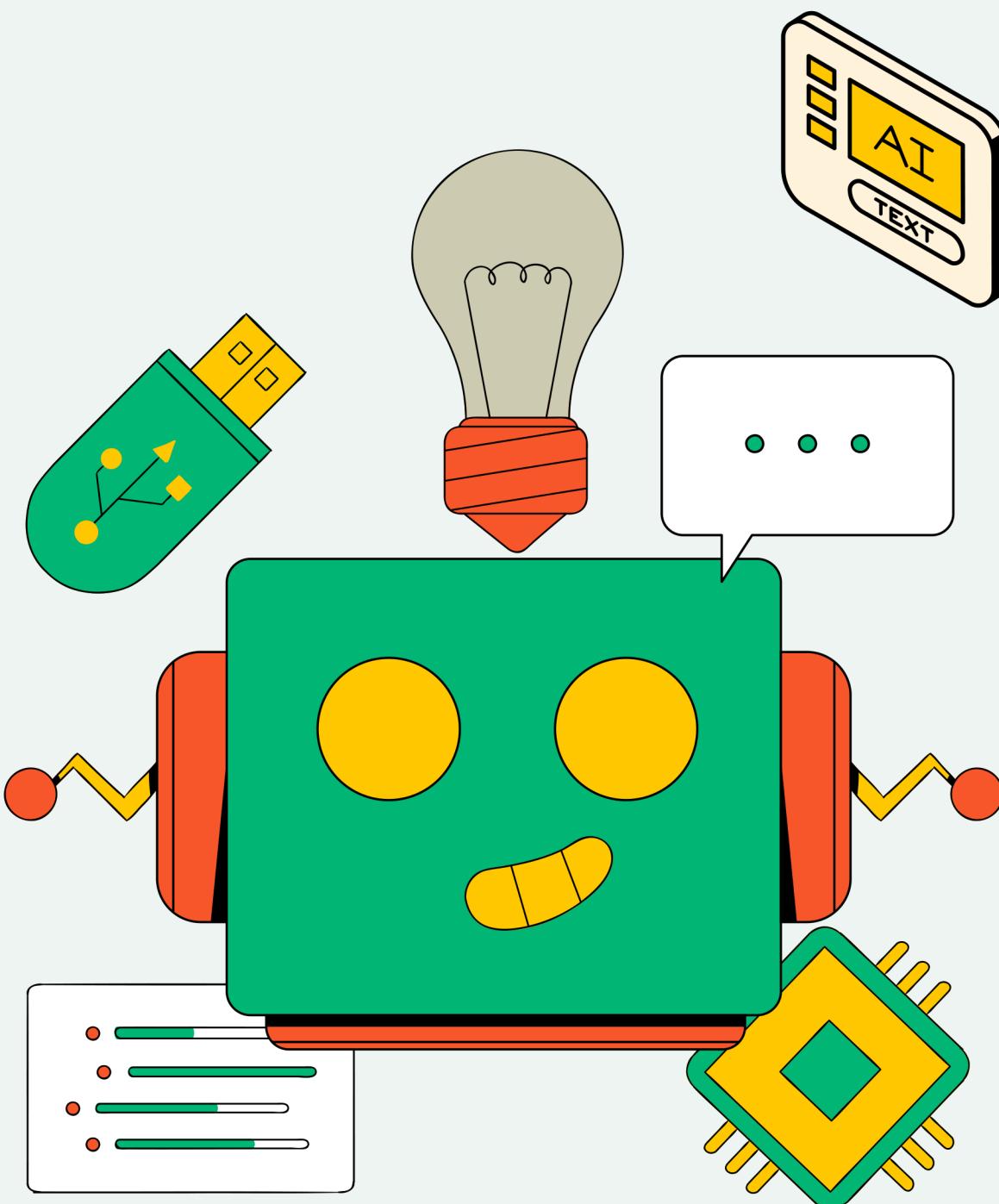


WE LEARN FOR THE FUTURE

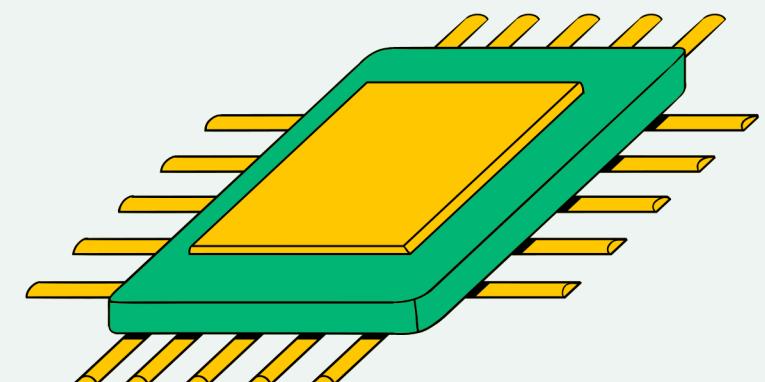


XGBOOST, METRICS, LOSS, AND ONE-HOT ENCODING ALTERNATIVES

PRESENTATION

PRESENTED BY:

SARA METAWEA

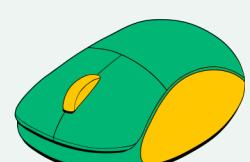
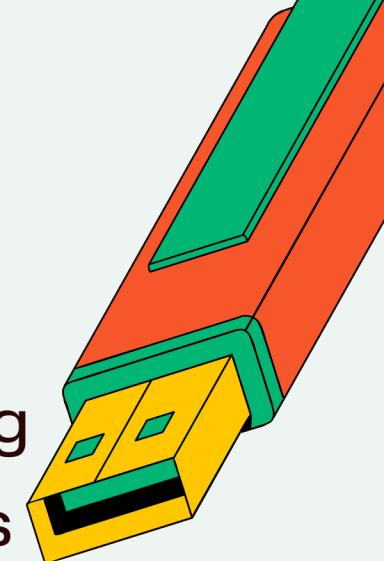


XGBOOST

- XGBoost (Extreme Gradient Boosting) is a powerful and efficient implementation of the gradient boosting framework. It is widely used for regression, classification, and ranking problems. XGBoost is known for its performance and speed, utilizing advanced features like tree pruning, parallel processing, and regularization to improve model accuracy and prevent overfitting.

- Example Code:

```
• import xgboost as xgb  
• from sklearn.datasets import load_iris  
• from sklearn.model_selection import train_test_split  
• data = load_iris()  
• X_train, X_test, y_train, y_test = train_test_split(data.data, data.target, test_size=0.2, random_state=42)  
• model = xgb.XGBClassifier(objective='multi:softprob', num_class=3)  
• model.fit(X_train, y_train)  
• preds = model.predict(X_test)
```



