

---

# STA 160 Final Project

---

## Group Members

Laura Valenzuela (919202591)  
Yamini Sharma (915286207)  
Thomas Venner (916857208)

## Abstract

The Shot Marilyns were created by Andy Warhol. There are five paintings. All of them have a different colored background and between the five of them, the colors of Marilyn's hair, collar, skin, and earring vary slightly from one to another. The purpose of the project is to extract the major colors from each picture and compare them to each other, using the fact that images are a collection of pixels which hold positional and color information in a 3-dimensional matrix. Color information can be codified in many different ways, and we will make use of the RGB and HSV color systems. Color dependence can be explored by calculating the relative entropy that displays the shared entropy of each color from each painting. On the painting level, we can compare color distributions of pixels in both the RGB and HSV color systems to determine subtle and nuanced differences between the paintings, in an effort to explain the difference in their elicited effects on the viewers. We can also extract the main colors using the k-means clustering algorithm to determine the number and type of the main colors used in each painting. Finally, we can focus on each ROI (region of interest) and use RGB and HSV representations to reveal subtle differences at the ROI level.

## 1 Introduction

The Shot Marilyns is a series of silkscreen paintings produced in 1964 by Andy Warhol. Each painting is a portrait of Marilyn Monroe. After her death in 1962, Warhol, who had a fascination with Hollywood and fame, began to use her in his work. He created the portraits of Monroe based on a publicity photo for her 1953 film Niagara. He painted five Marilyn portraits with different colored backgrounds: red, orange, blue, sage blue, and turquoise[3].

The story behind why they are called Shot Marilyns goes back to 1964. Warhol was in his Factory when two visitors arrived. Billy Name and a photographer named Dorothy Podber. When they looked at the paintings Podber asked if she could "shoot" the paintings. Assuming she meant photographing his paintings, Warhol accepted. Then, she revealed a revolver from her purse, and pulled the trigger, striking Marilyn right between the eyes. Four out of five paintings were wounded, known today as the Shot Marilyns[4]. Through EDA, entropy, color extraction and regions of interest, a person can appreciate how different each painting really is even though they are made from the same photograph.

## 2 The Shot Marilyns

The pictures of Warhol's paintings were obtained from 'The Interior Review' website[2]. They are 500 x 500 pixels. The photos were uploaded from the website to the project, where they were converted to a png file so it would add transparency. They were uploaded in an RGBA format, which stands for Red, Green, Blue, and Alpha: transparency. And then converted them into arrays for easier manipulation. In the arrays, all of the colors can range from 0 to 255. 0 meaning there is no trace of

that color and 255 meaning it mostly has that color. Most pictures in the data have a combination of colors, which means it is not likely for one color to be 255 or 0 but something in between; except for transparency which is almost always 255.

Figure 1: The Shot Marylins



### 3 Color Analysis

#### 3.1 RGB: Red, Green, Blue

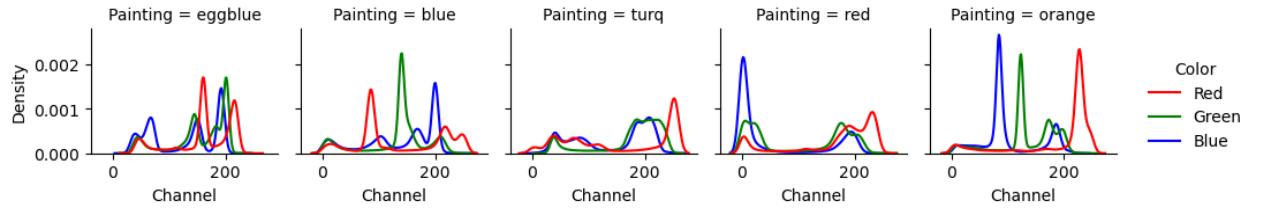


Figure 2: Painting RGB Representation

RGB is based on the theory that all visible colors can be made using the additive primary colors: red, green, and blue. As you adjust the amount of color used you can create different colors from the first three[5]. The distribution of the values of the color coordinates differ between the five paintings in Fig. 2. Since the background of each painting has a different color and it encompasses a lot of the painting, it should not be expected to see a lot of resemblance between the graphs. Nevertheless, there are some comparisons that can be made. In the orange and blue paintings, green has a lot of density in the range of [100, 150]. And the distribution of blue and red is almost inverted between the paintings because blue is more prominent in the blue painting and red is more important in the orange painting. Something odd is that the color red has more density in the orange painting than in the red one, but the color red has more distribution along the red painting.

In the graph for the red painting, we see a spike of blue hovering around 0 due to the fact that is barely used in that picture. And the instances where blue is used is usually mixed with the other colors and in low amounts. Between the five graphs, the color red is always present in the range [200,

255]. Meaning red is always used in the painting even if its density is low that is because the lipstick is red in all of the paintings and the yellow hair uses red to create the color.

