26%
Jaka część osób odrzuci ofertę	$2734 / (2734 + 7056) = 1367 / 4895 = \sim 0.279$	27.9%

Wykres przewidywanego zysku

Profit Chart

Profit Chart Settings...

Profit Chart Settings

Population:

50000

Fixed cost:

5000

Individual cost:

10

Revenue per individual:

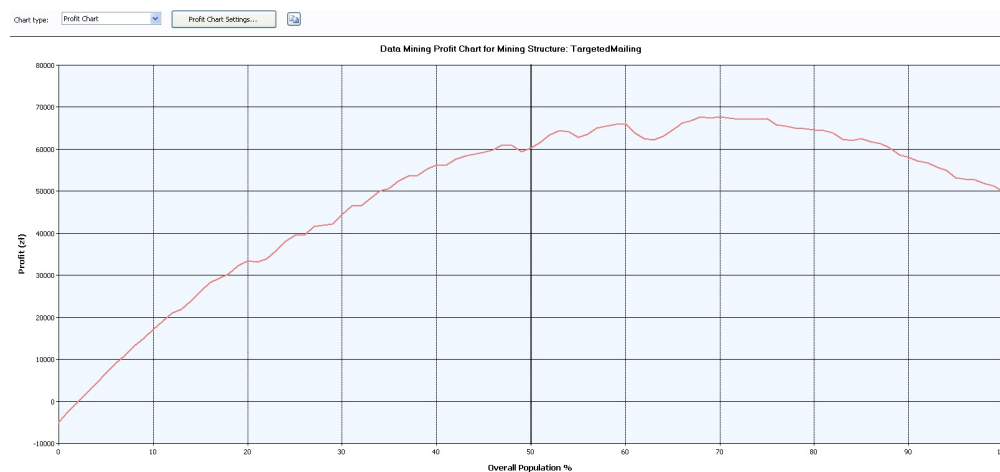
15

OK

Cancel

Help

Wybór parametrów wykresu.

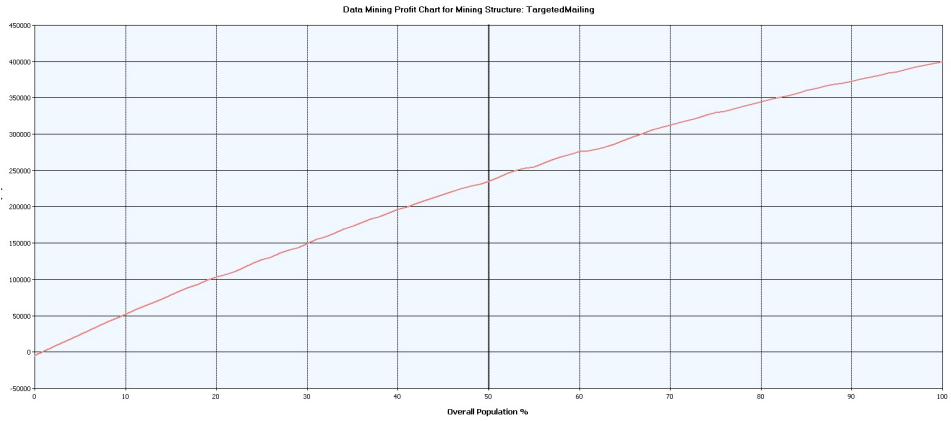


Wykres dla indywidualnego kosztu równego 10.

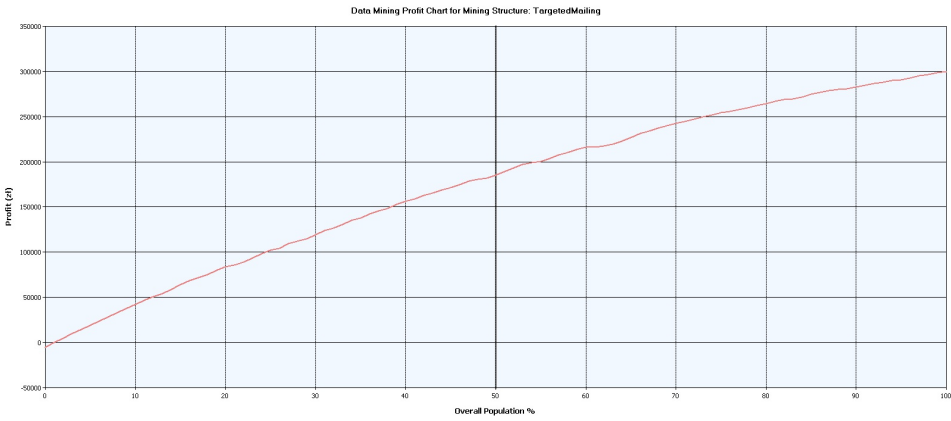
Jak można zinterpretować otrzymany wykres?

