

algorithm for email spam detection

Data Collection: Gather a dataset of emails labeled as spam or not spam (ham). You need a substantial amount of both spam and non-spam emails for training.

Preprocessing: Preprocess the emails to extract useful features. This may include:

- Tokenization: Splitting the text into words or phrases.
- Removing stopwords: Commonly occurring words like "the," "and," etc.
- Stemming or Lemmatization: Normalizing words to their base form.
- Feature Extraction: Converting text into numerical features that algorithms can understand. This could be bag-of-words representation, TF-IDF (Term Frequency-Inverse Document Frequency), or word embeddings.

Splitting Data: Divide the dataset into training and testing sets. Typically, around 70-80% of the data is used for training and the rest for testing.

Model Selection: Choose a suitable machine learning algorithm. Common choices include:

1. Naive Bayes
2. Support Vector Machines (SVM)
3. Logistic Regression
4. Random Forest
5. Gradient Boosting Machines (GBM)
6. Neural Networks

Training: Train the chosen model using the training data. The model learns to distinguish between spam and non-spam emails based on the features extracted from the emails.

Evaluation: Evaluate the trained model using the testing data to assess its performance. Common evaluation metrics include accuracy, precision, recall, and F1-score.

Tuning: Fine-tune the model hyperparameters to improve performance. This may involve techniques like cross-validation and grid search.

Deployment: Once satisfied with the performance, deploy the trained model to detect spam emails in real-time.

Monitoring and Updates: Monitor the model's performance over time and update it as needed with new data or retraining if the performance degrades.