

Etapa 1

Extragerea de attribute (Feature Extraction)

Am implementat un flux de extragere a atributelor folosind două metode principale: descriptori ORB și reducerea dimensionalității cu PCA.

1. Preprocesarea imaginilor

Am redimensionat imaginile la o dimensiune standard de **32x32 pixeli** pentru a uniformiza setul de date, ceea ce reduce complexitatea calculului și asigură consistență între imagini.

2. Extragerea descriptorilor ORB

Am utilizat ORB pentru a extrage caracteristici semnificative din imagini, cum ar fi colțuri și modele locale.

- **Parametrii aleși:**
 - `max_keypoints=300`: permite detectarea a până la 300 de puncte de interes pentru fiecare imagine.
 - `edge_threshold=6` pentru controlul sensibilității la marginile imaginii
- **Justificare:** ORB este rapid și eficient pentru imagini în tonuri de gri, ceea ce îl face potrivit pentru seturi mari de date. Este rezistent la modificări de scalare sau rotație.

3. Clustering cu KMeans (Bag of Words)

Descriptorii extrași au fost grupați în **500 cluster** utilizând KMeans, pentru a construi un vocabular vizual.

- Descriptorii din toate imaginile au fost concatenați într-un set comun, iar KMeans a fost aplicat pentru a învăța centrele clusterelor.
- Fiecare imagine a fost reprezentată ca o histogramă bazată pe frecvența clusterelor corespunzătoare descriptorilor săi.
- **Justificare:** Bag of Words oferă o reprezentare compactă și numerică a imaginilor, reducând complexitatea analizei ulterioare.

4. Reducerea dimensionalității cu PCA

Am aplicat PCA pentru a reduce dimensiunea histogramelor la **20 de componente principale**.

- Am standardizat datele înainte de PCA pentru a avea o medie de 0 și o varianță unitară.
- **Justificare:** PCA reduce redundanța dintre caracteristici, îmbunătățind eficiența modelelor și explicând ~95% din variabilitatea datelor cu mai puține atribute.

5. Selecția atributelor

Am utilizat metoda **VarianceThreshold** pentru a elimina caracteristicile cu varianță scăzută (sub 0.01). Aceasta păstrează doar caracteristicile relevante pentru clasificare.

- Selecția atributelor ajută la reducerea zgomotului și a timpului de procesare, îmbunătățind performanța modelului final.

Rezultate și observații

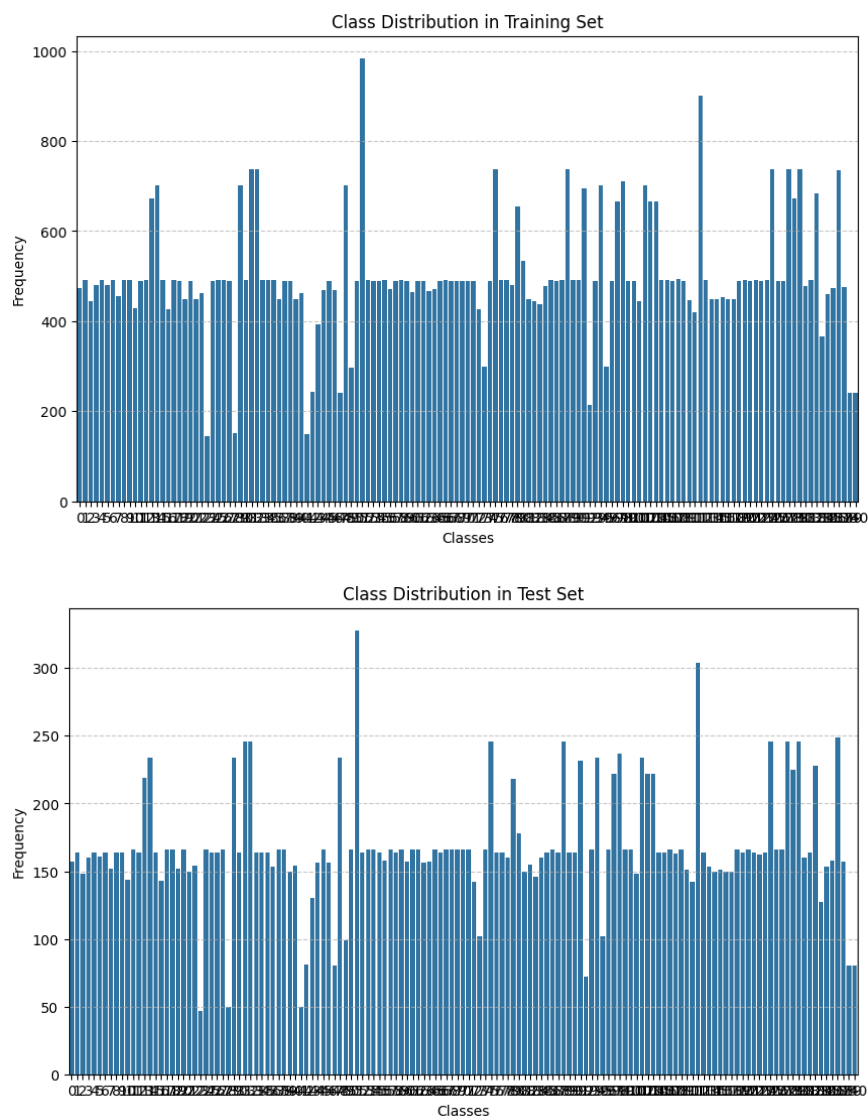
- **Dimensiuni finale:** După fluxul de preprocesare, dimensiunea setului de antrenament a fost redusă semnificativ la $(n_samples, 10)$, unde 10 este numărul de caracteristici finale.
- Metodele alese oferă o bună echilibrare între precizie și eficiență, reprezentând bine imaginile cu un model numeric compact și optimizat.

