

POLITECHNIKA BIAŁOSTOCKA

WYDZIAŁ MECHANICZNY

PRACOWNIA SPECJALISTYCZNA

Sztuczne sieci neuronowe i systemy ekspertowe

Ćwiczenie nr 3

Projektowanie optymalnego zespołu klasyfikatorów

KOD PRZEDMIOTU: MYAR2S22003M

Autor: Ostaszewicz Dawid

Kierunek: Automatyka i Robotyka, II stopień

Prowadzący: dr inż. Marcin Derlatka

Cel ćwiczenia

Celem ćwiczenia jest zaprojektowanie i implementacja optymalnego zespołu klasyfikatorów dla przesłanego zbioru danych. Zespół ma obejmować co najmniej jeden klasyfikator realizujący każdą z technik: bagging, boosting oraz stacking. Zadanie obejmuje:

- Wykorzystanie klasyfikatorów bazowych opartych na sieciach neuronowych, takich jak MLP (Multilayer Perceptron).
- Implementację klasyfikatorów bagging i boosting, takich jak BaggingClassifier, AdaBoostClassifier oraz GradientBoostingClassifier.
- Zastosowanie metody stacking, która łączy wyniki klasyfikatorów bazowych przy pomocy modelu meta-learnera, np. XGBoost.
- Eksperymentowanie z różnymi hiperparametrami klasyfikatorów oraz technikami optymalizacji, aby uzyskać jak najlepsze wyniki klasyfikacji.

1 Wczytanie danych

Dane wejściowe w formacie ARFF są wczytywane przy użyciu funkcji `arff.loadarff`. Zbiory danych są następnie konwertowane na format `pandas DataFrame` dla łatwiejszej manipulacji danymi.

```
1 Dry_Bean_Dataset_data, Dry_Bean_Dataset_meta = arff.loadarff('
   Dry_Bean_Dataset.arff')
2 testing_beans_data, testing_beans_meta = arff.loadarff('
   testing_beans.arff')
3 training_beans_data, training_beans_meta = arff.loadarff('
   training_beans.arff')
4
5 Dry_Bean_Dataset = pd.DataFrame(Dry_Bean_Dataset_data)
6 testing_beans = pd.DataFrame(testing_beans_data)
7 training_beans = pd.DataFrame(training_beans_data)
```

2 Przygotowanie danych

W tej części skryptu wydzielamy cechy (X) i etykiety (y) zarówno z danych treningowych, jak i testowych. Etykiety są konwertowane z typu `byte` na `string`.

```

1 X_train = training_beans.drop('Class', axis=1)
2 y_train = training_beans['Class']
3 X_test = testing_beans.drop('Class', axis=1)
4 y_test = testing_beans['Class']
5
6 y_train = y_train.str.decode('utf-8')
7 y_test = y_test.str.decode('utf-8')

```

3 Kodowanie etykiet

Użyto `LabelEncoder` do kodowania etykiet w postaci liczbowej, co jest wymagane do trenowania klasyfikatorów.

```

1 le = LabelEncoder()
2
3 y_train = le.fit_transform(y_train)
4 y_test = le.transform(y_test)

```

4 Inicjalizacja klasyfikatorów

Inicjalizowane są klasyfikatory bazowe: sieć neuronowa (`MLPClassifier`) oraz klasyfikator XGBoost (`XGBClassifier`). Sieć neuronowa jest używana w zespole bagging.

```

1 base_nn = MLPClassifier(hidden_layer_sizes=(100,), max_iter
    =500, random_state=42)
2 xgb_clf = XGBClassifier(random_state=42)

```

5 Bagging Classifier

Tworzy się klasyfikator `BaggingClassifier`, który stosuje metodę bagging z wykorzystaniem sieci neuronowej jako klasyfikatora bazowego. Klasyfikator jest trenowany na danych treningowych, a następnie dokonujemy predykcji na danych testowych.

```

1 bagging_clf = BaggingClassifier(estimator=base_nn, n_estimators
    =10, random_state=42)
2 bagging_clf.fit(X_train, y_train)

