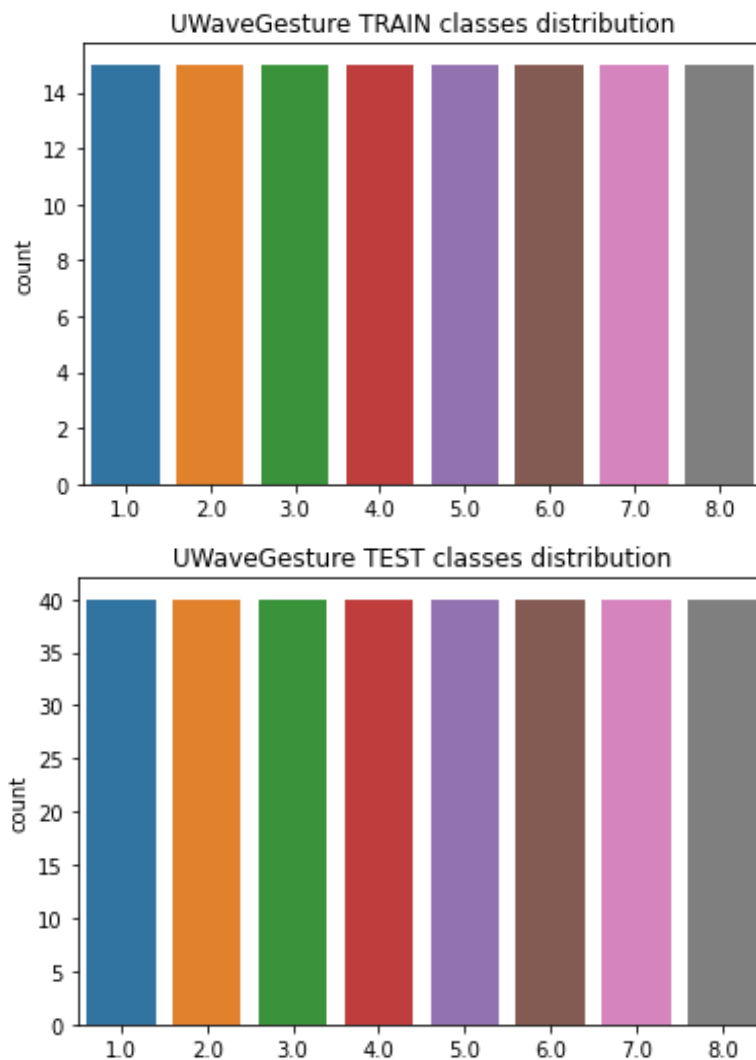


**Invatare Automata**  
**Tema - Etapa 1**  
**Lucian-Florin Grigore 343C4**  
*Facultatea de Automatica si Calculatoare*  
*Universitatea Politehnica, Bucuresti*

*Codul sursa al acestei lucrari poate fi gasit in Google Colab la [acest link](#).*

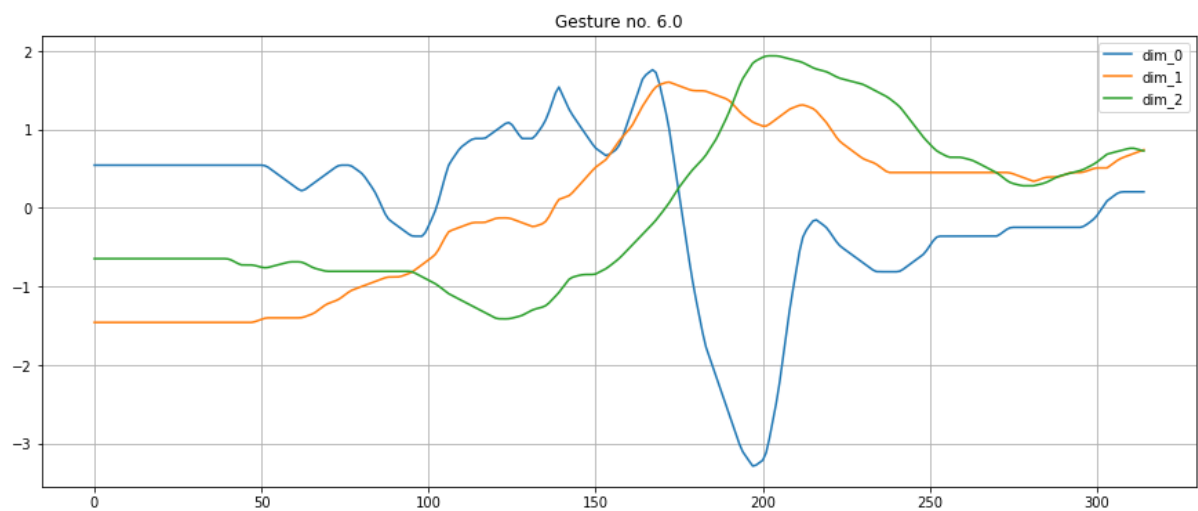
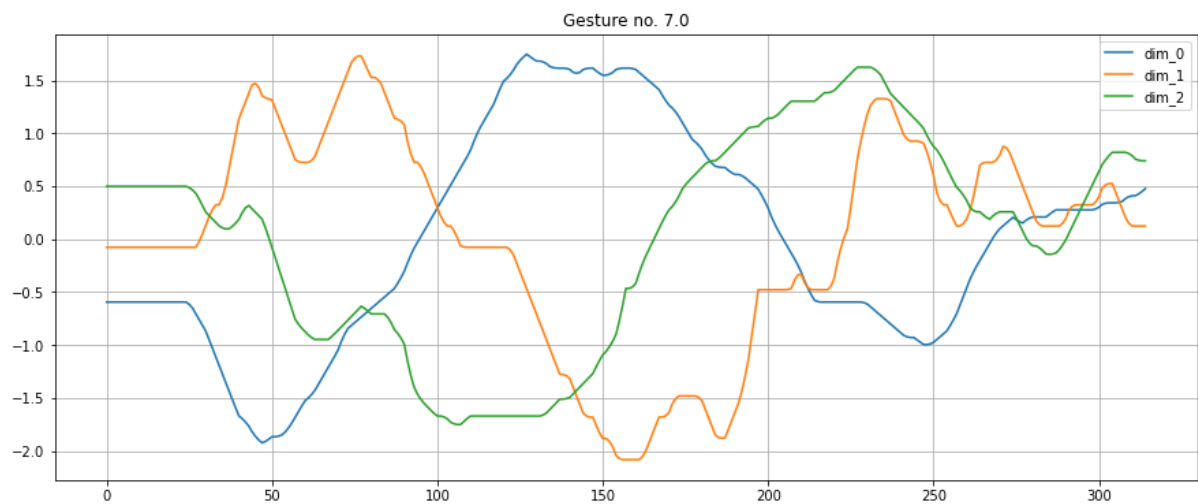
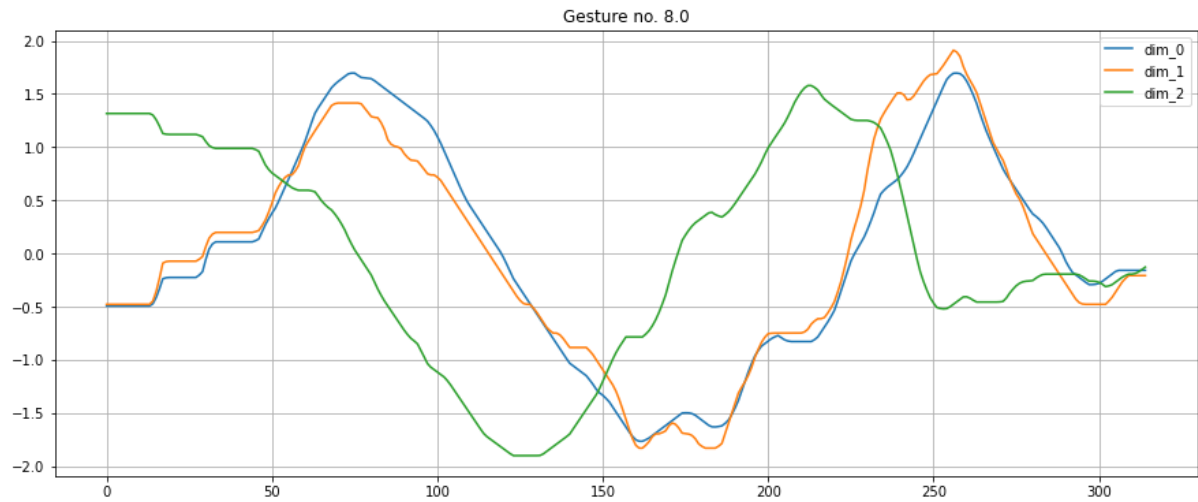
**Cerinta 1. Exploratory Data Analysis**