### 3.2 HSV: Hue Saturation Value

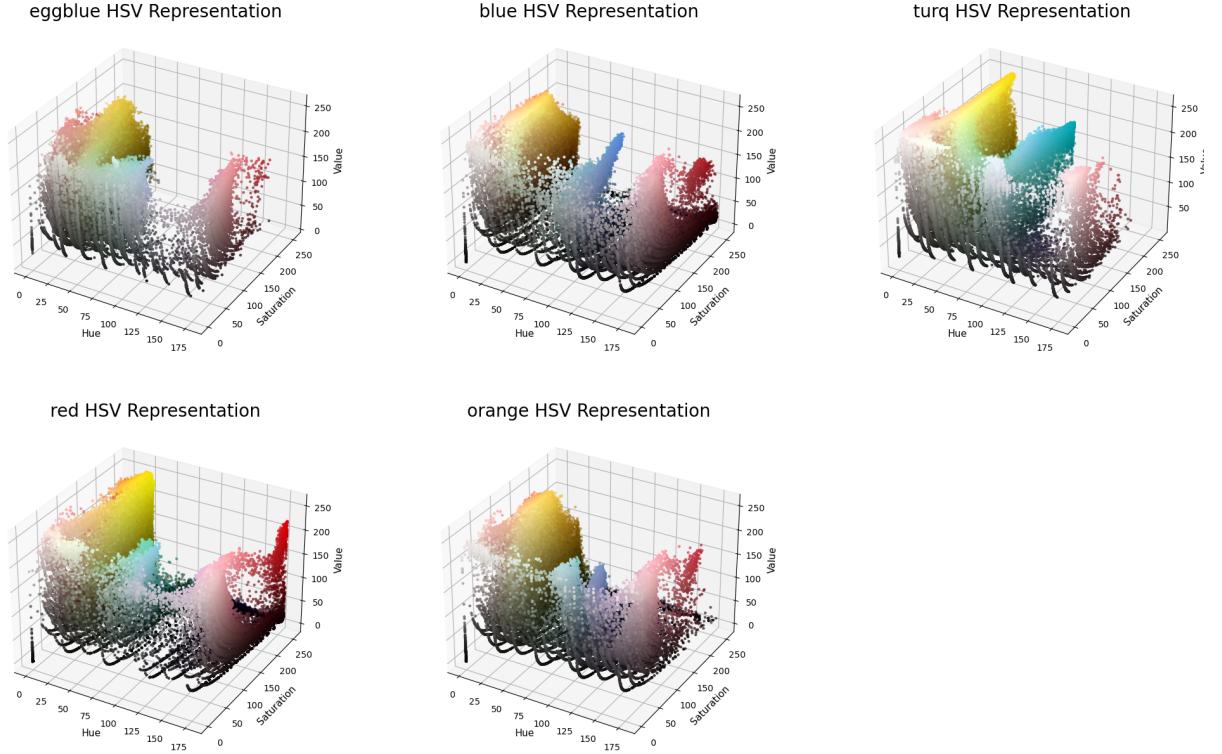


Figure 3: HSV Representation

The HSV scale is a numerical readout of the images that correspond to the colors contained inside[5]. Here, hue is measured from 0 to 180, saturation goes from 0 to 255 , and value is measured from 0 to 255. Hue is essentially the type of color, saturation is how much white is added to dilute the color (255 is maximum saturation), and value is how bright the color is (255 is maximum brightness). The HSV plots in Fig. 3 are going to help with obtaining the colors from each part of the painting and separating the different parts of the face and background. For example, looking at the turquoise HSV plot, one see a spike of turquoise in the middle of the graph. So, a very strong assumption can be made by saying that the middle part of the plot is the background. And the same thing happens with the blue HSV.

Another observation can be made by looking at the right end of the plots. All of them have some red in them. Meaning those are the lips since it is known that the red lips are present in all of the five paintings. Looking at the left end we see the color yellow and that is the hair of each painting. By knowing all of these observations, the separation of the lips, the hair, the background, etc. is going to be much easier.

### 3.3 Entropy

In this section, we will use the concept of entropy to explore color dependence. Color dependence is a good measure of lack of diversity in color composition in a painting. We explain entropy in the following equations.

$$S(c, d) = \{p \in S(I) | C(p) = c \& D(p) = d\} \quad (1)$$

$$P(C = c, D = d) = \frac{|S(c, d)|}{N} \quad (2)$$

$$P(C = c|D = d) = \frac{P(C = c, D = d)}{P(d)} \quad (3)$$

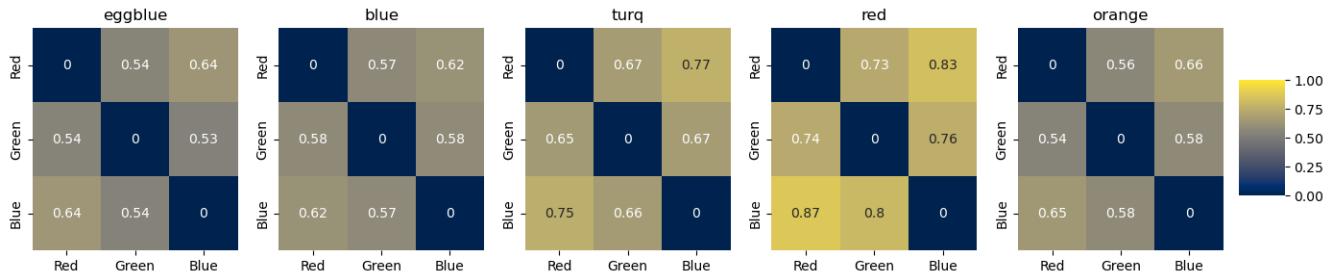
$$\begin{aligned} H(C|D = d) &= - \sum_{c=0}^{255} P(C = c|D = d) \log(P(C = c|D = d)) \\ &= - \sum_{c=0}^{255} \frac{P(C = c, D = d)}{P(d)} \log\left(\frac{P(C = c, D = d)}{P(d)}\right) \end{aligned} \quad (4)$$

$$H(C|D) = \sum_{d=0}^{255} P(D = d) H(C|D = d) \quad (5)$$

$$HR(C|D) = \frac{H(C|D)}{H(C)} \quad (6)$$

In Eq. 1 We let  $D$  and  $C$  represent 2 of the 3 main RGB colors, and  $d$  and  $c$  represent their respective values at a certain pixel  $p$ . In the first equation, we define the set  $S(c, d)$  which includes pixels which have values  $D = d$  and  $C = c$ . Then, we specify  $P(C = c, D = d)$  as being the cardinality of  $S(c, d)$  divided by the total number of pixels in Eq. 2. In other words,  $P(C = c, D = d)$  is the fraction of pixels that have  $C = c$  and  $D = d$ . Next, we make use of the definition of conditioning in Eq. 3 in order to make the substitution inside Eq. 4. Eq. 4 represents the conditional entropy of color  $C$  conditioned on  $D = d$ . We use the value of Eq. 4 for each  $d$  and weight it by the probability of  $D = d$ , and then sum all together to get the global conditional entropy of  $H(C|D)$  in Eq. 5. Finally in Eq. 6 we compute the relative entropy by the ratio of  $H(C|D)$  to  $H(C)$ . This ratio represents the fraction of information (in entropy) explained about  $C$  by conditioning on  $D$  as a whole. Note that a value of 0 indicates that color  $C$  is determined by color  $D$ , and a value of 1 indicates that  $C$  is independent of  $D$ . Proceeding with this understanding of relative entropy, we view Fig. 4 below to see the relative conditional entropies between each permutation pair of colors for each painting.

Figure 4: Painting Color Relative Entropies



We notice that for all the paintings, most of the entropy of a given color is not explained by conditioning on another color. Another way of putting this is that there is generally a large variety of color combinations present in the paintings; visually, this means that the paintings are all relatively color diverse in composition. Looking closely and comparing the paintings further, we note that the 'turq' and 'red' Shot Marilyns tend to have the highest relative entropy between the blue and red colors (given that the matrices are approximately symmetric), suggesting that these paintings have the widest variety of combinations of levels of blue and red. Conversely, we note that the 'egg blue' painting tends to have the lowest relative entropy between the red and green colors suggesting the lowest color diversity of combinations in red and green. We speculate that this is because of all the

paintings, the 'egg blue' painting is the most mono-chromatic, since the background, the eye shadow, and the head shadow are all approximately the same color, and Marilyn's face seemed almost blended in to these colors.

## 4 Color Extraction Using K-Means

### 4.1 Function Analysis

K-Means is a centroid-based clustering algorithm. It uses an iterative technique to group unlabeled data into K clusters based on cluster centers (centroids). The data in each cluster are chosen such that their average distance to their respective centroid is minimized. K-means clustering is a nondeterministic algorithm, meaning that the resultant clusters can change based on optimization criteria. The number of cluster is a hyper parameter which can be selected for using an elbow plot. The results of the K-Means algorithm are:

- K number of cluster centroids
- Data points classified into the clusters

Assuming we have input data points which correspond to pixel RGB color data from a particular painting  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n$  such that  $\mathbf{x}_i \in \mathbb{R}^3$ , we implement the following procedure for each  $k$  number of groups from  $k = 2, 3, \dots, K$  for a  $K$  large enough to be past the optimal  $k^*$ :

1. Randomly pick  $k$  points as the initial centroids from the dataset.
2. Find the Euclidean distance of each point in the dataset with the identified  $k$  points (cluster centroids).
3. Assign each data point to the closest centroid using the distance found in Step 2.
4. Find the new centroid by taking the average of the points in each cluster group.
5. Repeat 2 to 4 for a fixed number of iterations or until all observations are assigned to the same centroid as the previous step.

