

Mushroom Edibility Classification

YANFEIYUN WU (KATHY)

DATA SCIENCE INSTITUTE AT BROWN UNIVERSITY

OCT 19TH, 2023

[HTTPS://GITHUB.COM/KATHYWU1201/MUSHROOM_EDIBILITY_CLASSIFICATION_ML_PROJECT](https://github.com/kathywu1201/mushroom_edibility_classification_ml_project)

A photograph of two red mushrooms with white spots, likely Amanita muscaria, growing on a bed of green moss in a forest. The background is blurred, showing tree trunks and foliage. The image is partially obscured by a white, torn-edge graphic that separates it from the text on the right.

I want to know...

- If there exists a mushroom that has never seen before, will it be edible or it is poisonous?
- As the nature is evolving, new type of fungi come out, but how do we know if we are safe to eat it?
- Help them to identify if the mushroom is safe to consume.
- Reduce the risk of death and intoxication due to poisonous mushrooms.

Data Descriptions

- Classification problem identifying if a new fungi is poisonous or not.
- 'class' feature as Target Variable.
- Data from [Kaggle](https://www.kaggle.com/datasets/devzohaib/mushroom-edibility-classification/data), collected from Patrick Hardin's Mushrooms & Toadstools, and inspired by Jeff Schlimmer's Mushroom Data Set.
- Original Data: 61069 rows, 21 columns.

<https://www.kaggle.com/datasets/devzohaib/mushroom-edibility-classification/data>



