

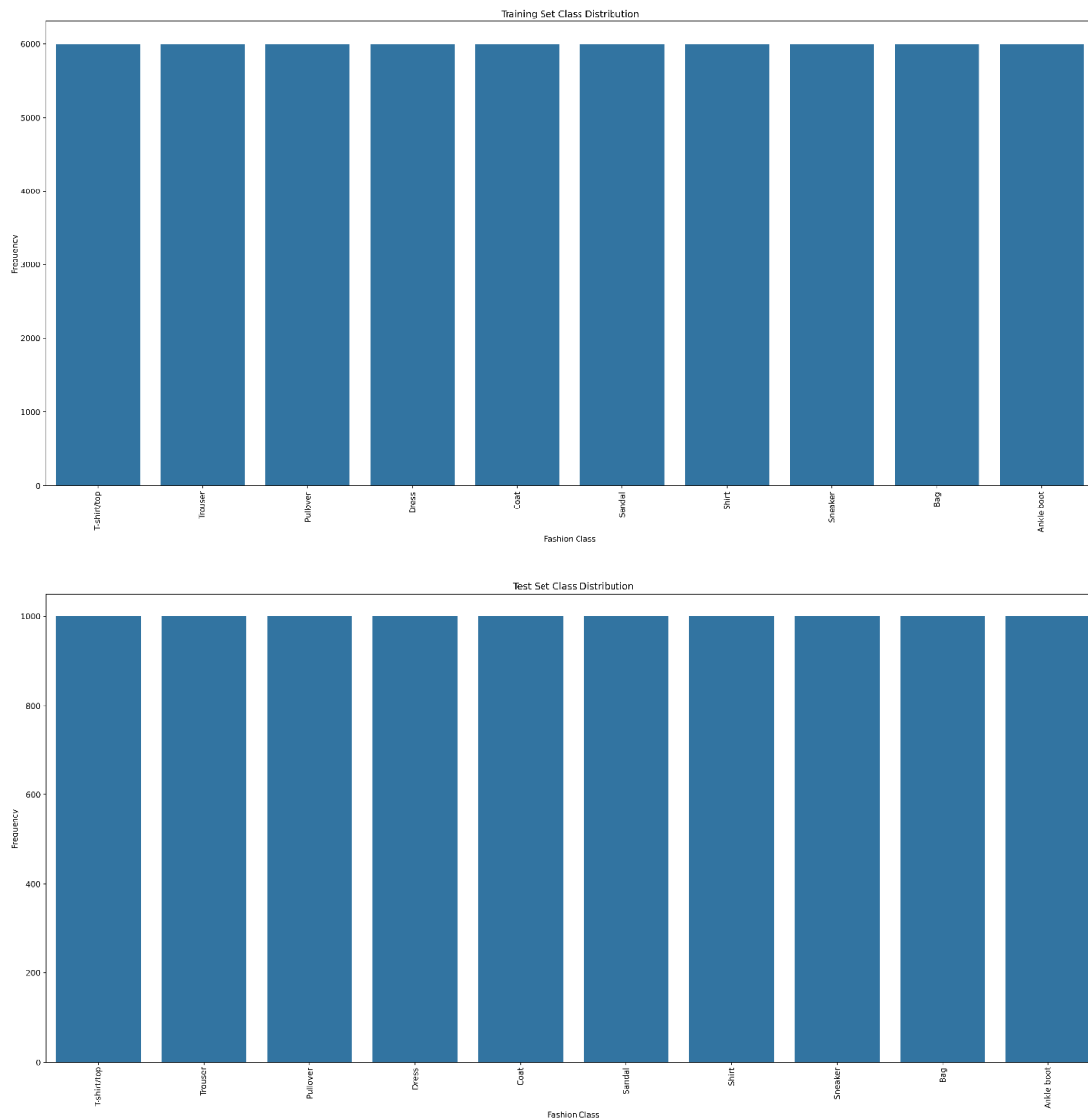
Tema 1 ML

Oana Maria Bacaran, 342C5

Pe datasetul fashion, in etapa de preprocesare, realizam un resize pe toate imaginile din dataset pentru a captura mai multe detalii/features la folosirea metodei HOG. Am folosit metodele de extragere de attribute HOG si PCA in pipeline. Am utilizat HOG deoarece acesta gaseste eficient directiile marginilor, texture si forme si reduce zgomotul pixelilor. In continuare am folosit PCA cu parametrul ncomponents = 50 pentru a pastra cat mai multe features din imaginea initiala si pentru a avea o varianta cat mai mare.

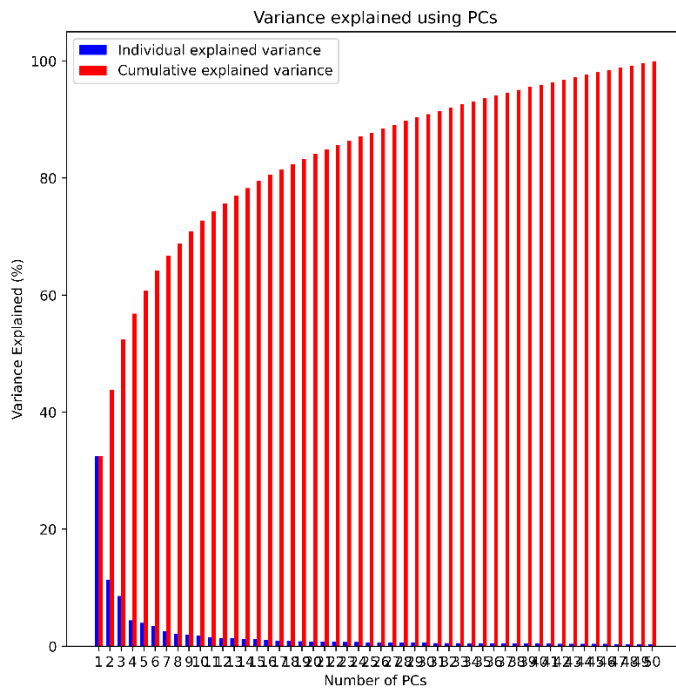
Pe datasetul fashion avem urmatoarele rezultate:

4.2.1. Analiza echilibrului de clase



În imaginile de mai sus a fost realizat analiza echilibrului de clase cu ajutorul unui countplot, folosind biblioteca seaborn. Se observă că numărul de elemente/imagini, atât în training cât și test, este egal per clasă.

