

# STATS 604 Project 3 Preanalysis Plan

Khoa Do, Surtai Han, Jialin He, Carlyle Morgan

2025-10-23

## Introduction

In this project, we attempt to analyze various procedures for preserving the freshness of cilantro. Namely, we will attempt to measure if the following have any effect:

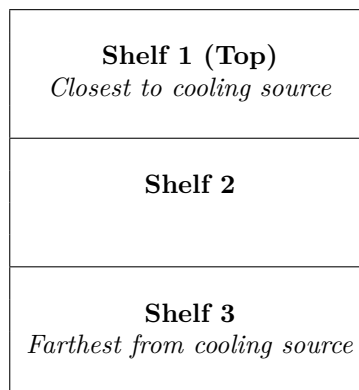
- Keeping cilantro in the fridge versus outside
- Keeping cilantro in a plastic bag versus not
- For refrigerated cilantro, keeping cilantro close to the cooling source at the top of the fridge versus farther away.

We hypothesize that cilantro kept in a fridge in a bag and far away from the cooling source will remain fresh for longer.

## Methodology

### Experimental setup

For convenience and budgetary concerns, we decided to use a group mate's mini-fridge to apply the treatment. A diagram of this mini-fridge is shown below:



*Mini-Fridge layout*

## Randomization Scheme

From our 40 cilantro stalks, we assigned a group of 5 cilantro stalks to one of eight treatment groups at random. The treatment plans for each of the eight groups is as follows:

Group	Location	Bagged
1	Outside the fridge	No
2	Outside the fridge	Yes
3	Fridge shelf 1	No
4	Fridge shelf 1	Yes
5	Fridge shelf 2	No
6	Fridge shelf 2	Yes
7	Fridge shelf 3	No
8	Fridge shelf 3	Yes

## Data Collection

The weight of each set of sprigs is measured in grams daily (thus, 8 measurements per day) using an AmazonBasic scale with dubious precision. One image of each set is taken daily under as similar lighting conditions as possible using the same smart phone camera.

## Data Processing

To quantify cilantro freshness from daily images, we use computer vision techniques with HSV (Hue, Saturation, Value) color space analysis. Each image is converted from RGB to HSV color space, which separates color information (hue) from intensity and lighting effects. We apply a calibrated HSV threshold mask (H: 35-85, S: 40-255, V: 40-255) to segment cilantro pixels from the background, producing a binary segmentation.

After segmentation, we use contour detection to identify exactly 5 sprigs per image. Contours are identified using the RETR\_EXTERNAL method, which finds outer boundaries of connected regions. To ensure we detect exactly 5 sprigs, we use an adaptive threshold approach: starting with a minimum area threshold of 10,000 pixels, we iteratively relax this threshold (reducing by 30% each iteration) until at least 5 contours are detected. Once found, the 5 largest contours are selected and assigned unique identifiers (1-5).

From each individual sprig, we extract the following metrics as indicators of freshness:

- **Mean Hue:** Average green color (0-179 scale)
- **Mean Saturation:** Color intensity (0-255 scale)
- **Mean Value:** Brightness level (0-255 scale)

Figure 1 illustrates the HSV segmentation and sprig identification process applied to a sample cilantro image. The original image (left) is converted to HSV color space and thresholded to create a binary mask (center), where white pixels represent detected cilantro and black pixels represent background. The adaptive contour detection algorithm then identifies the 5 largest sprigs (right), with each sprig outlined in a different color and labeled with a unique identifier (1-5). HSV statistics are calculated separately for each of the 5 selected sprigs.



Figure 1: HSV segmentation and sprig identification pipeline showing: (left) original cilantro image, (center) binary mask from HSV thresholding, and (right) the 5 largest sprigs identified via adaptive contour detection, each outlined in a distinct color and labeled with a unique identifier (1-5) for per-sprig tracking.

## Testing Procedure

### EDA

Following the data processing, we will have a time-series dataset where each of the 40 replicates has 3 individual sprig measurements. To visualize the decay of freshness over time, we will first generate a heat map where the Y-axis is the 8 different experimental conditions, the X-axis is Day, and the cell's color indicates the average Mean Saturation (or Mean Hue) for that condition on that day. This graph can provide a clear comparison of which conditions lose their freshness more rapidly. Following this, to compare different conditions, we will generate summary box plots of the Mean Saturation values. These box plots will be organized by our 8 experimental conditions every three days after the start of our experiment.

### Regression Model

To test the main effects and their interaction, we will build a regression model. The model is as below:

$$Y = \beta_1 \times Location + \beta_2 \times Bag + \beta_3 \times Location : Bag$$

Here  $Y$  is the above metric as indicators of freshness. And both Location and Bag are factor variables. From the summary of the model, we can get a view of whether the two variables and their interaction influences the freshness. First, we test the Location:Bag interaction term. If  $p < 0.05$ , we will conclude that the effect of bagging depends on the location. We will then proceed to use log-rank tests to compare “In Bag” vs. “Out of Bag” within each location separately to analyze the effects. If  $p > 0.05$ , we will conclude the effects are independent. We will remove the interaction term and use the simpler model:

$$Y = \beta_1 \times Location + \beta_2 \times Bag$$

And the p-values for the Location and Bag terms will test our primary hypotheses.

Things to test:

1. Fridge layer versus each other

- Compute mean of differences and compare
- Compute F statistic, permutation test on F statistic for the three levels
  - Sensitive to outliers

Possible issue: when we weigh cilantros collectively, we can't really apply a permutation test. So do we continue to use weight as a metric?

If no significance found, pool the fridge samples

2. Fridge versus outside
3. Out of bag vs in-bag

- If Test 1 was significant, there is some necessary stratification.

Should pre-analysis plan specify response to common issues? (i.e. If my data has a long tail, should I pre-specify that I plan to use medians. A: yes)

More versatile possibility: using ranks, which works with F statistics.

## Permutation Test Approach

- Permutation test on whether different layers of the fridge (which presumably has different temperature moisture) has an effect on preserving freshness, potentially pulling data from different layers together after this test.
- Permutation test on whether putting in a ziplock bag has an effect toward freshness: test between 30 sprigs in the fridge. Use conditioning if data from 3 layers differ and could not be pulled together.
- Permutation test on whether in fridge and out of fridge sprigs differ in freshness. E.g. normal permutation test between 20 out-of-ziplock-bag sprigs.

## Model Based Scheme

- ANOVA: use an F test to test significant differences in means between different layers of the fridge.
- Fits linear model (e.g. weight) as a function of time, use the slope as an indicator of weight loss speed.

## Preliminary Results

- Conclusion on whether putting cilantro sprigs in different layers affect freshness
- Conclusion on whether putting cilantro sprigs in the fridge helps preserving freshness
- Conclusion on whether putting cilantro sprigs in ziplocks helps preserving freshness