

VNP MERGED

Natural Language Processing (NLP) encompasses a wide range of tasks aimed at enabling machines to understand, interpret, and generate human language. Here are some of the most common NLP tasks along with their characteristics:

1. **Text Classification**

- **Description**: Assigning predefined labels or categories to a piece of text.
- **Examples**: Spam detection in emails, sentiment analysis (positive/negative/neutral), news categorization.
- **Key Characteristics**:
 - Supervised learning task.
 - Requires labeled data.
 - Utilizes models like Naive Bayes, SVMs, CNNs, RNNs, and Transformers.

2. **Named Entity Recognition (NER)**

- **Description**: Identifying and classifying named entities (e.g., people, organizations, locations, dates) within text.
- **Examples**: Extracting names of companies and dates from financial documents.
- **Key Characteristics**:
 - Sequence labeling task.
 - Requires labeled data with entities tagged.
 - Commonly uses models like Conditional Random Fields (CRFs) combined with LSTM or Transformer architectures.

3. **Part-of-Speech (POS) Tagging**

- **Description**: Assigning grammatical categories (e.g., noun, verb, adjective) to each word in a sentence.
- **Examples**: "The cat sat on the mat" -> "The/DET cat/NOUN sat/VERB on/ADP the/DET mat/NOUN".
- **Key Characteristics**:
 - Sequence labeling task.
 - Requires labeled corpora with POS tags.
 - Often uses models like HMMs, CRFs, or neural networks.

4. **Sentiment Analysis**

- **Description**: Determining the sentiment or emotional tone expressed in a piece of text (e.g., positive, negative, neutral).
- **Examples**: Analyzing customer reviews to gauge satisfaction.
- **Key Characteristics**:

- Can be treated as a text classification problem.
- Requires labeled data for supervised learning.
- Commonly uses models like Naive Bayes, LSTM, CNNs, and Transformers.

5. **Machine Translation**

- **Description**: Automatically translating text from one language to another.
- **Examples**: Translating a document from English to Spanish.
- **Key Characteristics**:
 - Sequence-to-sequence (seq2seq) task.
 - Requires large parallel corpora (bilingual text pairs).
 - Uses models like RNNs, LSTMs, and Transformers (e.g., Google's Transformer model in Google Translate).

6. **Text Summarization**

- **Description**: Producing a concise and coherent summary of a longer text document.
- **Examples**: Generating a summary of a news article.
- **Key Characteristics**:
 - Can be extractive (selecting key sentences) or abstractive (generating new sentences).
 - Requires labeled data for supervised learning.
 - Commonly uses seq2seq models, Transformers, and attention mechanisms.

7. **Question Answering (QA)**

- **Description**: Automatically answering questions posed in natural language by extracting information from a given text or dataset.
- **Examples**: Answering "Who is the president of the USA?" based on a document.
- **Key Characteristics**:
 - Can be open-domain (unrestricted topics) or closed-domain (restricted to a specific topic).
 - Uses models like BERT, GPT, and other Transformer-based architectures.

8. **Speech Recognition (Automatic Speech Recognition - ASR)**

- **Description**: Converting spoken language into text.
- **Examples**: Transcribing a conversation or a speech.
- **Key Characteristics**:
 - Sequence-to-sequence task.
 - Requires audio data paired with text transcripts.
 - Uses models like RNNs, LSTMs, and end-to-end architectures like DeepSpeech.

9. **Language Modeling**

- **Description**: Predicting the next word in a sentence or estimating the probability of a sequence of words.

- **Examples**: Autocomplete features in search engines, text generation.
- **Key Characteristics**:
 - Unsupervised learning task.
 - Can be generative (creating text) or discriminative (ranking sequences).
 - Uses models like N-grams, RNNs, LSTMs, and Transformers (e.g., GPT).

10. **Text Generation**

- **Description**: Generating coherent and contextually relevant text.
- **Examples**: Creating dialogues for chatbots, writing articles.
- **Key Characteristics**:
 - Sequence-to-sequence task.
 - Often builds on language models.
 - Uses RNNs, LSTMs, and Transformers (e.g., GPT).

11. **Coreference Resolution**

- **Description**: Determining which words in a sentence refer to the same entity.
- **Examples**: In "Alice went to the park. She enjoyed it," resolving "She" to "Alice."
- **Key Characteristics**:
 - Involves understanding the context and semantics.
 - Can be part of larger NLP systems like QA or summarization.
 - Often uses CRFs, neural networks, and attention mechanisms.

12. **Semantic Role Labeling (SRL)**

- **Description**: Identifying the predicate-argument structure in a sentence, i.e., who did what to whom, when, where, etc.
- **Examples**: In "John gave a book to Mary," "John" is the giver, "book" is the thing given, and "Mary" is the receiver.
- **Key Characteristics**:
 - Requires understanding of syntax and semantics.
 - Often uses models like CRFs, LSTMs, and Transformers.

13. **Dependency Parsing**

- **Description**: Analyzing the grammatical structure of a sentence and identifying relationships between words.
- **Examples**: Identifying the subject, object, and verb in a sentence.
- **Key Characteristics**:
 - Involves syntactic analysis.
 - Requires labeled corpora with syntactic trees.
 - Uses models like shift-reduce parsers, graph-based models, and neural networks.

