

Train-Test Split & Data Partitioning in ML: Complete Guide

1. What is Train-Test Split?

Train-Test Split is the standard procedure for dividing your dataset into "training" and "testing" subsets.

- **Training set:** Used by the ML algorithm to learn patterns.
- **Test set:** Used to evaluate how well the model performs on *unseen* data—shows how well it generalizes[1][2][3].

This fundamental technique ensures that models are tested on data they haven't seen during training, providing an honest assessment of their real-world performance.

2. Why Perform a Train-Test Split?

Train-test splitting is essential for several critical reasons:

- **Prevents overfitting:** Stops the model from learning the dataset too perfectly, which would cause it to fail on new data.
- **Simulates real-world performance:** By withholding test data during training, you get realistic performance estimates.
- **Ensures objective evaluation:** Enables unbiased measurement of model accuracy, precision, recall, and other metrics[2][3].

Without proper train-test splitting, you cannot trust your model's performance metrics, as they would be artificially inflated by memorization rather than genuine learning.

3. Common Ratios

Industry-standard splitting ratios include:

- **80:20** - Most popular for train:test splits
- **70:30** - Common when you have sufficient data
- **75:25** - Balanced middle ground
- **70:15:15** - Three-way split: train, validation, test
- **60:20:20** - Three-way split with larger validation set

The choice depends on your dataset size—larger datasets can afford smaller test percentages while maintaining statistical significance[4][5].

4. Key Industry Terms

Understanding these terms is crucial for professional ML work:

Training Set

The portion of data used to fit and learn the model parameters. This is where the algorithm identifies patterns and relationships.

Test Set

Data used only to evaluate the trained model. This set must remain completely unseen during training to provide honest performance metrics.

Validation Set

Optional subset used during model tuning, such as selecting optimal hyperparameters. Helps prevent overfitting to the test set through repeated evaluation.

Cross-Validation

Technique to further partition training data into multiple folds, rotating which subset acts as validation. Common methods include k-fold CV and stratified k-fold CV, which provide more robust performance estimates[3][6].

Random State / Seed

Parameter ensuring reproducibility of splits. Using the same random state guarantees identical splits across different runs.

Stratified Split

Maintains class distributions balanced across splits, critical for imbalanced datasets. Ensures minority classes are adequately represented in both train and test sets[7][8].

5. How to Do It in Python

Basic Usage

```
from sklearn.model_selection import train_test_split  
  
X_train, X_test, y_train, y_test = train_test_split(  
    X, y, test_size=0.2, random_state=42  
)
```

Parameters:

- `test_size`: Controls the percentage assigned to the test set (0.2 = 20%)
- `random_state`: Ensures repeated splits are identical (use 42 by convention)

For Imbalanced Classes

```
X_train, X_test, y_train, y_test = train_test_split(  
    X, y, test_size=0.2, stratify=y, random_state=42  
)
```

The `stratify` parameter ensures class proportions are maintained in both sets, preventing scenarios where rare classes might be missing from train or test sets.

6. Advanced Splitting: Validation & Cross-Validation

Train/Validation/Test Split

First split: separate test set (20%)

```
X_temp, X_test, y_temp, y_test = train_test_split(  
X, y, test_size=0.2, random_state=42  
)
```

Second split: separate validation from training (25% of 80% = 20% overall)

```
X_train, X_val, y_train, y_val = train_test_split(  
X_temp, y_temp, test_size=0.25, random_state=42  
)
```

Result: 60% train, 20% validation, 20% test

K-Fold Cross-Validation

```
from sklearn.model_selection import KFold  
  
kf = KFold(n_splits=5, shuffle=True, random_state=42)  
for train_idx, test_idx in kf.split(X):  
    X_train, X_test = X[train_idx], X[test_idx]  
    y_train, y_test = y[train_idx], y[test_idx]  
    # Train and evaluate model
```

Stratified K-Fold

```
from sklearn.model_selection import StratifiedKFold  
  
skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)  
for train_idx, test_idx in skf.split(X, y):  
    X_train, X_test = X[train_idx], X[test_idx]  
    y_train, y_test = y[train_idx], y[test_idx]  
    # Ensures class proportions in each fold
```

Stratified K-Fold is particularly important for classification tasks with imbalanced classes[6][8].

7. Industry Best Practices

Shuffle Data

Always shuffle before splitting to ensure randomness, unless working with time series data. Prevents systematic biases from data ordering.

```
train_test_split(X, y, test_size=0.2, shuffle=True, random_state=42)
```

For Time Series: Sequential Split

NEVER randomly split time series data. Use sequential splitting to respect temporal ordering:

First 80% for training, last 20% for testing

```
split_point = int(0.8 * len(data))
train_data = data[:split_point]
test_data = data[split_point:]
```

Stratify for Classification

For classification problems, especially with imbalanced classes, always use stratification to maintain class proportions.

