

Machine Learning Fundamentals

GenAI Learning Sprint

Iván Moreno

What is Machine Learning?

Definition

A field of study that gives computers the ability to learn from data without being explicitly programmed.

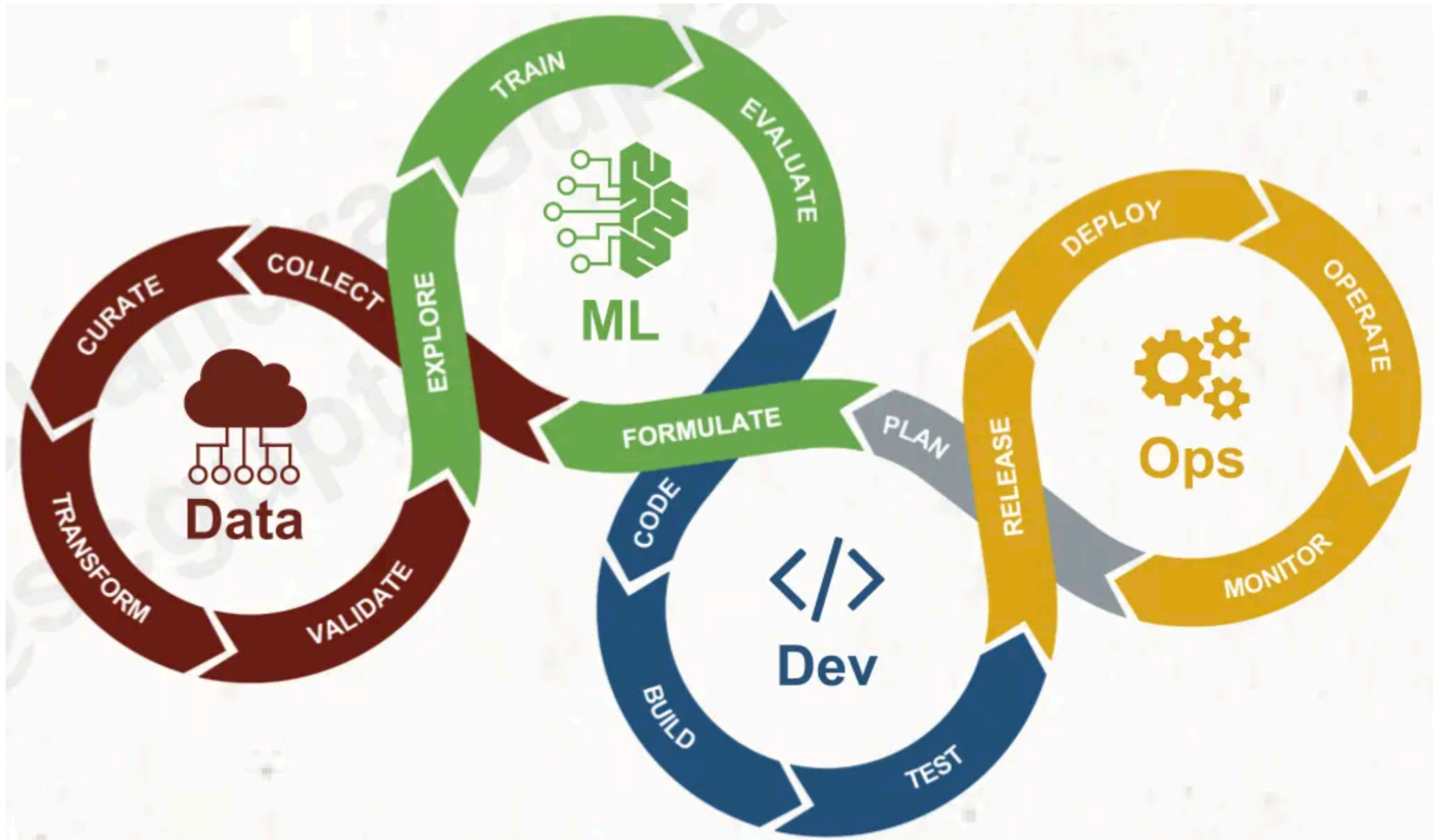
Why Machine Learning Matters

- Scalability in decision-making.
- Automating repetitive tasks.
- Unlocking insights from large datasets.