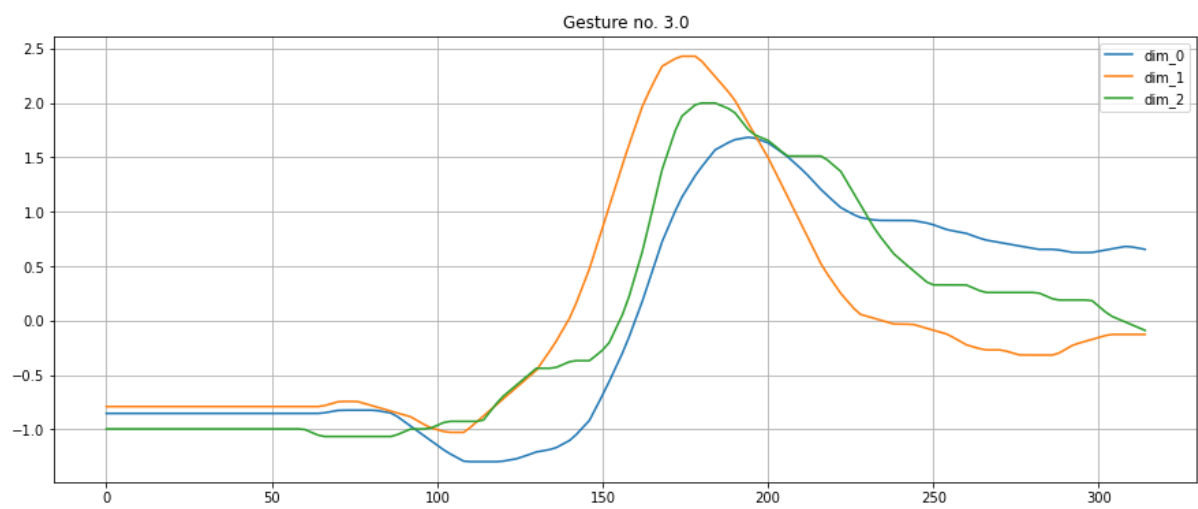
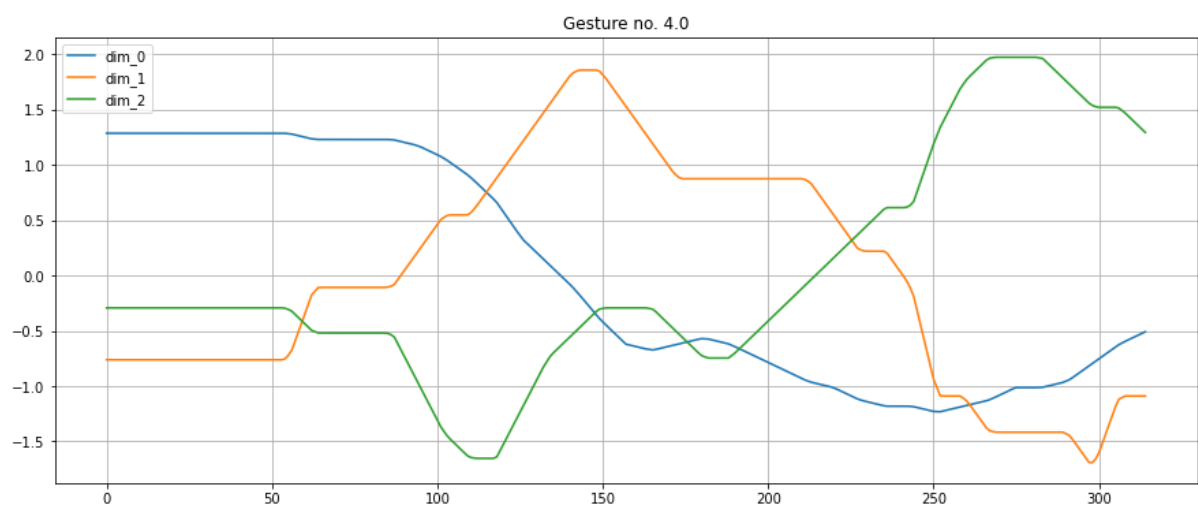
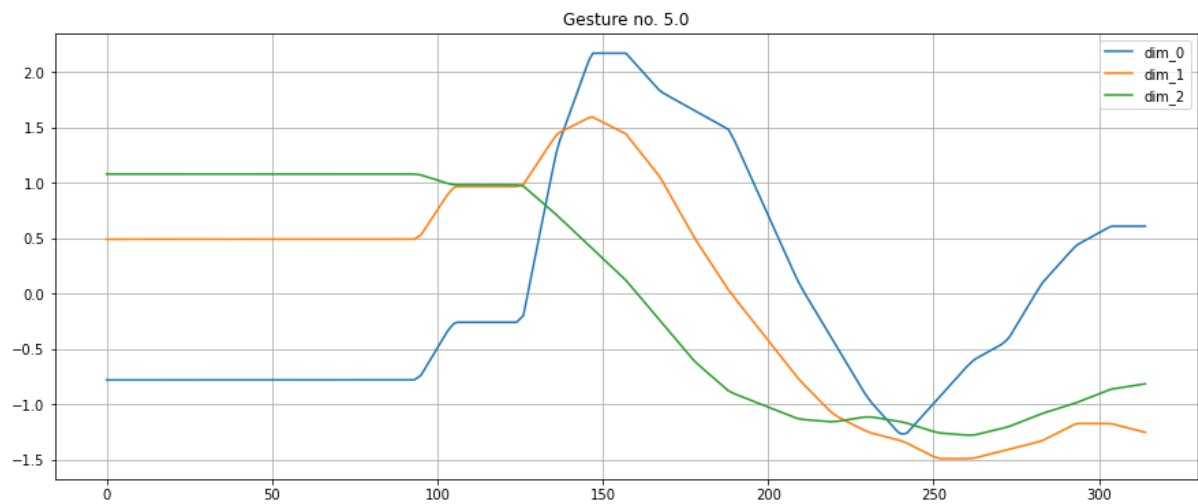
**Frecventa de aparitie a claselor in setul de date pentru UWaveGesture**

