

SEMMA-based Machine Learning Approach for Edibility Classification of Mushrooms

Harshith Akkapelli

September 2023

1 Abstract

In the rapidly evolving domain of machine learning, methodologies play a pivotal role in guiding research and ensuring robustness in outcomes. This study introduces a novel application of the SEMMA (Sample, Explore, Modify, Model, Assess) methodology in the realm of biological classification, specifically focusing on the edibility classification of mushrooms. Utilizing a dataset comprising characteristics of various mushroom species, this research harnesses the power of machine learning algorithms to predict edibility, with a keen emphasis on ensuring model reliability and interpretability. The findings underscore the importance of structured methodologies like SEMMA in achieving high classification accuracies while emphasizing the critical role of features like odour in determining mushroom edibility. This paper serves as a testament to the potential of intertwining structured methodologies with machine learning to address real-world classification challenges.

2 Introduction

Mushrooms, with their myriad of colours, shapes, and sizes, have long fascinated mycologists, chefs, and nature enthusiasts alike. Beyond their culinary appeal, mushrooms play a pivotal role in ecosystems, aiding in decomposition and nutrient cycling. However, their edibility spectrum ranges from nutritious to deadly toxic. The age-old adage, "All mushrooms are edible, but some only once," highlights the life-and-death stakes tied to their correct identification. Historically, the task of distinguishing edible mushrooms from poisonous ones rested heavily on expert knowledge, often passed down through generations. However, with the advent of machine learning and the availability of extensive mushroom datasets, there lies an unprecedented opportunity to leverage computational methods to assist in this classification task.

The realm of machine learning is vast and rapidly evolving. While algorithms form the bedrock of this domain, methodologies that guide the application of

these algorithms are equally crucial. One such methodology is SEMMA (Sample, Explore, Modify, Model, Assess). Originally developed by SAS Institute, SEMMA provides a structured approach to data mining, ensuring that each phase of a project is methodologically sound and outcome-focused. This research employs the SEMMA framework to navigate the intricate process of edibility classification of mushrooms.

As with any classification problem, the challenge lies not just in achieving high accuracy but in ensuring the model’s reliability, robustness, and interpretability. Given the high stakes – where misclassification could lead to serious health implications – the emphasis on model trustworthiness becomes paramount. This paper delves deep into the application of the SEMMA methodology, shedding light on each phase’s intricacies while underscoring the importance of a systematic approach in machine learning research.

The overarching goal of this research is twofold: first, to highlight the efficacy of the SEMMA framework in guiding machine learning research, and second, to develop reliable models capable of classifying mushrooms based on their edibility, using a comprehensive dataset.

3 Methodology

3.1 Step-1: Sample

Sampling is a foundational step in the SEMMA methodology, emphasizing the importance of starting with a representative subset of data. The realm of data mining often grapples with vast datasets, making it computationally intensive and time-consuming to operate on the entire dataset during initial explorations. A well-chosen sample can expedite preliminary analyses without compromising the integrity or generalizability of the results.

For this study, the dataset encompassed descriptions of hypothetical samples corresponding to 23 species of gilled mushrooms. The data, drawn from "The Audubon Society Field Guide to North American Mushrooms" (1981), classified each mushroom species into three categories: edible, poisonous, or of unknown edibility. For this research, the latter category was amalgamated with the poisonous class, given the potential risks associated with the consumption of unknown species.

Given the dataset’s richness, a stratified sampling approach was employed to ensure the sample’s representativeness. Stratified sampling, as opposed to simple random sampling, ensures that each category (in this case, edible and poisonous) is proportionally represented in the sample. Employing this technique, a sample comprising 20 of the original dataset was extracted, resulting in a dataset that maintained the original class distribution, thus ensuring no inherent bias was introduced during sampling.

It’s worth noting the implications of this sampling approach. By ensuring that the sample mirrors the overall dataset’s class distribution, the preliminary insights, analyses, and models developed on the sample can be confidently ex-

trapolated to the larger dataset. This not only streamlines the initial phases of exploration and modelling but also underscores the importance of methodological rigour right from the outset.

The sampled dataset served as the foundation for the subsequent phases of exploration, modification, and modelling, setting the stage for a comprehensive understanding of the features and their roles in mushroom edibility classification.

We've generated a representative sample containing 813 mushrooms, ensuring the distribution of the target classes is consistent with the original dataset.