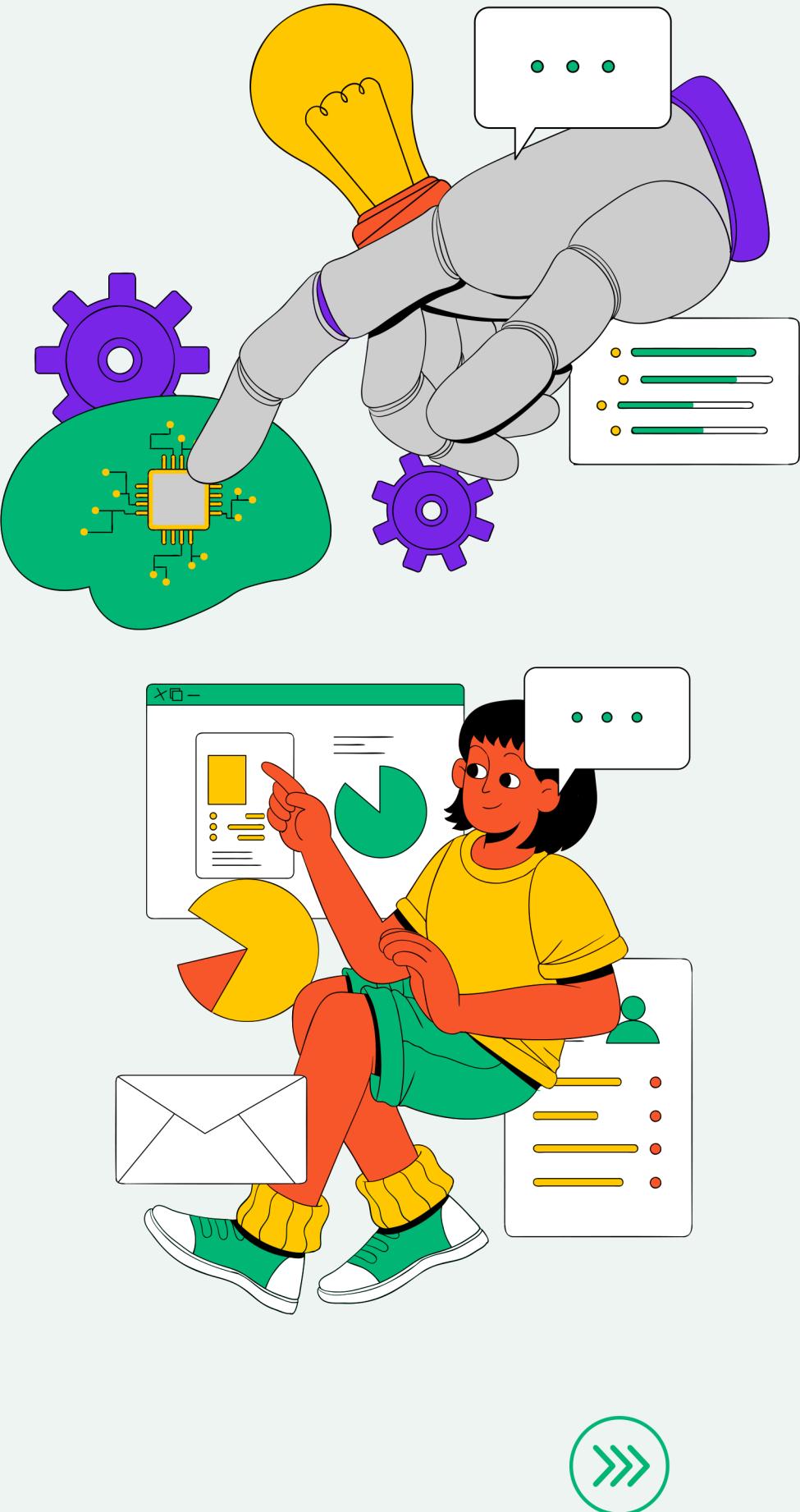
OPTIMIZATION

Optimization in machine learning involves the process of adjusting model parameters to minimize or maximize a specific objective function. This can involve various algorithms such as gradient descent, stochastic gradient descent, and more sophisticated techniques like Adam or RMSprop. The goal is to find the optimal set of parameters that reduce the error or loss function of the model.