4.2.2 Vizualizare cantitativă



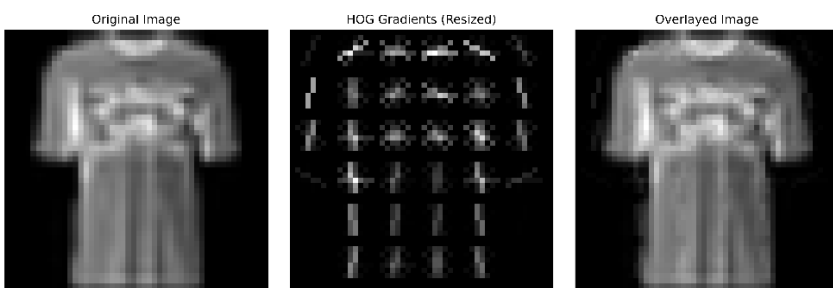
Graficul arată cât de multă informație este explicată de fiecare componentă principală (barele albastre) și cât de multă informație totală este explicată dacă luăm mai multe componente împreună (barele roșii). Observăm că primele componente principale explică cea mai mare parte din varianță, iar după aproximativ 30 de componente, graficul roșu aproape că se stabilizează, ceea ce înseamnă că restul componentelor contribuie foarte puțin.

Acest lucru înseamnă că putem reduce dimensiunea datelor păstrând doar primele 30 de componente, fără să pierdem prea multă informație. Acest lucru ajută la simplificarea procesării datelor și face modelele mai rapide, dar fără să afecteze prea mult performanța. PCA este util aici pentru că elimină redundanțele și ne permite să lucrăm cu mai puține caracteristici.

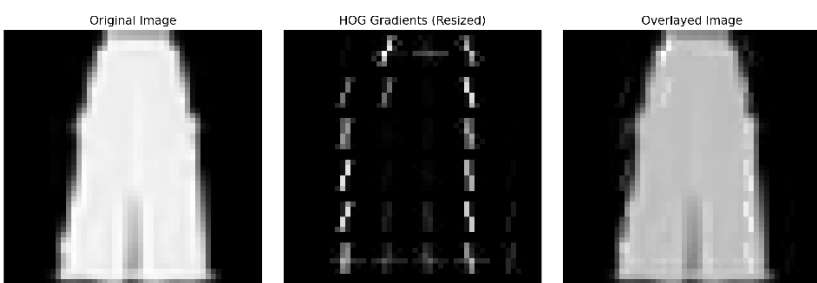
Vizualizare calitativa

Plot-area gradientilor hog peste imaginile initiale (vizualizare intermediara)