Core Principles of Machine Learning

- **Data-Driven Decisions:** Leveraging data for predictive modeling.
- **Model Representation:** Choice of models (e.g., linear models, decision trees, neural networks).
- **Generalization:** Balancing model complexity and performance on unseen data.
- **Optimization:** Loss functions, cost functions, and the gradient descent method.
- **Evaluation:** Accuracy, precision, recall, F1-score, ROC-AUC.

The Machine Learning Pipeline



Key Stages in the ML Pipeline

1. Problem Definition:

- Identify the objective (e.g., classification, regression).

2. Data Collection:

- Gather data from reliable sources.

3. Data Preprocessing:

- Handling missing values, data normalization, feature engineering.

4. Model Selection:

- Choose appropriate algorithms based on the problem and data characteristics.

Key Stages in the ML Pipeline (cont.)

5. Training:

- Split data into training/validation/test sets. Train the model.

6. Evaluation:

- Use performance metrics to assess model quality.

7. Deployment:

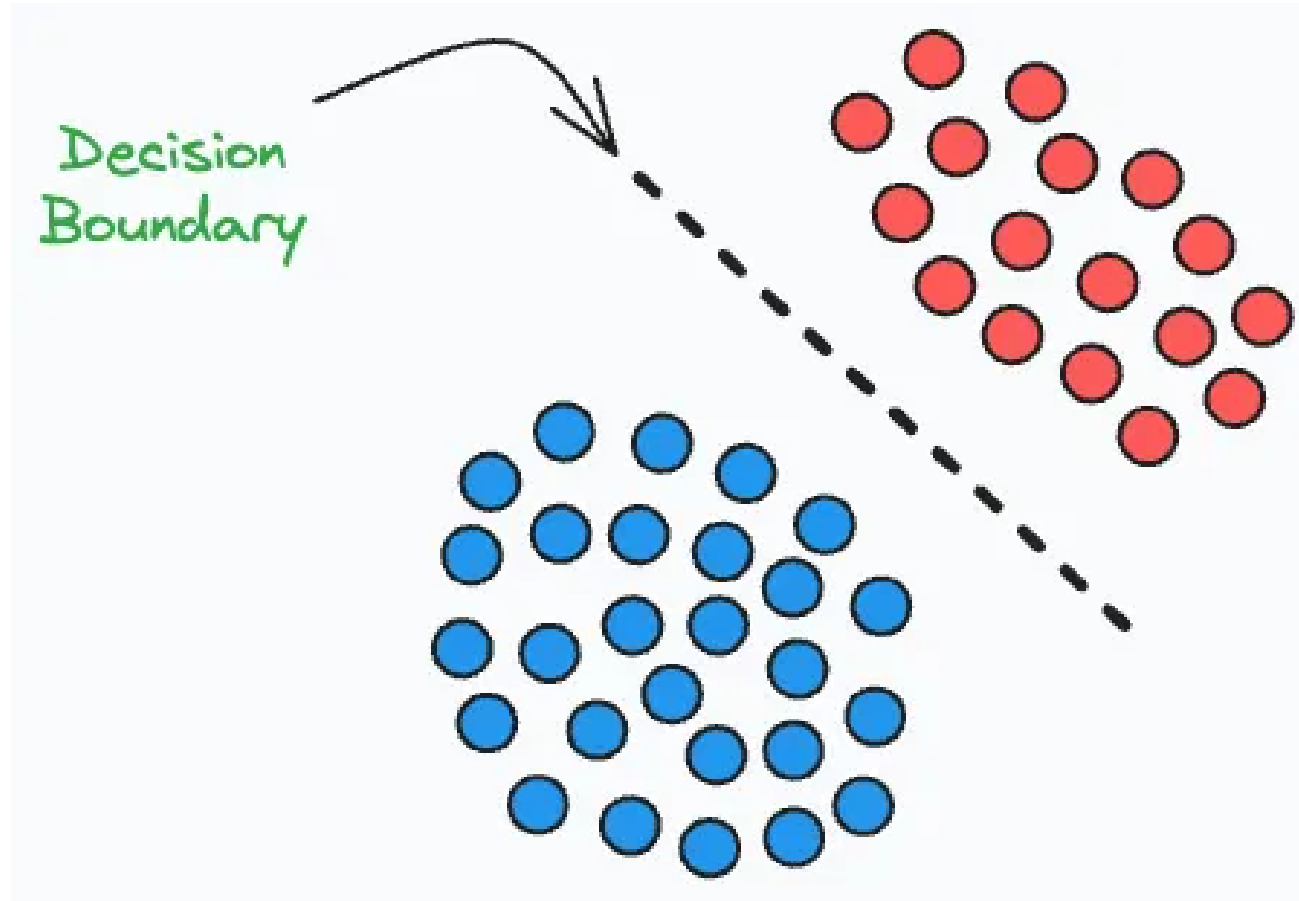
- Integrate the model into production environments.