These tasks often overlap and are combined in various NLP applications, such as chatbots, sentiment analysis tools, translation systems, and more. The choice of task, model, and preprocessing steps largely depends on the specific application and the nature of the data.

Whether **GPT or BERT** is better for a specific NLP task depends on the nature of the task and the characteristics of the models. Here's a comparison of GPT and BERT across the listed tasks:

1. Text Classification

- **Better Model: BERT**
- **Reason:** BERT is designed for tasks that require understanding the context of a word in both directions (bidirectional), which is highly beneficial for text classification tasks. It excels in extracting contextual information from text for classification purposes.

2. Named Entity Recognition (NER)

- **Better Model: BERT**
- **Reason:** BERT's bidirectional nature allows it to better understand the context surrounding entities, making it more effective for NER tasks, where understanding the full context is crucial.

3. Part-of-Speech (POS) Tagging

- **Better Model: BERT**
- **Reason:** BERT's ability to capture the contextual meaning of words makes it highly effective for POS tagging, which relies on understanding the role of each word in a sentence.

4. Sentiment Analysis

- **Better Model: BERT**
- **Reason:** BERT's fine-tuned models for sentiment analysis can capture the subtleties in language that indicate sentiment. Its bidirectional context understanding is especially useful in detecting sentiment that depends on the surrounding words.

5. Machine Translation

- **Better Model: GPT**

- **Reason:** GPT, especially when used in sequence-to-sequence models (like in some machine translation tasks), is powerful due to its autoregressive nature, which is beneficial for generating coherent translations.

6. Text Summarization

- **Better Model: GPT**
- **Reason:** GPT is more suited for generative tasks like text summarization, where creating new text that is coherent and contextually relevant is key. Its autoregressive approach helps in producing fluent summaries.

7. Question Answering (QA)

- **Better Model: BERT**
- **Reason:** BERT is specifically pre-trained on QA tasks (like in SQuAD), and its bidirectional nature makes it excellent at understanding and extracting relevant information from a given text to answer questions.

8. Speech Recognition (ASR)

- **Better Model: Neither BERT nor GPT is typically used for ASR.**
- **Reason:** Speech recognition tasks generally require specialized models like RNNs, LSTMs, or Transformer-based models specifically designed for audio processing, such as Wav2Vec.

9. Language Modeling

- **Better Model: GPT**
- **Reason:** GPT is a language model trained to predict the next word in a sequence, making it inherently better at tasks requiring generation of coherent text sequences.

10. Text Generation

- **Better Model: GPT**
- **Reason:** GPT is designed for generating text, and its autoregressive nature allows it to produce coherent and contextually relevant text, making it ideal for tasks like dialogue generation and creative writing.

11. Coreference Resolution

- **Better Model: BERT**
- **Reason:** BERT's bidirectional context understanding allows it to effectively resolve references by considering the full context of a text, which is crucial for accurate coreference resolution.

12. Semantic Role Labeling (SRL)

- **Better Model: BERT**
- **Reason:** BERT's ability to capture detailed context makes it better suited for SRL, where understanding the relationships between different parts of a sentence is essential.

13. Dependency Parsing

- **Better Model: BERT**
- **Reason:** BERT's strong contextual understanding helps in accurately parsing the grammatical structure of sentences, which is essential for dependency parsing.

Summary:

- **BERT** is generally better for tasks that require deep understanding of the context within a sentence, such as text classification, NER, POS tagging, sentiment analysis, QA, coreference resolution, SRL, and dependency parsing.
- **GPT** shines in generative tasks, such as text summarization, text generation, and language modeling, due to its autoregressive nature.

For most NLP tasks involving understanding or classifying text, BERT is superior, while GPT is preferable for tasks involving text generation.

Type of NLP Task – FROM EXAM

The task you are describing falls under **Text Classification.** Specifically, it involves classifying news articles into different categories such as sports, politics, finance, economy, and entertainment. This is a supervised learning problem where each news article (text input) is assigned a label representing its category.

Required Preprocessing Steps

1. **Data Collection**:

- Gather a labeled dataset of news articles categorized by topic (sports, politics, finance, etc.).

2. **Text Cleaning**:

- **Lowercasing**: Convert all text to lowercase to ensure uniformity.
- **Remove Punctuation**: Strip out punctuation marks.
- **Remove Stop Words**: Eliminate common words that do not carry significant meaning (e.g., "the," "and," "is").
- **Lemmatization/Stemming**: Reduce words to their base or root form (e.g., "running" becomes "run").
- **Remove Special Characters/Numbers**: Strip out any special characters, digits, or non-alphanumeric symbols.

3. **Tokenization**:

- Split the text into individual words or tokens.

4. **Text Sequencing**:

- Convert the sequence of words (tokens) into a sequence of integers where each unique word corresponds to a unique integer.

5. **Padding**:

- Pad the sequences to ensure all input sequences are of the same length. This is essential for feeding into the LSTM model.

6. **Word Embeddings**:

- Use pre-trained word embeddings (e.g., GloVe, Word2Vec) or train your own embeddings. This step converts words into dense vector representations that capture semantic meaning.

Model Architecture

A typical LSTM-based model architecture for text classification might look like this:

1. **Embedding Layer**:

- Input: Sequences of integers (padded).
- