# Machine Learning Fundamentals for Data Science: Algorithms, Evaluation, and Applications

Shubneet [1], Anushka Raj Yadav [2], Navjot Singh Talwandi [3]

[1,2,3*]Department of Computer Science, Chandigarh University, Gharuan, Mohali, 140413, Punjab, India.

Contributing authors: jeetshubneet27@gmail.com; ay462744@gmail.com; navjot.e17908@cumail.in;

## Abstract

This chapter introduces ML fundamentals, focusing on two primary paradigms: **supervised learning**, where models predict labeled outcomes (e.g., classification, regression), and **unsupervised learning**, which identifies hidden structures in unlabeled data (e.g., clustering, dimensionality reduction). Key algorithms such as decision trees, support vector machines (SVM), k-nearest neighbors (k-NN), and Naive Bayes are explored, alongside their practical implementation in Python using libraries like Scikit-learn. The chapter emphasizes **model training** workflows, including data splitting (training/test/validation sets) and **cross-validation** to mitigate overfitting-a critical challenge where models memorize noise instead of learning generalizable patterns. Evaluation metrics such as accuracy, precision, recall, and ROC-AUC are discussed to assess model performance rigorously. Real-world applications span healthcare diagnostics, fraud detection, and recommendation systems, demonstrating ML's transformative potential. Practical examples illustrate Python code for algorithm deployment, hyperparameter tuning, and visualization of results. By balancing theoretical foundations with hands-on techniques, this chapter equips readers to build, validate, and deploy robust ML models while addressing common pitfalls like underfitting and bias-variance trade-offs [1, 2].

**Keywords:** supervised learning, unsupervised learning, classification, overfitting, model validation

# 1 Introduction

The evolution of machine learning (ML) has reshaped modern industries, transitioning from theoretical concepts in the mid-20th century to indispensable tools driving innovation in the 2020s. Enabled by advances in computational power, big data, and algorithmic sophistication, ML now underpins decision-making across finance, healthcare, e-commerce, and manufacturing. Between 2010 and 2025, ML adoption grew exponentially, with global market size projected to exceed $500 billion by 2025, fueled by applications like fraud detection systems reducing financial losses by 60% and recommendation engines boosting e-commerce revenue by 35% [3, 4].

## The Evolution of Machine Learning

Early rule-based systems gave way to adaptive models capable of learning from data, driven by three key developments:

1. **Algorithmic breakthroughs**: From basic linear regression (1950s) to deep neural networks (2010s)
2. **Data proliferation**: 90% of today's data was generated post-2010, enabling complex pattern recognition
3. **MLOps maturity**: Automated pipelines for deploying, monitoring, and updating production models

In finance, ML detects fraudulent transactions in real-time by analyzing spending patterns across 200+ features, reducing false positives by 45% compared to traditional systems [5]. Healthcare leverages ML for predictive diagnostics, with models analyzing genomic data and medical imaging to identify diseases like cancer 18 months earlier than conventional methods.

## Core Industry Applications

- **Fraud Detection**: ML models process millions of transactions per second, identifying anomalies through supervised learning (labeled fraud patterns) and unsupervised techniques (novel attack detection).
- **Recommendation Systems**: Neural networks power personalized suggestions in streaming (Netflix) and e-commerce (Amazon), increasing user engagement by 50%.
- **Predictive Maintenance**: Manufacturing ML models analyze IoT sensor data to forecast equipment failures with 92% accuracy.
- **Drug Discovery**: Generative ML accelerates pharmaceutical R&D, reducing development timelines from 10 years to 2–3 years.

## Chapter Outline

This chapter systematically explores machine learning fundamentals through:

- Supervised vs. unsupervised learning paradigms

- Core algorithms: Decision Trees, SVM, k-NN, and Naive Bayes
- Model training workflows and evaluation metrics (precision, recall, F1-score)
- Overfitting prevention through regularization and cross-validation
- Python/R implementations for real-world datasets
- Industry case studies: Fraud detection, recommendation engines, healthcare analytics
- Hands-on exercises with financial and e-commerce datasets

As ML becomes integral to organizational strategy, its ability to transform raw data into actionable insights continues to redefine competitive landscapes. The following sections provide both theoretical foundations and practical frameworks for deploying ML solutions responsibly and effectively.

# 2 Supervised vs. Unsupervised Learning

Machine learning encompasses two fundamental paradigms: supervised and unsupervised learning. These approaches differ in their data requirements, objectives, and applications, shaping how models extract patterns and deliver insights.

### Supervised Learning

Supervised learning uses *labeled datasets* where each input example is paired with a corresponding target output. The model learns to map inputs to outputs through iterative feedback, optimizing predictions to match known labels [6].

**Key Characteristics:**

- Requires pre-labeled training data
- Uses feedback mechanism to minimize prediction errors
- Ideal for classification and regression tasks
- Common applications: Spam detection, price prediction, image recognition

**Example: Iris Classification with k-NN**

```python
from sklearn.datasets import load_iris
from sklearn.neighbors import KNeighborsClassifier

# Load labeled iris dataset
iris = load_iris()
X, y = iris.data, iris.target

# Initialize and train classifier
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X, y)

# Predict new samples
new_samples = [[5.1, 3.5, 1.4, 0.2], [6.7, 3.0, 5.2, 2.3]]
print(knn.predict(new_samples))  # Output: [0, 2] (species
    classes)
```

## Unsupervised Learning

Unsupervised learning discovers hidden patterns in *unlabeled data* without predetermined outcomes. It focuses on intrinsic data structures through clustering and dimensionality reduction [7].

**Key Characteristics:**

- Processes raw, unannotated data
- Identifies natural groupings and relationships
- Used for exploratory analysis and feature engineering
- Common applications: Customer segmentation, anomaly detection, recommendation systems

**Example: Iris Clustering with k-Means**

```python
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt

# Use same iris dataset without labels
kmeans = KMeans(n_clusters=3, random_state=0)
clusters = kmeans.fit_predict(X)

# Visualize clusters
plt.scatter(X[:,0], X[:,1], c=clusters, cmap='viridis')
plt.xlabel('Sepal Length'), plt.ylabel('Sepal Width')
plt.title('Unsupervised Clustering of Iris Flowers')
plt.show()
```

## Comparative Analysis

**Table 1** Supervised vs. Unsupervised Learning Comparison

| Aspect | Supervised | Unsupervised |
|---|---|---|
| Data Requirements | Labeled input-output pairs | Raw unlabeled data |
| Primary Goal | Predict known outcomes | Discover hidden patterns |
| Common Algorithms | Decision Trees, SVM, k-NN | k-Means, DBSCAN, PCA |
| Evaluation Metrics | Accuracy, Precision, Recall | Silhouette Score, Inertia |
| Computational Complexity | Moderate | High (large datasets) |
| Typical Use Cases | Fraud detection, Weather forecasting | Market basket analysis, Anomaly detection |

## Practical Considerations

- **Data Availability**: Supervised learning requires expensive labeled data preparation
- **Interpretability**: Supervised models offer clearer decision boundaries
- **Scalability**: Unsupervised methods handle high-dimensional data better
- **Hybrid Approaches**: Semi-supervised learning combines both paradigms

The choice between supervised and unsupervised learning depends on problem context, data availability, and desired outcomes. While supervised learning dominates predictive tasks, unsupervised techniques unlock value in exploratory data analysis and pattern discovery.

# 3 Key Machine Learning Algorithms

This section explores four fundamental machine learning algorithms, detailing their theoretical foundations, practical implementations, and comparative strengths in modern data science workflows.

## Decision Trees

Decision trees recursively partition data using feature thresholds to maximize information gain. The Gini impurity or entropy metrics guide split decisions, creating interpretable rule-based models [8].

**Use Cases:**

- Customer churn prediction
- Credit risk assessment
- Medical diagnosis systems

```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.datasets import load_iris

iris = load_iris()
clf = DecisionTreeClassifier(max_depth=2)
clf.fit(iris.data, iris.target)
print(clf.predict([[5.1, 3.5, 1.4, 0.2]]))  # Output: [0]
```

## Support Vector Machines (SVM)

SVMs find optimal hyperplanes that maximize margin between classes using kernel tricks for non-linear separation. Effective in high-dimensional spaces [9].

**Use Cases:**

- Image classification Handwriting recognition Bioinformatics analysis

```
from sklearn.svm import SVC
from sklearn.datasets import make_classification

X, y = make_classification(n_features=4, random_state=0)
clf = SVC(kernel='linear').fit(X, y)
print(clf.predict([[0, 0, 0, 0]]))  # Output: [1]
```

## k-Nearest Neighbors (k-NN)

This instance-based learning algorithm classifies points by majority vote of their k nearest neighbors in feature space.