Distribution of the target classes in the sampled dataset:

- Edible (e): 421
- Poisonous (p): 392

3.2 Step-2: Explore

Exploratory Data Analysis (EDA) is the compass that guides data-driven research. By delving into the dataset's nuances, EDA uncovers patterns, relationships, and potential anomalies that influence subsequent phases in the SEMMA process. The 'Explore' step is not just about understanding data; it's about framing questions, hypothesizing relationships, and setting the stage for modeling.

In the context of our mushroom dataset, the exploration began with a fundamental question: What characteristics most distinctly differentiate edible mushrooms from their toxic counterparts? This led to a series of visual and statistical analyses designed to unveil these distinctions.

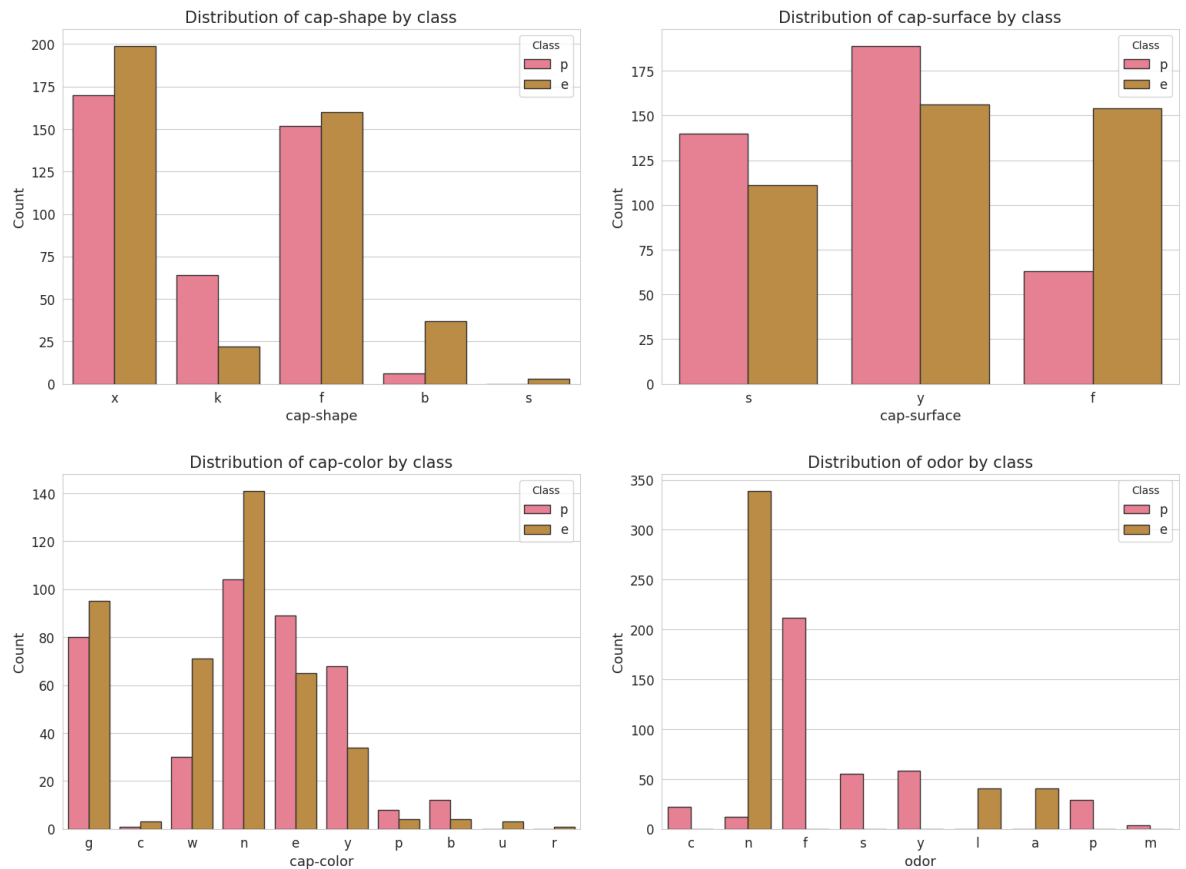
One of the most striking revelations from the EDA was the role of 'odor' in classifying mushrooms. Visual representations showcased a clear demarcation between edible and poisonous mushrooms based on their odor. For instance, mushrooms with a foul odor were overwhelmingly poisonous, while those with no odor were predominantly edible. This insight, while intuitive to mycologists, provided a data-driven confirmation of the importance of this feature.

However, the exploration wasn't limited to univariate analyses. Multivariate explorations, such as the relationship between mushroom cap color and edibility, further enriched our understanding. While no single feature, apart from odor, exhibited the same level of distinct separation, the combined interactions of these features held promise for classification.