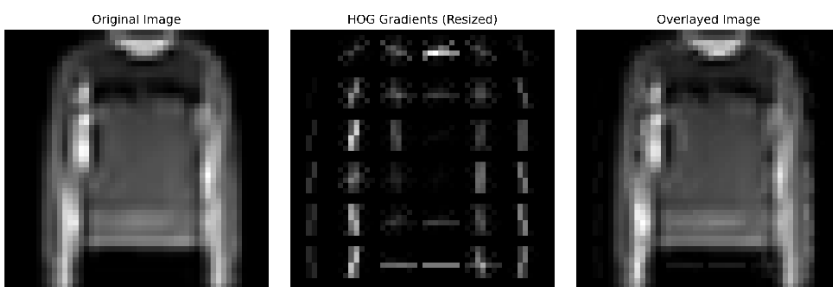
Clasa 0/T-shirt/top



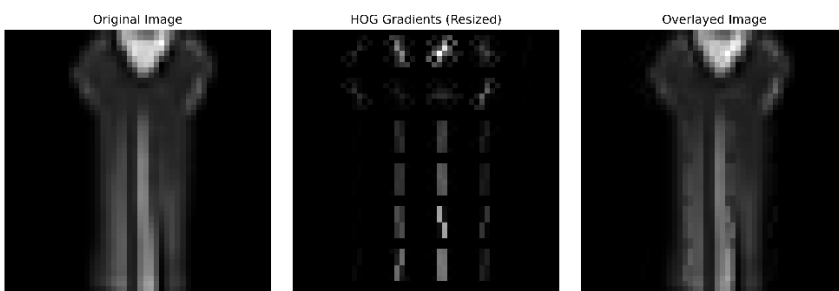
Clasa 1/Trouser



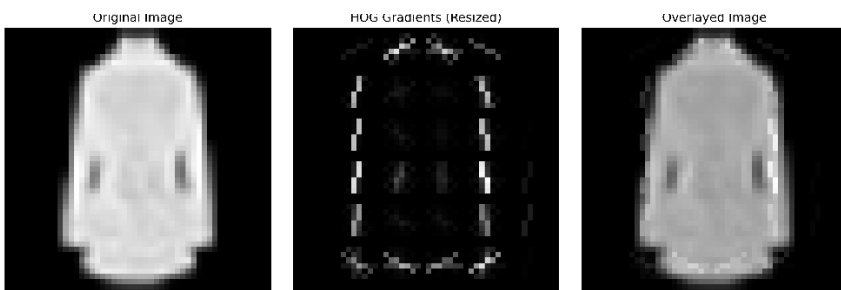
Clasa 2/Pullover



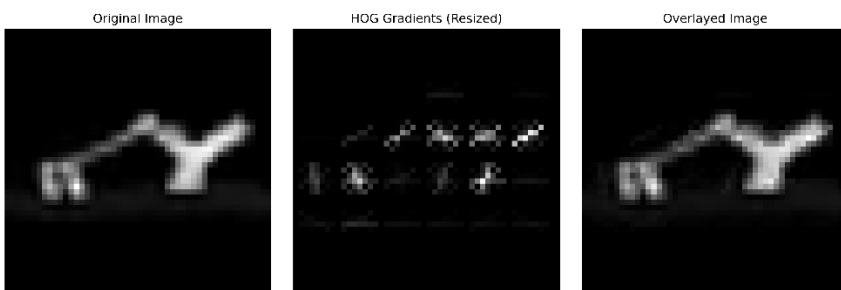
Clasa 3/Dress



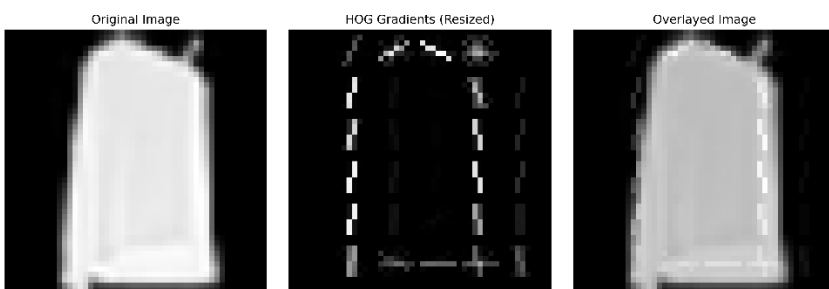
Clasa 4/Coat



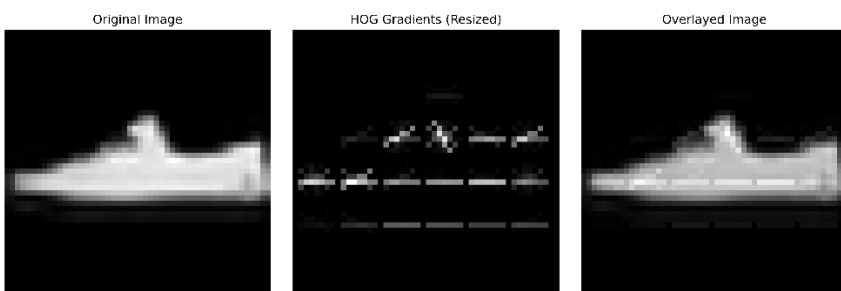
Clasa 5/Sandal



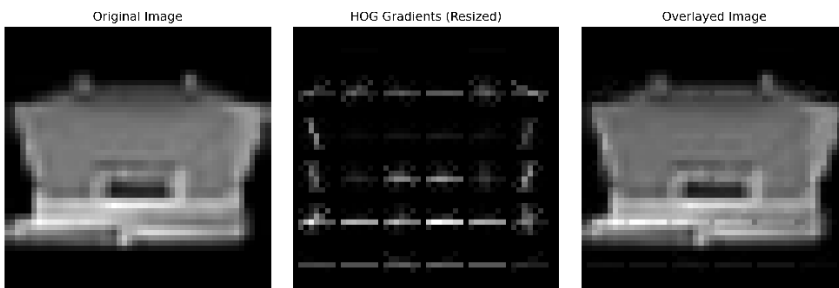
Clasa 6/Shirt



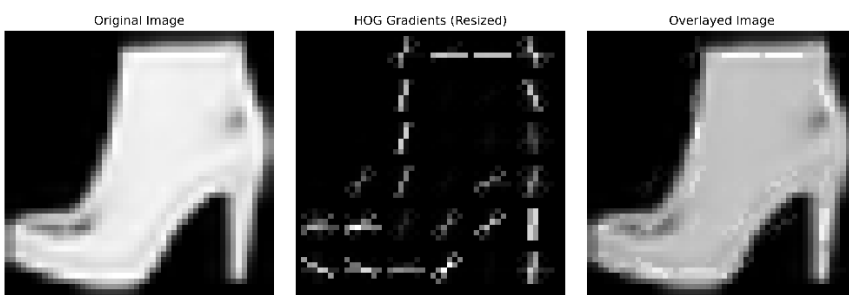
Clasa 7/Sneaker



Clasa 8/Bag



Clasa 9/Ankle boot



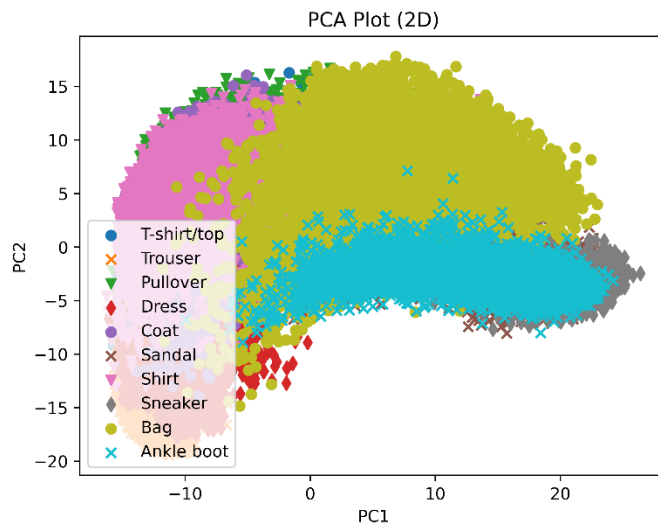
Vizualizare intermediara a aplicarii metodei HOG pe imaginile initiale

Acesta este rezultatul vizualizarii calitative in urma aplicarii HOG pe imagini selectate random din fiecare clasa. Este o modalitate utila de a observa daca cele mai importante margini, contururi si texturi au fost identificate.

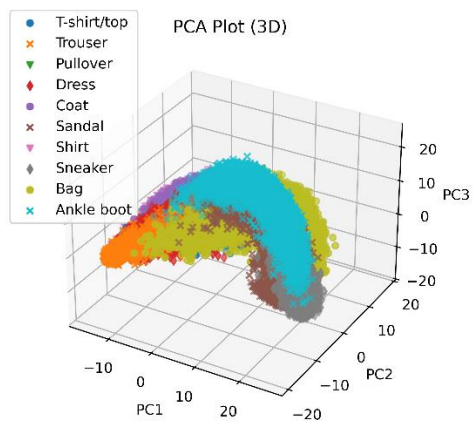
Plot-area componentelor principale

Aceasta metoda de vizualizare calitativa peste hog+pca este utila deoarece putem observa zgomotul si suprapunerea claselor, de asemenea, putem vedea majoritatea informatiei din datele principale intr-un mod mai usor de inteles deoarece avem informatiile comprimate in dimensiuni mai mici.

Pe 2 componente:



Pe 3 componente:



Atat pe vizualizarea bidimensională cât și pe cea tridimensională se observă că clasele sunt bine grupate dar există suprapuneri semnificative între ele. Acest lucru sugerează că unele clase pot fi greu de separate doar cu parametrii dați la PCA și HOG

4.3

Total features before selection: 900

Total features after selection: 675

Setul initial, inaintea aplicarii SelectPercentile, care elimina atributele cu varianta mica, continea toate informatiile calculate in urma aplicarii HOG si a StandardScaler. Dupa aplicarea selectiei, au fost eliminate caracteristicile mai putin relevante, ramanand doar 75% din cele initiale.