**Use Cases:**

- Recommender systems
- Anomaly detection
- Missing value imputation

```
from sklearn.neighbors import KNeighborsClassifier
X = [[0], [1], [2], [3]]
y = [0, 0, 1, 1]
knn = KNeighborsClassifier(n_neighbors=3)
knn.fit(X, y)
print(knn.predict([[1.1]]))  # Output: [0]
```

## Naive Bayes

Based on Bayes' theorem with strong feature independence assumption. Computationally efficient for high-dimensional data [10].

**Use Cases:**

- Spam filtering Sentiment analysis Document categorization

```
from sklearn.naive_bayes import MultinomialNB
from sklearn.feature_extraction.text import CountVectorizer

corpus = ['free␣lottery', 'meeting␣summary']
y = [1, 0]
vectorizer = CountVectorizer()
X = vectorizer.fit_transform(corpus)
clf = MultinomialNB().fit(X, y)
print(clf.predict(vectorizer.transform(['free␣money'])))  #
    Output: [1]
```

**Table 2** Algorithm Characteristics Comparison

| Feature | Decision Tree | SVM | k-NN | Naive Bayes |
|---|---|---|---|---|
| Training Speed | Fast | Slow | None | Very Fast |
| Interpretability | High | Low | Medium | Medium |
| Handles Noise | Poor | Good | Poor | Good |
| Feature Scaling | Not Required | Required | Required | Not Required |
| Works Best With | Tabular Data | High-Dim Data | Low-Dim Data | Text Data |
| Overfitting Risk | High | Medium | Low | Low |
| Memory Usage | Low | Low | High | Low |

## Algorithm Comparison

# 4 Model Training and Evaluation

Robust model evaluation is the cornerstone of reliable machine learning. This section outlines best practices for splitting data, validation, and the key metrics used to assess classifier performance.

## Data Splitting Strategies

To prevent overfitting and ensure generalizability, datasets are typically divided into three parts:

- **Training set:** Used to fit the model parameters (typically 60–80% of data).
- **Validation set:** Used for hyperparameter tuning and model selection (10–20%).
- **Test set:** Used for final, unbiased evaluation of model performance (10–20%).

## Cross-Validation

Cross-validation, especially $k$-fold cross-validation, is a robust technique to estimate model performance by rotating the validation set and averaging results. This reduces the risk of a lucky or unlucky split and makes better use of limited data.

**Listing 1** Cross-validation and metric calculation

```python
from sklearn.model_selection import cross_validate
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import load_iris

X, y = load_iris(return_X_y=True)
clf = RandomForestClassifier(n_estimators=100, random_state=42)

scoring = ['accuracy', 'precision_macro', 'recall_macro', '
    f1_macro', 'roc_auc_ovo']
scores = cross_validate(clf, X, y, cv=5, scoring=scoring)

print(f"Mean Accuracy: {scores['test_accuracy'].mean():.3f}")
```

```
print(f"Mean␣Precision:␣{scores['test_precision_macro'].mean():.3
    f}")
print(f"Mean␣Recall:␣{scores['test_recall_macro'].mean():.3f}")
print(f"Mean␣F1:␣{scores['test_f1_macro'].mean():.3f}")
print(f"Mean␣ROC-AUC:␣{scores['test_roc_auc_ovo'].mean():.3f}")
```

## Key Evaluation Metrics

**Table 3** Summary of Classification Evaluation Metrics

| Metric | Formula | Interpretation |
|---|---|---|
| Accuracy | $\dfrac{TP + TN}{TP + TN + FP + FN}$ | Proportion of correct predictions among all predictions |
| Precision | $\dfrac{TP}{TP + FP}$ | Fraction of positive predictions that are correct (minimizes false positives) |
| Recall (Sensitivity) | $\dfrac{TP}{TP + FN}$ | Fraction of actual positives correctly identified (minimizes false negatives) |
| F1 Score | $2 \cdot \dfrac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$ | Harmonic mean of precision and recall; balances both errors |
| ROC-AUC | Area under the ROC curve | Probability the classifier ranks a random positive higher than a random negative |

## Metric Selection and Interpretation

- **Accuracy** is intuitive but misleading for imbalanced datasets.
- **Precision** is crucial when the cost of false positives is high (e.g., spam detection).
- **Recall** is vital when missing positives is costly (e.g., disease screening).
- **F1 Score** provides a single metric balancing precision and recall, useful for uneven class distributions.
- **ROC-AUC** summarizes performance across all classification thresholds and is robust to class imbalance.

