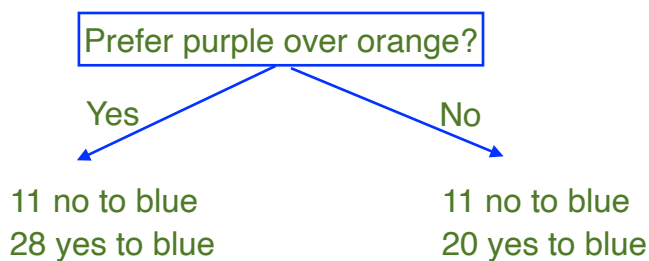


Please refer to the  
 Color preferences.csv  
 file posted on Canvas to find the probability values needed for this first page of questions.

**Q1 (9pts):** Create a stump decision tree by splitting the data based on purple preference, where the left child holds all the samples that answered Yes, and the right child holds all the samples that answered No.

**Q1a:** Given the task of predicting blue preference, what's the Total Weighted Gini Impurity of this split?



$$GI_{yes} = 1 - [(11/39)^2 + (28/39)^2] = 0.405$$

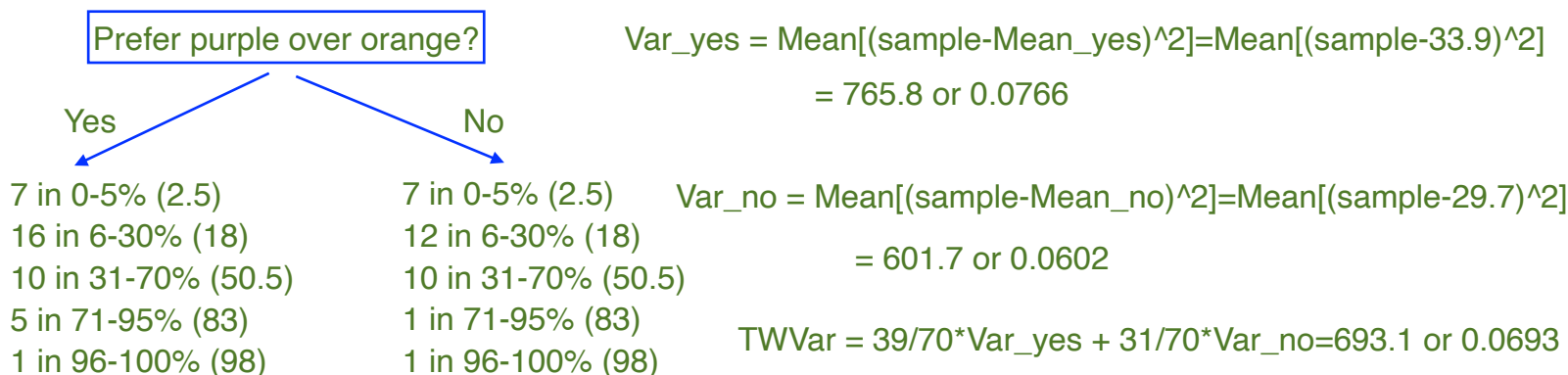
$$GI_{no} = 1 - [(11/31)^2 + (20/31)^2] = 0.458$$

$$TWGI = 39/70 * GI_{yes} + 31/70 * GI_{no} = 0.428$$

**Q1b:** What would be the predicted labels of the left and right leaves?

Both leaves predict yes in response to preferring blue over red

**Q1c:** Given the task of predicting percentage of white t-shirts, what's the Total Weighted Variance of this split?



$$\begin{aligned} Var_{yes} &= \text{Mean}[(\text{sample} - \text{Mean}_{yes})^2] = \text{Mean}[(\text{sample} - 33.9)^2] \\ &= 765.8 \text{ or } 0.0766 \end{aligned}$$

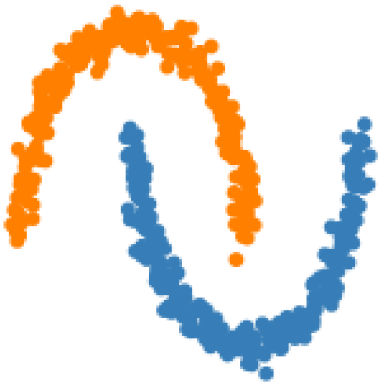
$$\begin{aligned} Var_{no} &= \text{Mean}[(\text{sample} - \text{Mean}_{no})^2] = \text{Mean}[(\text{sample} - 29.7)^2] \\ &= 601.7 \text{ or } 0.0602 \end{aligned}$$

$$TWVar = 39/70 * Var_{yes} + 31/70 * Var_{no} = 693.1 \text{ or } 0.0693$$

**Q1d:** What would be the predicted values of the left and right leaves?

Left predicts 33.9%, Right predicts 29.7%

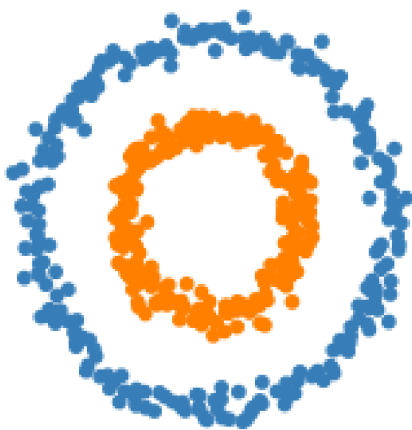
**Q2 (6pts):** For each of the following clusters, give a clustering method you would NOT recommend, and a reason why you think it would fail.  
Don't repeat clustering methods



**Q2a:**

k-means  
centroid/complete-linkage hierarchical/agglomerative  
Mean Shift  
EM with GMM

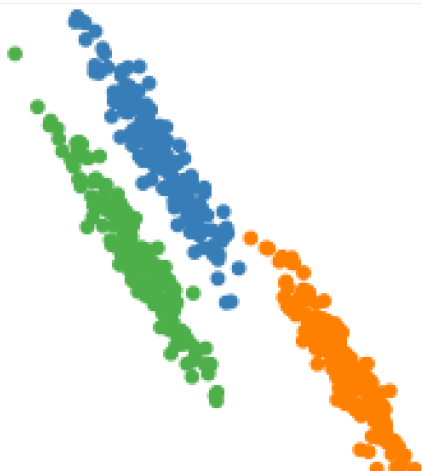
Will all fail because the data is not spherical and/or not linearly separable



**Q2b:**

k-means  
centroid/complete-linkage hierarchical/agglomerative  
Mean Shift  
EM with GMM

Will all fail because the data is not linearly separable



**Q2c:**

k-means  
centroid/complete-linkage hierarchical/agglomerative  
MeanShift

Will all fail because the clusters are not spherical

DBScan could fail because the low density areas should not be outliers