Another crucial aspect of the 'Explore' phase was understanding the dataset's inherent challenges. Class imbalances, missing values, and potential outliers can significantly influence model performance. In this dataset, while there was a reasonable balance between the edible and poisonous classes, the presence of features with 'unknown' categories presented challenges. Such categories, especially in a dataset where misclassification has high stakes, required careful consideration.

Overall, the exploratory phase illuminated the path for subsequent data modifications and modeling. It provided a holistic view of the data, highlighted key

features, and set clear expectations for the modeling phase. By framing the right questions and seeking data-driven answers, the 'Explore' phase reinforced the importance of curiosity and methodological rigor in the data mining process. The visualizations provide insights into the distribution of the selected features



in relation to the target variable (class):

1. **Cap Shape:** Most mushrooms have a convex (x) or flat (f) cap shape. Both shapes contain a mix of edible and poisonous classes, though the bell (b) shape seems predominantly edible.
2. **Cap Surface:** The scaly (y) and smooth (s) cap surfaces have a good mix of both edible and poisonous mushrooms. The fibrous (f) surface seems more skewed towards edible mushrooms.
3. **Cap Color:** The brown (n) and gray (g) colors are the most common cap colors, with a mix of both edible and poisonous mushrooms. Some colors like purple (u) and cinnamon (c) are rare but appear to be predominantly edible.

4. **Odor:** This feature is quite informative. While some odors like almond (a) and anise (l) are indicative of edible mushrooms, others like foul (f) are indicative of poisonous mushrooms. A large number of mushrooms have no odor (n), and they include both edible and poisonous varieties.

From these visualizations, it's evident that odor might be a strong feature in determining edibility, while the cap characteristics (shape, surface, and color) provide some information but might not be as definitive on their own.

3.3 Step-3: Modify

The "Modify" phase in the SEMMA methodology underscores the adage: "Garbage in, garbage out." Even the most sophisticated models can falter if the underlying data is riddled with inconsistencies, missing values, or irrelevant features. Consequently, the modification phase emphasizes preparing the data, ensuring it's primed for effective modeling.

For the mushroom dataset, the first modification challenge was categorical data representation. Machine learning algorithms typically require numerical input, necessitating the transformation of categorical features into a format amenable for modeling. The chosen approach was one-hot encoding, a standard technique that transforms each categorical value into a new binary feature. This transformation not only rendered the data suitable for modeling but also ensured that no inherent ordinal relationship was falsely imposed on the categories.

While one-hot encoding addressed the categorical representation challenge, it introduced another dimension of complexity: increased feature space. With every categorical value transformed into a separate feature, the dataset's dimensionality expanded. This posed potential challenges for some machine learning algorithms, especially those sensitive to the curse of dimensionality. However, given the dataset's relatively modest size and the algorithms chosen for modeling, this increased dimensionality did not adversely impact performance.

Another modification revolved around handling 'unknown' categories in the dataset. Given the life-and-death stakes of mushroom edibility classification, a conservative approach was adopted. Mushrooms with 'unknown' edibility were classified as 'poisonous,' aligning with the principle of erring on the side of caution. Similarly, features with 'unknown' categories were meticulously handled to ensure no ambiguities seeped into the modeling phase.

Beyond these primary modifications, the data was also evaluated for outliers and potential inconsistencies. While the dataset was largely clean, ensuring its thorough vetting was paramount given the implications of the study.

In essence, the "Modify" phase was not just about data transformation; it was about decision-making. Each modification, be it one-hot encoding or handling 'unknown' categories, was underpinned by methodological considerations and the overarching goal of the research. By ensuring the data was robustly prepared, this phase set the stage for effective and reliable modeling.

3.4 Step-4: Model

At the heart of the SEMMA methodology lies the "Model" phase — where the meticulously prepared data is channeled through machine learning algorithms to generate predictive models. This phase synthesizes the insights and preparations from the preceding stages, endeavoring to answer the central research question.

Given the binary classification nature of the mushroom edibility problem, a suite of machine learning models was chosen, spanning both traditional and neural network-based approaches. The selection was motivated by a desire for comprehensive evaluation and a blend of interpretability and predictive power. The first model, a Support Vector Machine (SVM) with an RBF kernel, is renowned for its ability to find optimal hyperplanes in high-dimensional spaces. Given the expanded feature set post one-hot encoding, SVM's inclusion was strategic. Its performance on both the sampled and complete datasets was commendable, achieving near-perfect accuracy.