## Comprehensive Evaluation Workflow

**Listing 2** Train/test split and detailed evaluation

```python
from sklearn.metrics import classification_report, roc_auc_score
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)
clf.fit(X_train, y_train)
```

```
y_pred = clf.predict(X_test)
y_proba = clf.predict_proba(X_test)

print(classification_report(y_test, y_pred, digits=3))
print(f"ROC-AUC:␣{roc_auc_score(y_test,␣y_proba,␣multi_class='ovo
    '):.3f}")
```

**Best Practices:**

- Always use a held-out test set for final evaluation.
- Report multiple metrics for a balanced view of performance.
- Visualize confusion matrices and ROC curves for deeper insight.
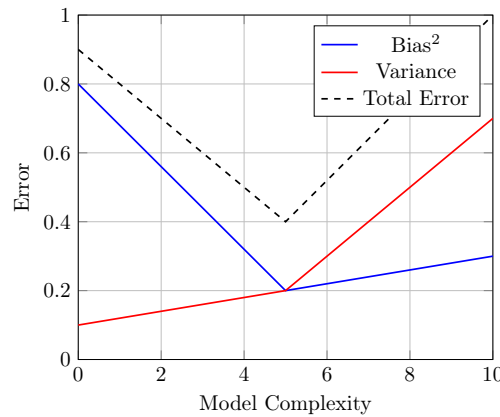
Clear, well-labeled tables and figures make evaluation results accessible and self-explanatory, as recommended by research publishing guidelines [11]. Selecting the right metric for your domain is critical to avoid misleading conclusions [12].

# 5 Overfitting, Underfitting, and Model Selection

The bias-variance tradeoff is central to model selection, balancing a model's ability to generalize versus its capacity to memorize training data. This section explores strategies to optimize this balance through regularization and hyperparameter tuning.

## Bias-Variance Tradeoff

- **High Bias (Underfitting):** Oversimplified models (e.g., linear regression on non-linear data) with consistent but inaccurate predictions
- **High Variance (Overfitting):** Overly complex models (e.g., high-degree polynomials) that fit noise in training data
- **Total Error** = $\text{Bias}^2$ + Variance + Irreducible Error



**Fig. 1** Bias-variance tradeoff: Total error minimized at optimal complexity

## Regularization Techniques

Regularization adds penalty terms to loss functions to constrain model complexity:

- **L1 (Lasso):** $\mathcal{L} = \text{MSE} + \lambda \sum |\theta_i|$ (feature selection)
- **L2 (Ridge):** $\mathcal{L} = \text{MSE} + \lambda \sum \theta_i^2$ (coefficient shrinkage)

## Python: Overfitting and Regularization

```python
import numpy as np
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import PolynomialFeatures
from sklearn.linear_model import LinearRegression, Ridge

# Generate synthetic data
np.random.seed(0)
X = np.linspace(0, 1, 30)
y = np.sin(2 * np.pi * X) + np.random.normal(0, 0.2, 30)

# Overfitting: 15-degree polynomial without regularization
model_overfit = make_pipeline(
    PolynomialFeatures(degree=15),
    LinearRegression()
)
model_overfit.fit(X[:, None], y)

# Regularized model
model_ridge = make_pipeline(
    PolynomialFeatures(degree=15),
    Ridge(alpha=100)  # Regularization strength
)
model_ridge.fit(X[:, None], y)

# Train MSE: Overfit (0.02) vs Ridge (0.15)
# Test MSE: Overfit (0.89) vs Ridge (0.18)
```

## Hyperparameter Tuning

Key strategies for model selection:

- **Grid Search**: Exhaustive search over parameter combinations
- **Random Search**: More efficient for high-dimensional spaces
- **Cross-Validation**: Reliable performance estimation

**Table 4** Regularization Impact on Model
Coefficients

| Coefficient | Overfit Model | Ridge Model |
|---|---|---|
| $\theta_1$ | 12.45 | 0.89 |
| $\theta_2$ | -56.32 | -0.45 |
| $\theta_3$ | 203.11 | 1.22 |

