



Rockery : Color Clustering Module

Research

Key Concepts

- **Color Clustering** \Rightarrow a technique used to group similar colors in an image.
- **K-means** \Rightarrow an unsupervised ML algorithm
- image loading and processing libraries
 - PIL** - for opening & processing images
 - OpenCV** - processing images
 - matplotlib** - for visualizing color clusters
- Clustering algorithms
 - K-means
 - **DBSCAN** \Rightarrow density based clustering
- **Feature Extraction** \Rightarrow Extracting the most important colors from the image (eg. Pixel Sampling or histograms)

K Means

- alg iteratively divides data points into k clusters by minimizing the variance in each cluster
- how it works:
 1. each point randomly assigned to a k cluster
 2. Compute the centroid (functionally the center) for each cluster
 3. reassign each data point to the cluster w/ the closest center
 4. Repeat until cluster assignments stop changing
- need to select k number of clusters to group the data into
- Elbow method lets us graph the inertia & visualize the point it decreases linearly
 - point referred to as **elbow** \Rightarrow good estimate for the best k value

Image Color Extraction Using K means Clustering

- two types of graphs: raster & vector

- **vector graphs** \Rightarrow based on geometry shapes

- Zoom in, shape is always clear

- drawn by clear geometric formulas

- **Raster graphs** \Rightarrow array like data mix of pixels, smallest unit of a dot containing color information

- number of pixels \Rightarrow resolution

- Color Spaces

- Every pixel in raster is given a color

- **image = data set of colors**

- define the colors w/ numbers

- color map \Rightarrow one dimensional color information

- palette & color swatches \Rightarrow 2D

- default for ML is RGB

- deconstruct the photo into pixels w/ RGB data points

- K means Clustering

- Similarities recognized as Euclidean Distance

- dynamic process to find the centroid

- other points distance from the centroid

- non deterministic alg \Rightarrow outputs can be adjusted

- can stop when we find K clusters consistent w/ intuition

- use the Elbow test or Silhouette to find the optimal K value

- Elbow not always helpful

- when the K value is larger than 8 \Rightarrow not always visually apparent

Kick off

Project: functional Color analyzer written in python

- Published to a custom domain
 - need to wrap it in an API \Rightarrow probably fast API
 - where would it be hosted? github actions?
 - Easiest would be to test locally

Getting Started

- if over a certain % of the image is white Remove the background?
- for the inertia \Rightarrow list the graph points and find the one where it dips

Day 1 Summary

- Conducted trial runs of finding k means
- Cluster quality
 - Elbow method \Rightarrow identifies the "elbow" point where the SSE starts to level off
 - Silhouette Score \Rightarrow how well separated the clusters are
 - Davies-Bouldin \Rightarrow Balances Cluster Compactness & Separation
 - Calinski-Harabasz \Rightarrow Emphasize between cluster dispersion
- Next Steps
 - find a way to average out the k values
 - try on a more colorful image
 - maybe all 4 graphs w/ the different k values

Day 3

- Created a function that saves the scatter plots
- CH is consistently overkill
- made either right on the money or slightly under
- run again comparing median, mean, mode
- which one to go with
 - mode is still mostly useless \Rightarrow Removing
 - if they agree on the # of clusters go w/ that
 - find the Correlation Coefficient
 - if it is extremely close to a line (the abs is close to 1)
 - lowest number

- intermediately close.
- middle number
- all over the place
- highest number

- low, med, high choice	if 2 # are the same (low or high)
- low: $0.99 - 0.97$	} lower } higher
- med: $0.96 - 0.6$	
- low: $0.6 - 0$	

- all the same # \Rightarrow go w/ that number

2nd Try Alg

- linear Corr pretty much useless
- Also Ch method pretty much useless \Rightarrow overkill
- Solid Color Background Skewing Results
- maybe look at main colors and how distinctly different they are?
 - distance formula

~~Not Gonna
Work~~

3rd Shot

- difference from the mean is typically between 15 and 53 when picked

- distance ranges 110 - 200

- if no # are in range pick the one closest

- if more than one w/i range \Rightarrow pick the lowest

} possible Alg

Fast API

- validation for all data types
- Security already integrated → if we need it

End points

Render Probably the best idea to host this

/home

- Erases any images that
- back to the home page
- GET

/upload

- POST
- upload the image

/results/{image_id}

- GET
- endpoint to retrieve processed image
- all graphs

/Results/{image_id}/Pie

- GET
- original image & Pie chart

/Results/{image_id}/Polar

- GET
- original image & Circular polar graph

/Results/{image_id}/bar

- GET
- original image & 3D bar graph