4.4

Class 0: T-shirt/top Class 5: Sandal

Class 1: Trouser Class 6: Shirt

Class 2: Pullover Class 7: Sneaker

Class 3: Dress Class 8: Bag

Class 4: Coat Class 9: Ankle boot

Model	Accuracy	Class 0 Precision	Class 0 Recall	Class 0 F1-Score	Class 1 Precision	Class 1 Recall	Class 1 F1-Score	Class 2 Precision	Class 2 Recall	Class 2 F1-Score	Class 3 Precision	Class 3 Recall	Class 3 F1-Score	Class 4 Precision	Class 4 Recall	Class 4 F1-Score	Class 5 Precision	Class 5 Recall	Class 5 F1-Score	Class 6 Precision	Class 6 Recall	Class 6 F1-Score	Class 7 Precision	Class 7 Recall	Class 7 F1-Score	Class 8 Precision	Class 8 Recall	Class 8 F1-Score	Class 9 Precision	Class 9 Recall	Class 9 F1-Score
0 Logistic Regression	0.737000	0.773000	0.740000	0.756000	0.836000	0.899000	0.866000	0.555000	0.586000	0.570000	0.649000	0.670000	0.659000	0.628000	0.575000	0.600000	0.832000	0.840000	0.836000	0.470000	0.446000	0.458000	0.853000	0.832000	0.843000	0.915000	0.868000	0.891000	0.846000	0.914000	0.878000
1 SVM	0.764000	0.730000	0.776000	0.752000	0.927000	0.916000	0.922000	0.606000	0.648000	0.626000	0.680000	0.757000	0.716000	0.662000	0.709000	0.685000	0.848000	0.862000	0.855000	0.569000	0.442000	0.497000	0.857000	0.703000	0.773000	0.968000	0.885000	0.925000	0.811000	0.947000	0.874000
2 Random Forest	0.788000	0.814000	0.753000	0.782000	0.960000	0.922000	0.941000	0.608000	0.670000	0.637000	0.696000	0.864000	0.771000	0.603000	0.654000	0.627000	0.917000	0.907000	0.912000	0.525000	0.340000	0.413000	0.890000	0.875000	0.883000	0.924000	0.960000	0.942000	0.912000	0.938000	0.925000
3 XGBoost	0.773000	0.798000	0.764000	0.781000	0.950000	0.930000	0.940000	0.561000	0.619000	0.589000	0.703000	0.836000	0.764000	0.563000	0.515000	0.538000	0.927000	0.908000	0.917000	0.485000	0.408000	0.443000	0.889000	0.868000	0.879000	0.936000	0.948000	0.942000	0.897000	0.938000	0.917000

```
Best Logistic Regression Parameters: {'C': np.float64(1.8619556021711627), 'multi_class': 'multinomial'}, Best Score: 0.8601166666666668
Best SVM Parameters: {'C': np.float64(3.6946411021099497), 'kernel': 'poly'}, Best Score: 0.9080999999999999
Best Random Forest Parameters: {'n_estimators': 200, 'max_features': 'log2', 'max_depth': None}, Best Score: 0.87985
Best XGBoost Parameters: {'n_estimators': 300, 'max_depth': 9, 'learning_rate': 0.1}, Best Score: 0.8946666666666667
```

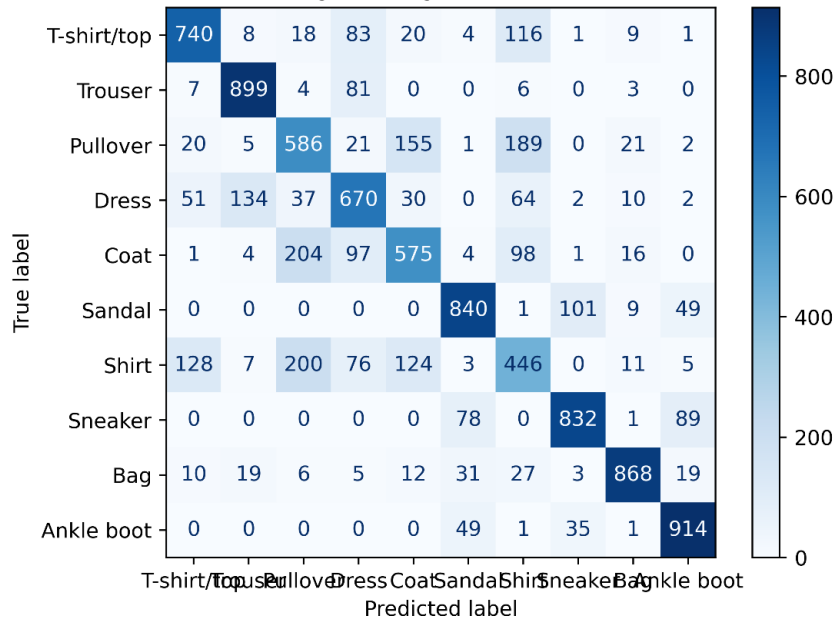
Pentru partea de antrenare, am folosit functia RandomizedSearchCV(cross-validation) pentru a afla cei mai buni parametri, iar apoi am antrenat si prezis pe fiecare model in parte. In urma rularii programului, desi SVM a avut cel mai bun scor, acuratetea generala mai mare a oferit o Random Forest.

In general, clasele au rezultate bune cu acuratete in jur de 90% (este posibil sa se fi realizat overfit pe unele clase).Clasele cu cele mai bune predictii sunt clasele Trouser, Bag si Ankle boot, deoarece au cea mai mare precizie. Un motiv pentru acest lucru este conturul usor de distins al acestora in urma HOG, si faptul ca au componentele bine grupate in analiza componentelor principale.

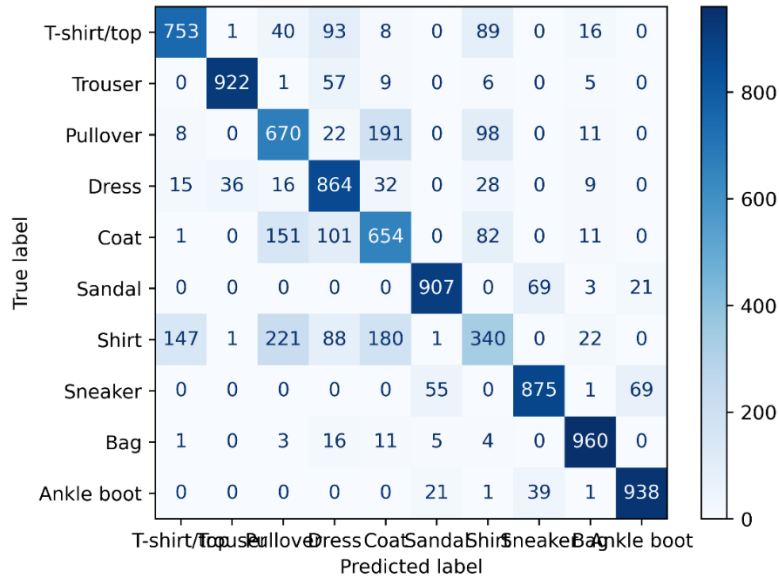
La coada clasamentului se afla clasele Shirt, Coat si Pullover, acestea sunt usor confundabile cand vine vorba de forma lor. Acestea necesita metode mai eficiente pentru o separare eficienta.