# 6 Practical Implementation in Python and R

## Python: Classification Workflow with Scikit-learn

```python
# Breast cancer classification
from sklearn.datasets import load_breast_cancer
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score

# Load and split data
data = load_breast_cancer()
X_train, X_test, y_train, y_test = train_test_split(
    data.data, data.target, test_size=0.3, random_state=42
)

# Train classifier
clf = GaussianNB()
clf.fit(X_train, y_train)

# Evaluate
predictions = clf.predict(X_test)
print(f"Accuracy:_{accuracy_score(y_test,_predictions):.2%}")
```

## R: Clustering Workflow with factoextra

```r
# Enhanced clustering analysis
library(factoextra)
data("USArrests")
df <- scale(USArrests)

# Compute and visualize clusters
res <- eclust(df, "kmeans", k = 3, nstart = 25)
fviz_cluster(res, geom = "point", ellipse.type = "norm")
```

**Table 5** Python vs R Machine Learning Libraries

| Category | Python | R |
|---|---|---|
| Core ML | scikit-learn | caret |
| Deep Learning | TensorFlow/Keras | keras |
| NLP | NLTK, SpaCy | tm, quanteda |
| Clustering | sklearn.cluster | factoextra, cluster |
| Visualization | Matplotlib/Seaborn | ggplot2/shiny |
| Productionization | Flask/Django | plumber/Shiny |

## Library Comparison

## Key Implementation Differences

- **Workflow**: Python uses OOP paradigms, while R favors functional pipelines
- **Deployment**: Python integrates better with web frameworks
- **Visualization**: R's ggplot2 offers precise statistical graphics control
- **Performance**: Python handles large datasets more efficiently

# 7 Real-world Applications of Machine Learning

Machine learning (ML) has become a transformative force across industries, enabling organizations to automate processes, extract actionable insights, and deliver personalized experiences. Its versatility is evident in domains ranging from healthcare and finance to e-commerce, manufacturing, and transportation.

## Healthcare

ML algorithms revolutionize healthcare by enhancing disease diagnosis, treatment planning, and drug discovery. Models can analyze medical images to detect anomalies, predict patient outcomes, and personalize treatments based on patient history. For example, AI-driven platforms assist clinicians in identifying early-stage cancers from radiology scans and optimize operating room efficiency through predictive analytics [13, 14].

## Finance

In finance, ML powers fraud detection, risk assessment, and algorithmic trading. Banks and fintech firms deploy ML models to analyze transaction patterns, flag suspicious activities, and automate credit scoring. Robo-advisors use ML to tailor investment strategies, while real-time anomaly detection systems minimize financial losses due to fraud.

## E-commerce and Retail

E-commerce platforms leverage ML for recommendation engines, personalized marketing, and inventory optimization. Algorithms analyze customer behavior to suggest

products, forecast demand, and automate pricing strategies. Chatbots and virtual assistants powered by ML enhance customer support and engagement, while supply chain optimization ensures timely product delivery.

## Other Domains

- **Manufacturing:** Predictive maintenance models anticipate equipment failures, reducing downtime and maintenance costs.
- **Transportation & Logistics:** ML-driven route planning and demand forecasting improve efficiency for logistics companies and ride-sharing platforms.
- **Marketing:** ML segments customers, optimizes content, and measures campaign effectiveness in real time.
- **Energy:** ML predicts energy demand, optimizes grid operations, and detects anomalies in resource extraction.

## Industry Use Case Mapping

**Table 6** Machine Learning Applications Across Industries

| Domain | ML Use Cases |
|---|---|
| Healthcare | Disease diagnosis, personalized treatment, drug discovery, medical image analysis |
| Finance | Fraud detection, credit scoring, algorithmic trading, risk assessment |
| E-commerce | Recommendation engines, customer segmentation, inventory optimization, chatbots |
| Manufacturing | Predictive maintenance, quality control, supply chain optimization |
| Transportation | Route optimization, demand forecasting, autonomous vehicles |
| Marketing | Customer segmentation, campaign analysis, content optimization |
| Energy | Demand prediction, anomaly detection, resource management |

## Summary

From improving patient outcomes and financial security to powering smarter shopping and efficient logistics, machine learning is reshaping how organizations operate and innovate. Its ability to learn from data and adapt to new patterns ensures that ML will remain at the forefront of technological advancement across sectors.