```

```
3 y_pred_bagging = bagging_clf.predict(X_test)
```

6 Boosting Classifier

Następnie tworzymy klasyfikator boostingowy `GradientBoostingClassifier`, który jest trenowany na danych treningowych.

```
1 boosting_clf = GradientBoostingClassifier(n_estimators=50,
      random_state=42)
2 boosting_clf.fit(X_train, y_train)
3 y_pred_boosting = boosting_clf.predict(X_test)
```

7 Stacking Classifier

W przedstawionym fragmencie kodu zastosowano metodę stacking do łączenia predykcji trzech klasyfikatorów bazowych. Proces rozpoczyna się od zdefiniowania trzech modeli bazowych: `bagging_clf`, `boosting_clf` oraz `xgb_clf`, które odpowiadają za generowanie predykcji na podstawie danych treningowych. Następnie, przy użyciu funkcji `stacking`, wykonywane jest połączenie wyników predykcji tych klasyfikatorów. Dla każdego modelu bazowego obliczane są predykcje "out-of-fold", które stanowią dane wejściowe do meta-learnera.

Funkcja `stacking` przyjmuje następujące parametry: dane treningowe `X_train`, `y_train`, dane testowe `X_test`, oraz liczbę foldów krosvalidacji (`n_folds=4`). Ustawienie `regression=False` wskazuje, że problem jest klasyfikacyjny, a parametr `metric=accuracy_score` określa miarę oceny modelu. Dodatkowo, parametr `stratified=True` zapewnia zachowanie proporcji klas w trakcie krosvalidacji.

```
1 models = [bagging_clf, boosting_clf, xgb_clf]
2 S_train, S_test = stacking(models,
3 X_train, y_train, X_test,
4 regression=False,
5 mode='oof_pred_bag',
6 needs_proba=False,
7 save_dir=None,
8 metric=accuracy_score,
9 n_folds=4,
10 stratified=True,
```

```

11  shuffle=True,
12  random_state=42,
13  verbose=2)
14
15  final_estimator = XGBClassifier(random_state=42)
16
17  final_estimator = final_estimator.fit(S_train, y_train)
18  y_pred_stacking = final_estimator.predict(S_test)

```

8 Optymalizacja hiperparametrów dla MLPClassifier

Poniższy kod przedstawia konfigurację siatki hiperparametrów (`param_grid`) do wykorzystania z klasą `GridSearchCV`. Celem jest optymalizacja klasyfikatora `MLPClassifier` z różnymi ustawieniami warstw ukrytych, funkcji aktywacji, solverów oraz parametru regularyzacji `alpha`.

```

1  param_grid = {
2      'hidden_layer_sizes': [(50,), (100,), (150,)],
3      'activation': ['relu', 'tanh', 'logistic'],
4      'solver': ['adam', 'sgd'],
5      'alpha': [0.0001, 0.001, 0.01]
6  }

```

Program automatycznie sprawdza konfiguracje modeli dla parametrów zadanych w tym listingu. Tabela z wynikami została wstawiona na końcu sprawozdania.

