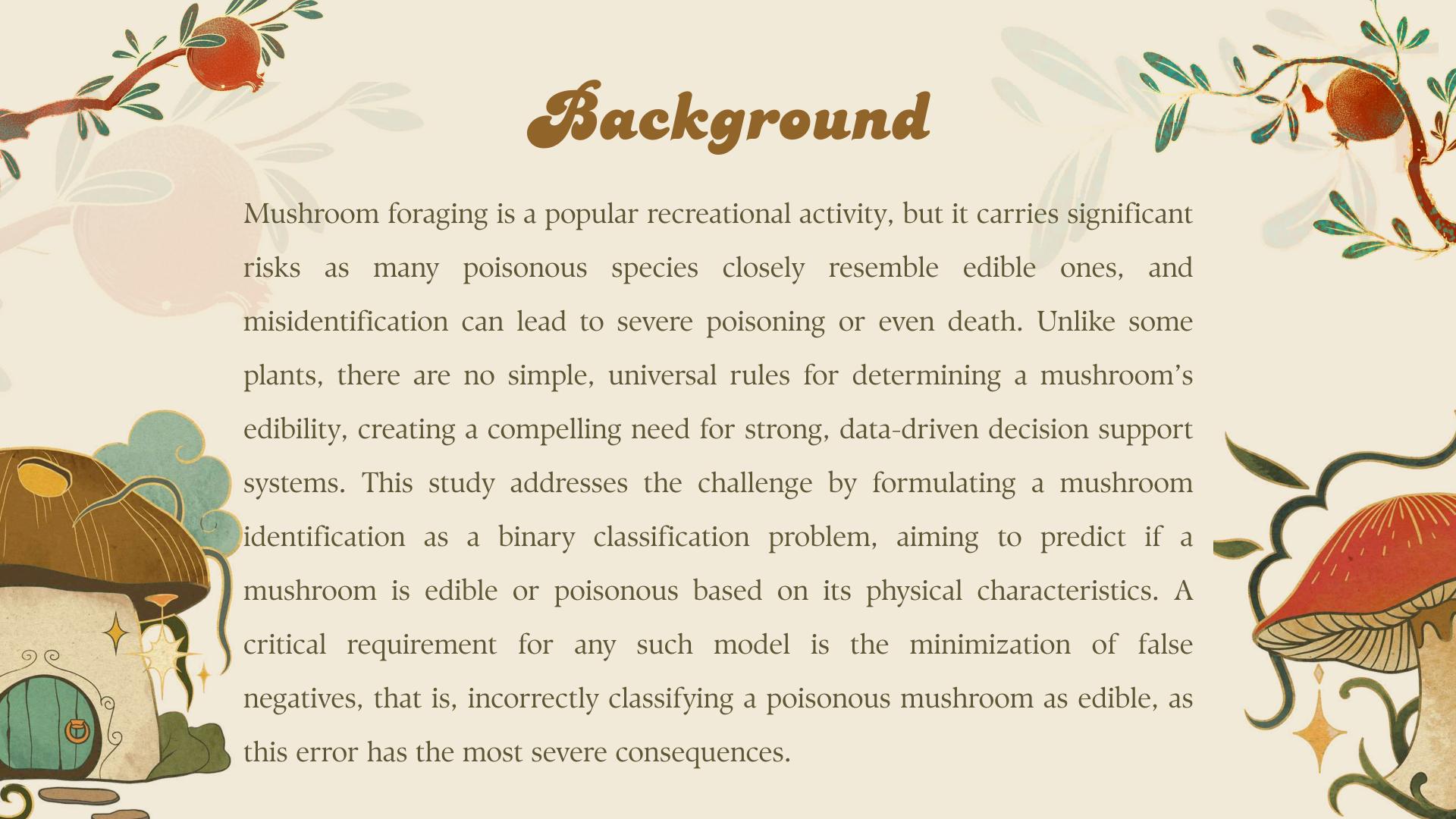