Euclidean Distance between two points in space is given by:

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2} \quad (7)$$

If each cluster centroid is denoted by  $\mathbf{c}_j$  where  $j = 1, 2, \dots, k$ , then each data point  $\mathbf{x}_i$  is assigned to a cluster based on:

$$\operatorname{argmin}_{\mathbf{c}_j \in C} d(\mathbf{c}_j, \mathbf{x}_i)^2 \quad (8)$$

The new centroid from the clustered group of points can be found through:

$$\mathbf{c}_j = \frac{1}{|S_j|} \sum \mathbf{x}_i \quad (9)$$

where  $\mathbf{x}_i \in S_j$  is the set of all points assigned to  $\mathbf{c}_j$ . Additionally,  $|S_j|$  denotes the cardinality of the set.

### 4.2 Evaluation

We implemented K-means to extract the main colors from each of the 5 images.

We further evaluated the images for an optimal value of  $k$  using distortion and the Elbow method. Distortion (Eq. 10) is calculated as the average of the squared distances from the cluster centers of the respective clusters to each data point, and is a measure of the fit of the algorithm[7].

$$\text{Distortion} = \frac{1}{K} \sum_{j=1}^K \frac{1}{|S_j|} \sum_{i=1}^{|S_j|} d(\mathbf{x}_i, \mathbf{c}_j)^2 \quad (10)$$

Finally, to obtain the optimal  $k$ , we iterate the values of  $k$  from  $1, \dots, n$ , and calculate the distortion for each value of  $k$  in the given range. Figure 5 shows the decrease in distortion begins to slow at  $K = 4$  which gives us the “elbow” of this graph.

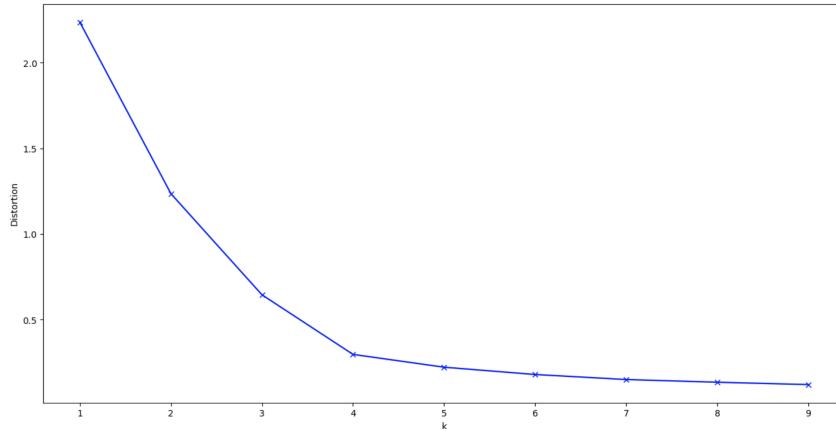


Figure 5: Elbow Method to Find Optimal "k"

Upon fixing the optimal value for K, we get the extracted colors as seen in Figure 6.