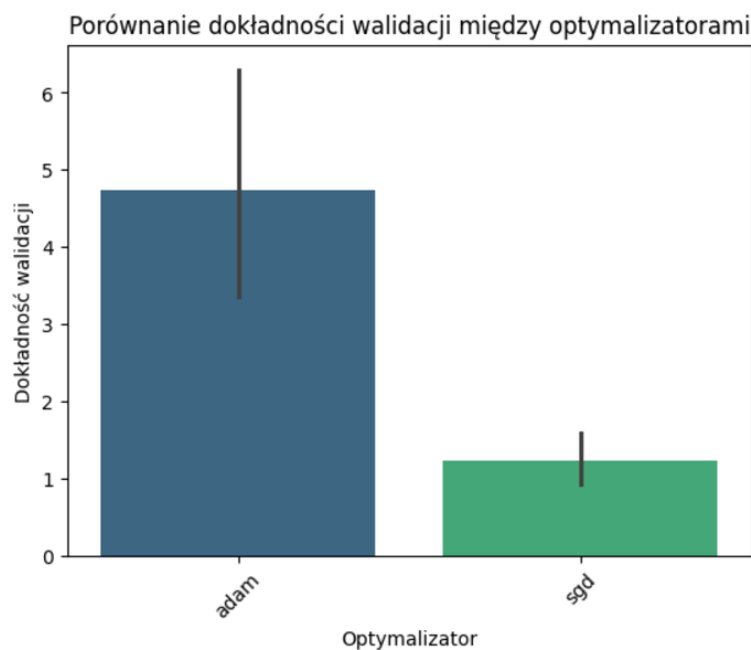
9 Wyniki

Rysunek 1 przedstawia porównanie dokładności walidacji uzyskanej przez modele optymalizatorów: Adam i SGD. Na osi X umieszczone zostały poszczególne optymalizatory, a oś Y przedstawia średnią dokładność walidacji, którą osiągnęły modele w trakcie procesu uczenia. Wysokość słupków pokazuje skuteczność poszczególnych optymalizatorów w poprawie wyników walidacji, umożliwiając ich bezpośrednie porównanie.

Z wykresu wynika, że optymalizator Adam osiąga wyraźnie lepsze wyniki niż SGD. Adam przewyższa SGD pod względem średniej dokładności walidacji, co sugeruje jego większą efektywność w optymalizacji modelu. Powodem, dla którego Adam wypada lepiej, jest jego zdolność do dynamicznego dostosowywania współczynnika uczenia do każdej wagi w modelu. Algorytm

ten korzysta z adaptacyjnych współczynników uczenia, które automatycznie modyfikują tempo nauki na podstawie historii gradientów, co pozwala na szybszą i bardziej stabilną zbieżność, szczególnie w przypadku trudnych funkcji kosztu.

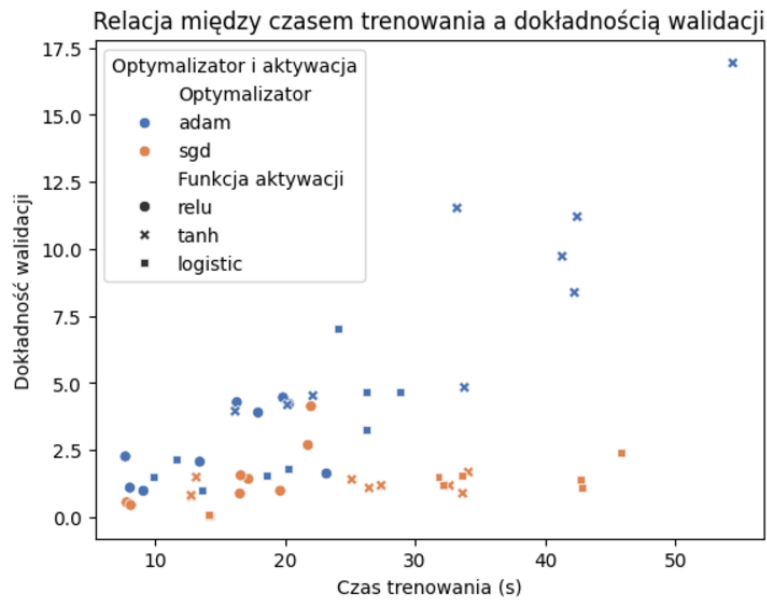
Z kolei optymalizator SGD, mimo że jest prosty i szeroko stosowany, może wykazywać wolniejszą zbieżność i być bardziej podatny na wpadanie w minima lokalne, szczególnie gdy tempo uczenia nie jest odpowiednio dobrane. W związku z tym, dla omawianych modeli, Adam okazuje się bardziej efektywny, co znajduje odzwierciedlenie w wyższej średniej dokładności walidacji w porównaniu do SGD. Rysunek 2 typu scatterplot ilustruje zależność między czasem trenowania



