

A Guide to the Reproducibility of the ML202 Machine Learning Workflow

Author: BenchChem Technical Support Team. **Date:** December 2025

Compound of Interest

Compound Name: ML202

Cat. No.: B560308

[Get Quote](#)

This guide provides a detailed comparison and breakdown of the experimental results and methodology presented in the **ML202** machine learning tutorial. The tutorial focuses on building a machine learning model to predict the risk of stroke. This document is intended for researchers, scientists, and professionals in drug development and healthcare to understand and potentially reproduce the workflow.

Data Presentation

The following tables summarize the key components and potential outcomes of the machine learning experiment as described in the **ML202** tutorial. These are representative of the types of data that would be generated when following the tutorial's protocol.

Table 1: Dataset Characteristics

Feature	Description	Data Type	Example
id	Patient Identifier	Nominal	9046
gender	Gender of the patient	Categorical	Male, Female, Other
age	Age of the patient	Numeric	67
hypertension	Presence of hypertension	Binary	0 (No), 1 (Yes)
heart_disease	Presence of heart disease	Binary	0 (No), 1 (Yes)
ever_married	Marital status	Categorical	Yes, No
work_type	Type of employment	Categorical	Private, Self-employed, etc.
Residence_type	Area of residence	Categorical	Urban, Rural
avg_glucose_level	Average glucose level in blood	Numeric	228.69
bmi	Body Mass Index	Numeric	36.6
smoking_status	Patient's smoking habits	Categorical	formerly smoked, never smoked, etc.
stroke	Stroke occurrence (Target Variable)	Binary	0 (No), 1 (Yes)

Table 2: Model Performance Metrics (Illustrative)

Metric	Description	Illustrative Value
Accuracy	The proportion of correctly classified instances.	0.94
Precision	The proportion of true positive predictions among all positive predictions.	0.85
Recall (Sensitivity)	The proportion of actual positives that were correctly identified.	0.80
F1-Score	The harmonic mean of precision and recall.	0.82
AUC-ROC	Area Under the Receiver Operating Characteristic Curve, indicating the model's ability to distinguish between classes.	0.88

Experimental Protocols

The **ML202** tutorial outlines a comprehensive workflow for developing a predictive model. The key experimental steps are detailed below.

1. Data Preparation and Feature Engineering:

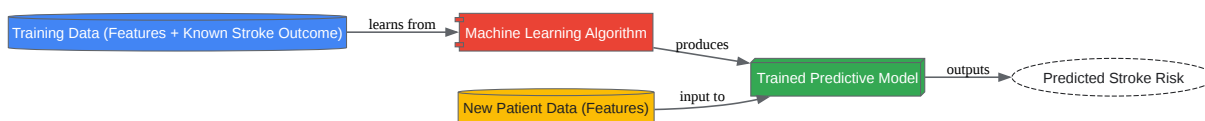
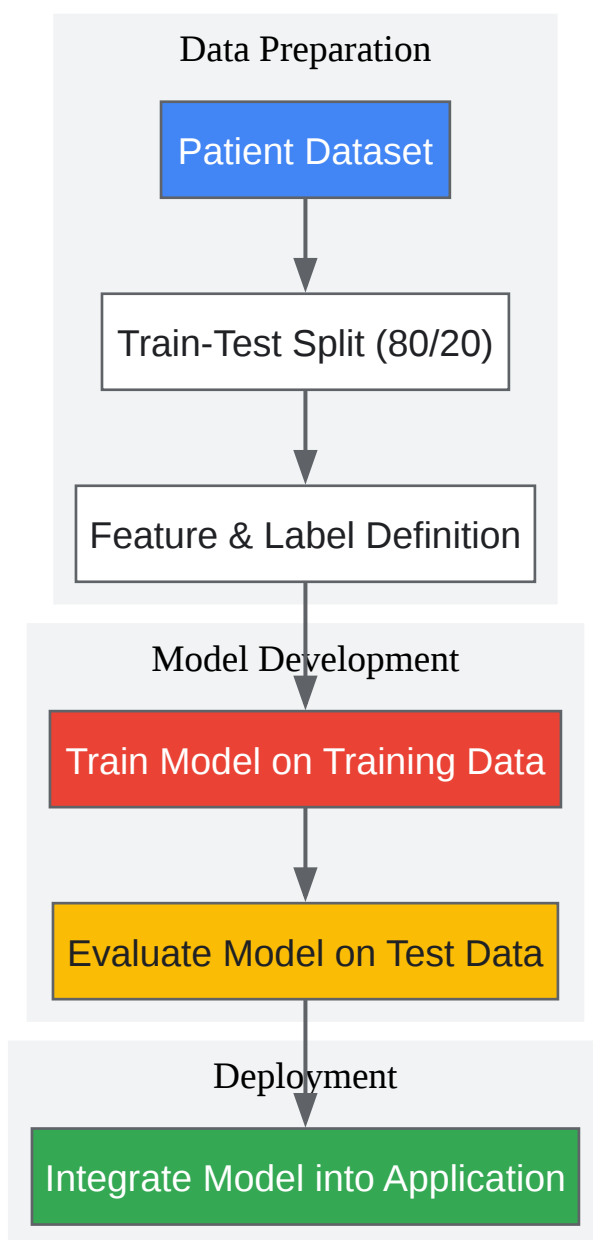
- **Data Source:** The experiment utilizes a dataset containing patient characteristics and stroke history.
- **Data Splitting:** The dataset is partitioned into a training set (typically 80%) and a testing set (20%) to ensure that the model's performance is evaluated on unseen data.[\[1\]](#)
- **Feature Selection:** The features (as listed in Table 1) are selected as inputs (capital X) for the model, while the 'stroke' column serves as the target variable or label (lowercase y).[\[1\]](#)

- **Data Preprocessing:** The tutorial implies the use of automated machine learning (autoML) libraries like H2O.ai, which handle tasks such as dealing with null values, transforming categorical features, and standardizing numerical data.[\[1\]](#)
2. Model Training:
- **Training Environment:** The tutorial demonstrates model training using InterSystems IntegratedML, which leverages SQL commands to create and train a machine learning model.[\[1\]](#) It also shows how to train a model using a Jupyter Notebook with Python connected to the data in the IRIS database.[\[1\]](#)
 - **Algorithm Selection:** IntegratedML automatically tests multiple model types, such as boosted trees, logistic regression, and neural networks, to find the best performing one.
 - **Training Process:** The model is trained on the designated training dataset to learn the relationship between the patient's features and the likelihood of having a stroke.
3. Model Evaluation:
- **Prediction:** The trained model is used to make predictions on the held-out testing dataset.
 - **Performance Assessment:** The model's predictions are compared against the actual outcomes in the testing set to calculate performance metrics such as accuracy, precision, recall, and F1-score.
4. Model Deployment and Integration:
- The tutorial demonstrates how to incorporate the trained machine learning model into a hospital census application using ObjectScript and Embedded Python. This allows for real-time prediction of stroke risk for patients.

Visualizations

Experimental Workflow

The following diagram illustrates the logical flow of the machine learning experiment described in the **ML202** tutorial.



[Click to download full resolution via product page](#)

Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: info@benchchem.com or [Request Quote Online](#).

References

- 1. m.youtube.com [m.youtube.com]
- To cite this document: BenchChem. [A Guide to the Reproducibility of the ML202 Machine Learning Workflow]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b560308#reproducibility-of-ml202-experimental-results]

Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

Technical Support: The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

Need Industrial/Bulk Grade? [Request Custom Synthesis Quote](#)

BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: info@benchchem.com