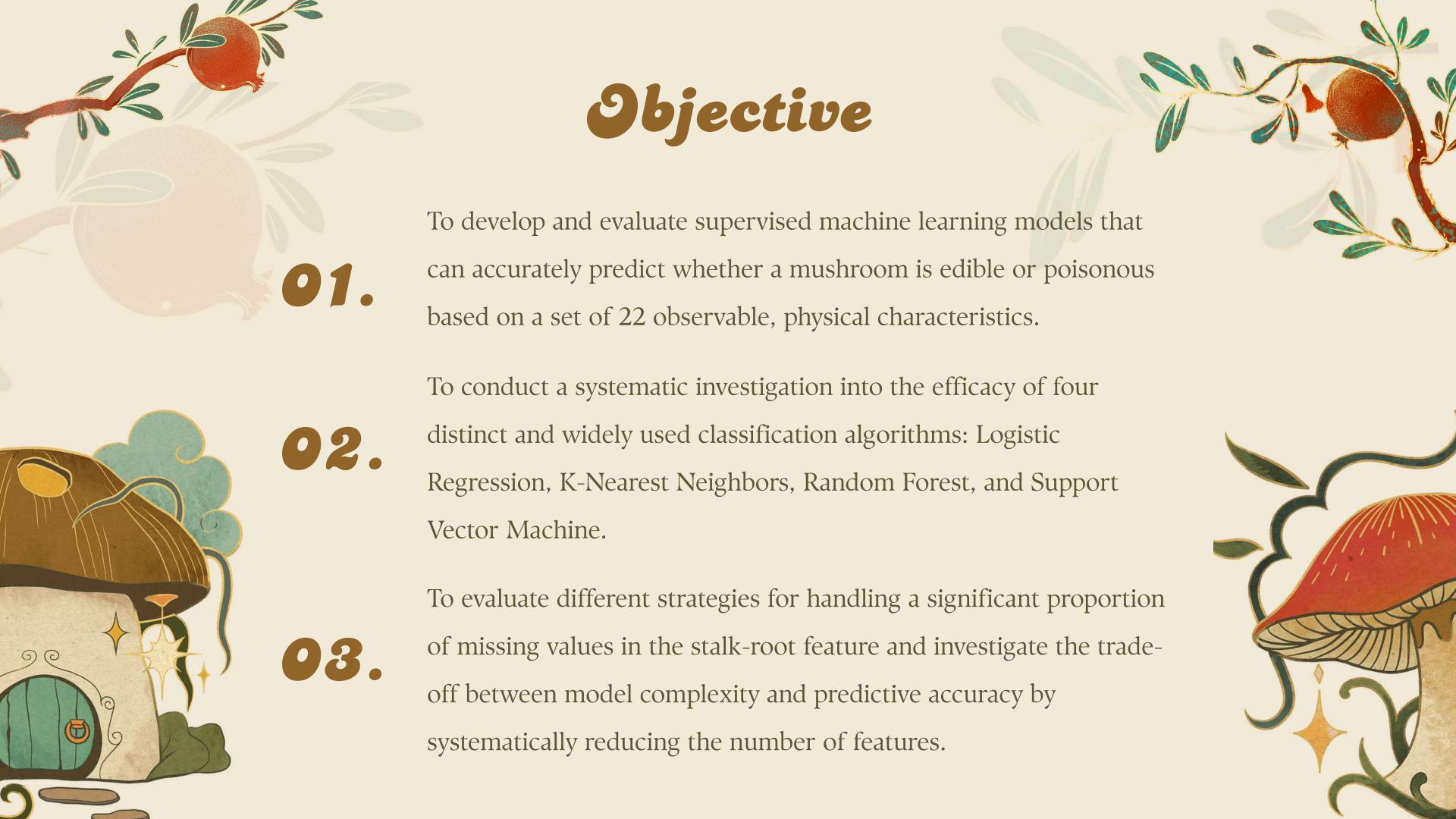


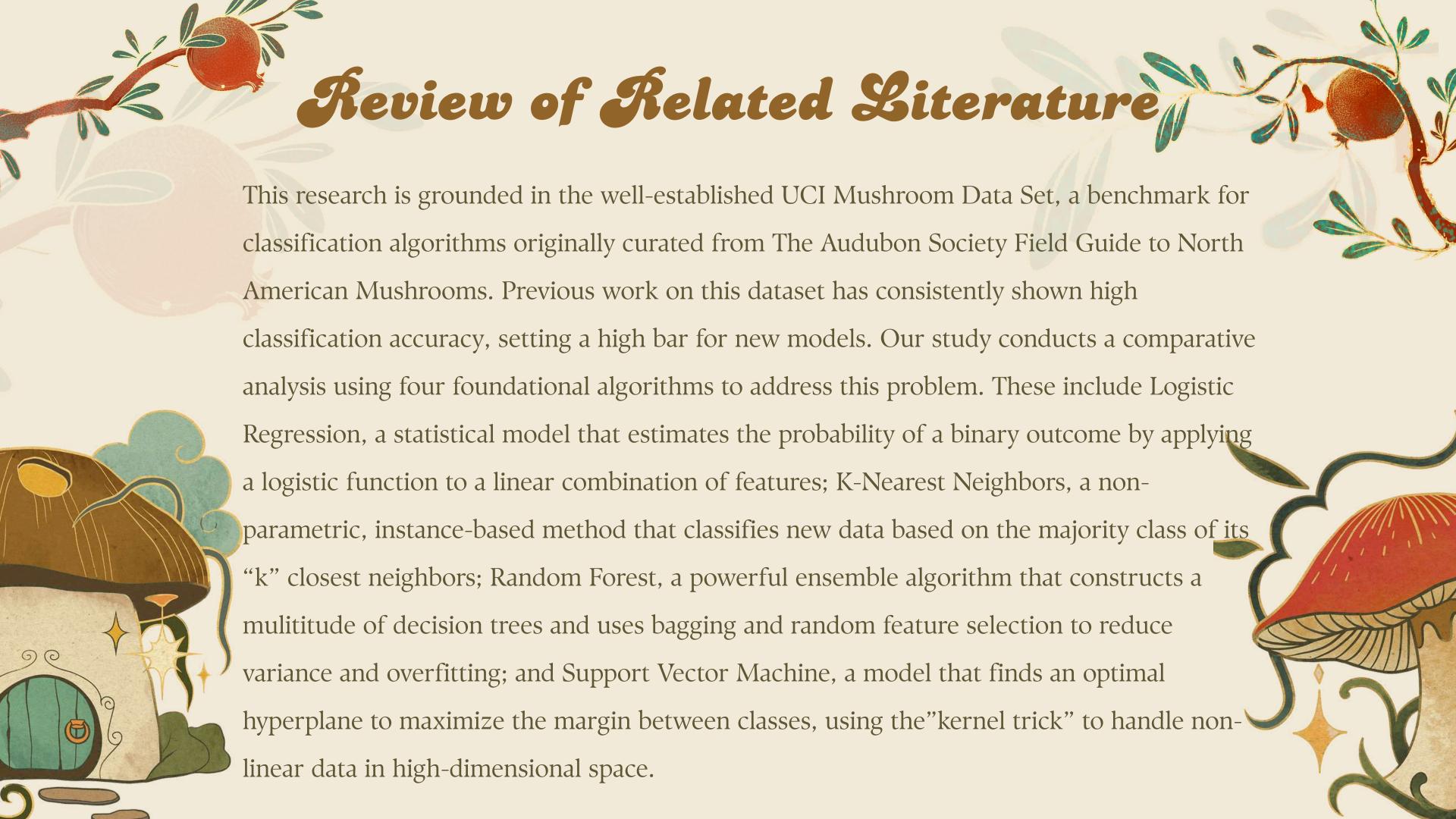