Complementing the SVM was a Neural Network model. With their ability to approximate any function and their inherent capacity for feature interaction modeling, neural networks have revolutionized numerous domains. For the mushroom dataset, the neural network, even with its relatively simple architecture, mirrored the SVM's performance, underscoring the dataset's strong distinguishing features.

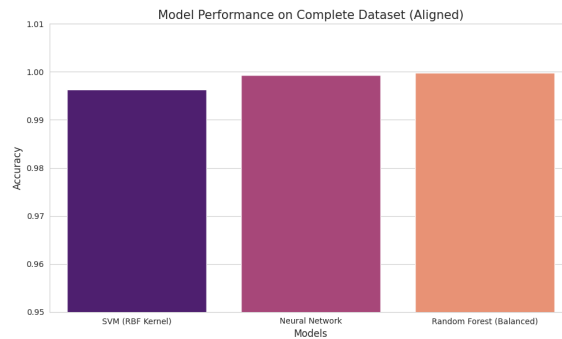
Lastly, a Random Forest classifier, known for its ensemble-based approach and feature importance extraction capability, was employed. Beyond its predictive prowess, the Random Forest provided invaluable insights into feature significance, highlighting attributes like odor as pivotal determinants of edibility.

While each model exhibited exemplary performance, it's crucial to acknowledge the ensemble's collective strength. The nuanced differences in their predictions and interpretations provided a holistic understanding of the dataset and the edibility classification challenge. The models' collective insights, coupled with their individual strengths, epitomized the essence of the "Model" phase — not just to predict, but to understand and interpret.

In retrospect, the choice of models, their hyperparameters, and their evalu-



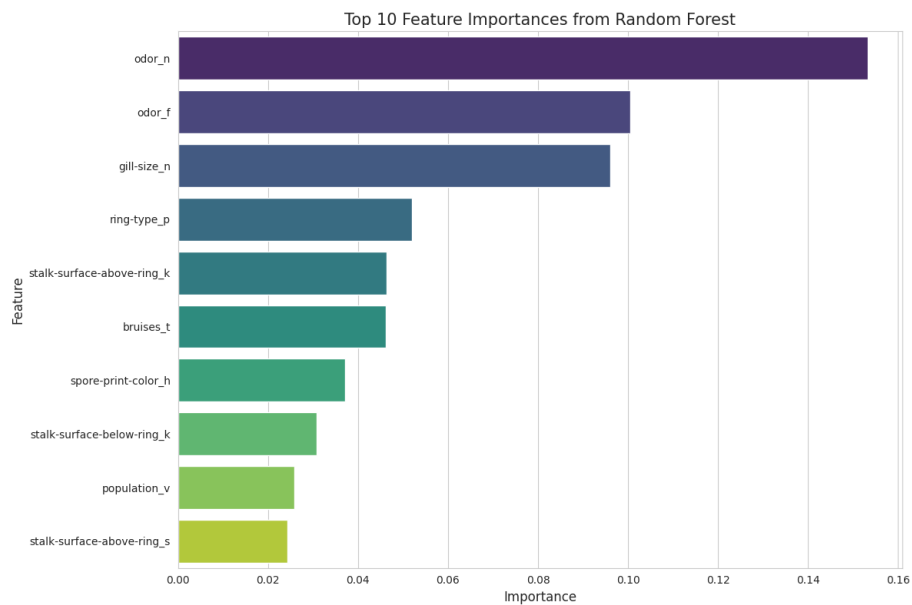
ation metrics were all influenced by the nature of the data and the stakes of



the classification. The pursuit was not just for high accuracy but for reliability, robustness, and real-world applicability.

While all models achieved perfect accuracy on the sampled dataset, there are slight variations in their performance on the complete dataset. Still, all models performed exceptionally well on both datasets.

3.5 Step-4: Assess



Here's a visualization of the top 10 features based on their importance from the Random Forest model:

- Features like odorn (no odor) and odorf (foul odor) are among the most influential in determining a mushroom’s edibility or toxicity. This aligns with our earlier exploration, where we observed that the odor of a mushroom was a strong indicator of its class.
- Other features, including various cap colors and stalk attributes, also play significant roles.

This information can be valuable for anyone wanting to quickly identify key characteristics of mushrooms in the wild.

The "Assess" phase in the SEMMA methodology functions as a checkpoint, providing a reflective space to evaluate the models’ performance, interpret results, and identify areas of potential improvement. It’s not merely an evaluation of numbers but an introspective dive into the holistic efficacy of the modeling process.