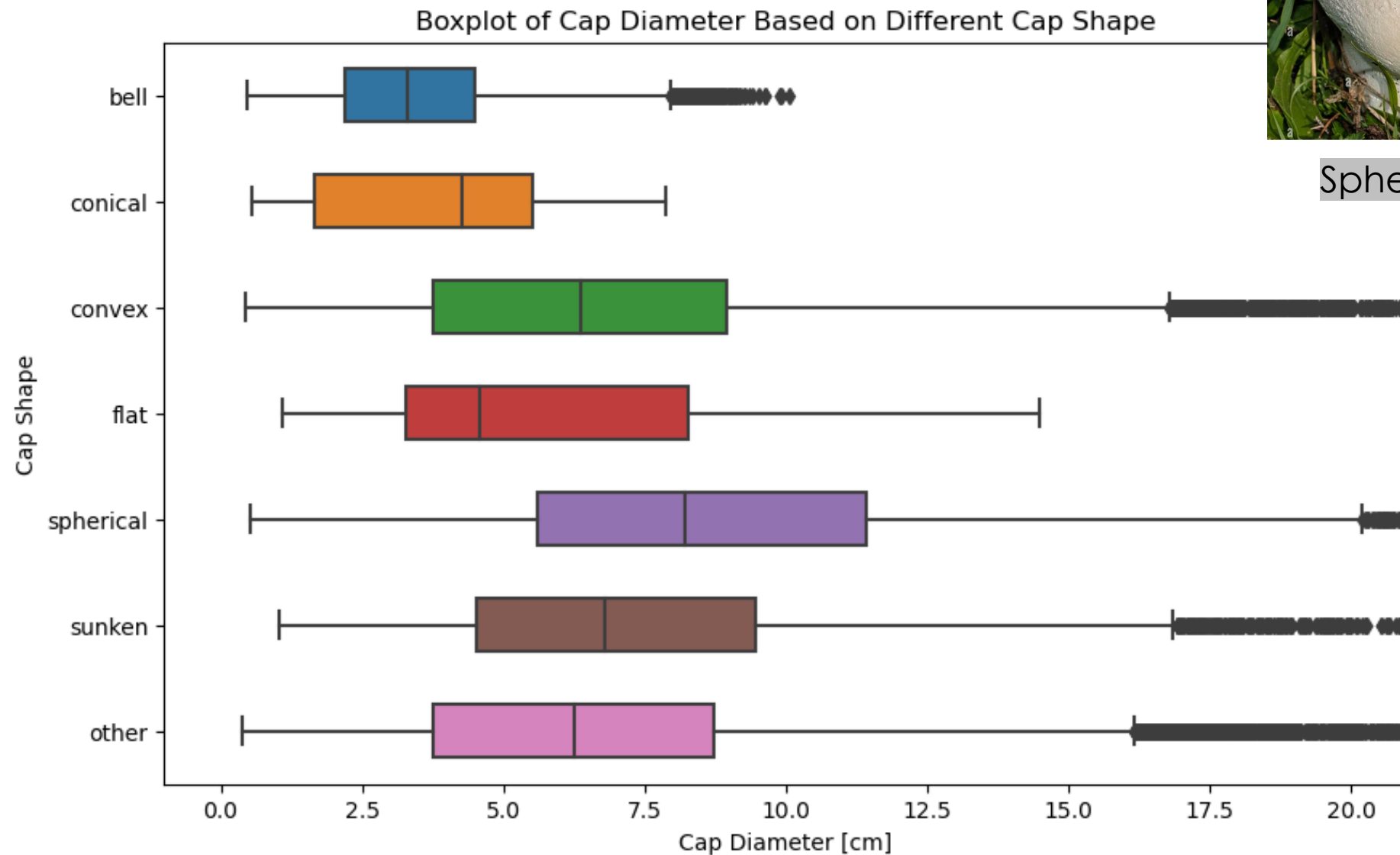
Boxplot of Cap Diameter



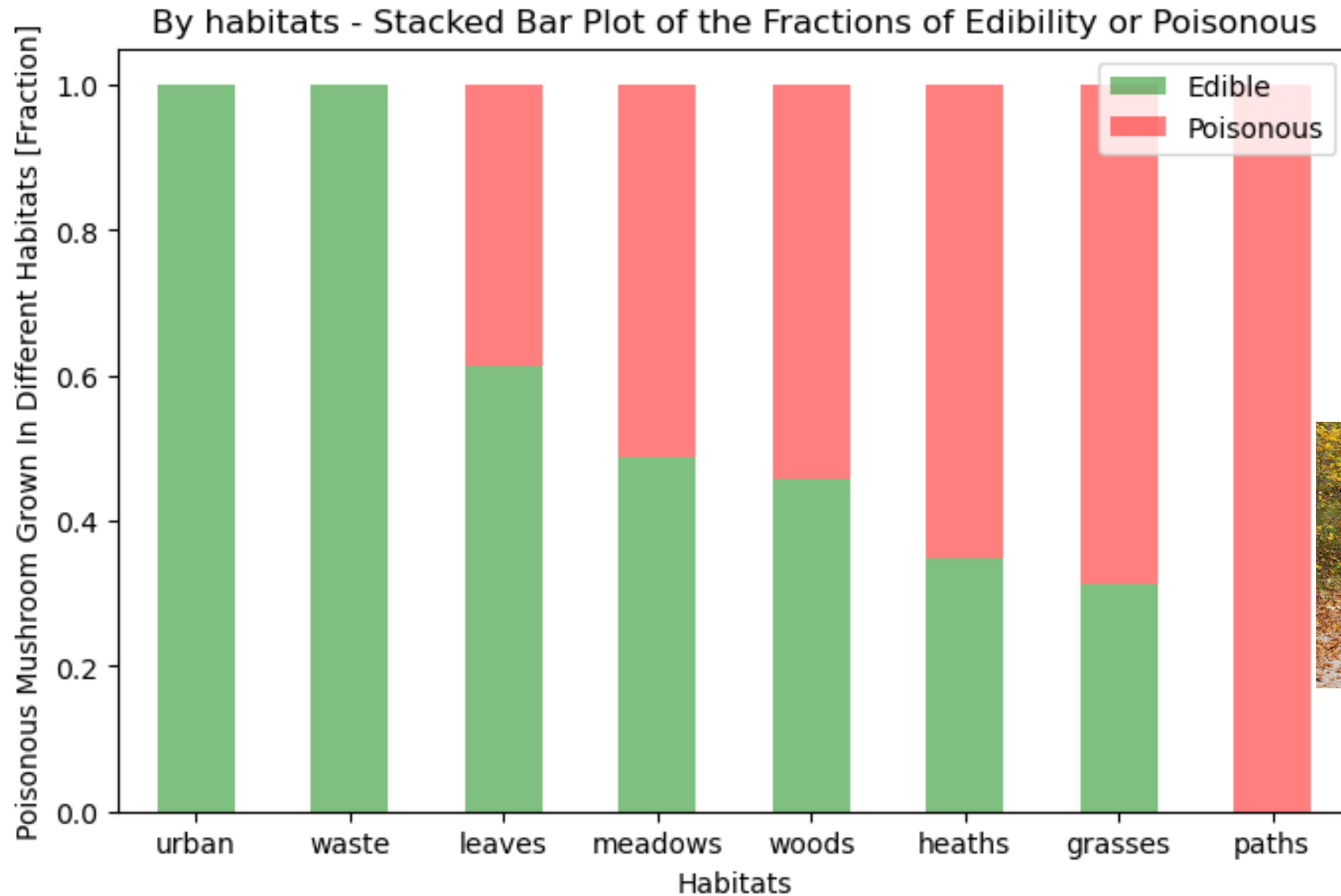