Presenter

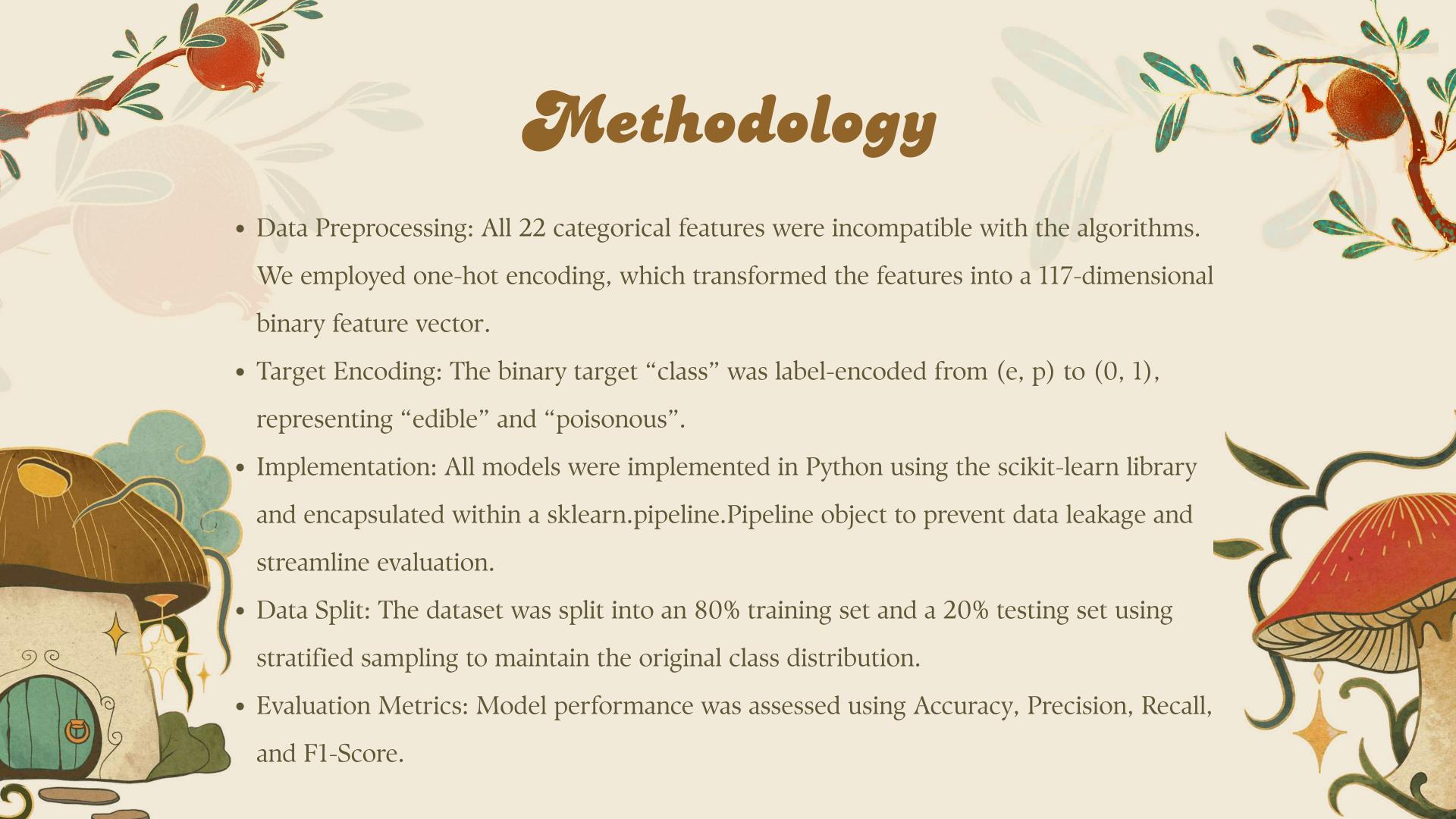
- Gabriel Angelo Vinas
- Ma. Mae Sid Santos