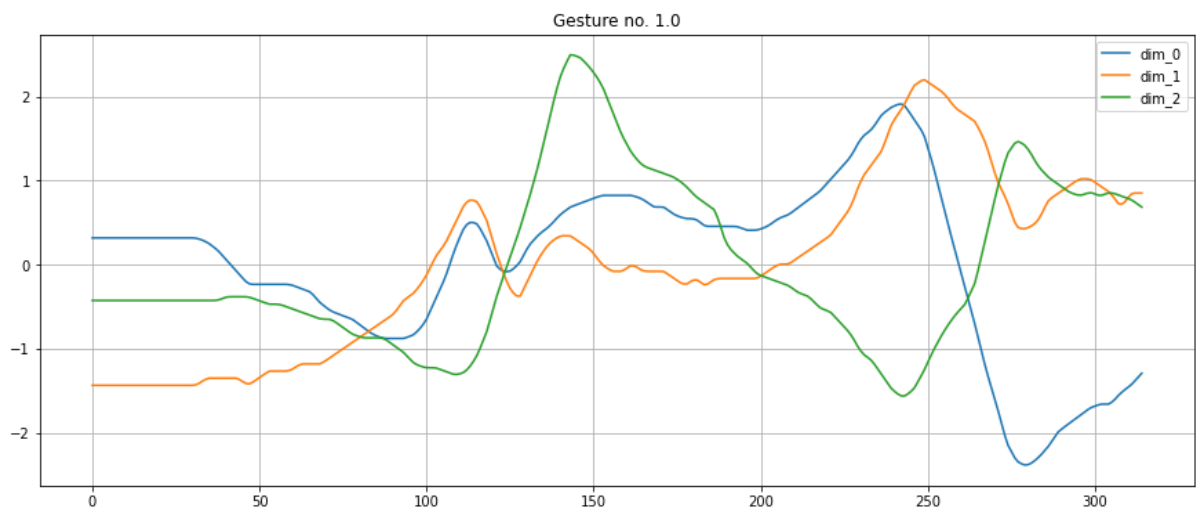
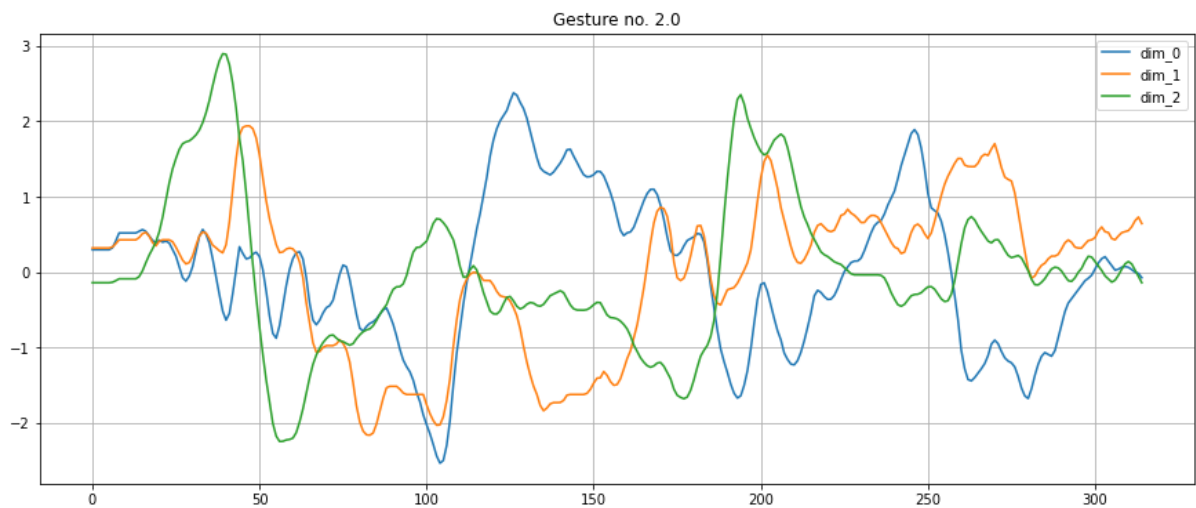


Observam ca setul de date UWaveGesture contine un numar egal de clase atat in setul de antrenare, cat si in cel de testare.

## Afisarea unei serii temporale pentru fiecare gest

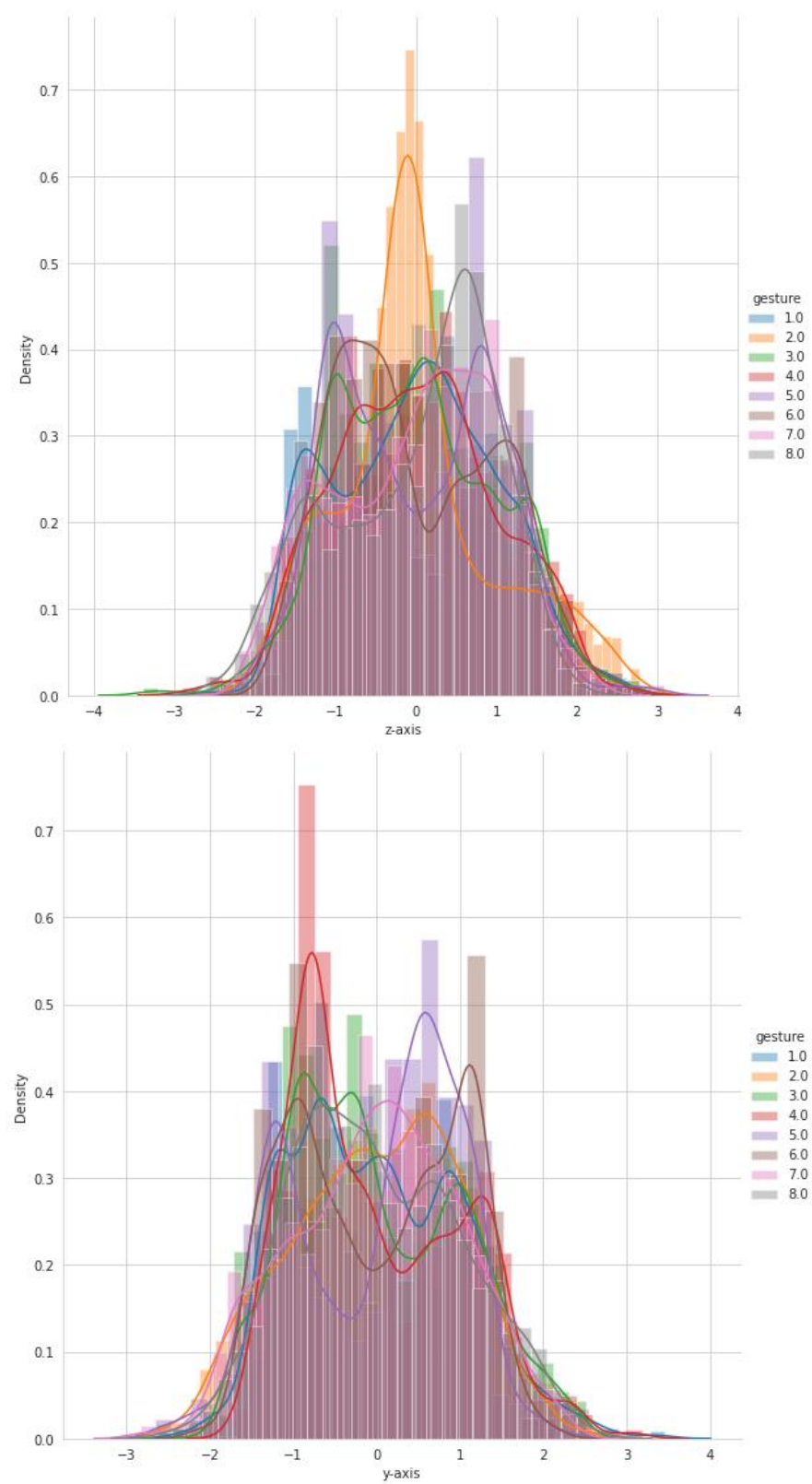


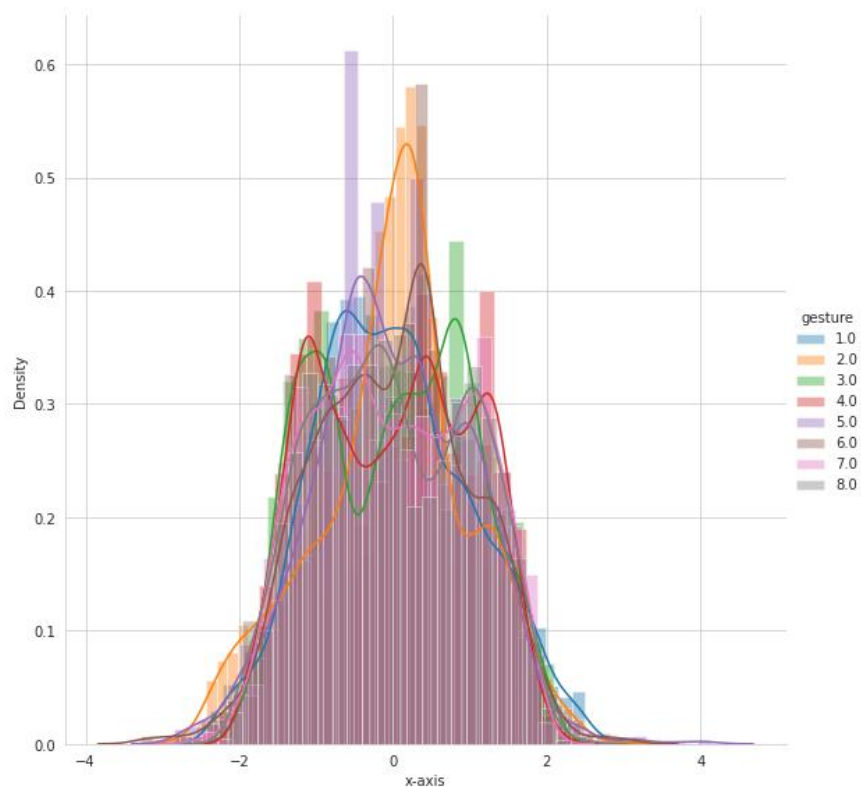




Pentru fiecare serie temporală am luat din setul de antrenare primul exemplu din clasa respectivă.