Bell shape



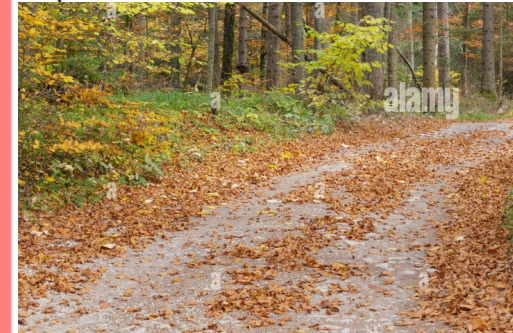
Spherical shape



By Habitats - Stacked Bar Plot

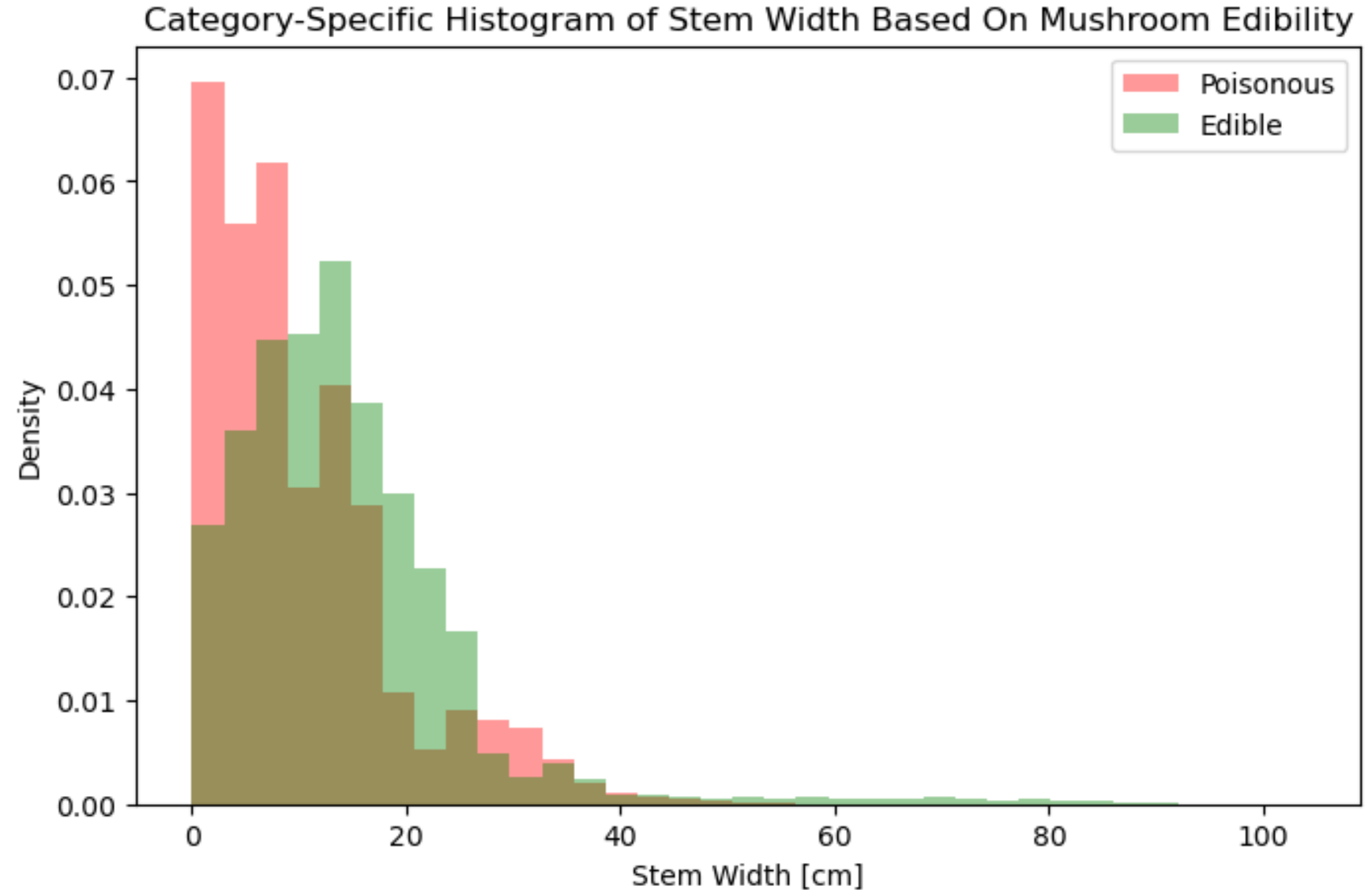
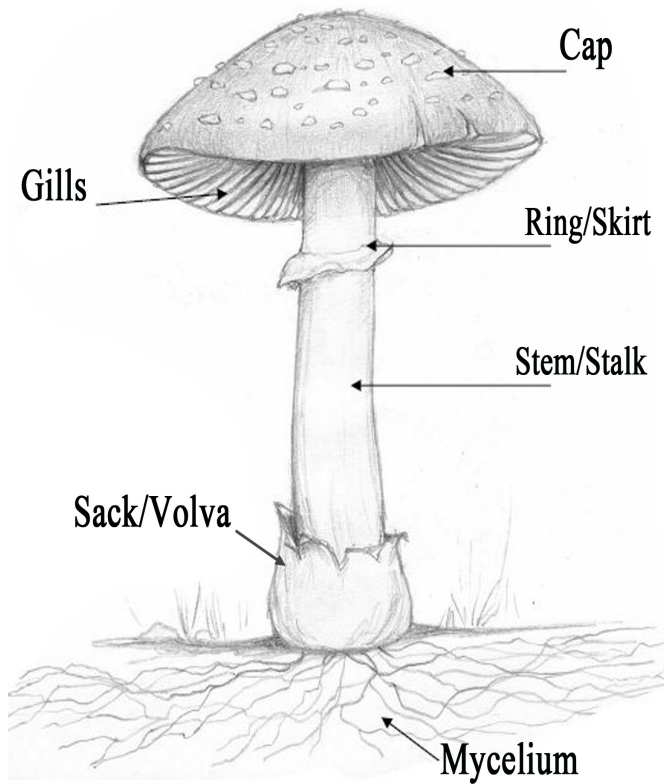


