

Classificação de Cogumelos (Mushroom UCI) usando modelos de aprendizado de máquina: proposta experimental e comparação com a literatura

Matheus Henrique, João Guilherme, Renan Aprígio, Arthur Tomé

Curso Integrado de Desenvolvimento de Sistemas – Instituto Federal de Educação, Ciência e Tecnologia de Pernambuco (IFPE - Campus Jaboatão dos Guararapes)

jglf@discente.ifpe.edu.br, atx@discente.ifpe.edu.br, radm@discente.ifpe.edu.br,
mham4@discente.ifpe.edu.br

Orientador: Luciano de Souza

Resumo.

Este trabalho apresenta uma solução experimental para a tarefa de classificação binária (comestível vs. venenoso) sobre o dataset Mushroom (UCI). Implementamos um pipeline reproduzível de pré-processamento (One-Hot Encoding, escalonamento quando necessário) e avaliamos sete modelos: GaussianNB, Logistic Regression, Decision Tree, KNN, SVM, RandomForest e XGBoost. As medições foram obtidas via validação estratificada (Stratified K-Fold) e as métricas comparadas com resultados reportados na literatura recente. Observou-se que modelos baseados em árvore/ensemble (Decision Tree, RandomForest, XGBoost) atingiram precisão perfeita (1.000), enquanto modelos probabilísticos e lineares apresentaram desempenho inferior (GaussianNB ≈ 0.922 ; LogisticRegression ≈ 0.948). Discutimos por que esses resultados concordam com trabalhos anteriores, implicações sobre viés/overfitting e sugestões para validação adicional.

Abstract.

This paper reports an experimental study on the Mushroom (UCI) dataset. We implemented a reproducible pipeline and evaluated seven classifiers (GaussianNB, Logistic Regression, Decision Tree, KNN, SVM, RandomForest, XGBoost) using stratified cross-validation. Tree-based and ensemble methods reached perfect precision (1.000). We compare our findings to the literature and discuss dataset characteristics that explain the results and propose further robustness checks.

1. Introdução

A identificação automática de cogumelos como comestíveis ou venenosos tem sido um problema padrão em estudos introdutórios e comparativos de aprendizado de máquina devido à disponibilidade do Mushroom dataset (UCI). Diversos trabalhos relatam desempenhos muito altos, especialmente para classificadores baseados em árvores e ensembles. Este trabalho apresenta uma solução experimental completa (pré-processamento, pipeline, sete modelos), reporta métricas relevantes e compara os resultados empíricos obtidos com os achados da literatura recente.

2. Trabalhos Relacionados

Estudos comparativos e relatórios aplicados sobre o Mushroom frequentemente mostram que modelos de árvore e ensembles alcançam acurácias em torno de 98–100%; modelos lineares e probabilísticos normalmente ficam abaixo desses valores (90–96%); KNN e SVM costumam apresentar resultados próximos aos ensembles quando o pré-processamento é adequado.

3. Materiais e Métodos

3.1. Dataset

Usou-se o dataset Mushroom da UCI (8.124 instâncias, 22 atributos categóricos + rótulo class), com classes edible e poisonous. Atributos notáveis incluem odor, spore-print-color e gill-size.

3.2. Pré-processamento

Leitura do CSV bruto; tratamento de valores faltantes (dataset padrão não tem valores ausentes significativos); One-Hot Encoding para todas as variáveis categóricas; escalonamento aplicado quando o classificador exige (SVM, LogisticRegression); validação com StratifiedKFold (n_splits=5, shuffle=True, random_state=42).

3.3. Modelos avaliados

Gaussian Naive Bayes; Logistic Regression; Decision Tree; K-Nearest Neighbors; Support Vector Machine; Random Forest (n_estimators=200); XGBoost.

3.4. Métricas

precisão, Precisão (macro), Recall (macro) e F1-Score (macro). Métricas calculadas a partir das pontuações médias em validação cruzada estratificada (5 folds).

4. Resultados

Os resultados médios obtidos nos experimentos foram os seguintes (ordenados pelos valores fornecidos):

Tabela 1. Resultados médios por modelo (Precisão, Recall e F1-Score)

Modelo	precisã o	Precisã o	Recall	F1-Sco re
Decisio nTree	1.0000 00	1.0000 00	1.0000 00	1.0000 00
Rando mForest	1.0000 00	1.0000 00	1.0000 00	1.0000 00
XGBoo st	1.0000 00	1.0000 00	1.0000 00	1.0000 00
KNN	0.9963 08	0.9923 86	1.0000 00	0.9961 78
SVM	0.9926 15	0.9987 05	0.9859 34	0.9922 78
Logisti cRegression	0.9476 92	0.9439 49	0.9475 70	0.9457 56
Gaussia nNB	0.9218 46	0.9098 87	0.9296 68	0.9196 71

5. Discussão — Comparação com a Literatura

As principais concordâncias com a literatura: árvores e ensembles alcançaram precisão perfeita, como relatado em diversos estudos; KNN e SVM obtiveram desempenho muito alto ($\geq 99\%$) quando o pré-processamento foi adequado; logistic regression e gaussianNB apresentaram desempenho inferior ($\approx 92\text{--}95\%$).

Interpretação crítica: o dataset é altamente solucionável, o que explica os resultados perfeitos; recomenda-se cautela quanto a overfitting e verificação de data leakage. Sugestões para robustez: validação externa, remoção de atributos altamente informativos (ex.: odor), análise de importância de features (SHAP), experimentos com ruído sintético.

6. Conclusões

Implementamos um pipeline reprodutível e avaliamos sete classificadores. Decision Tree, RandomForest e XGBoost atingiram 100% de precisão, corroborando a literatura. KNN e SVM performaram $>99\%$ e modelos lineares/probabilísticos ficaram em 92–95%. Recomendamos análises adicionais para testar robustez e generalização.

7. Referências (seleção)

UCI Machine Learning Repository — Mushroom Dataset.

Wagner, D.; Heider, D.; Hattab, G. (2021). Mushroom data creation, curation and simulation to support classification tasks.

Tarawneh, O.; Tarawneh, M.; Sharrab, Y.; Husni, M. (2022). Mushroom classification using machine-learning techniques.

Paudel, S.; Bhatta, R. (2022). Comparative study: Random Forest and other classifiers on Mushroom dataset.