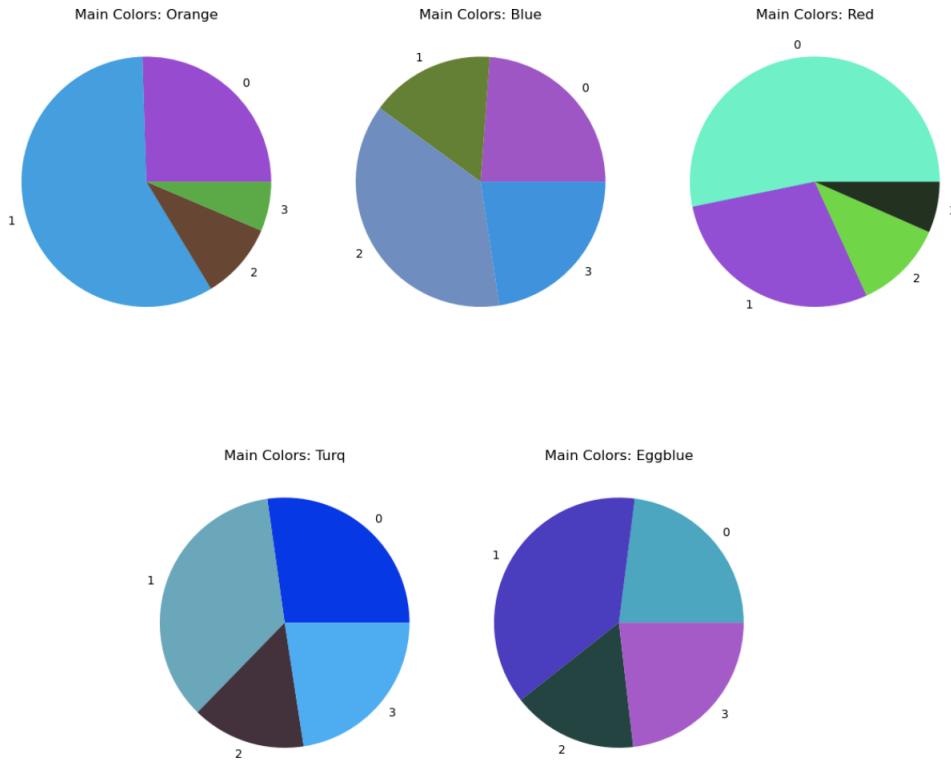


Figure 6: Color Extraction

## 5 Regions of interest

We used our 3-d visualizations of HSV measures for each of the paintings to guide our decisions in picking relevant color thresholds to extract the ROIs (Region of Interest). The ROIs that we

considered were her hair, her skin, and her eye shadow, as we believe numeric analysis can illuminate the subtle differences between these centerpiece ROIs to enrich interpretation of the painting. We used color thresholding and different masking techniques to extract the relevant ROIs. For the hair ROI, we just used basic color threshold masking. For the skin ROI, we used double masking (explained more in the skin section). Finally, for the eye shadow ROI, we used rectangular masking first and then applied color thresholding to more accurately extract the eye shadow color from the other similar colors (more details in the eye shadow section).