Wykres mówi o tym jakiego zysku można się spodziewać wysyłając kampanie reklamową do wybranego procentu populacji uwzględniając całkowity koszt kampanii, koszt kampanii na osobę, możliwy zysk od osoby oraz prawdopodobieństwo, czy dana osoba zdecyduje się na zakupy.

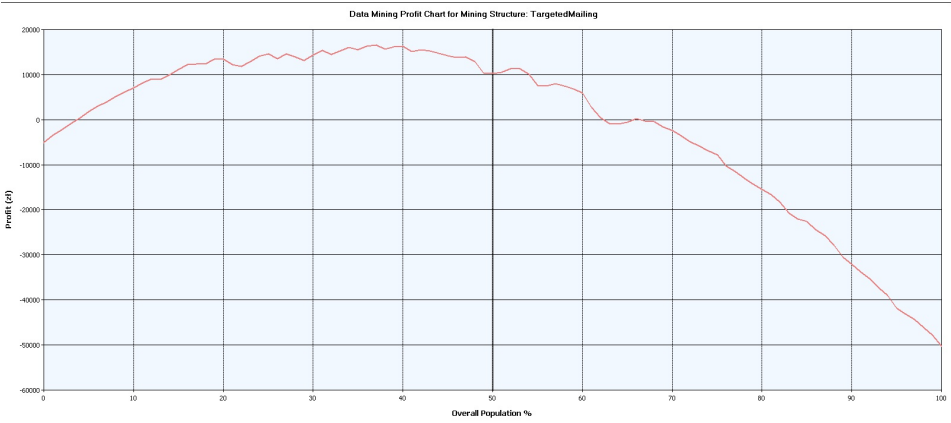
Zbadaj kształt wykresu dla różnych wartości



Wykres dla indywidualnego kosztu równego 3.



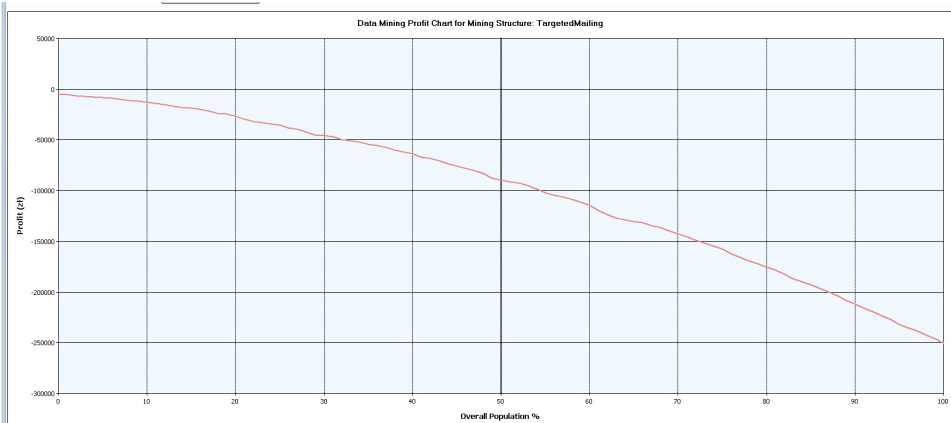
Wykres dla indywidualnego kosztu równego 5.



Wykres dla indywidualnego kosztu równego 12.

Im wyższy koszt wykonania kampanii na osobę (przy takim samym potencjalnym zysku), tym bardziej się opłaca wysłać kampanie do mniejszej ilości, lepiej wyselekcjonowanych ludzi, przy których prawdopodobieństwo sukcesu jest większe.

Dlaczego przy ustawieniu Individual Cost powyżej wartości 15 wykres zysku nie posiada wartości dodatnich?



Wykres dla indywidualnego kosztu równego 16.

Wykres nie posiada wartości dodatnich, ponieważ potencjalny zysk z osoby, decydującej się na zakupy, jest niższy niż koszt wykonania kampanii.

Zmiany kolumn w modelu

Counts for TM_Decision_Tree on [Bike Buyer]:

Predicted	0 (Actual)	1 (Actual)
0	5938	2239
1	3414	6893

Nowa wartość macierzy klasyfikacji.