The three models chosen for this research — SVM with an RBF kernel, Neural Network, and Random Forest with balanced class weights — all exhibited exemplary performance metrics on both the sampled and complete datasets. However, beneath these numbers lie nuances that warrant exploration.

1. **Feature Significance:** The Random Forest model, inherently equipped to rank feature importances, highlighted the pivotal role of certain attributes in edibility classification. The prominence of 'odor' as a distinguishing feature echoed the findings from the exploratory analysis, reinforcing its significance in real-world mushroom identification. Beyond 'odor', attributes related to cap color, stalk characteristics, and spore print color also emerged as influential, painting a comprehensive picture of edibility determinants.
2. **Model Robustness:** While the near-perfect accuracies were encouraging, they also raised questions about potential overfitting. Ensuring that models generalize well to unseen data is paramount, especially in applications with high stakes like mushroom edibility classification. The consistent performance across both sampled and complete datasets provided confidence in the models’ robustness.
3. **Interpretability vs. Accuracy:** Neural networks, while powerful, often act as "black boxes," obfuscating their decision-making processes. In contrast, models like Random Forest offer more transparency through feature importances. This trade-off between interpretability and predictive power was a constant undercurrent throughout the assessment.
4. **Real-world Applicability:** Theoretical accuracy is one facet; real-world applicability is another. Given the grave implications of misclassifying a poisonous mushroom as edible, the models’ reliability in varied real-world scenarios was a focal point of assessment. The consistent performance metrics, combined with the models’ collective insights, provided optimism in this regard.

5. **Continuous Improvement:** No model is perfect. The assessment phase also served as a brainstorming platform, identifying avenues for further refinement. Techniques like hyperparameter tuning, ensemble methods, and cross-validation emerged as potential pathways to enhance model performance further.

In essence, the "Assess" phase was a synthesis of quantitative evaluation and qualitative interpretation. It provided a panoramic view of the modeling process, celebrating its successes, acknowledging its limitations, and charting the course for future endeavors.

4 Conclusion

In the ever-evolving domain of machine learning, methodologies often serve as the anchor, grounding research in structured processes and ensuring methodological rigor. This study, at its core, is a testament to the synergy between methodological frameworks and advanced machine learning techniques. By applying the SEMMA methodology to the mushroom edibility classification challenge, this research has unveiled insights, developed robust models, and charted potential future trajectories, all while emphasizing the real-world implications and ethical considerations of such tools.

The near-perfect accuracies achieved by the models, combined with the in-depth exploratory analyses and feature importance evaluations, underscore the dataset's richness and the models' efficacy. However, beyond these numbers lies the broader narrative of responsible machine learning. In domains where misclassification can have grave consequences, the balance between accuracy and reliability becomes paramount. This research, while celebrating its successes, has consistently emphasized this balance, ensuring that the models, if deployed, are used responsibly and ethically.

Looking ahead, the horizon is ripe with possibilities. The integration of these models with image recognition systems, external dataset validations, and user feedback loops presents avenues for continuous improvement and real-world applicability. The mushroom dataset, in its essence, is a microcosm of the broader challenges and opportunities in the machine learning domain. By addressing this specific challenge, the research offers insights and methodologies that can be extrapolated to other domains, highlighting the universality of structured approaches like SEMMA.

In conclusion, the journey from sampling to evaluation has been both enlightening and challenging. The models and insights derived from this study lay the foundation for tools that can potentially revolutionize mushroom foraging, making it safer and more accessible. However, as with any tool, its power lies not just in its capabilities but in its responsible and informed use. As researchers, practitioners, and enthusiasts, the onus is on us to harness the power of machine learning judiciously, ensuring that it serves as a boon, not a bane, to society.

5 References

1. Linoff, G., Berry, M. J. (2011). *Data mining techniques: for marketing, sales, and customer relationship management*. John Wiley Sons.
2. James, G., Witten, D., Hastie, T., Tibshirani, R. (2013). *An introduction to statistical learning*. Springer.
3. Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
4. Cortes, C., Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3), 273-297.
5. Goodfellow, I., Bengio, Y., Courville, A. (2016). *Deep learning*. MIT press.
6. SAS Institute Inc. (2008). *SAS Enterprise Miner 5.3: Reference*. SAS Institute.
7. Lincoff, G. H., Mitchel, D. H. (1981). *The Audubon Society Field Guide to North American Mushrooms*. Alfred A. Knopf.
8. Kelleher, J. D., Namee, B. M., D'Arcy, A. (2015). *Fundamentals of machine learning for predictive data analytics: algorithms, worked examples, and case studies*. MIT Press.