Wataset and the second second

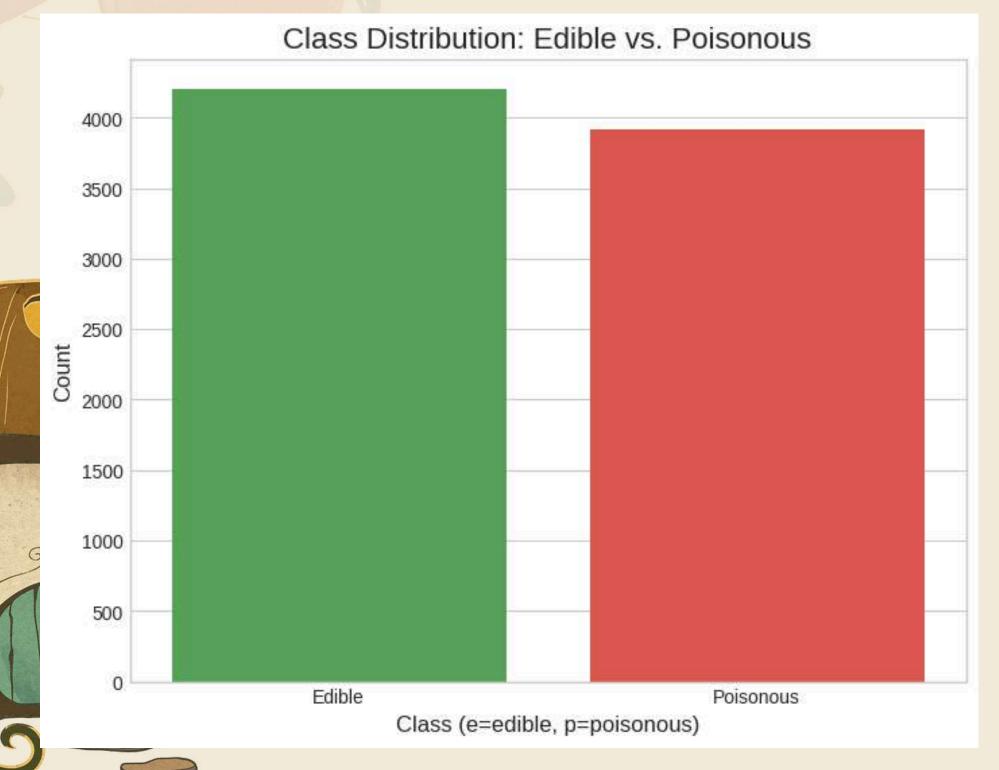
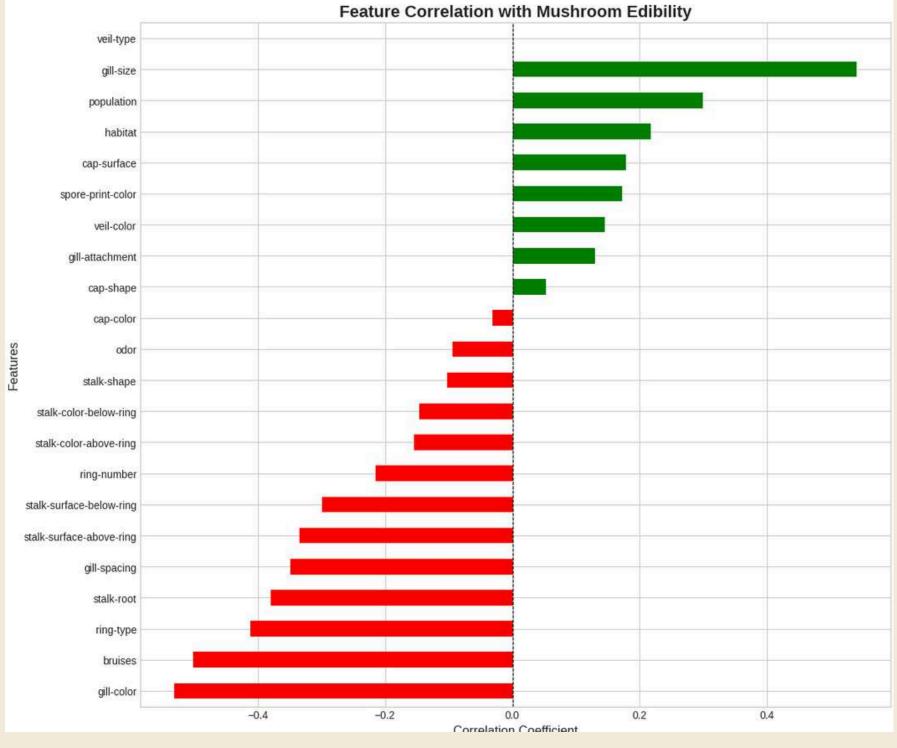


TABLE I: Selected Dataset Attribute Descriptions

Feature	Codes	Description	
cap-shape	b, c, x, f, k, s	bell, conical, convex, flat, knobbed, sunken	
odor	a, l, c, y, f, m, n, p, s	almond, anise, creosote, fishy, foul, etc.	
gill-color	k, n, b, h, g, r, o, p, u, e, w, y	black, brown, buff, chocolate, gray, etc.	
bruises	t, f	bruises, no bruises	
habitat	g, l, m, p, u, w, d	grasses, leaves, meadows, paths, etc.	





The Problem: A significant data quality challenge was the 'stalk-root' feature, where 2,480 instances (30.5% of the data) had a missing value denoted by '?'.

Objective: To determine the optimal strategy for handling these missing values and quantify its impact on model performance.