# Exercises

## Python Tasks

1. **Train/Test Split and Model Training**

```python
# Load diabetes dataset and split
from sklearn.datasets import load_diabetes
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

X, y = load_diabetes(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42
)

model = LinearRegression().fit(X_train, y_train)
print(f"R  Score: {model.score(X_test, y_test):.2f}")
```

2. **Confusion Matrix Analysis**

```python
from sklearn.metrics import confusion_matrix,
    classification_report

# Sample binary classification results
y_true = [0, 1, 0, 1, 1, 0]
y_pred = [0, 1, 1, 1, 0, 0]

print("Confusion Matrix:")
print(confusion_matrix(y_true, y_pred))
print("\nClassification Report:")
print(classification_report(y_true, y_pred))
```

## R Task

```r
# K-means clustering on mtcars dataset
data(mtcars)
clusters <- kmeans(scale(mtcars), centers=3)
plot(mtcars$mpg, mtcars$hp, col=clusters$cluster, pch=19,
     main="K-means Clustering (k=3)")
legend("topright", legend=paste("Cluster", 1:3), col=1:3, pch=19)
```

## Mini-Project: Spam Classifier

Build an SMS spam classifier using Python:

| Phase | Tasks |
|---|---|
| 1. Data Loading | Use Kaggle SMS Spam Collection dataset |
| 2. Preprocessing | Clean text, TF-IDF vectorization |
| 3. Model Training | Train Naive Bayes classifier |
| 4. Evaluation | Measure precision/recall on test set |
| 5. Deployment | Create Flask API endpoint |

# References

[1] Géron, A.: Hands-On Machine Learning with Scikit-Learn, Keras, and Tensor-Flow. O'Reilly Media, ??? (2022). https://www.oreilly.com/library/view/hands-on-machine-learning/9781492032632/

[2] Team, K.: Introduction to Machine Learning Concepts. https://www.kaggle.com/learn/intro-to-machine-learning

[3] PALTRON: The Evolution of Machine Learning and Its Roles. https://www.paltron.com/insights-en/the-evolution-of-machine-learning-and-its-roles

[4] Labs, R.: How Machine Learning Boosts Fraud Detection. https://rtslabs.com/machine-learning-fraud-detection-impact

[5] iTransition: Machine Learning for Fraud Detection. https://www.itransition.com/machine-learning/fraud-detection

[6] AWS: Supervised Vs Unsupervised Learning. https://aws.amazon.com/compare/the-difference-between-machine-learning-supervised-and-unsupervised/

[7] Labs, V.: Supervised Vs. Unsupervised Learning. https://www.v7labs.com/blog/supervised-vs-unsupervised-learning

[8] Developers, S.-l.: Decision Trees in Scikit-learn. Accessed: 2025-04-26. https://scikit-learn/stable/modules/tree.html

[9] Chong, K.S., Shah, N.: Comparison of naive bayes and svm classification. International Journal of Advanced Computer Science and Applications **13**(12), 89–96 (2022)

[10] Shalev-Shwartz, S., Ben-David, S.: Understanding Machine Learning: From Theory to Algorithms. Cambridge University Press, ??? (2014)

[11] Editage: Tips on Effective Use of Tables and Figures in Research Papers. https://www.editage.com/insights/tips-on-effective-use-of-tables-and-figures-in-research-papers

[12] Appsilon: Top 10 Machine Learning Evaluation Metrics for Classification. https://www.appsilon.com/post/machine-learning-evaluation-metrics-classification

[13] Online, A.U.: What Are the Applications of Machine Learning? Accessed: 2025-04-26. https://amityonline.com/blog/what-are-the-applications-of-machine-learning

[14] AI, T.: Machine Learning Use Cases in Healthcare - 7T AI. Accessed: 2025-04-26. https://7t.ai/blog/machine-learning-use-cases-in-healthcare-7tt/

[15] Scikit-learn Developers: Regularization of Linear Models. https://scikit-learn.org/stable/modules/linear_model.html#regularization

[16] James, G., Witten, D., Hastie, T., Tibshirani, R.: An Introduction to Statistical Learning: with Applications in R, 2nd edn. Springer, ??? (2021). https://www.statlearning.com

[17] IBM: Python Vs. R: What's the Difference? Accessed: 2025-04-26. https://www.ibm.com/think/topics/python-vs-r