Keep Splits Consistent

Use the same split for all model comparison experiments to ensure fair comparisons. Save train/test indices for reproducibility[9][10].

Don't Peek at Test Data

Never tune your model after seeing test set results. Use the validation set for hyperparameter tuning, and reserve the test set for final evaluation only.

8. Pitfalls & Extra Concepts

Overfitting

Evaluating on training data produces artificially high scores. Always use held-out test data for honest performance assessment.

Data Leakage

Ensure no information about the test set "leaks" into the training phase. Common sources:

- Feature engineering using statistics from entire dataset
- Preprocessing fitted on combined train+test data
- Using future information in time series

Imbalanced Data

Without stratification, classes may be distributed unevenly across sets. A minority class might be completely absent from the test set, making evaluation impossible.

Small Dataset Considerations

For very small datasets ($n < 1000$), consider:

- Using cross-validation instead of simple train-test split
- Leave-one-out cross-validation for extremely small datasets
- Bootstrapping methods for variance estimation

9. Real Industry Applications

Finance: Fraud Detection

Ensure the "fraud" class is present in both train and test sets via stratified splitting. Fraud is typically rare (1-5% of transactions), making stratification critical[11].

Healthcare: Disease Prediction

Split patient records carefully to avoid data leakage. Ensure multiple records from the same patient don't appear in both train and test sets.

Customer Analytics: Churn Prediction

Careful validation prevents overzealous tuning. Use time-based splits if you want to predict future churn based on historical data.

E-commerce: Product Recommendations

Time-aware splits ensure you're predicting future purchases based on past behavior, not the reverse.

10. Summary Table

Term	Use in ML	Typical Ratio
Training Set	Model learning	60–80%
Validation Set	Tuning/hyperparameters	10–20%
Test Set	Final evaluation	10–20%

Table 1: Standard data partition ratios in machine learning workflows

11. Complete Example: Full Workflow

```
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report
import pandas as pd
```

Load data

```
df = pd.read_csv('data.csv')
X = df.drop('target', axis=1)
y = df['target']
```

Split 1: Separate test set (20%)

```
X_temp, X_test, y_temp, y_test = train_test_split(
    X, y, test_size=0.2, stratify=y, random_state=42
)
```

Split 2: Separate validation set (20% of original)

```
X_train, X_val, y_train, y_val = train_test_split(
    X_temp, y_temp, test_size=0.25, stratify=y_temp, random_state=42
)
```

Train model

```
model = RandomForestClassifier(random_state=42)
model.fit(X_train, y_train)
```

Validate and tune

```
val_score = model.score(X_val, y_val)
print(f"Validation Accuracy: {val_score:.3f}")
```

Cross-validation on training set for robust estimation

```
cv_scores = cross_val_score(model, X_train, y_train, cv=5, scoring='accuracy')
print(f"CV Accuracy: {cv_scores.mean():.3f} (+/- {cv_scores.std():.3f})")
```

Final test evaluation (only once!)

```
test_score = model.score(X_test, y_test)
y_pred = model.predict(X_test)
print(f"\nTest Accuracy: {test_score:.3f}")
print("\nClassification Report:")
print(classification_report(y_test, y_pred))
```

References

- [1] Built In. (2025). Train Test Split: What it Means and How to Use It. <https://builtin.com/data-science/train-test-split>
- [2] Machine Learning Mastery. (2020). Train-Test Split for Evaluating Machine Learning Algorithms. <https://www.machinelearningmastery.com/train-test-split-for-evaluating-machine-learning-algorithms/>
- [3] GeeksforGeeks. (2020). Splitting Data for Machine Learning Models. <https://www.geeksforgeeks.org/machine-learning/splitting-data-for-machine-learning-models/>
- [4] Encord. (2025). Training, Validation, Test Split for Machine Learning Datasets. <https://encord.com/blog/train-val-test-split/>
- [5] V7 Labs. (2024). Train Test Validation Split: How To & Best Practices. <https://www.v7labs.com/blog/train-validation-test-set>
- [6] Real Python. (2025). Split Your Dataset With scikit-learn's train_test_split(). <https://realpython.com/train-test-split-python-data/>
- [7] scikit-learn. (2008). train_test_split documentation. https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html
- [8] GeeksforGeeks. (2022). How To Do Train Test Split Using Sklearn In Python. <https://www.geeksforgeeks.org/machine-learning/how-to-do-train-test-split-using-sklearn-in-python/>
- [9] Stanford CS230. Splitting into train, dev and test sets. <https://cs230.stanford.edu/blog/split/>
- [10] KDnuggets. (2020). Dataset Splitting Best Practices in Python. <https://www.kdnuggets.com/2020/05/dataset-splitting-best-practices-python.html>
- [11] Google Developers. (2025). Dividing the original dataset. <https://developers.google.com/machine-learning/crash-course/overfitting/dividing-datasets>