Three Strategies Tested:

- 1. Strategy A (Categorical Treatment): Treat '?' as a distinct category, 'not observable', testing the hypothesis that its absence is informative.
- 2. Strategy B (Modal Imputation): Replace '?' with the most frequent (modal) value in the column, which is 'bulbous' ('b').
- 3. Strategy C (Feature Removal): Remove the entire 'stalk-root' feature from the dataset







TABLE II: Performance Comparison of Stalk-Root Handling Strategies (Accuracy / F1-Score)

Model	Strategy A (Keep '?')	Strategy B (Impute)	Strategy C (Drop)
LR	99.94% / 99.94%	99.88% / 99.87%	99.88% / 99.87%
KNN	99.82% / 99.81%	99.82% / 99.81%	99.94% / 99.94%
RF	100.0% / 100.0%	100.0% / 100.0%	100.0% / 100.0%
SVM	100.0% / 100.0%	100.0% / 100.0%	100.0% / 100.0%

Result: As shown in Table II, performance was "remarkably consistent" across all three strategies. The average accuracy for all models differed by less than 0.1%. Random Forest and SVM achieved 100% accuracy in all three scenarios.

It suggests that stalk-root is not a critical feature for this dataset. and this experiment may contribute to practical implication in data gathering of mushrooms.

Final Decision: We selected Strategy A (Keep '?') to use for our final model. This "prioritizes data integrity" and is a "more scientifically robust approach" as it avoids introducing artificial certainty through imputation or discarding data (like removal).



Experiment 2 Impact of feature selection

Objective: To investigate the relationship between the number of features and model performance, and to determine if a simpler, more parsimonious model could achieve comparable accuracy.

Ranked all 117 one-hot encoded features using a Random Forest model's feature importance scores (Gini impurity reduction).

Trained and evaluated all four models on four distinct feature sets:

All 117, Top 20, Top 10, and Top 5.



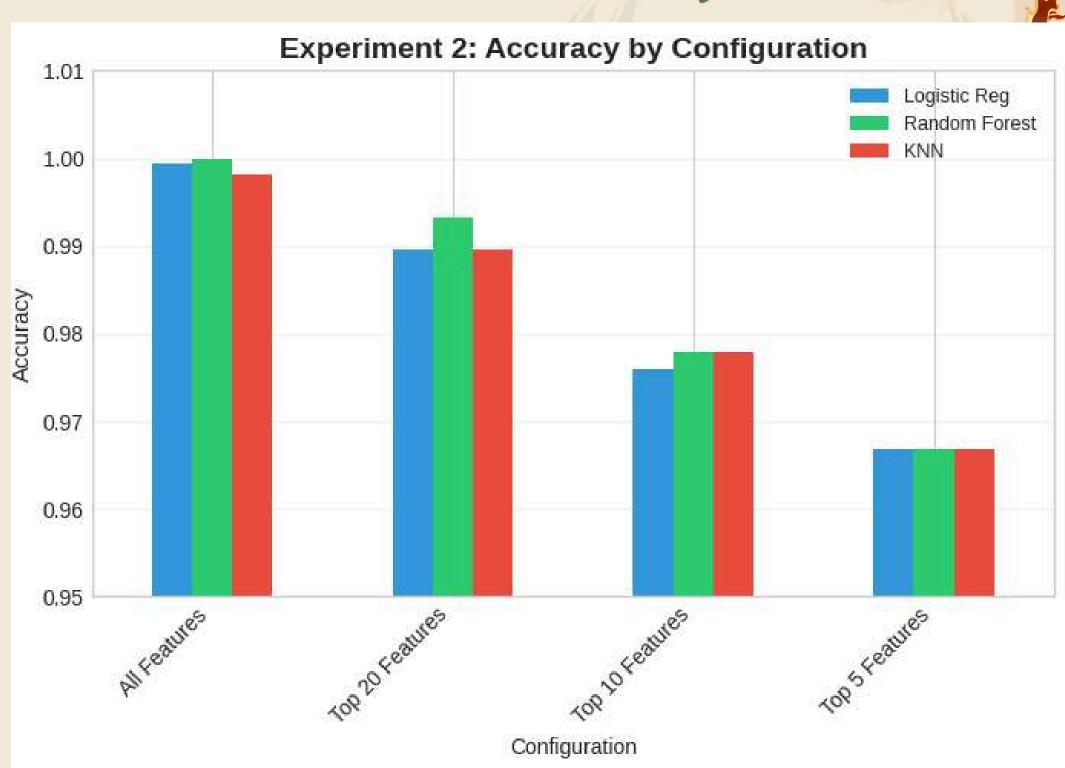




Configuration	Average Accuracy
All Features	0.999179
Top 20 Features	0.990769
Top 10 Features	0.977231
Top 5 Features	0.966769