8. Monitoring & Maintenance:

- Continuously monitor model performance and update as needed.

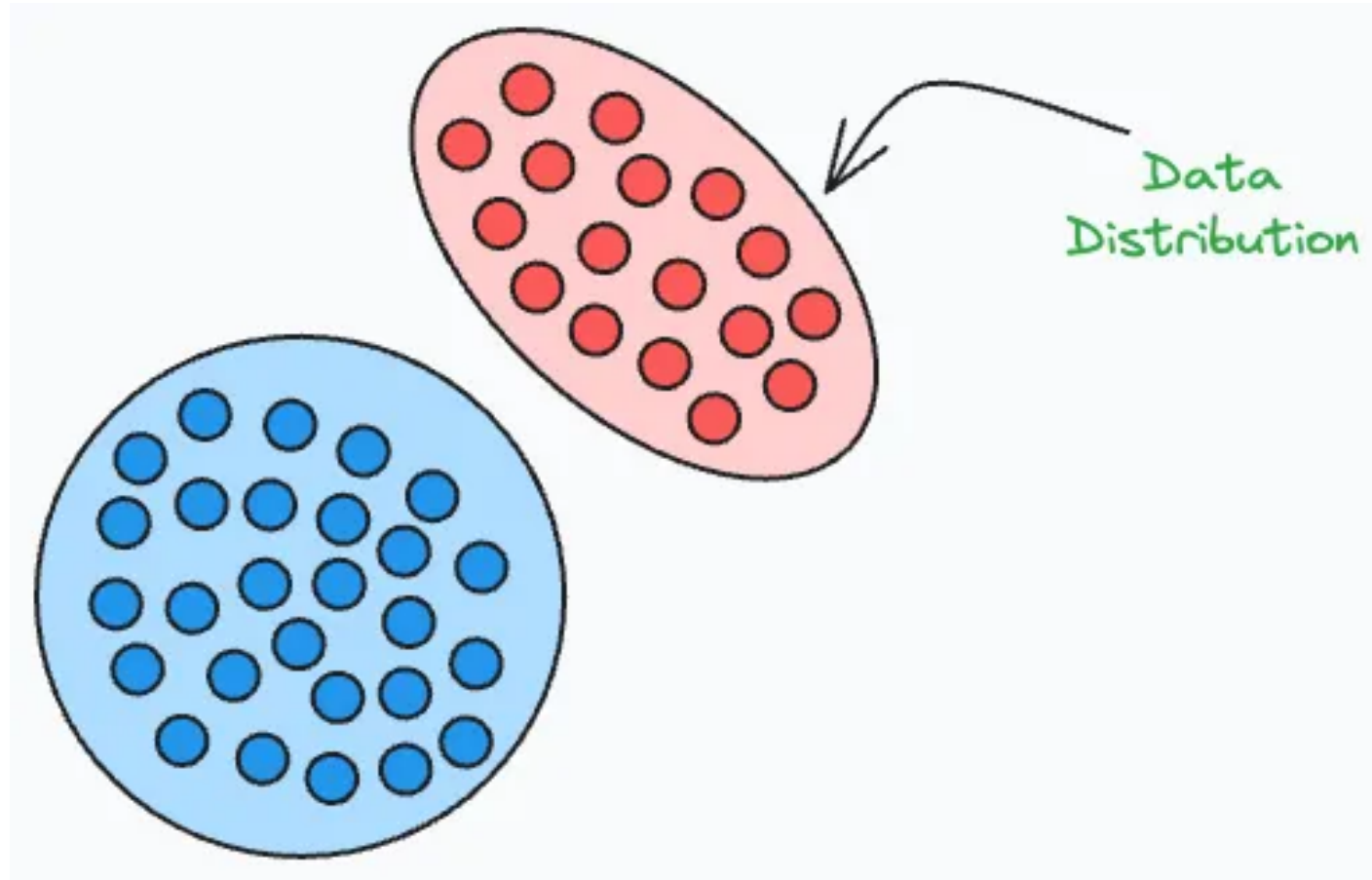
Discriminative vs. Generative Models

Discriminative Models



- Focus on predicting the target variable directly (e.g., $P(y|X)$).
- **Examples:** Logistic Regression, Support Vector Machines (SVMs), Neural Networks.
- **Advantages:** Often provide better predictive accuracy.