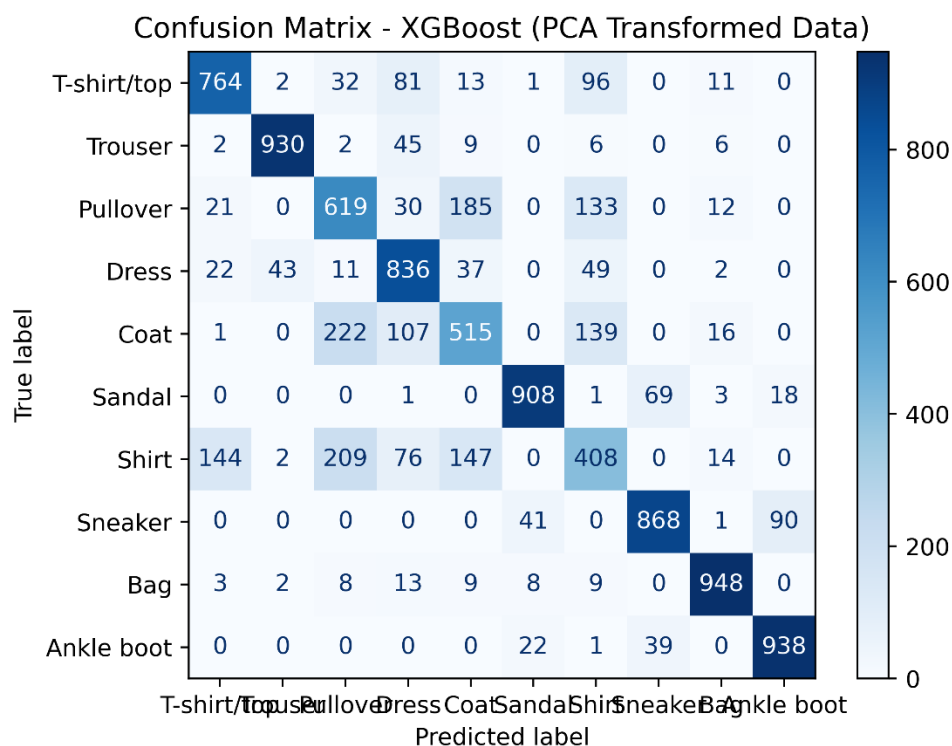
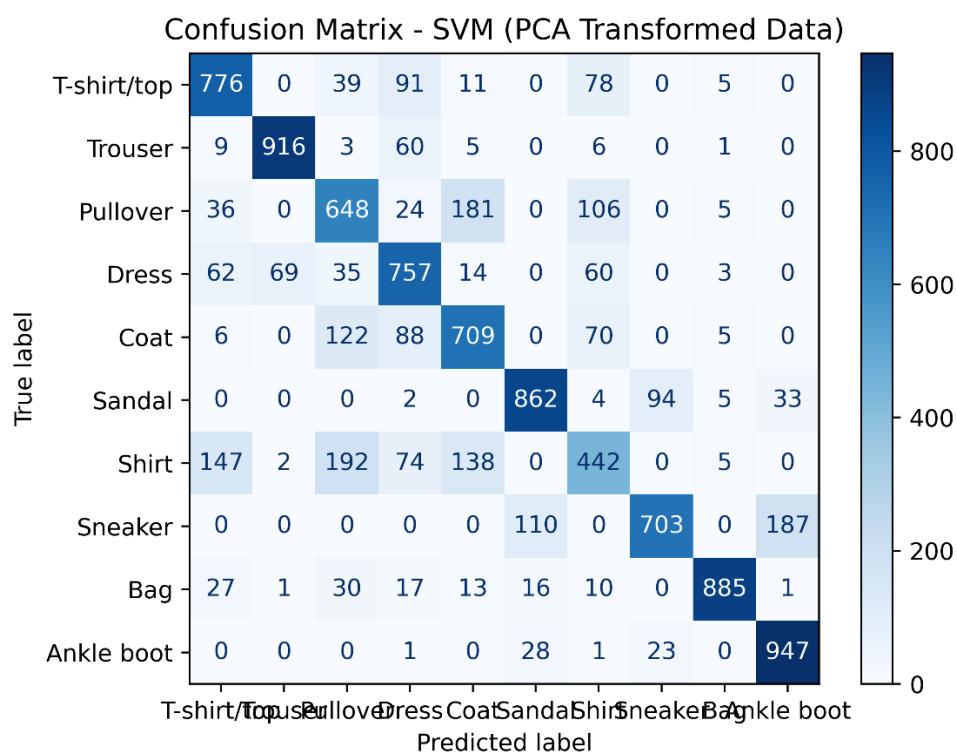
Confuzia intre aceste 3 clase poate fi observata si in matricile de confuzie de mai jos, unde exista valori mai mici pe diagonala.

Confusion Matrix - Logistic Regression (PCA Transformed Data)



Confusion Matrix - Random Forest (PCA Transformed Data)



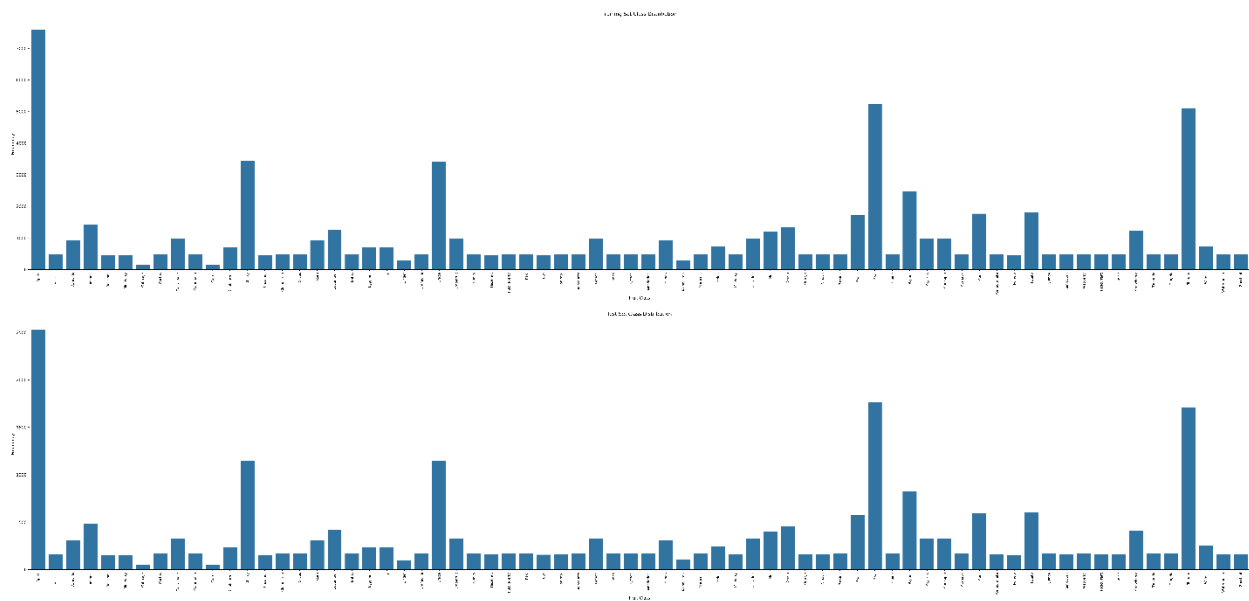


Pe datasetul fruits, am ales să aplic algoritmul ORB datorită vitezei sale ridicate și a acurateței mari în determinarea orientării componentelor imaginii. Acesta este util în captarea detaliilor esențiale ale caracteristicilor imaginii, fiind un descriptor eficient pentru detectarea și descrierea punctelor de interes. În paralel, am folosit HOG (Histogram of Oriented Gradients) pentru că oferă o descriere robustă a formelor și conturilor din imagini, fiind eficient în analiza texturii și a structurii vizuale a obiectelor, ceea ce este deosebit de util pentru recunoașterea fructelor, care adesea au forme complexe și texturi variate.