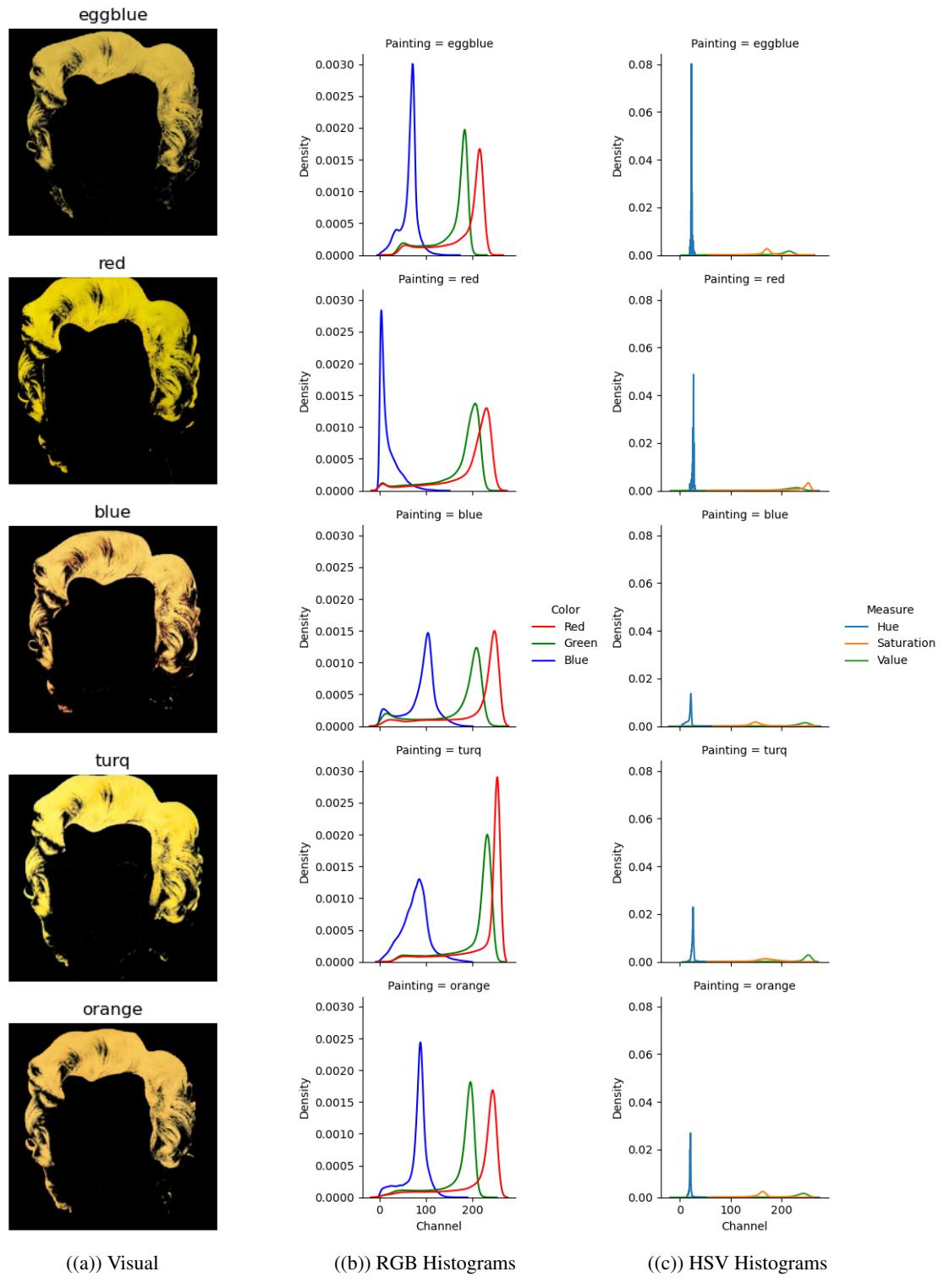


Figure 7: Hair Analysis