Urban habitat



Path habitat

Category-Specific Histogram of Stem Width



Splitting

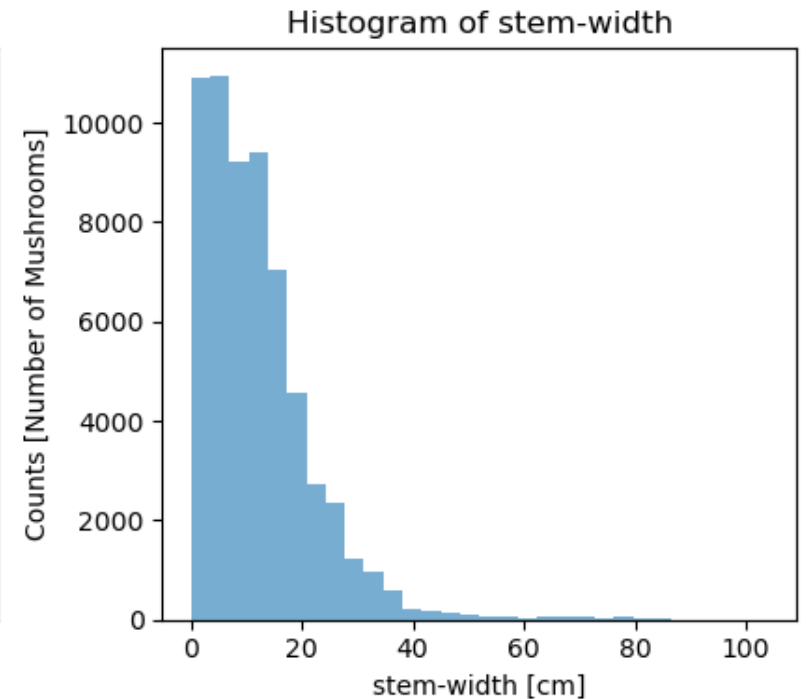
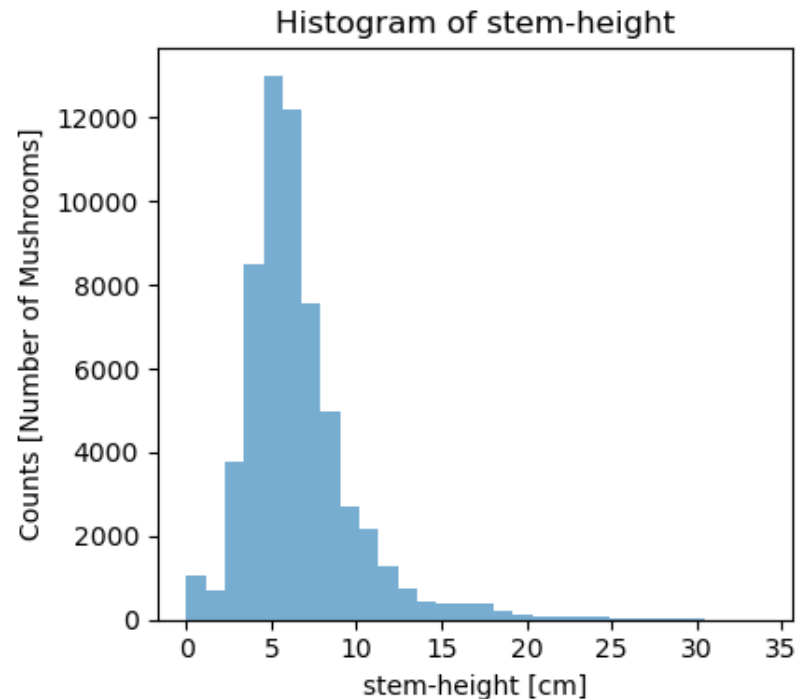
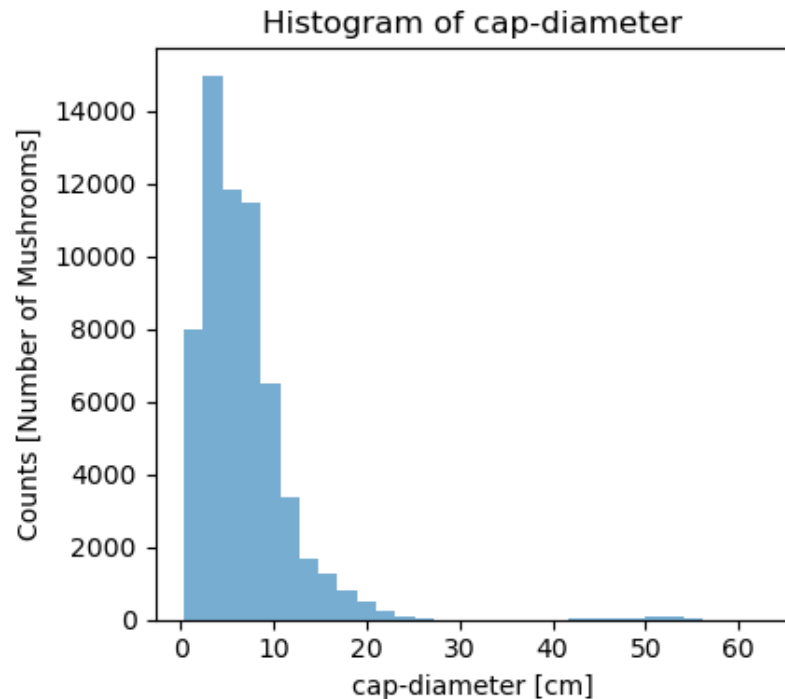
- The ratio of target variable (classification) is about 4:6 for edible and poisonous respectively.
- Use general **train_test_split** to set 20% of the dataset as test set and the rest 80% as other for **Kfold Cross Validation**.
- The shape of each other and test set are the following:
X_other: (48855, 20), **X_test:** (12214, 20)
y_other: (48855,) , **y_test:** (12214,)

Preprocessing

- 2 preprocessors are used:
 - (i) Unordered categorical data: **one-hot encoder**
 - (ii) Continuous features: **Standard Scalar**
- Most of the features are categorical data and unordered:
e.g. color, shape, surface, etc.
- Two binary features - **does-bruise-or-bleed** and **has-ring** – will not perform any encoders.

Standard Scalar

- **cap-diameter, stem-height, stem-width**
- The features are heavy tailed, thus choose standard scalar



Missing Values

- Among **20 features** (exclude the target variable), **9 features** contains missing values.
- All the missing values are categorical features.
- **stem-root, veil-type, veil-color, spore-print-color** have more than 80% of missing values

```
fraction of missing values in features:  
cap-surface      0.231214  
gill-attachment  0.161850  
gill-spacing     0.410405  
stem-root        0.843931  
stem-surface     0.624277  
veil-type        0.947977  
veil-color       0.878613  
ring-type        0.040462  
spore-print-color 0.895954
```

After Splitting and Preprocessing...

- Number of features **before** transformation: 20
Number of features **after** transformation: 126
- Most of the features are categorical data, thus the number of features after preprocessing **increases**
- **No features** contain missing values after transformation



Thanks for listening!