Generative Models



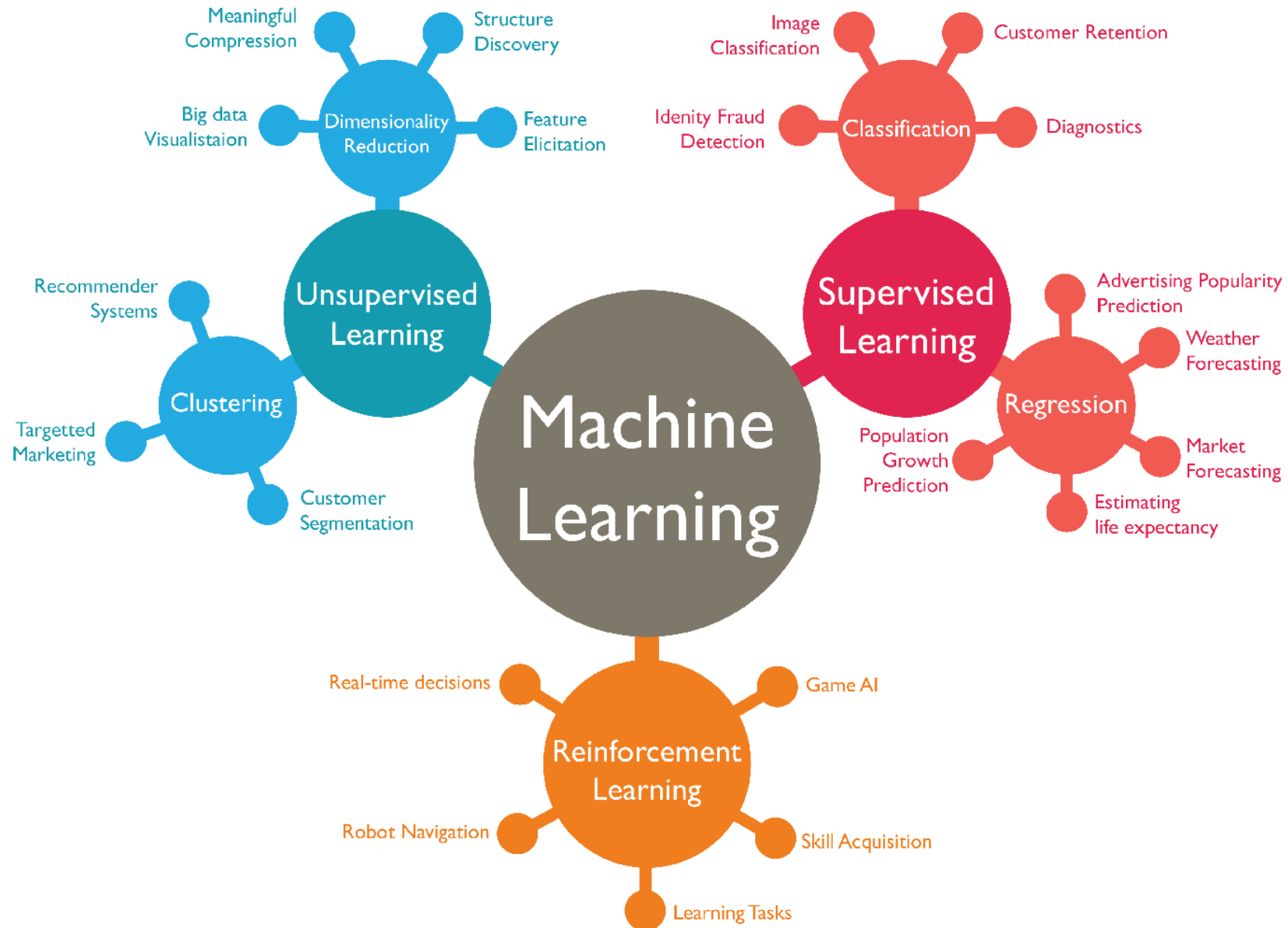
When to Use Each Type of Model



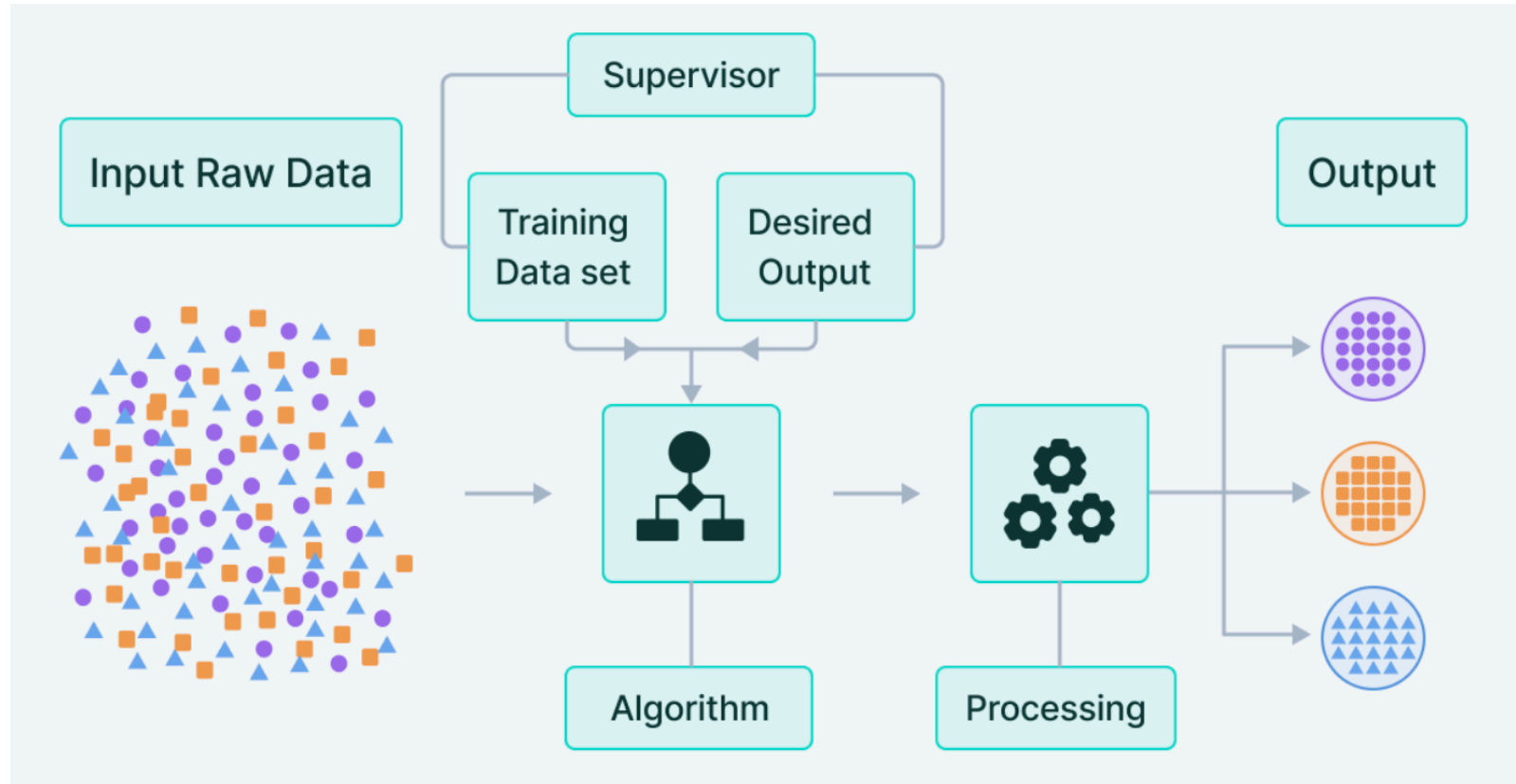
Tip

Choose based on the problem context—discriminative models for classification accuracy, generative models for data understanding and synthesis.