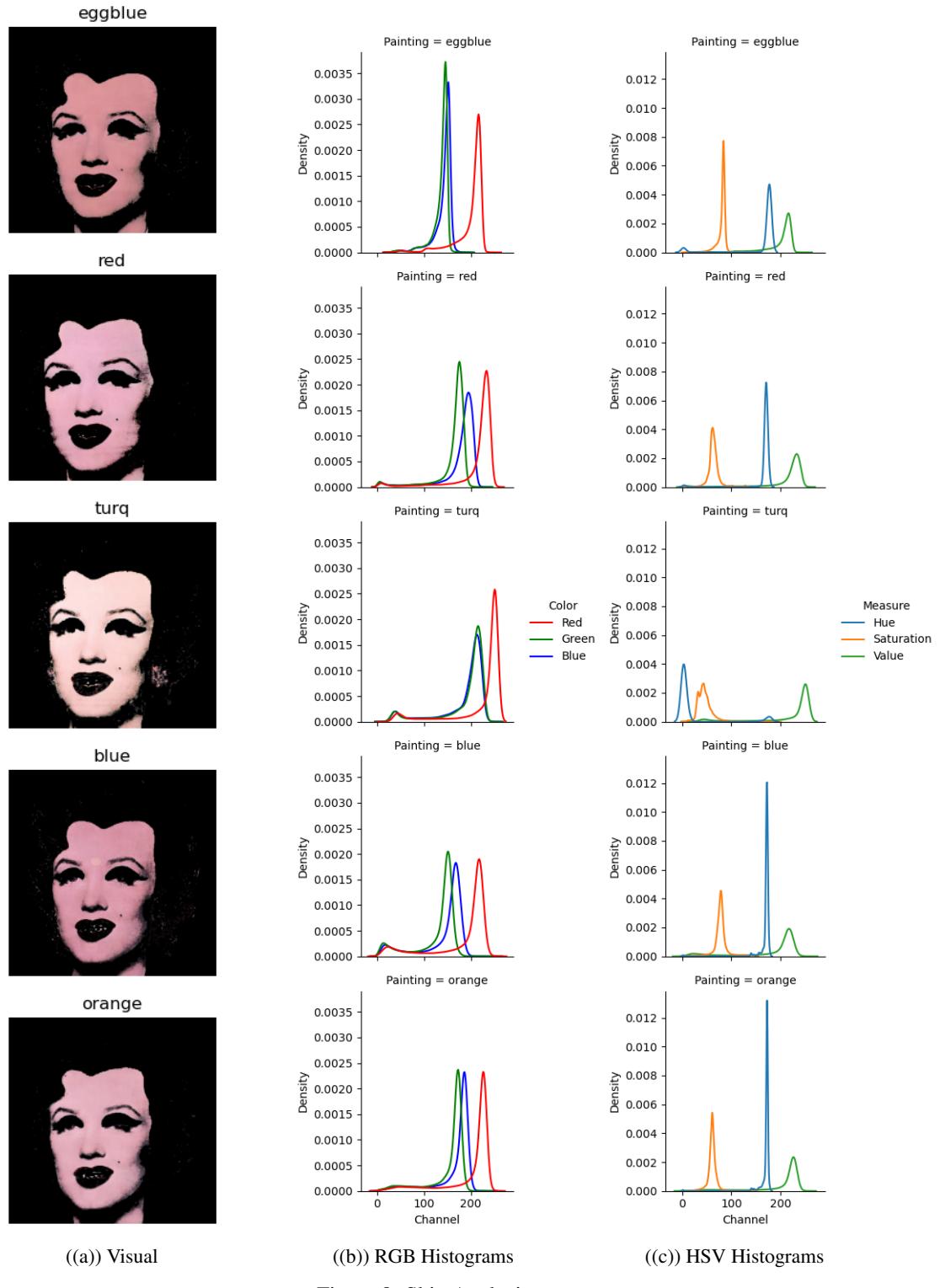
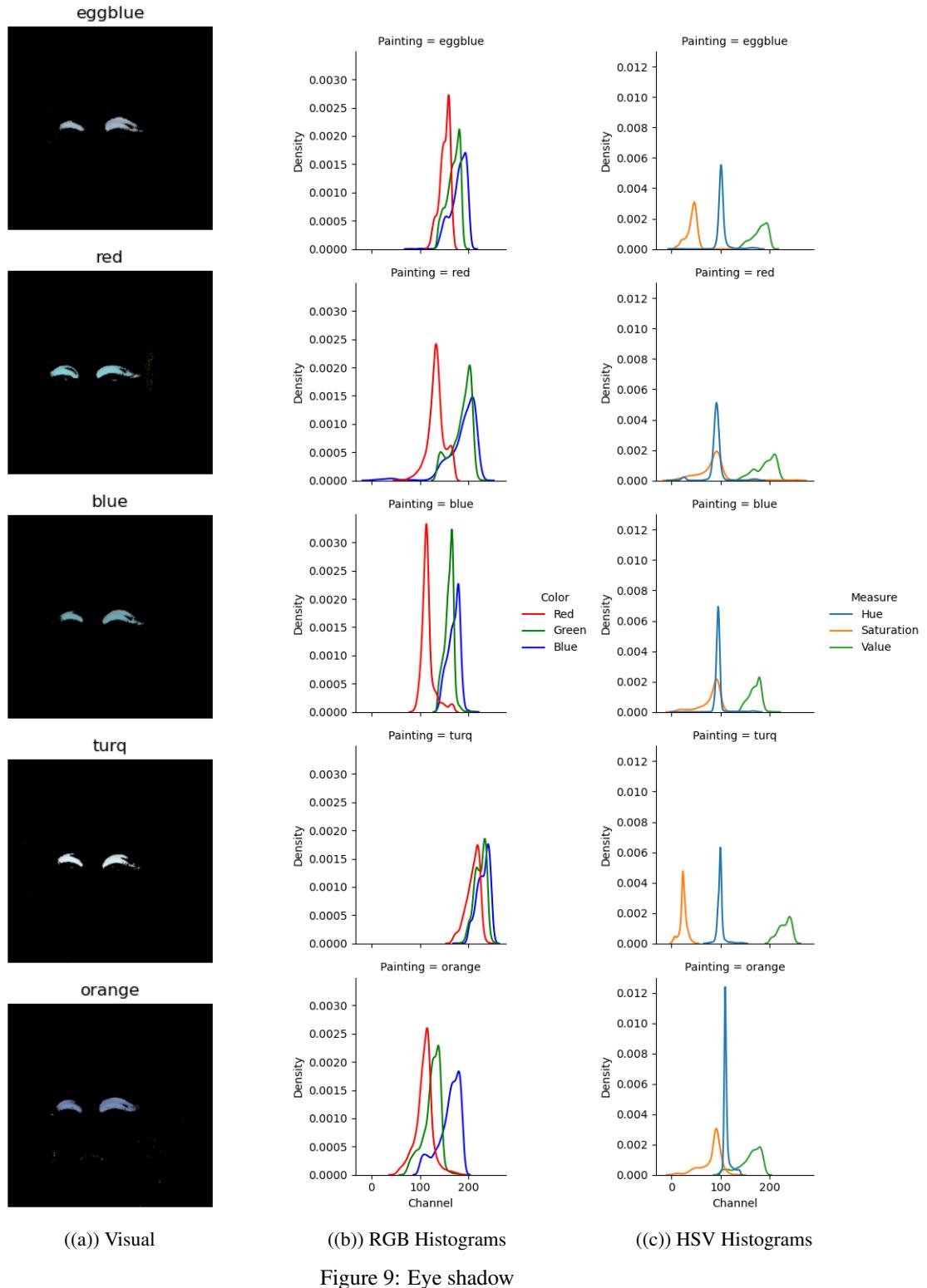