- Reducing features causes significant accuracy loss (3.24%)
- All features contribute meaningfully to predictions
- Model complexity is necessary for optimal performance

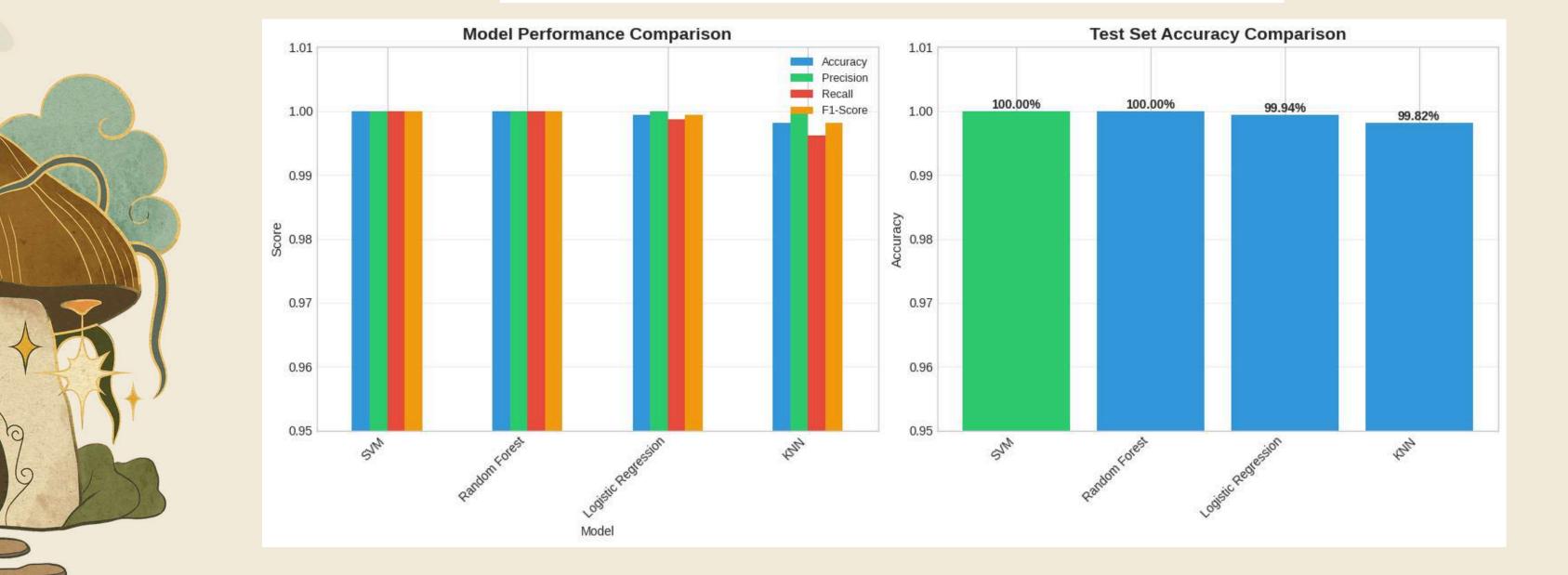
Final Decision: Use all features including the keeping of '?' in the stalk-root for production model to maintain accuracy and to keep the integrity of the original data.