Acest flux poate fi extins sau ajustat în funcție de performanțele modelelor de clasificare ulterioare.

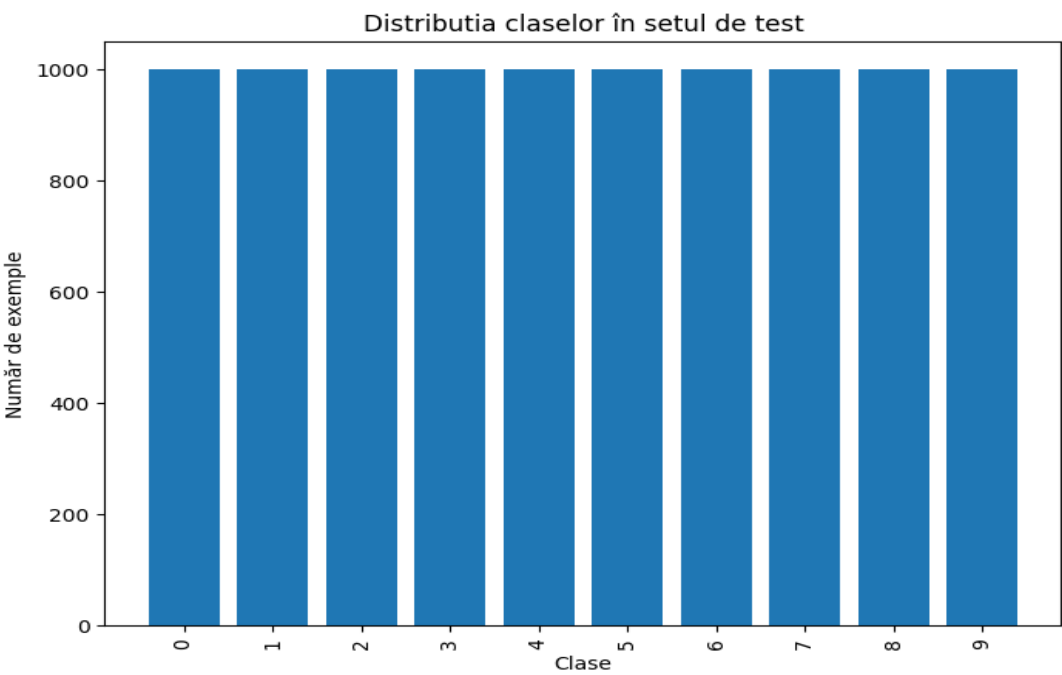
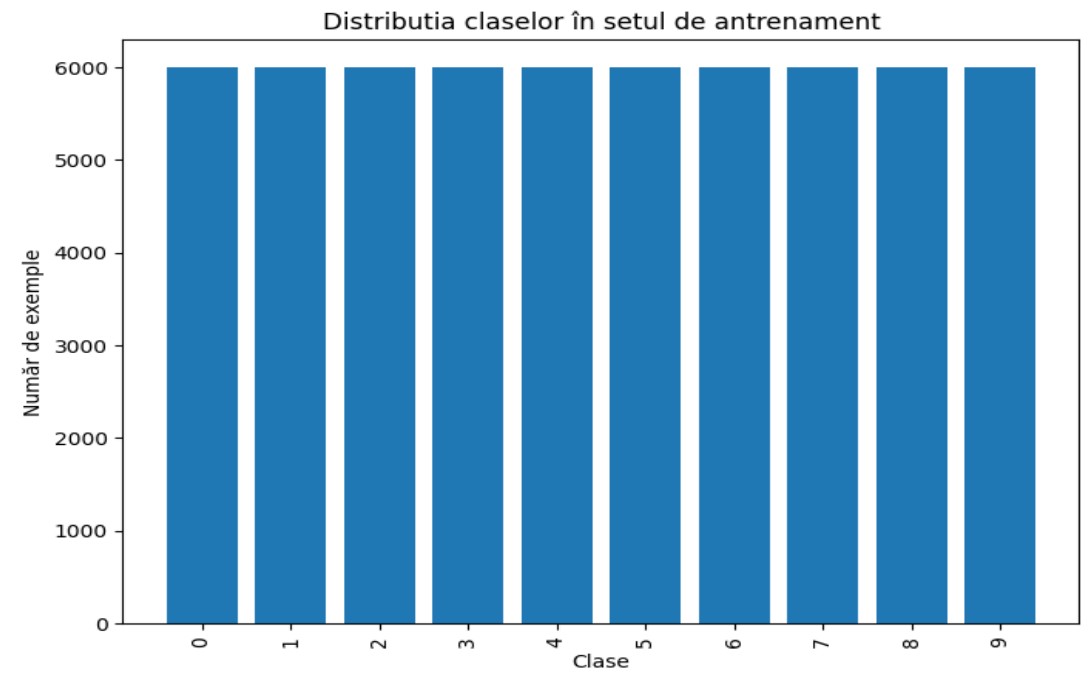
Vizualizarea atributelor extrase