Learning Paradigms

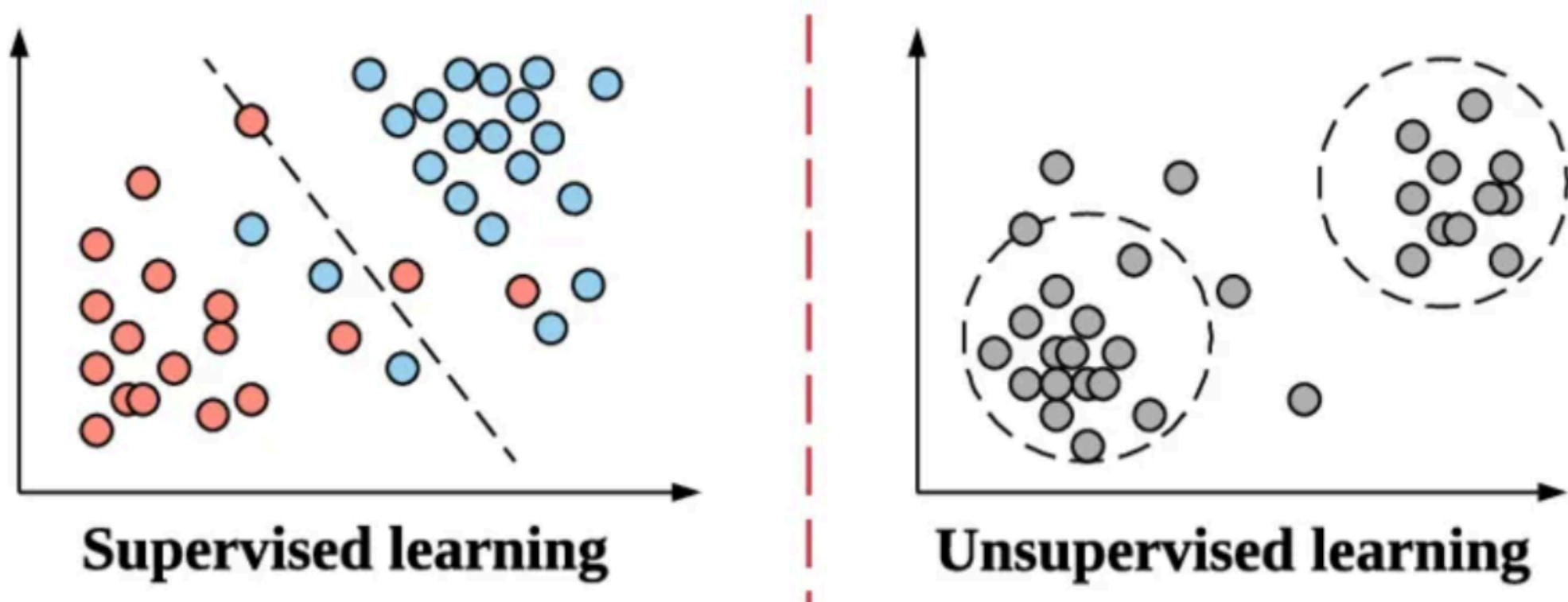


Supervised Learning



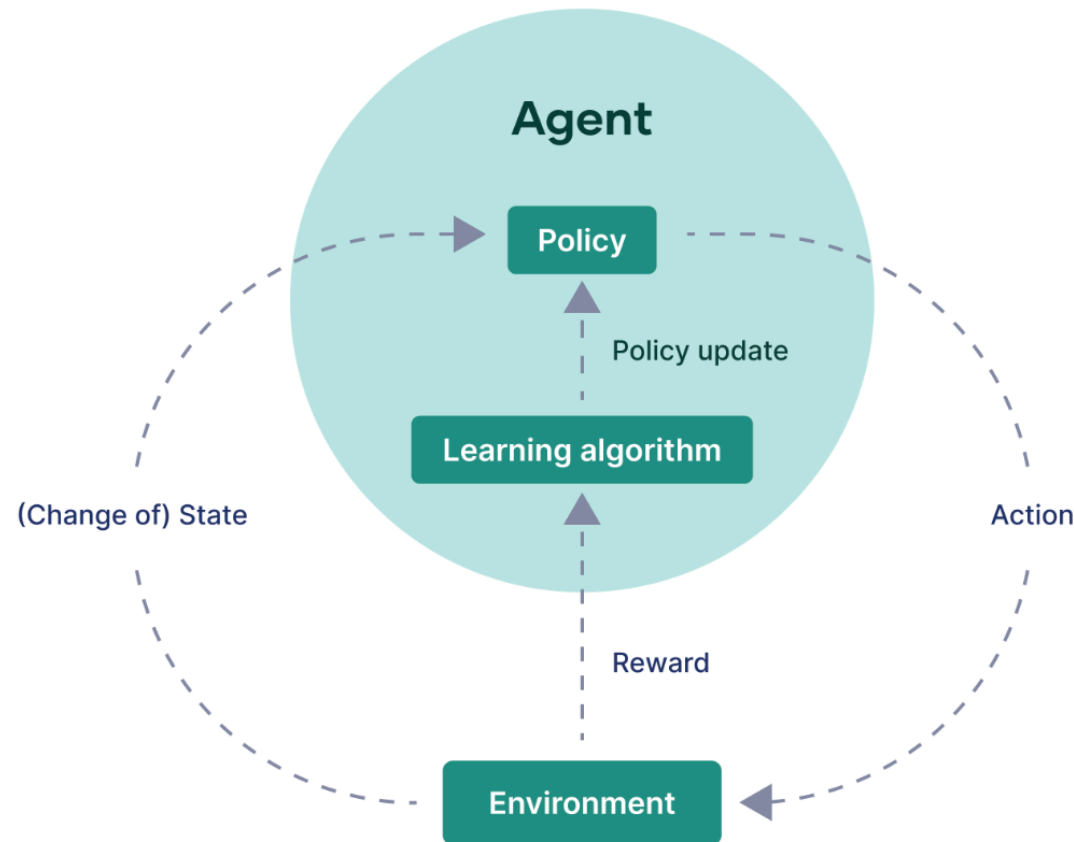
- **Goal:** Learn a function that maps inputs to outputs using labeled data.
- **Common Algorithms:** Linear Regression, Decision Trees, Support Vector Machines (SVMs), k-Nearest Neighbors (k-NN).
- **Applications:** Spam detection, medical diagnosis, fraud detection.