Figure 8: Skin Analysis



((b)) RGB Histograms

((c)) HSV Histograms

Figure 9: Eye shadow

## 5.1 Hair

The first region of interest taken from the five paintings is Marilyn Monroe's hair. As shown in Fig. 7, we see that the hair's color is yellow but after further inspection it is noticeable how different each

shade of yellow is from one another. Some are brighter and have an electric tone to them while others look more opaque almost looking a little bit orange in color.

When looking at the RGB plots right next to the pictures of her hair, we notice how the green and red lines are always very close together. This makes a lot of sense since these two colors make yellow. And whenever the color yellow is brighter you see higher density in the graph. The color blue is almost always around the range between [0, 100], since not a lot of blue is used except for when the hair has an opaque look with not a lot of shininess. From the HSV plots, we see how the saturation and the value are almost always around 0 in density while the hue varies from one painting to the other but is mainly in the range of [0, 100].

## 5.2 Skin

When we take a look at Fig 8, the next region of interest taken from the paintings is the skin which includes her face and part of her neck. The colors of her face range from kind of red to a light cream color. Taking a look at the RGB plots we see that the three colors, red, blue, and green are always together. This is because the face is using a combination of the primary colors to create new colors and in most of the cases the density is similar for the three colors.

In the HSV plots, hue ranges higher than the plots taken from the hair except when the skin is kind of white looking. The saturation is low because most of the colors of her face are lighter than the color of her hair. The highest saturation appears in the painting where the face has a darker color. Note that for extraction of her skin, we had to employ more than one mask. This is because her skin is composed of hues that are at the tail ends of the HSV 'H' spectrum and thus could not be thresholded for by a single mask threshold.

## 5.3 Eye shadow

Note that for extraction of the eye shadow, we had to apply a rectangular mask which blocked out any other colors that were similar to the eye shadow. This served as a shield such that our thresholding would only pick up the colors associated with her eye shadow. See Fig. 10 below for a visual of rectangular masking.



Figure 10: Rectangular Masking

In Fig. 9, the region of interest is the eye shadow that Marilyn Monroe is wearing. The colors are variations of blue but none of them is a pure blue. That is why the RGB plots show combinations of the three colors being very close together with the color blue being the highest in range most of the time. The hue has a big density when we analyze the eye shadow of the painting with the orange background and when we look at it we see that this eye shadow is the most different from the rest. Red and blue are similar in color and we can see that in their plots. Very similar densities and color range. This is a very small region to analyze but as you can see, something so small contains a whole lot of information and it gives us insight about the painting. Warhol did not paint each eye shadow the same which would have been easier. He went into the thought process of choosing each color for each painting.

## 6 Conclusion

The Shot Marilyns are an iconic series of paintings. They have been relevant for almost 60 years and have become a part of pop culture. Andy Warhol was a famous artist adept at using color to evoke meaning and emotions. Analyzing these pictures through algorithms and statistical methods made this project insightful. The project delved into studying the composition of different colors, and the meaning of their combination in a single painting. The first step taken after uploading the images was to do an exploratory data analysis where we saw the color composition of each images and how the colors ranged from [0, 255]. We implemented the same process with HSV plots to get a sense of the individual colors that were in the regions of interest. The ROIs gave us a lot of information about each painting that wasn't very clear to the naked eye. We were able to approach the paintings in different ways and noticed the slightly different colors Warhol used for each painting. How something we thought was just yellow, was a specific shade of yellow that worked the best with the skin tone and the background. Warhol took that into account with each of the five paintings he did.

As the audience, we used statistical methods and software packages to disintegrate the shots and extract colors and shadows. These elements, across the different shots, resulted in an interpretation of emotional states. Although his subject is same across all of his shots, Warhol's experiments with different colors imbued a set of unique attitudes and atmospheres. In certain instances, Monroe's hair looked grayed, and her smile looked unkempt and stale. Her blue laden eyes made her look tired, and different hues enhanced this nuance of the shot. As the pictures became more disintegrated and visible to us, we sensed a real person behind the camera, just like the rest of us. Despite having a team dedicated to make her look flawless, Marilyn's eyes, lips, and hair portrayed otherwise. Warhol was able to create an emotionally diverse selection of Marilyn's from a single print by leveraging the viewers' ability to assign meaning to certain colors and shapes.

## 7 References

- [1] Hsieh, F (2023), *Data Analysis from Scientific Perspective.*, UC Davis.
- [2] The Interior Review. "A Visual Critique of Warhol's Shot Sage Blue Marilyn, 1964." The Interior Review, 10 May 2022, [www.theinteriorreview.com/story/2022/5/10/critically-assessing-warhols-shot-sage-blue-marilyn](http://www.theinteriorreview.com/story/2022/5/10/critically-assessing-warhols-shot-sage-blue-marilyn). Accessed 20 May 2023.
- [3] "Shot Marilyns." Wikipedia, 12 May 2023, [www.en.wikipedia.org/wiki/ShotMarilyns](http://www.en.wikipedia.org/wiki/ShotMarilyns). Accessed 8 Jun. 2023.
- [4] Ghighi, Emma . "Andy Warhol, the Shot Marilyns, and His Early Silkscreens." Revolvery Gallery, 10 Jan. 2022, [www.revolverwarholgallery.com/andy-warhol-the-shot-marilyns-and-his-early-silkscreens/](http://www.revolverwarholgallery.com/andy-warhol-the-shot-marilyns-and-his-early-silkscreens/). Accessed 8 Jun. 2023.
- [5] MasterClass. "Hue, Saturation, Value: How to Use HSV Color Model in Photography." Masterclass, 29 Sept. 2021, [www.masterclass.com/articles/how-to-use-hsv-color-model-in-photography](http://www.masterclass.com/articles/how-to-use-hsv-color-model-in-photography).
- [6] Stone, Rebecca. "Image Segmentation Using Color Spaces in OpenCV + Python." Real Python, [www.realpython.com/python-opencv-color-spaces/](http://www.realpython.com/python-opencv-color-spaces/). Accessed 9 June 2023.
- [7] Ng, Andrew. "The k-means clustering algorithm". CS 229.