Rysunek 1: Czas treningu

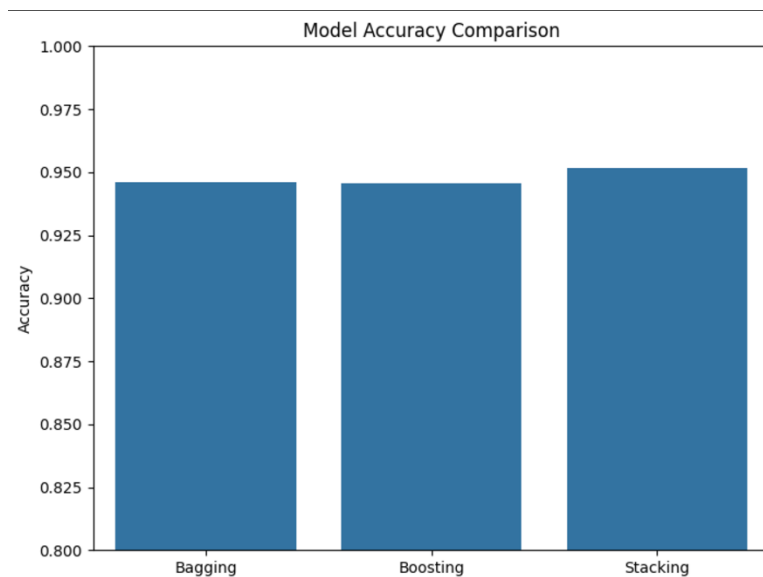
modelu (na osi X) a jego dokładnością walidacji (na osi Y). Punkty na wykresie reprezentują różne optymalizatory i funkcje aktywacji, a ich rozmieszczenie wskazuje, jak czas treningu wpływa na osiąganą dokładność. Każdy punkt jest oznaczony różnym kolorem w zależności od używanego optymalizatora, co pozwala na rozróżnienie wpływu różnych optymalizatorów na wyniki. Dodatkowo, styl punktów wskazuje, która funkcja aktywacji została zastosowana w modelu, co umożliwia ocenę, jak ta zmienna również wpływa na wydajność modeli. Wykres zawiera legendę, która wskazuje, który kolor odpowiada któremu optymalizatorowi oraz jaką funkcję aktywacji reprezentuje dany styl punktu. Dzięki temu wykresowi można szybko zobaczyć, czy istnieje korelacja między czasem trenowania a dokładnością, a także jak różne optymalizatory i funkcje aktywacji wpływają na efektywność modeli.

Ostatecznie, porównanie wyników uzyskanych przez różne metody jest przedstawione na Rysunku 4. Zgodnie z wykresem, metoda stacking osiąga średnio lepsze wyniki w porównaniu do baggingu i boosting. Stacking, który polega na łączeniu predykcji różnych klasyfikatorów



Rysunek 2: Czas treningu

bazowych przy użyciu modelu drugiego poziomu (meta-learnera), wydaje się skuteczniejszy w wykorzystywaniu różnorodności wyników uzyskanych przez modele bazowe.



Rysunek 3: Czas treningu

10 Wnioski

Porównywanie wyników uzyskanych przez różne modele klasyfikacyjne, bazujących na odmiennych założeniach teoretycznych i strukturach wewnętrznych, jest z metodologicznego punktu widzenia nieadekwatne. Brak pełnej znajomości struktury przestrzeni problemowej, w której operują te modele, stanowi kluczowy czynnik uniemożliwiający rzetelną interpretację wyników.

Niejednorodność struktur wewnętrznych sprawia, że porównanie wyników tych modeli jest

nieadekwatne bez uwzględnienia potencjalnych różnic w przestrzeni cech. Przestrzeń problemowa, w której modele te operują, jest często nieliniowa, skomplikowana i zależna od specyfiki zbioru danych. Bez pełnej wiedzy o jej strukturze, nawet dokładne wyniki jednego modelu mogą być nieodpowiednie lub źle zinterpretowane w kontekście innego, co prowadzi do błędnych wniosków.

W związku z tym, brak pełnej świadomości tego, jak dany algorytm przetwarza dane, w jaki sposób interpretuje zależności pomiędzy cechami, a także w jaki sposób przebiega optymalizacja, sprawia, że próby porównań pomiędzy algorytmami mogą prowadzić do mylnych wniosków, które zniekształcają rzeczywisty obraz efektywności modeli w kontekście złożoności problemu.