Oblicz na nowo całkowity błąd predykcji i błąd dla kupna roweru. Czy otrzymane wyniki są lepsze czy gorsze od poprzednio uzyskanych?

Wskaźnik	Działanie	Wynik
Całkowity błąd predykcji	$(3414 + 2239) / (5938 + 3414 + 2239 + 6893) = 5653 / 18484 = \sim 0.3058$	30.6%
Jaka część osób odrzuci ofertę	$3414 / (3414 + 6893) = 3414 / 10307 = \sim 0.33$	33%

Wyniki są gorsze od poprzednich.

Problem 2 – kontynuacja modułu 1

Algorytm naiwny Bayesowski

Structure	TM_Decision_Tree	TM_Bayes
	Microsoft_Decision_Trees	Microsoft_Naive_Bayes
Age	Input	Ignore
Bike Buyer	Predict	Predict
Commute Distance	Input	Input
Customer Key	Key	Key
English Education	Input	Input
English Occupation	Input	Input
First Name	Ignore	Ignore
Gender	Input	Input
Geography Key	Input	Ignore
House Owner Flag	Input	Input
Last Name	Ignore	Ignore
Marital Status	Input	Input
Number Cars Owned	Input	Input
Number Children At Home	Input	Input
Region	Ignore	Ignore
Total Children	Input	Input
Yearly Income	Input	Ignore

Nowy algorytm Bayesowski.

Attribute profiles					
Attributes	States	Populatio... Size: 17484	0 Size: 8875	1 Size: 8609	missing Size: 0
Commute Distance	<ul style="list-style-type: none"> 0-1 Miles 2-5 Miles 1-2 Miles 5-10 Miles Other 				
English Education	<ul style="list-style-type: none"> Bachelors Partial College High School Graduate Degree Other 				
English Occupation	<ul style="list-style-type: none"> Professional Skilled Manual Management Clerical Other 				
Marital Status	<ul style="list-style-type: none"> M S Missing 				
Number Cars Owned	<ul style="list-style-type: none"> 2 1 0 3 Other 				
Number Children At Home	<ul style="list-style-type: none"> 0 1 2 3 Other 				
Total Children	<ul style="list-style-type: none"> 0 2 1 4 Other 				

Wpływ poszczególnych atrybutów na wyniki.

Characteristics for 1		
Attributes	Values	Probability
Number Children At Home	0	
Marital Status	M	
Marital Status	S	
Commute Distance	0-1 Miles	
English Education	Bachelors	
English Occupation	Professional	
Number Cars Owned	1	
Number Cars Owned	0	
Total Children	0	
Number Cars Owned	2	
English Education	Partial College	
English Occupation	Skilled Manual	
Total Children	1	
Total Children	2	
Commute Distance	2-5 Miles	
English Education	Graduate Degree	
English Occupation	Clerical	
Commute Distance	1-2 Miles	
English Occupation	Management	
English Education	High School	
Commute Distance	5-10 Miles	
Number Children At Home	1	
English Occupation	Manual	
Total Children	3	
Number Children At Home	2	
Commute Distance	10+ Miles	
Total Children	4	
Number Cars Owned	3	
English Education	Partial High School	
Number Cars Owned	4	
Total Children	5	
Number Children At Home	3	
Number Children At Home	4	

Wpływ poszczególnych wartości atrybutów na wyniki.

Czy aktualnie stworzony model jest lepszy czy gorszy od drzew decyzyjnych?

Counts for TM_Bayes on [Bike Buyer]:		
Predicted	0 (Actual)	1 (Actual)
0	5912	3621
1	3440	5511

Wynik macierzy klasyfikacji.

Algorytm jest gorszy.

Zmienić kolumny wejściowe dla tego modelu, czy można w ten sposób poprawić wyniki dla tego algorytmu?

Counts for TM_Bayes on [Bike Buyer]:		
Predicted	0 (Actual)	1 (Actual)
0	6073	3484
1	3279	5648

Wynik macierzy klasyfikacji po dodaniu poprzednio wyłączonych atrybutów (imię, nazwisko, region).

Wynik delikatnie się poprawił, jednak wciąż jest słabszy niż algorytm na podstawie drzew decyzyjnych.