Distributia claselor

Fruits

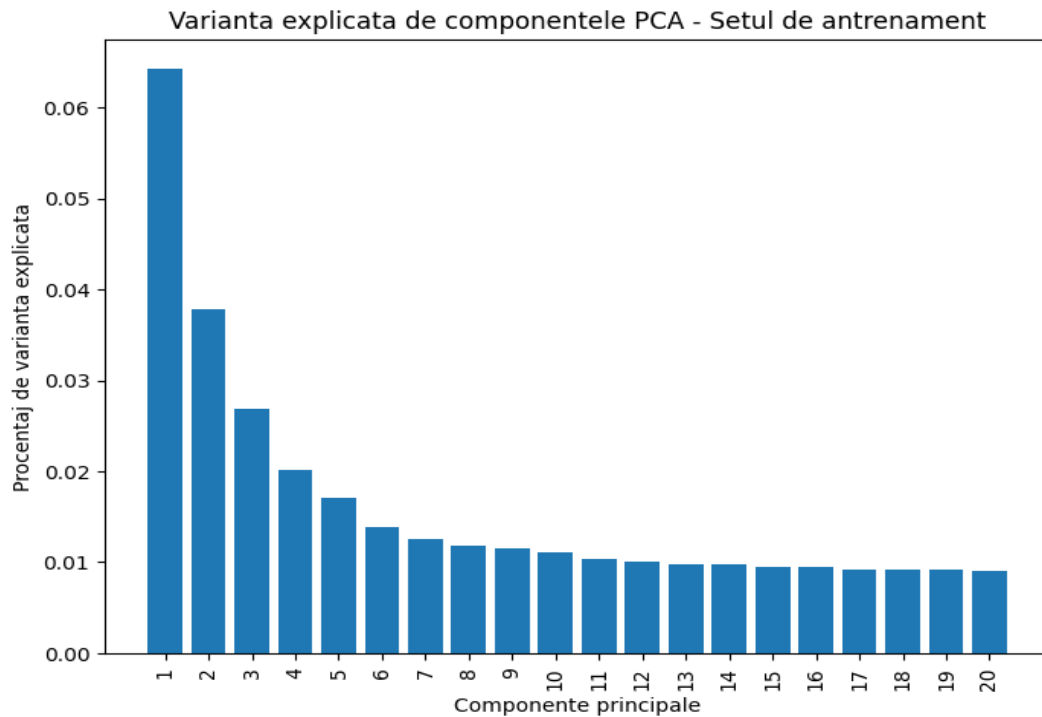


Fashion

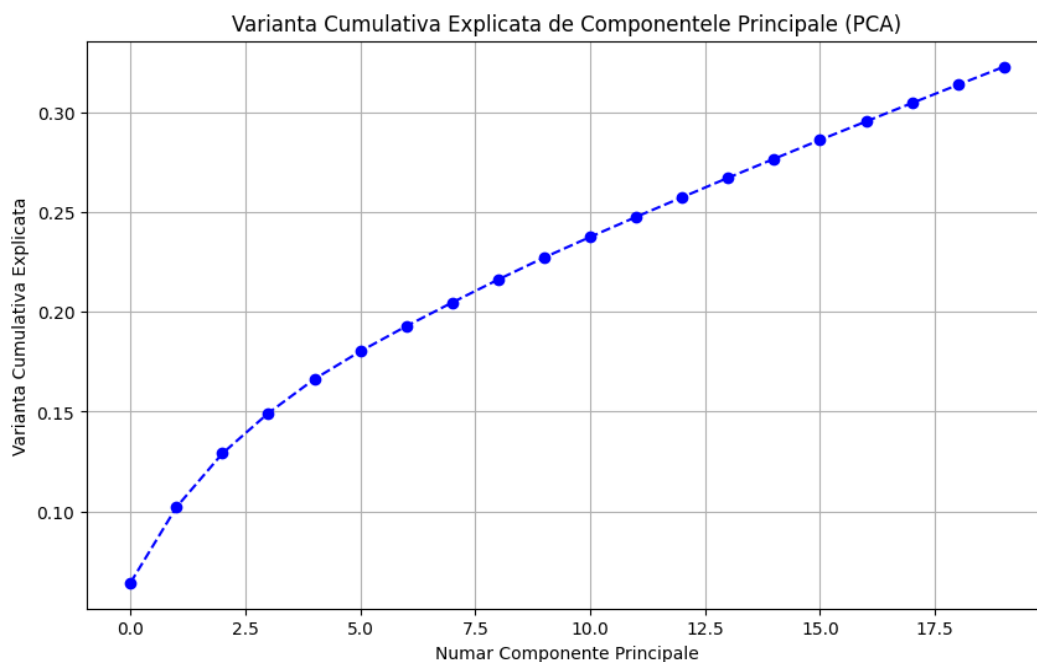


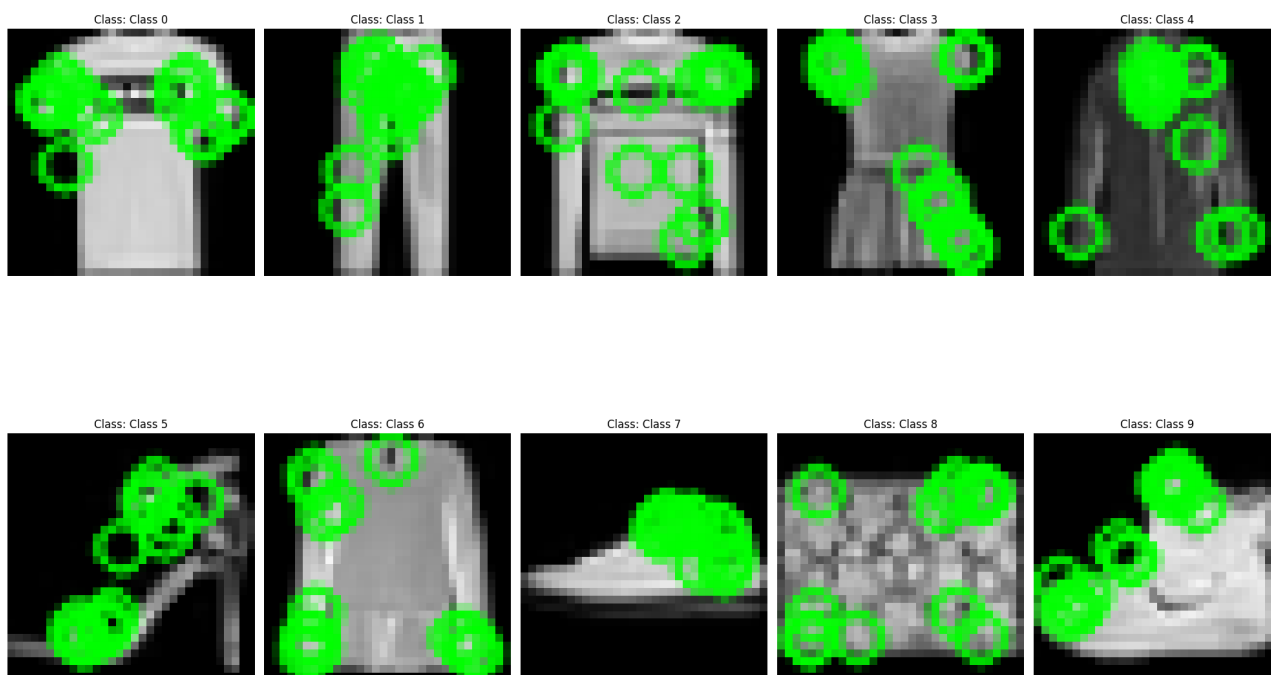
VIZUALIZARE CANTITATIVA SI CALITATIVA

FASHION



Observăm că prima componentă principală explică o mare parte din varianța datelor, urmată de componentele a doua și a treia care încă contribuie semnificativ, dar ponderea lor scade rapid. Restul componentelor explică doar o fracțiune foarte mică din varianța totală.





FRUITS