Pentru a îmbunătăți performanța clasificării, am combinat features extrase de ambele metode, ORB și HOG, folosind hstack. Astfel, am creat un set de caracteristici mai complet și mai informativ, care poate contribui la o performanță mai bună a modelului, prin captarea atât a detaliilor fine (cu ORB), cât și a informațiilor legate de formele și texturile globale (cu HOG).

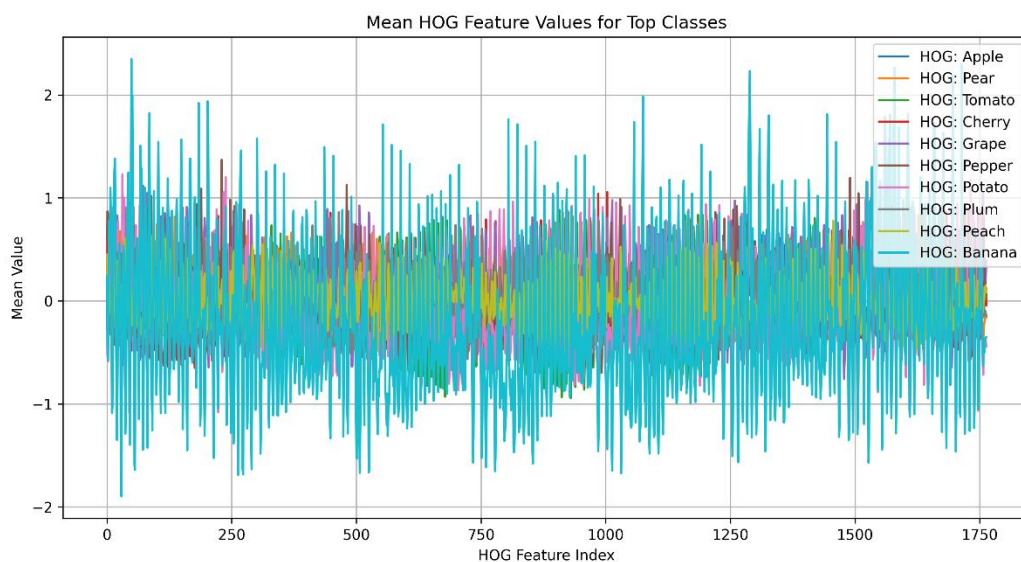
Pe datasetul fruits avem următoarele rezultate:

4.2.1. Analiza echilibrului de clase



Am grupat fructele din aceeași categorie (Apple, Apple 1 și Golden Apple sunt în aceeași categorie Apple) și au rezultat 70 de clase. Se observă faptul că există mai multe imagini la fructele mai comune precum Apple.

4.2.2 Vizualizare cantitativa



Vizualizarea cantitativa pe HOG de mai sus arata o distributie numerica a valorilor descriptorilor HOG pe top 10 clase, adică a intensității gradientelor într-o gamă de direcții specific.

Vizualizare calitativa

Acesta este rezultatul aplicării metodei ORB pe fructele ce se afla in top 10. Se observa ca pe unele imagini orb nu a detectat niciun keypoint. O posibilitate a acestui lucru este faptul ca fie imaginea este prea mica (trebuie redimensionata), fie nu identifica detalii indeajuns de unice.





4.3

Total features before selection: 1796

Total features after selection: 1347

Setul initial, inaintea aplicarii SelectPercentile, care elimina attributele cu varianta mica, continea toate informatiile calculate in urma aplicarii HOG+ORB si a StandardScaler. Dupa aplicarea selectiei, au fost eliminate caracteristicile mai putin relevante, ramanand doar 75% din cele initiale.

Reducerea numărului de caracteristici poate ajuta la prevenirea overfitting-ului a modelului și la îmbunătățirea generalizării, permițându-i să se concentreze pe cele mai relevante trăsături ale datelor.

4.4

La datasetul fruits, nu am mai aplicat cross-validation, ci am realizat un set de validare cu ajutorul functiei train_test_split. Am initializat modelele, apoi am generat combinatii de hiperparametrii predefiniti pentru a-l alege pe cei care dau o acuratete mai mare. Modelul era ulterior antrenat cu acesti hiperparametri. Am reusit sa rulez doar Logistic Regression si Random Forest, celelalte doua modele fiind mult prea de durata. Dintre acestea doua, Logistic Regression are acuratetea cea mai mare, la fel si weighted si max average.

Din cele doua matrice de confuzie, se poate observa ca fructul care a fost prezis cel mai bine este Apple, care are si cel mai mare numar de imagini de antrenare. Se poate observa ca in general clasele sunt distinse usor, neexistand valori foarte mari in afara diagonalei principale.

	0	1
Model	Random Forest_n_estimators=100_max_depth=20_max_features=log2	Logistic Regression_C=0.1_multi_class=multinomial
Accuracy	0.786000	0.851000
Apple Precision	0.497000	0.802000
Apple Recall	0.963000	0.895000
Apple F1-Score	0.656000	0.846000
Apricot Precision	0.987000	0.828000
Apricot Recall	0.463000	0.793000
Apricot F1-Score	0.631000	0.810000
Avocado Precision	0.888000	0.902000
Avocado Recall	0.793000	0.864000
Avocado F1-Score	0.838000	0.883000
Banana Precision	0.946000	0.973000
Banana Recall	0.682000	0.740000
Banana F1-Score	0.792000	0.840000
Beetroot Precision	0.779000	0.875000
Beetroot Recall	0.353000	0.420000
Beetroot F1-Score	0.486000	0.568000
Blueberry Precision	1.000000	0.974000
Blueberry Recall	0.675000	0.734000
Blueberry F1-Score	0.806000	0.837000
Cabbage Precision	0.979000	0.887000
Cabbage Recall	1.000000	1.000000
Cabbage F1-Score	0.989000	0.940000
Cactus Precision	0.879000	0.844000
Cactus Recall	0.789000	0.880000
Cactus F1-Score	0.832000	0.861000
Cantaloupe Precision	0.810000	0.804000
Cantaloupe Recall	0.991000	0.991000
Cantaloupe F1-Score	0.892000	0.888000
Carambola Precision	0.989000	0.975000
Carambola Recall	0.524000	0.711000
Carambola F1-Score	0.685000	0.822000
Carrot Precision	1.000000	1.000000
Carrot Recall	1.000000	1.000000
Carrot F1-Score	1.000000	1.000000
Cauliflower Precision	0.969000	1.000000
Cauliflower Recall	0.940000	1.000000
Cauliflower F1-Score	0.954000	1.000000
Cherry Precision	0.951000	0.981000
Cherry Recall	0.905000	0.937000
Cherry F1-Score	0.927000	0.959000