Algorytm oparty o Sztuczne Sieci Neuronowe

Structure	TM_Decision_Tree	TM_Bayes	Neural_Network
	Microsoft_Decision_Trees	Microsoft_Naive_Bayes	Microsoft_Neural_Network
Age	Input	Ignore	Input
Bike Buyer	Predict	Predict	Predict
Commute Distance	Input	Input	Input
Customer Key	Key	Key	Key
English Education	Input	Input	Input
English Occupation	Input	Input	Input
First Name	Input	Input	Input
Gender	Input	Input	Input
Geography Key	Input	Ignore	Input
House Owner Flag	Input	Input	Input
Last Name	Input	Input	Input
Marital Status	Input	Input	Input
Number Cars Owned	Input	Input	Input
Number Children At Home	Input	Input	Input
Region	Input	Input	Input
Total Children	Input	Input	Input
Yearly Income	Input	Ignore	Input

Nowy algorytm SSN.

Jakie kolumny nie mogą być użyte przez sieć neuronową?

Wszystkie dostępne kolumny mogą być użyte poprzez sieć neuronową.

Porównaj opracowane modele. Które kolumny można zignorować?

Counts for TM_Decision_Tree on [Bike Buyer]:

Predicted	0 (Actual)	1 (Actual)
0	6644	2076
1	2708	7056

Counts for TM_Bayes on [Bike Buyer]:

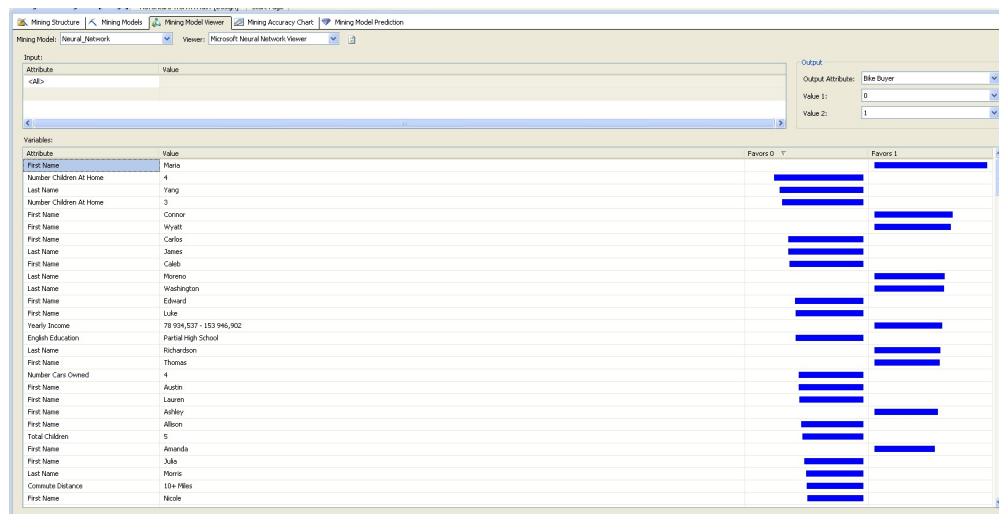
Predicted	0 (Actual)	1 (Actual)
0	6073	3484
1	3279	5648

Counts for Neural_Network on [Bike Buyer]:

Predicted	0 (Actual)	1 (Actual)
0	6085	3125
1	3267	6007

Wynik macierzy klasyfikacji dla wszystkich argumentów.

Powyżej widać, że algorytm SSN jest trochę lepszy niż Bayesowski, jednak gorszy niż drzewa decyzyjne.



Analiza modelu.

Analiza modelu pokazuje, że duży wpływ na dane mają kolumny First and Last Name. Oczywiście jest, że te kolumny nie prezentują żadnych użytecznych informacji, więc można je wyłączyć.

Counts for TM_Decision_Tree on [Bike Buyer]:

Predicted	0 (Actual)	1 (Actual)
0	6644	2076
1	2708	7056

Counts for TM_Bayes on [Bike Buyer]:

Predicted	0 (Actual)	1 (Actual)
0	6014	3476
1	3338	5656

Counts for Neural_Network on [Bike Buyer]:

Predicted	0 (Actual)	1 (Actual)
0	6521	3092
1	2831	6040

Wynik macierzy klasyfikacji po wyłączeniu parametrów First and Last Name.

Powyżej widać, że wyniki delikatnie uległy poprawie, jednak nie jest to znaczący wynik.

Zmień kolumny wejściowe dla tego modelu, czy można w ten sposób poprawić wyniki dla tego algorytmu?

Kilkukrotna próba zmiany kolumn wejściowych nie spowodowała poprawy wyników.