Jeżeli proces, nawet skomplikowany, można opisać analitycznie, zdecydowanie lepiej jest zastosować podejście analityczne. Każdy z modeli może również różnie reagować na różne dane wejściowe, co nie zostało sprawdzone, dlatego z przeprowadzonych badań nie można wyciągnąć poprawnych wniosków. Przeprowadzenie większej liczby badań na sprawdzenie skuteczności działania modeli, czasowo wykracza poza czas przeznaczony na realizację zajęć. (mam nadzieję, że to nic złego że próbowałem Pythona zamiast Weki, ale ciekawiło mnie jak to by działa :).

Ogólnie czas uzyskiwania tabelki na końcu to trochę ponad 2 godziny. Ciekawe jest też to, że można było wygenerować dane bezpośrednio do formuły Latex z poziomu Pythona. Tylko się nie mieściła na stornie, ale można ją przybliżyć w pdf i nie powinna tracić jakości.

Cena transakcji (t)	Dokładność: ułtkiej	Cena parafekla (t)	Spisat: kłd na rynek	Emalja aktywniej	Breakdown: czy	Wersyj: strzo	Ory: mialor	Shorik: kłrparametr	Spisat: dokłdność: CV	Weak: CV2	Weak: CV3	Weak: CV4	Weak: CV5	Okolnosc: CV	L: rze: syk
6416031	0495101	0409255	040213b	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402102	0402944	0405433	0409059	0409026	0403004	13
741463	0404062	0401389	0403148	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0401283	0402954	0402704	0402380	0402440	0401647	48
424964	0400175	0401096	0400175	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0401778	0402956	0402704	0401709	0402400	0401570	17
2503666	0400693	0401731	0400173	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402230	0402956	0402704	0402909	0402906	0401670	16
2461445	0401133	0401074	0401133	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402238	0402956	0402704	0402909	0402906	0401670	29
7228088	0402641	0400855	0402641	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0401877	0402956	0402704	0402909	0402906	0401670	40
19434223	0403088	0401853	0401853	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402408	0402956	0402704	0402909	0402906	0401670	7
16333072	0403389	0401732	0403389	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0401357	0402956	0402704	0402909	0402906	0401670	38
21693844	0402223	0401733	0402223	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0401333	0402956	0402704	0402909	0402906	0401670	35
8481467	040224	0401721	040224	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402140	0402956	0402704	0402909	0402906	0401670	15
1410258	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402171	0402956	0402704	0402909	0402906	0401670	2
217855	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402106	0402956	0402704	0402909	0402906	0401670	30
1601031	0400543	0400543	0400543	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402107	0402956	0402704	0402909	0402906	0401670	30
2435332	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	33
2510227	040525	0401000	040525	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	6
1418578	040525	0401000	040525	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	33
2541303	0401311	0401311	0401311	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	45
2646798	040525	0401000	040525	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	42
54413109	040224	040224	040224	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	18
1408212	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	36
8380662	0401772	0401772	0401772	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	9
33217859	0402947	0402947	0402947	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	32
148471	0408614	0408614	0408614	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	31
32414388	0402967	0402967	0402967	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	31
141759	0401319	0401319	0401319	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	34
33789020	0401764	0401764	0401764	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	5
6809066	0402984	0402984	0402984	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	39
1411755	0401409	0401409	0401409	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	22
2527875	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	22
4628289	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	49
33242171	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	20
7494596	0401144	0401144	0401144	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	19
2495739	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	53
11469150	0401445	0401445	0401445	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	19
14301660	0401265	0401265	0401265	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	53
25214616	0402284	0402284	0402284	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	53
2840608	0401108	0401108	0401108	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	24
4248601	0403233	0403233	0403233	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	46
1447625	0402597	0402597	0402597	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	54
18555469	0401627	0401627	0401627	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	26
145875	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	27
2421469	0401000	0401000	0401000	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	27
45456850	0401329	0401329	0401329	ru	0401000	(b)	adum	'verification': 'ru', 'value': '04011', 'hidden': 'ru', 'dms': '(04)', 'value': 'value'	0402108	0402956	0402704	0402909	0402906	0401670	47