Cauliflower F1-Score	0.954000	1.000000
Cherry Precision	0.951000	0.981000
Cherry Recall	0.905000	0.937000
Cherry F1-Score	0.927000	0.959000
Chestnut Precision	1.000000	0.828000
Chestnut Recall	0.516000	0.627000
Chestnut F1-Score	0.681000	0.714000
Clementine Precision	0.991000	0.948000
Clementine Recall	0.687000	0.886000
Clementine F1-Score	0.811000	0.916000
Cocos Precision	0.813000	0.860000
Cocos Recall	0.813000	0.928000
Cocos F1-Score	0.813000	0.893000
Corn Precision	0.730000	0.641000
Corn Recall	0.480000	0.487000
Corn F1-Score	0.579000	0.553000
Cucumber Precision	0.865000	0.902000
Cucumber Recall	0.767000	0.859000
Cucumber F1-Score	0.813000	0.880000
Dates Precision	0.949000	0.765000
Dates Recall	0.789000	0.765000
Dates F1-Score	0.862000	0.765000
Eggplant Precision	0.888000	0.911000
Eggplant Recall	0.708000	0.742000
Eggplant F1-Score	0.788000	0.818000
Fig Precision	1.000000	0.880000
Fig Recall	0.611000	0.752000
Fig F1-Score	0.759000	0.811000
Ginger Precision	0.969000	0.894000
Ginger Recall	0.636000	0.768000
Ginger F1-Score	0.768000	0.826000
Granadilla Precision	1.000000	1.000000
Granadilla Recall	1.000000	0.982000
Granadilla F1-Score	1.000000	0.991000
Grape Precision	0.790000	0.890000
Grape Recall	0.951000	0.991000
Grape F1-Score	0.863000	0.938000
Grapefruit Precision	0.994000	0.991000
Grapefruit Recall	0.952000	0.961000
Grapefruit F1-Score	0.972000	0.975000
Guava Precision	1.000000	0.982000

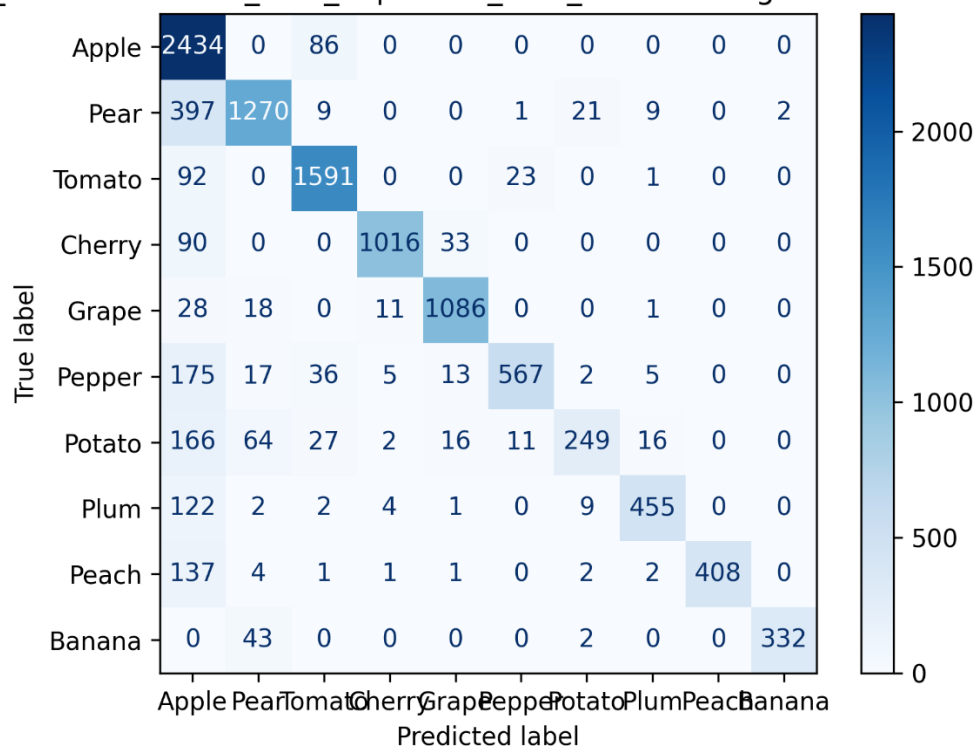
Eggplant Precision	0.888000	0.911000
Eggplant Recall	0.708000	0.742000
Eggplant F1-Score	0.788000	0.818000
Fig Precision	1.000000	0.880000
Fig Recall	0.611000	0.752000
Fig F1-Score	0.759000	0.811000
Ginger Precision	0.969000	0.894000
Ginger Recall	0.636000	0.768000
Ginger F1-Score	0.768000	0.826000
Granadilla Precision	1.000000	1.000000
Granadilla Recall	1.000000	0.982000
Granadilla F1-Score	1.000000	0.991000
Grape Precision	0.790000	0.890000
Grape Recall	0.951000	0.991000
Grape F1-Score	0.863000	0.938000
Grapefruit Precision	0.994000	0.991000
Grapefruit Recall	0.952000	0.961000
Grapefruit F1-Score	0.972000	0.975000
Guava Precision	1.000000	0.982000
Guava Recall	0.663000	0.976000
Guava F1-Score	0.797000	0.979000
Hazelnut Precision	0.994000	0.952000
Hazelnut Recall	1.000000	1.000000
Hazelnut F1-Score	0.997000	0.975000
Huckleberry Precision	0.982000	0.943000
Huckleberry Recall	0.657000	0.801000
Huckleberry F1-Score	0.787000	0.866000
Kaki Precision	1.000000	1.000000
Kaki Recall	0.783000	0.747000
Kaki F1-Score	0.878000	0.855000
Kiwi Precision	0.897000	0.927000
Kiwi Recall	0.667000	0.974000
Kiwi F1-Score	0.765000	0.950000
Kohlrabi Precision	1.000000	0.947000
Kohlrabi Recall	0.580000	0.573000
Kohlrabi F1-Score	0.734000	0.714000
Kumquats Precision	0.940000	0.947000
Kumquats Recall	0.952000	0.976000
Kumquats F1-Score	0.946000	0.961000
Lemon Precision	0.914000	0.805000
Lemon Recall	0.933000	1.000000
Lemon F1-Score	0.924000	0.892000
Limes Precision	1.000000	0.944000