Results & Discussions

TABLE III: Final Model Performance on Test Set

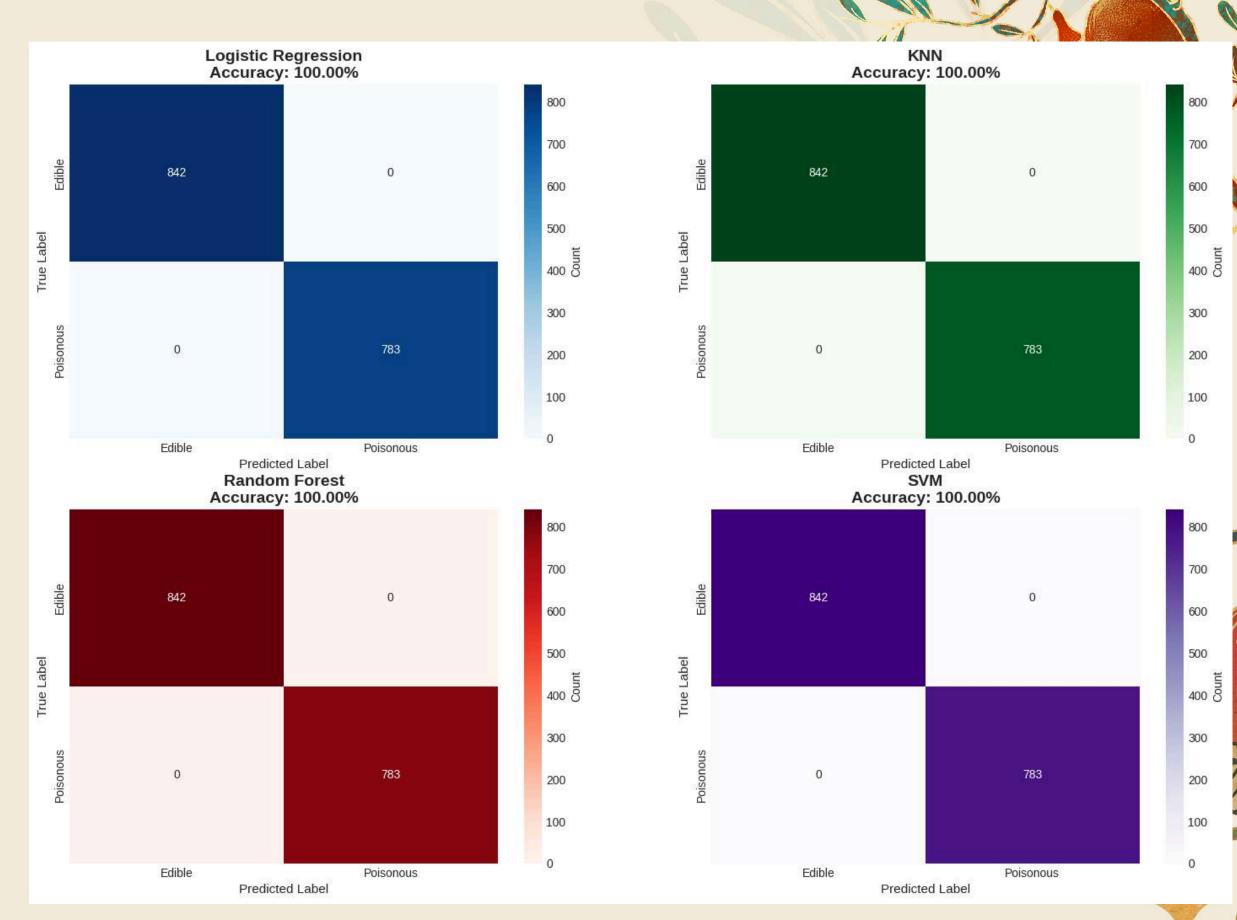
Model	Accuracy	Precision	Recall	F1-Score
LR	99.94%	100.0%	99.87%	99.94%
KNN	99.82%	100.0%	99.62%	99.81%
RF	100.0%	100.0%	100.0%	100.0%
SVM	100.0%	100.0%	100.0%	100.0%



Results & Discussions

Final Model Performance:

- Using all 117 features (Strategy A),
 both Random Forest and Support
 Vector Machine achieved perfect
 100% scores across all metrics
 (Accuracy, Precision, Recall, and
 F1-Score) on the held-out test set.
- Critically, as shown in the
 confusion matrices, there were
 zero false negatives (a poisonous
 mushroom classified as edible)
 and zero false positives.



Results & Discussions

Dominant Predictive Features:

- Feature importance analysis
 revealed that features related to
 'odor' are "overwhelmingly the
 most predictive".
- odor_n (odor=none) was the
 single most important feature,
 followed by odor_f (odor=foul).
 This aligns with anecdotal
 knowledge in mycology.