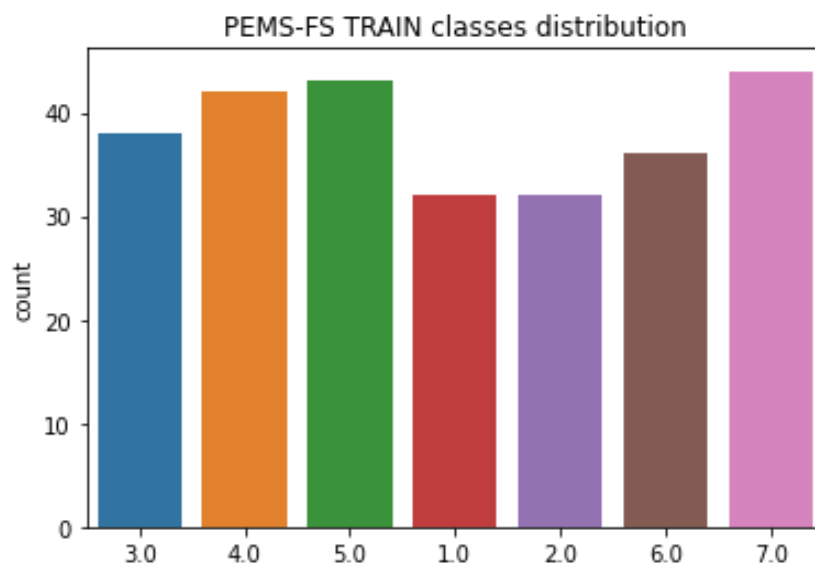
Distributia valorilor per fiecare axa, per gest

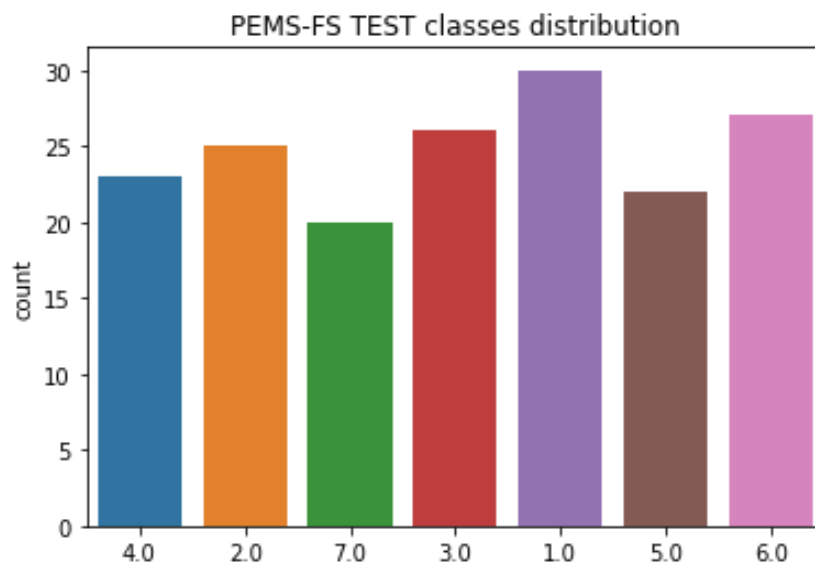




Se poate observa ca majoritatea datelor se afla intr-un interval relativ restrans, ceea ce ar putea reprezenta o dificultate in antrenare, atunci cand se doreste optimizarea modelului si obtinerea unei acuratete foarte buna.

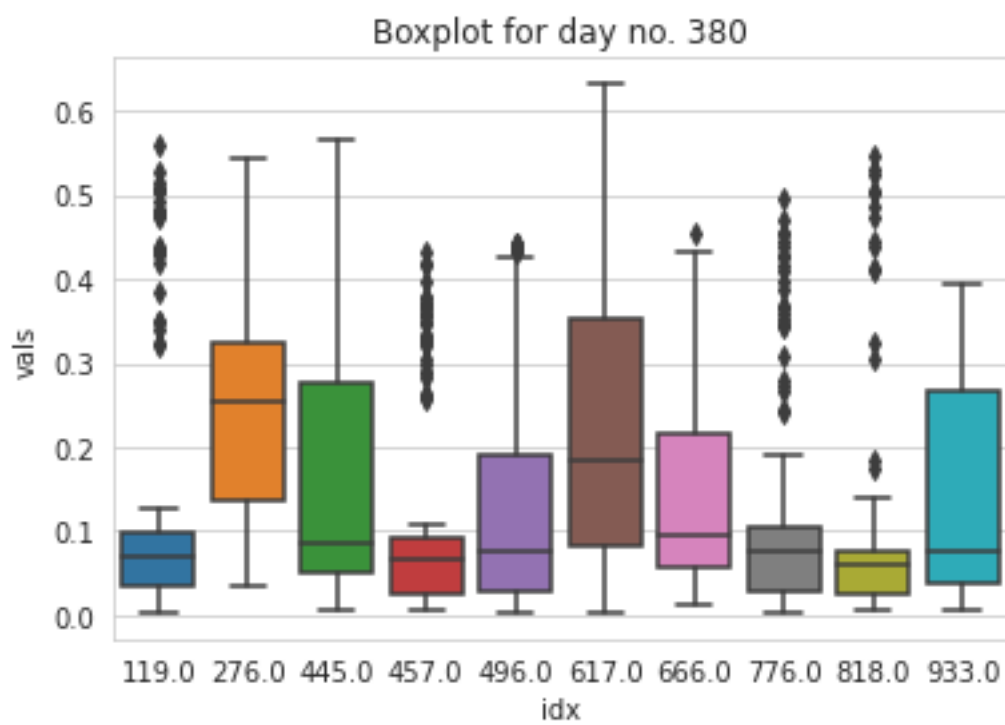
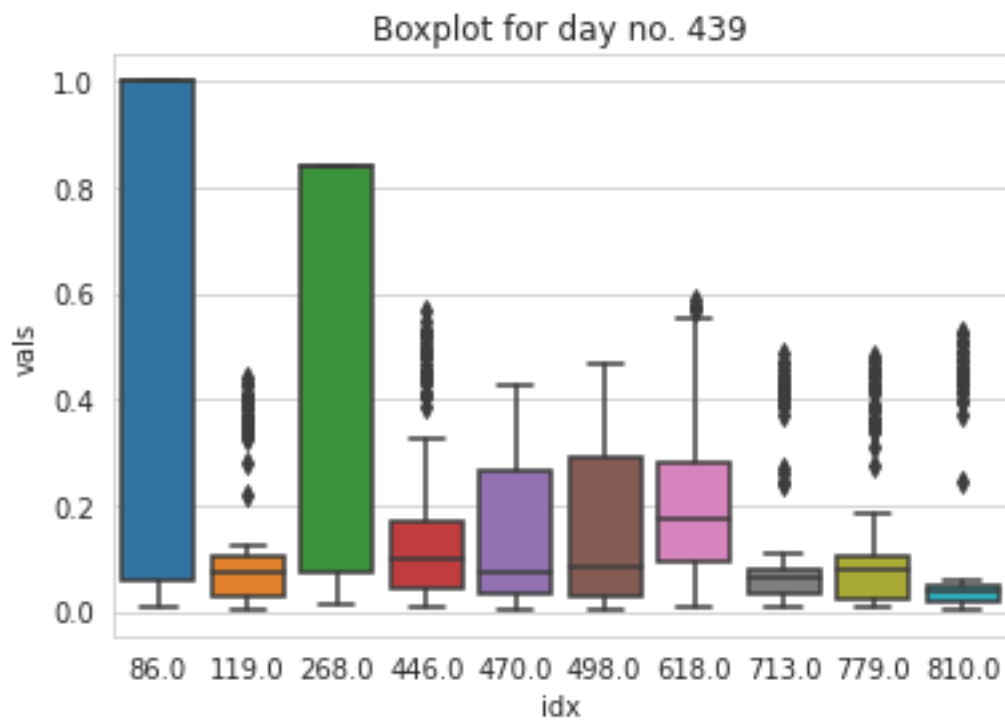
### Frecventa de aparitie a claselor in setul de date pentru UWaveGesture dataset