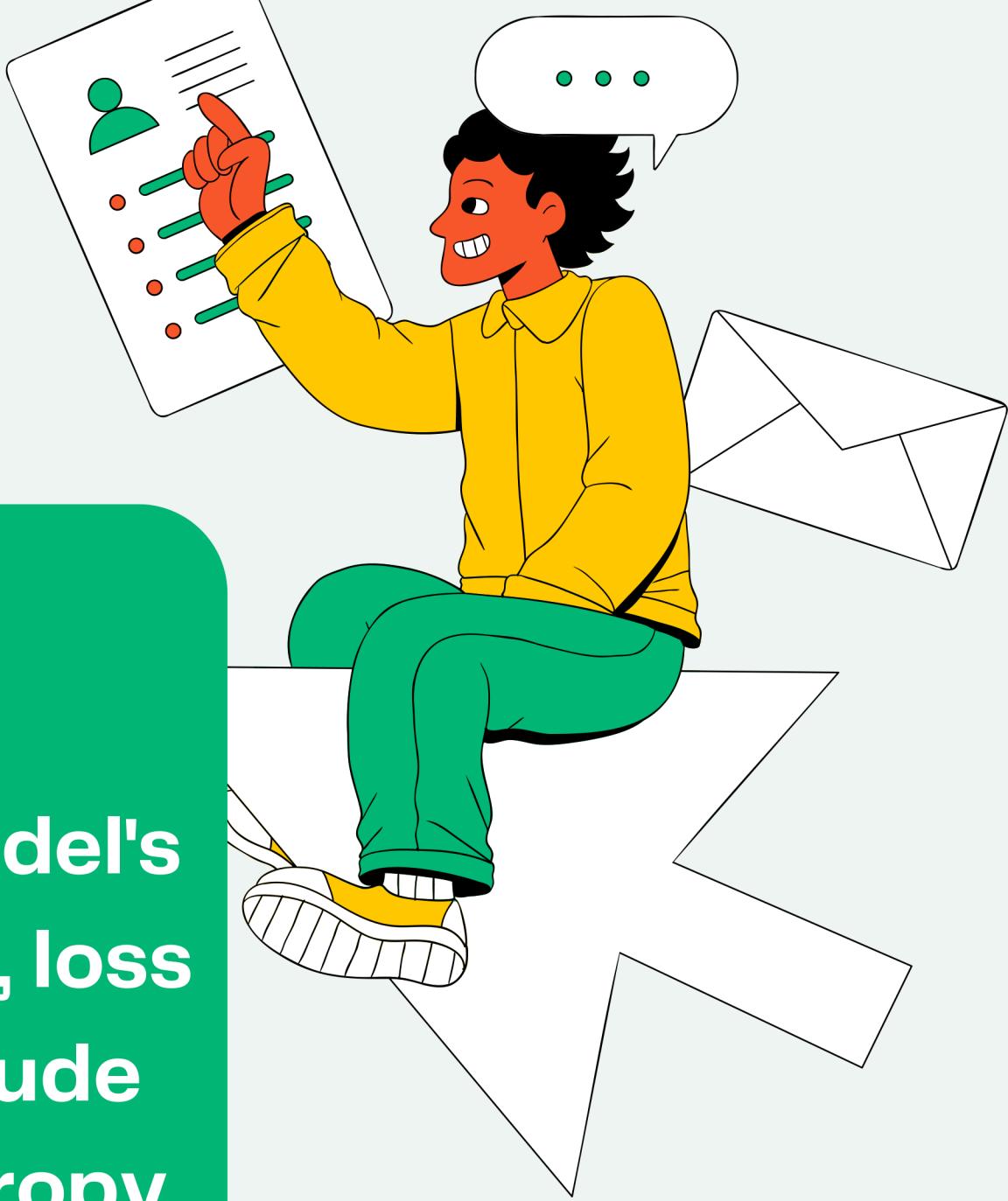
METRICS

Metrics are used to evaluate the performance of a machine learning model. Common metrics include accuracy, precision, recall, F1-score for classification problems, and Mean Squared Error (MSE), Mean Absolute Error (MAE), or R-squared for regression problems. Choosing the right metric depends on the problem and what aspect of the model's performance is most critical.



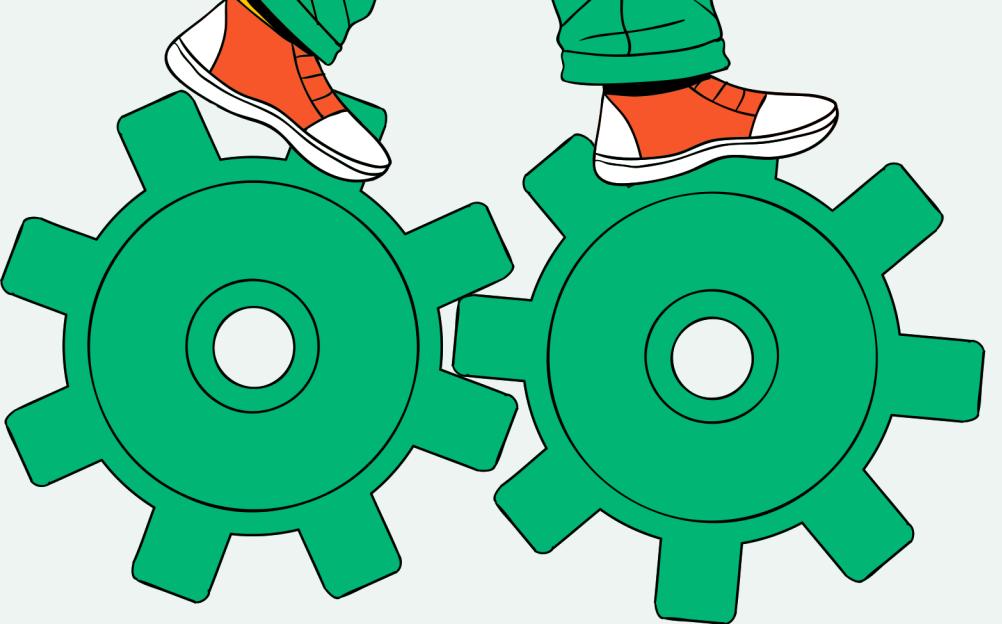
LOSS

Loss functions, or cost functions, measure how well a model's predictions match the actual data. In supervised learning, loss functions guide the optimization process. Examples include Mean Squared Error (MSE) for regression and Cross-Entropy Loss for classification. The choice of loss function impacts the training process and the performance of the model.



ONE-HOT ENCODING AND ITS ALTERNATIVES

One-hot encoding is a technique used to convert categorical variables into a binary vector representation. Each category is represented as a vector with a single high (1) value and all other values being low (0). This method is useful for algorithms that require numerical input data. However, there are alternative encoding techniques:



- Label Encoding: Converts categorical values into integer codes.
- Ordinal Encoding: Similar to label encoding but maintains order among categories.
- Target Encoding: Replaces categories with the average of the target variable.
- Binary Encoding: Converts integers to binary digits and then creates separate columns for each bit.
- Frequency Encoding: Replaces categories with the frequency of the category in the dataset.
- Hashing Encoding: Uses a hash function to convert categories into a fixed number of columns.

REFERENCES

Here are some resources for further reading:

- [XGBoost Documentation](#)
- Optimization Algorithms in Machine Learning
- Evaluation Metrics in Machine Learning
- [Loss Functions in Machine Learning](#)
- One-Hot Encoding and Alternatives