Unsupervised Learning



- **Goal:** Find hidden patterns or intrinsic structures in unlabeled data.
- **Techniques:** Clustering (K-means, Hierarchical Clustering), Dimensionality Reduction (PCA, t-SNE).
- **Applications:** Market segmentation, anomaly detection, gene expression analysis.

Reinforcement Learning



- **Goal:** Learn a policy to maximize cumulative reward in an environment.
- **Key Concepts:** Agent, Environment, Actions, Rewards, Policy, Value Functions.
- **Applications:** Robotics, game playing (e.g., AlphaGo), autonomous driving.

Evaluating Model Performance

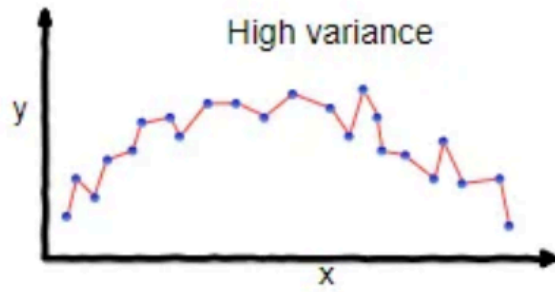
Evaluation Metrics

- Depending on the problem type, different metrics are used:
 - **Classification:** Accuracy, Precision, Recall, F1-Score, ROC-AUC.
 - **Regression:** Mean Squared Error (MSE), Mean Absolute Error (MAE), R-Squared.

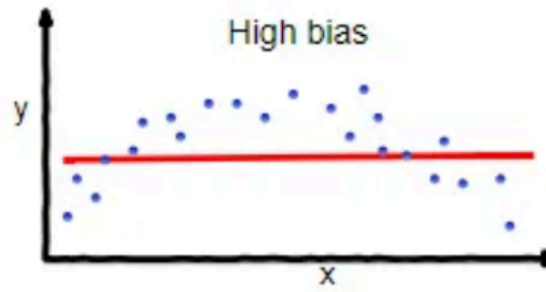
Cross-Validation

- **Purpose:** Estimate model performance on unseen data.
- **Why?** Prevent overfitting and assess generalization capabilities.
- **Techniques:** K-Fold Cross-Validation, Stratified Cross-Validation.

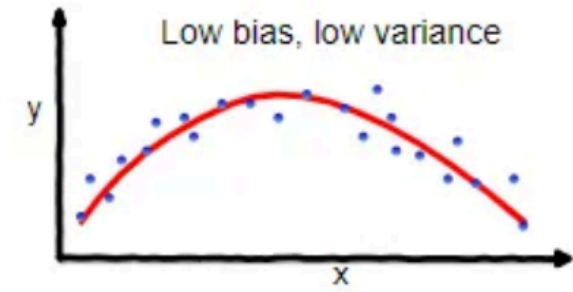
The Bias-Variance Tradeoff



overfitting



underfitting



Good balance

Represents the tradeoff between model complexity and generalization.

- **Overfitting:** Model captures noise instead of underlying patterns.
- **Underfitting:** Model is too simple to capture the data's complexity.



Goal

Find the right balance to minimize error on unseen data.

Industry Use Case: Classification Exercise

Scenario

- A financial institution needs to detect fraudulent transactions.

Steps

1. **Data Preprocessing:** Cleaning and transforming transaction data.
2. **Model Selection:** Evaluate Logistic Regression, Random Forests, and XGBoost.
3. **Model Training:** Split data into training and test sets. Train models.
4. **Evaluation:** Compare performance using **ROC–AUC** and precision-recall metrics.
5. **Interpretation:** Analyze feature importance for insights.

Conclusion and Next Steps

Summary:

- Machine learning bridges data and intelligent systems.
- Discriminative and generative models have different strengths.
- Diverse learning paradigms cater to various problem types.
- Deep learning pushes boundaries in AI capabilities.

Q&A

- Open floor for questions and discussion.