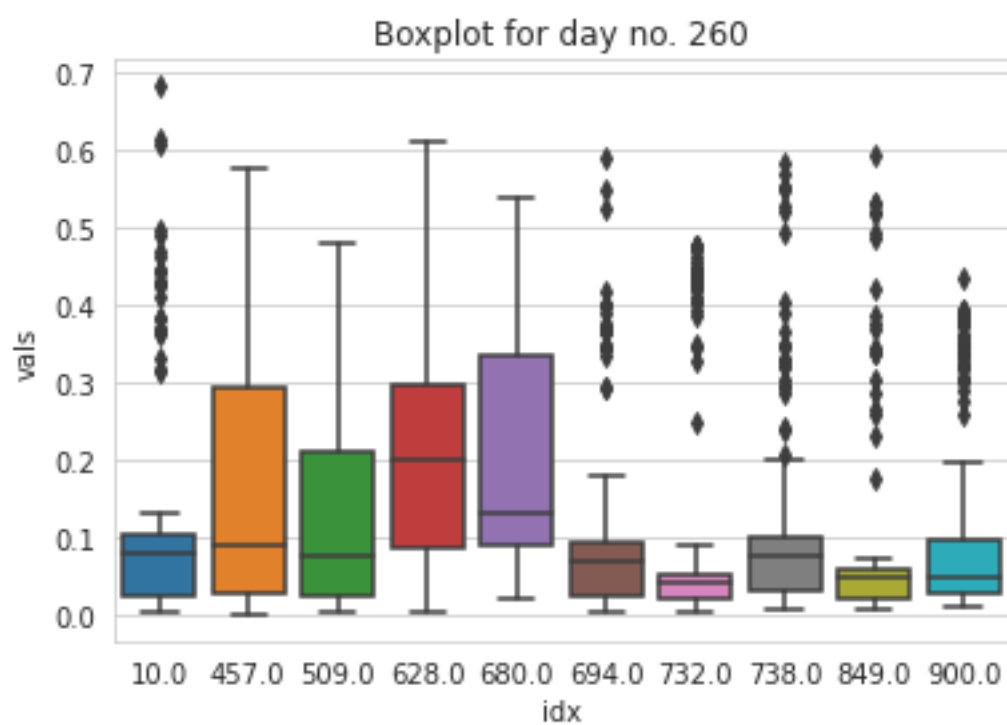
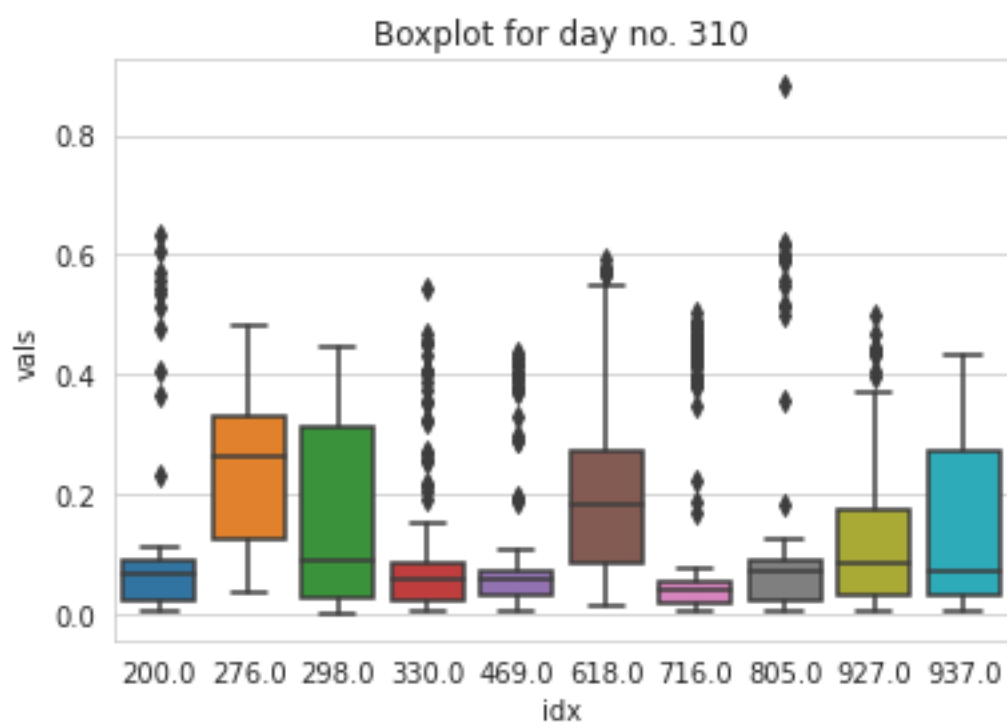


Observam ca setul de date PEMS-FS contine un numar inegal de clase atat in setul de antrenare, cat si in cel de testare. Cu toate acestea, nu se poate considera nicio clasa redundanta, toate avand un numar apropiat de exemple.

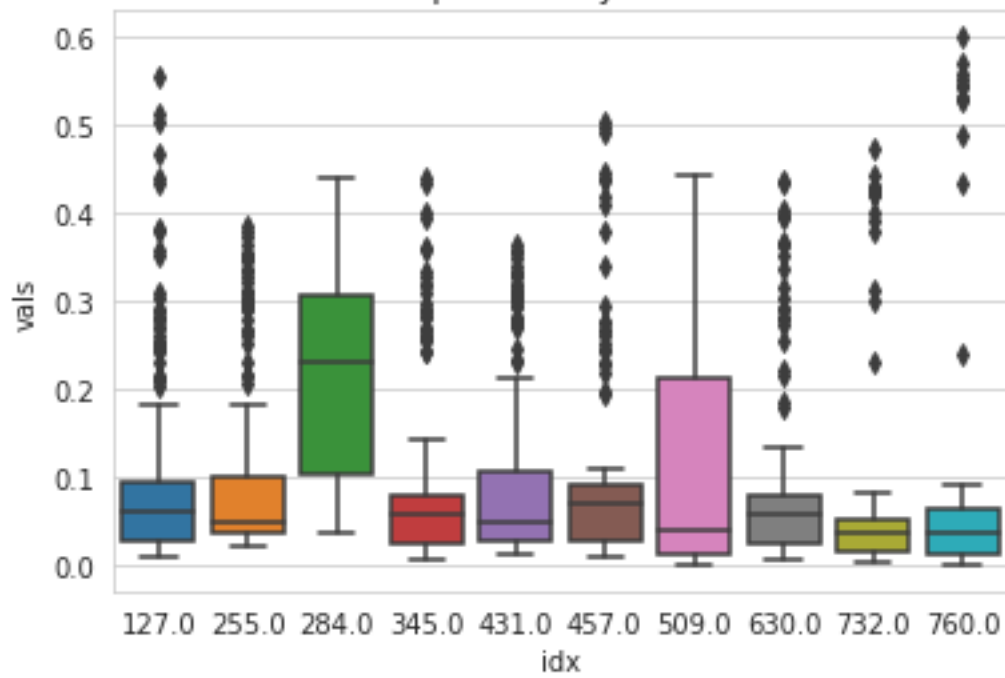
**Varierea ratei de ocupare pentru top 10 senzori cu deviatia cea mai mare pentru 8 zile selectate arbitrar uniform din totalul zilelor**