Limes Precision	1.000000	0.944000
Limes Recall	0.614000	0.813000
Limes F1-Score	0.761000	0.874000
Lychee Precision	0.941000	0.960000
Lychee Recall	0.861000	1.000000
Lychee F1-Score	0.899000	0.979000
Mandarine Precision	0.981000	0.985000
Mandarine Recall	0.633000	0.783000
Mandarine F1-Score	0.769000	0.872000
Mango Precision	0.958000	0.759000
Mango Recall	0.594000	0.727000
Mango F1-Score	0.733000	0.743000
Mangostan Precision	1.000000	1.000000
Mangostan Recall	0.627000	0.912000
Mangostan F1-Score	0.771000	0.954000
Maracuja Precision	0.713000	0.781000
Maracuja Recall	0.614000	0.711000
Maracuja F1-Score	0.660000	0.744000
Melon Precision	0.882000	0.838000
Melon Recall	0.854000	0.992000
Melon F1-Score	0.868000	0.909000
Mulberry Precision	0.924000	0.970000
Mulberry Recall	0.970000	0.994000
Mulberry F1-Score	0.946000	0.982000
Nectarine Precision	0.951000	0.465000
Nectarine Recall	0.299000	0.364000
Nectarine F1-Score	0.455000	0.408000
Nut Precision	0.995000	0.821000
Nut Recall	0.505000	0.720000
Nut F1-Score	0.670000	0.767000
Onion Precision	0.709000	0.603000
Onion Recall	0.643000	0.783000
Onion F1-Score	0.674000	0.681000
Orange Precision	1.000000	0.894000
Orange Recall	1.000000	1.000000
Orange F1-Score	1.000000	0.944000
Papaya Precision	0.736000	0.792000
Papaya Recall	0.561000	0.695000
Papaya F1-Score	0.637000	0.740000
Passion Precision	0.961000	0.986000
Passion Recall	0.596000	0.861000
Passion F1-Score	0.736000	0.920000

Papaya Recall	0.561000	0.695000
Papaya F1-Score	0.637000	0.740000
Passion Precision	0.961000	0.986000
Passion Recall	0.596000	0.861000
Passion F1-Score	0.736000	0.920000
Peach Precision	0.964000	0.818000
Peach Recall	0.693000	0.735000
Peach F1-Score	0.806000	0.774000
Pear Precision	0.648000	0.693000
Pear Recall	0.714000	0.837000
Pear F1-Score	0.680000	0.758000
Pepino Precision	0.953000	0.904000
Pepino Recall	0.614000	0.620000
Pepino F1-Score	0.747000	0.736000
Pepper Precision	0.925000	0.823000
Pepper Recall	0.714000	0.783000
Pepper F1-Score	0.806000	0.803000
Physalis Precision	1.000000	0.962000
Physalis Recall	0.930000	0.991000
Physalis F1-Score	0.964000	0.976000
Pineapple Precision	0.875000	0.945000
Pineapple Recall	0.982000	0.942000
Pineapple F1-Score	0.926000	0.944000
Pitahaya Precision	1.000000	0.988000
Pitahaya Recall	0.693000	0.982000
Pitahaya F1-Score	0.819000	0.985000
Plum Precision	0.828000	0.747000
Plum Recall	0.727000	0.846000
Plum F1-Score	0.774000	0.793000
Pomegranate Precision	0.969000	0.958000
Pomegranate Recall	0.384000	0.689000
Pomegranate F1-Score	0.550000	0.801000
Pomelo Precision	0.983000	0.929000
Pomelo Recall	0.739000	0.941000
Pomelo F1-Score	0.843000	0.935000
Potato Precision	0.449000	0.744000
Potato Recall	0.428000	0.581000
Potato F1-Score	0.438000	0.652000
Quince Precision	1.000000	1.000000
Quince Recall	0.681000	0.729000
Quince F1-Score	0.810000	0.843000
Rambutan Precision	0.994000	1.000000
Rambutan Recall	0.888000	0.888000
Rambutan F1-Score	0.941000	0.944000

Potato F1-Score	0.438000	0.652000
Quince Precision	1.000000	1.000000
Quince Recall	0.681000	0.729000
Quince F1-Score	0.810000	0.843000
Rambutan Precision	0.994000	1.000000
Rambutan Recall	0.970000	0.982000
Rambutan F1-Score	0.981000	0.991000
Raspberry Precision	0.835000	0.947000
Raspberry Recall	0.699000	0.747000
Raspberry F1-Score	0.761000	0.835000
Redcurrant Precision	1.000000	1.000000
Redcurrant Recall	1.000000	1.000000
Redcurrant F1-Score	1.000000	1.000000
Salak Precision	1.000000	0.920000
Salak Recall	0.994000	1.000000
Salak F1-Score	0.997000	0.959000
Strawberry Precision	0.852000	0.876000
Strawberry Recall	0.761000	0.934000
Strawberry F1-Score	0.804000	0.904000
Tamarillo Precision	0.982000	0.988000
Tamarillo Recall	0.970000	1.000000
Tamarillo F1-Score	0.976000	0.994000
Tangelo Precision	1.000000	0.878000
Tangelo Recall	0.717000	1.000000
Tangelo F1-Score	0.835000	0.935000
Tomato Precision	0.876000	0.905000
Tomato Recall	0.941000	0.920000
Tomato F1-Score	0.908000	0.913000
Walnut Precision	0.827000	0.915000
Walnut Recall	1.000000	1.000000
Walnut F1-Score	0.905000	0.956000
Watermelon Precision	1.000000	0.903000
Watermelon Recall	0.675000	0.650000
Watermelon F1-Score	0.806000	0.756000
Zucchini Precision	0.847000	0.856000
Zucchini Recall	1.000000	1.000000
Zucchini F1-Score	0.917000	0.922000
Macro Avg Precision	0.914000	0.891000
Macro Avg Recall	0.752000	0.843000
Macro Avg F1-Score	0.809000	0.860000
Weighted Avg Precision	0.835000	0.856000
Weighted Avg Recall	0.786000	0.851000
Weighted Avg F1-Score	0.788000	0.849000

n_estimators=100_max_depth=20_max_features=log2 Confusion Matrix (HC)



ic Regression_C=0.1_multi_class=multinomial Confusion Matrix (HOG+ORB)