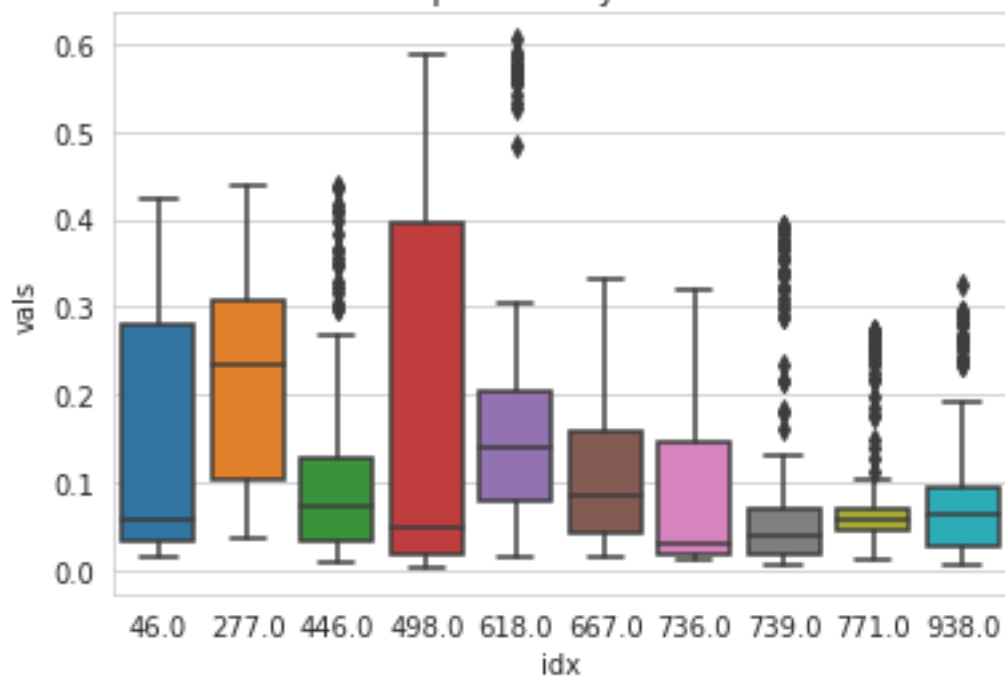


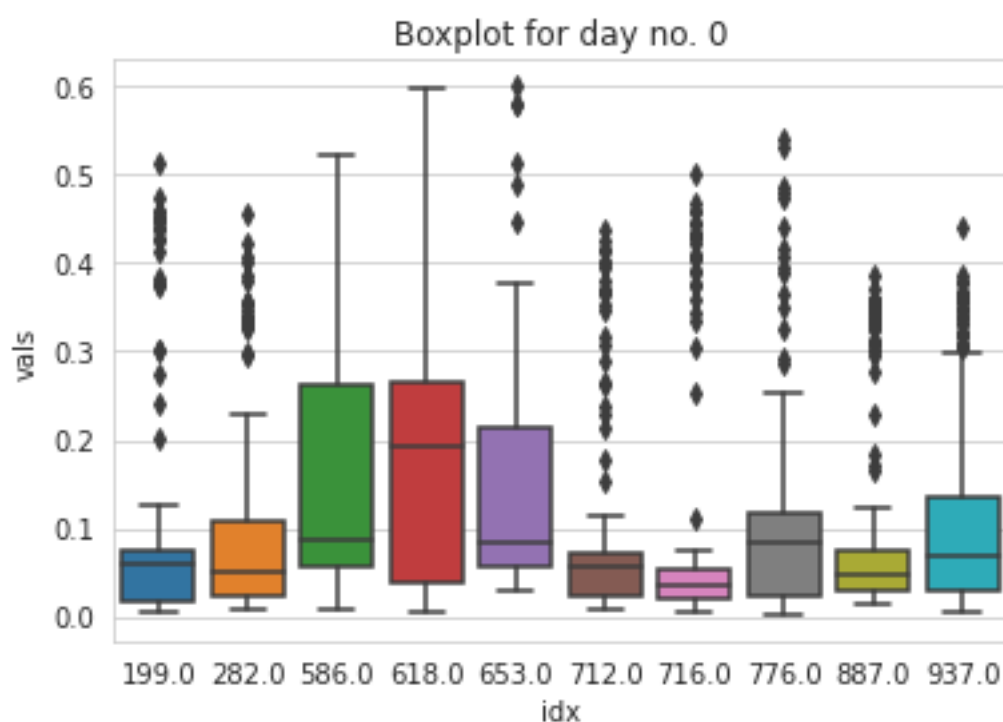
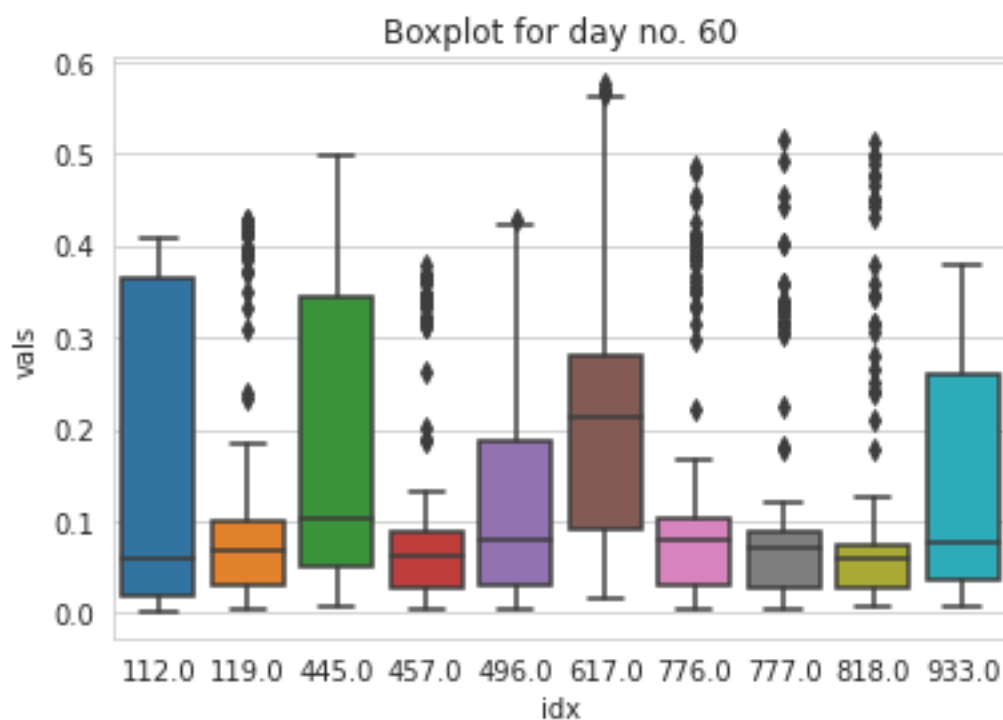


Boxplot for day no. 180



Boxplot for day no. 120



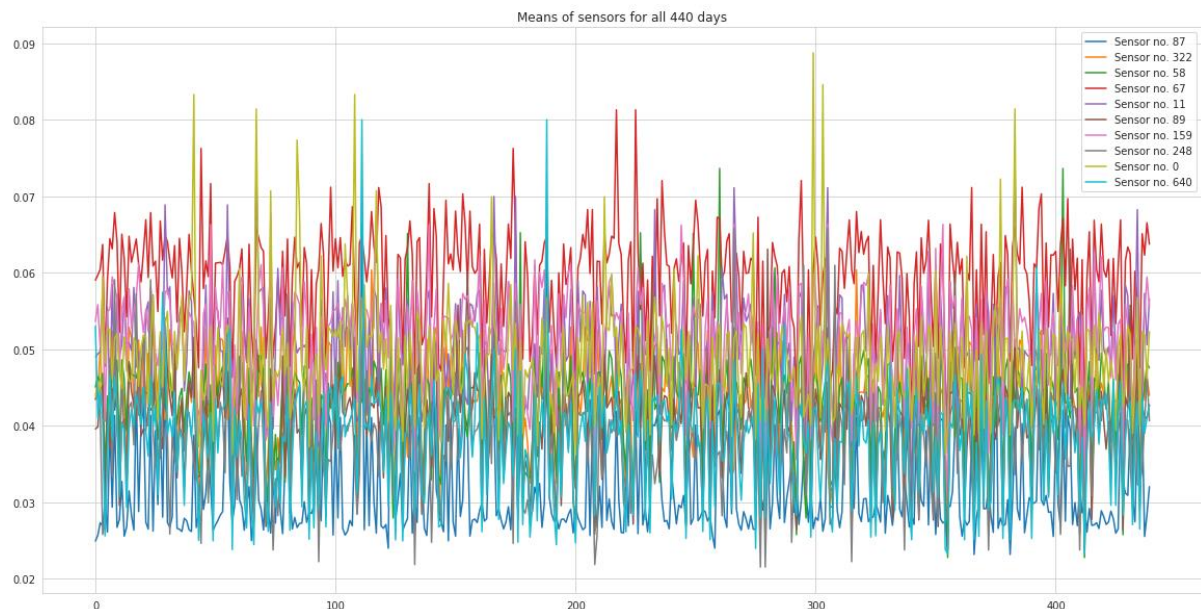


În urma analizei acestor grafice, cu mici excepții, se observă cum majoritatea valorilor înregistrate de senzorii selectați se află înspre limita inferioară a spectrului de valori. Asta poate reprezenta 2 lucruri:

1. Senzorii care au valori înregistrate mai mici (poate prin amplasarea lor în locația respectivă) au și o sensibilitate la variație, rezultând date mai împrăștiate comparativ cu ceilalți senzori.

2. Posibil ca toti senzorii sa aiba valori mai apropiate de limita inferioara si atunci e nevoie de un pas in plus la preprocesarea datelor, dupa confirmarea acestei supozitii.

### **Evolutia mediilor celor mai relevanti 10 senzori pe durata tuturor celor 440 de zile**



Putem remarca ca valorile in general sunt restranse intr-un interval mic de valori, apropiat destul de mult de origine.

**Cerinta 2.** Pentru cerinta a doua am folosit datasetul *UWaveGesture*

**Feature Selection:** Pentru a reduce datele de input la o dimensiune care poate fi gestionata si mai usor de analizat, am aplicat urmatoarele operatii:

- Am impartit fiecare axa (x, y si z) in ferestre de lungime 105 -> rezulta 3 ferestre per fiecare axa = 9 ferestre in total
- Pentru fiecare astfel de fereastră am facut media valorilor din seria de timp
- O intrare X din setul de date reprezinta aceste 9 valori obtinute in urma operatiilor de mai sus

In continuare, analiza atributelor si antrenarea modelelor este realizata pe aceasta noua reprezentare a datelor.

## Extragerea atributelor

Applying mean on x\_axis: -1.1325396825710079e-07  
Applying mean on y\_axis: -1.191991341994739e-07  
Applying mean on z\_axis: -2.6096681093963078e-08  
Applying std on x\_axis: 0.6537793630419304  
Applying std on y\_axis: 0.7495870421412952  
Applying std on z\_axis: 0.7088489128654653  
Applying avg absolute diff on x\_axis: 0.5317972085502645  
Applying avg absolute diff on y\_axis: 0.6496156217923061  
Applying avg absolute diff on z\_axis: 0.6049543094942061  
Applying min on x\_axis: -1.363360857142857  
Applying min on y\_axis: -1.3331605714285715  
Applying min on z\_axis: -1.3638615238095237  
Applying max on x\_axis: 1.360968380952381  
Applying max on y\_axis: 1.3764914000000001  
Applying max on z\_axis: 1.33198819047619  
Applying max-min diff on x\_axis: 2.7243292380952377  
Applying max-min diff on y\_axis: 2.7096519714285714  
Applying max-min diff on z\_axis: 2.695849714285714  
Applying median on x\_axis: -0.02353086190476192  
Applying median on y\_axis: 0.03759409047619046  
Applying median on z\_axis: 0.09206555714285715  
Applying median abs dev on x\_axis: 0.45422849999999999  
Applying median abs dev on y\_axis: 0.6274349095238094  
Applying median abs dev on z\_axis: 0.548894542857143  
Applying IQR on x\_axis: 0.912718819047619  
Applying IQR on y\_axis: 1.2576291285714287  
Applying IQR on z\_axis: 1.158704673809524  
Applying negative count on x\_axis: 632  
Applying negative count on y\_axis: 673  
Applying negative count on z\_axis: 716  
Applying positive count on x\_axis: 688  
Applying positive count on y\_axis: 647  
Applying positive count on z\_axis: 604  
Applying values above mean on x\_axis: 632  
Applying values above mean on y\_axis: 673  
Applying values above mean on z\_axis: 716  
Applying values below mean on x\_axis: 688  
Applying values below mean on y\_axis: 647  
Applying values below mean on z\_axis: 604  
Applying number of peaks on x\_axis: 451  
Applying number of peaks on y\_axis: 442  
Applying number of peaks on z\_axis: 441  
Applying skewness on x\_axis: 0.09000602953872863  
Applying skewness on y\_axis: -0.07910393643284531  
Applying skewness on z\_axis: -0.18151097976780978  
Applying kurtosis on x\_axis: -0.7231033556102662  
Applying kurtosis on y\_axis: -1.181820388478901  
Applying kurtosis on z\_axis: -1.0416019685057663  
Applying energy on x\_axis: 5.642042413121732  
Applying energy on y\_axis: 7.416825685449181  
Applying energy on z\_axis: 6.632561512771293  
Average resultant acc is 44.3750263226313  
Signal magnitude area is 1.786367140079365

De pe urma acestor metrici, valorile obtinute nu indica vreo anomalie evidenta.

## Antrenare de modele ML

Folosind percentile=10 (valoarea default din sklearn) am fi folosit doar un atribut din cele 9, ceea ce este destul de riscant intrucat se pierde foarte multa informatie pentru fiecare exemplu.

### Rezultate pentru folosirea Select Percentile cu percentile=50, adica folosirea a 4 din 9 attribute per fiecare intrare din dataset

```
--- Performance Analysis for Random Forest classifier ---
Best accuracy was 0.7833333333333333
    with params {'bootstrap': False, 'max_depth': 5, 'n_estimators': 100}

--- Performance Analysis for SVM classifier ---
Best accuracy was 0.8333333333333334
    with params {'C': 0.15, 'kernel': 'rbf'}

--- Performance Analysis for XGBoost classifier ---
Best accuracy was 0.725
    with params {'learning_rate': 0.15, 'max_depth': 2, 'n_estimators': 100}
```

### Rezultate pentru folosirea Select Percentile cu percentile=100, adica folosirea tuturor celor 9 din 9 attribute per fiecare intrare din dataset

```
--- Performance Analysis for Random Forest classifier ---
Best accuracy was 0.8
    with params {'bootstrap': False, 'max_depth': 50, 'n_estimators': 50}

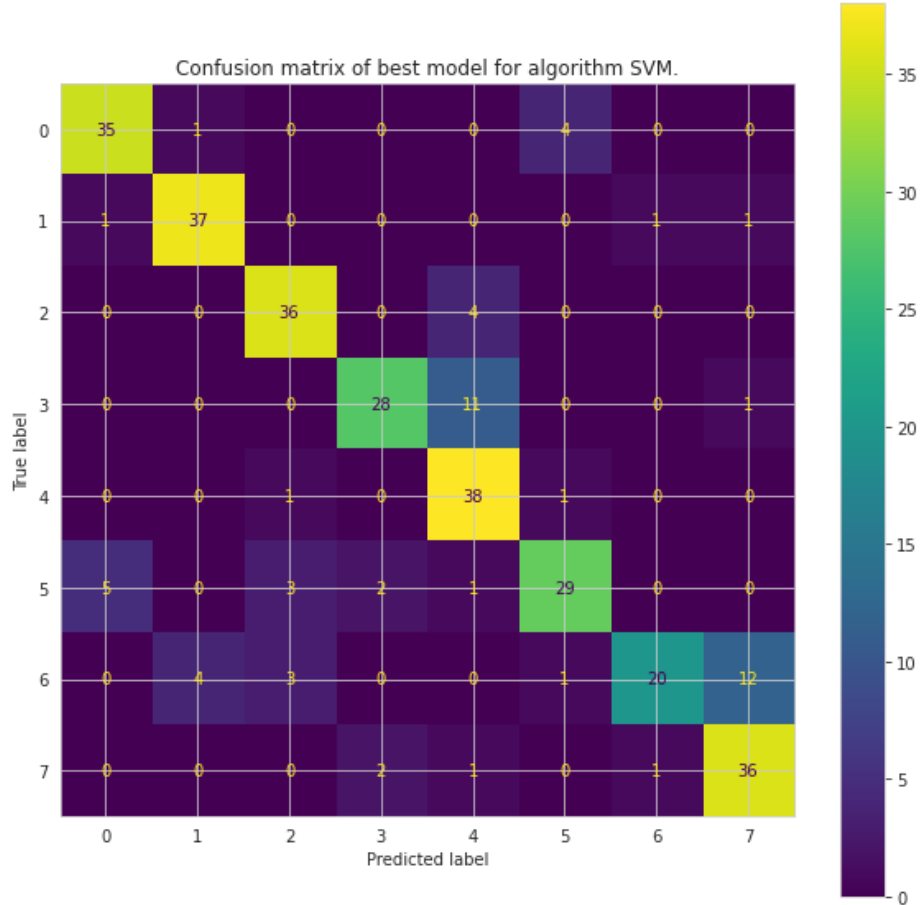
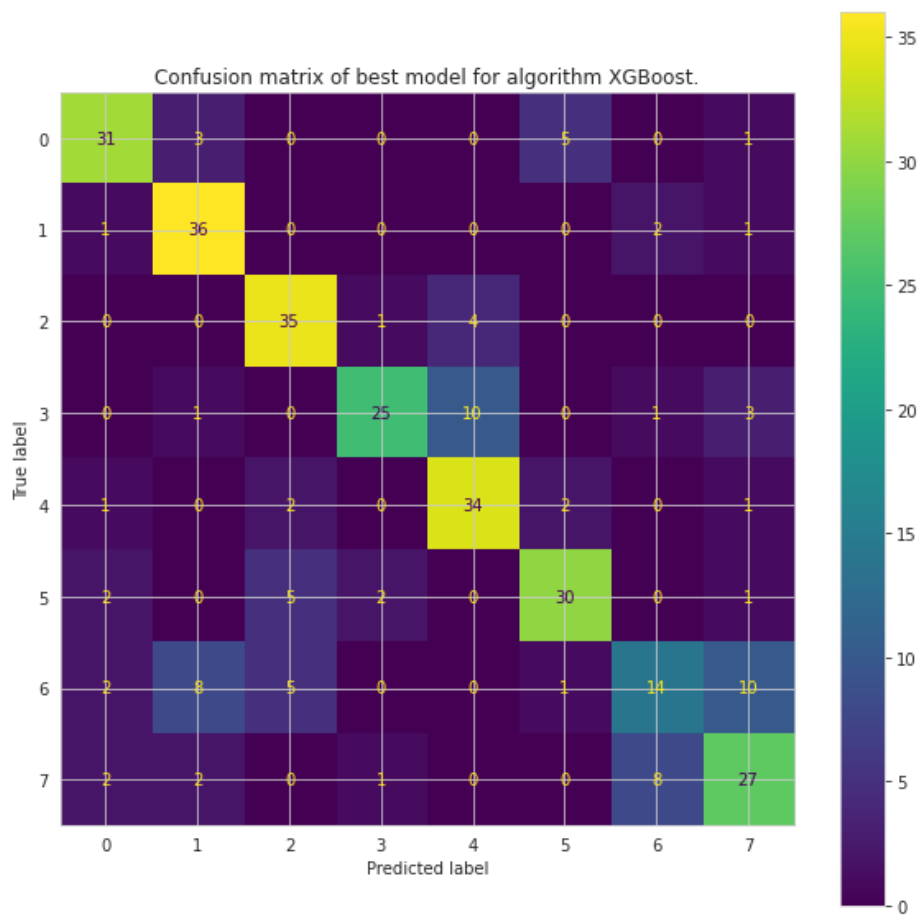
--- Performance Analysis for SVM classifier ---
Best accuracy was 0.8333333333333334
    with params {'C': 0.15, 'kernel': 'rbf'}

--- Performance Analysis for XGBoost classifier ---
Best accuracy was 0.7333333333333333
    with params {'learning_rate': 0.2, 'max_depth': 2, 'n_estimators': 150}
```

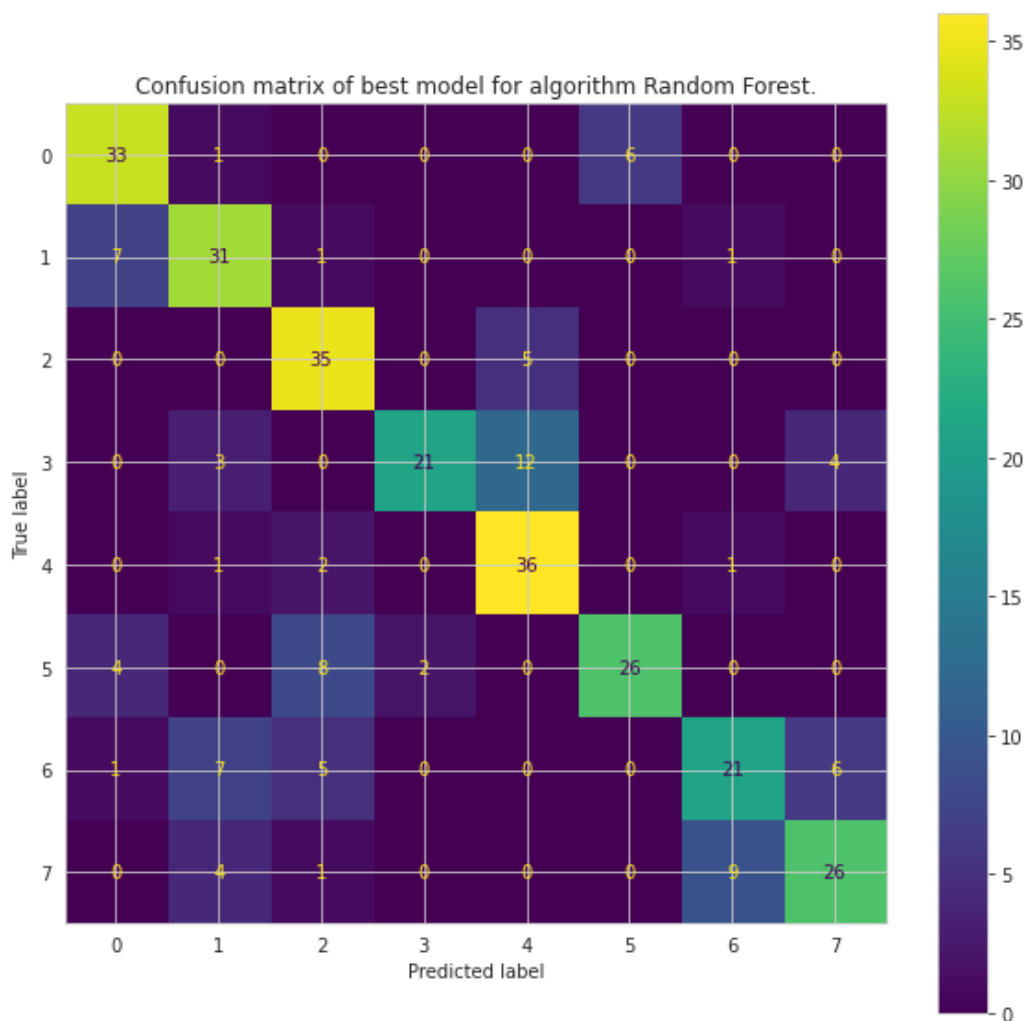
Pentru XGBoost learning\_rate joaca un rol foarte important. Pentru SVC, kernelul "rbf" pare a fi cel mai constant dpdv al performantei obtinute. De asemenea, pare ca valoarea de 0.15 pentru "C" este ideala intrucat ofera cele mai bune performante. Random Forest pare ca prefera un numar finit de estimatori si o adancime maxima care nu este infinita. Dar aici intervine si dimensiunea relativ scazuta a setului de date.

Observam ca in mod constant Support Vector Machine Classifier obtine cea mai buna acuratete pe setul de antrenare. Am considerat in continuare modelul antrenat folosind percentile=50, intrucat acuratetea la antrenare este aceeaasi, dar volumul de date este considerabil mai mic, imbunatatind astfel performanta.

	General Accuracy	Classes	1	2	3	4	5	6	7	8
Classifiers and Parameters (best performing)										
<b>Random Forest</b> <i>Bootstrap: False</i> <i>Max_depth: 50,</i> <i>N_estimators: 50</i>	Train: 0.8 Test: 0.72	Precision Recall F1	0.73 0.82 0.77	0.66 0.77 0.76	0.67 0.87 0.76	<b>0.91</b> 0.52 0.66	0.68 0.9 0.77	0.81 0.65 0.72	0.65 0.52 0.58	0.72 0.65 0.68
<b>Support Vector Machine</b> <i>C: 0.15</i> <i>Kernel: rbf</i>	Train: <b>0.83</b> Test: <b>0.81</b>	Precision Recall F1	0.85 0.87 0.86	0.88 <b>0.92</b> <b>0.9</b>	0.84 0.9 0.87	0.87 0.7 0.77	0.69 0.95 0.8	0.83 0.72 0.77	0.9 0.5 0.65	0.72 0.9 0.8
<b>XGBoost</b> <i>Learning_rate: 0.2</i> <i>Max_depth: 2</i> <i>N_estimators: 150</i>	Train: 0.73 Test: 0.72	Precision Recall F1	0.79 0.77 0.78	0.72 0.9 0.8	0.74 0.87 0.8	0.86 0.62 0.72	0.7 0.85 0.77	0.79 0.75 0.77	0.56 0.35 0.43	0.61 0.67 0.64







Mai sus sunt prezentate rezultatele pentru cea mai buna combinatie de parametrii pentru fiecare algoritm, urmarind: acuratetea generala, recall, precision si F1 (ultimele trei la nivel de clasa).

De asemenea, sunt afisate matricile de confuzie pentru acesti algoritmi.

Toate aceste date sunt obtinute de pe urma predictiilor pe setul de testare.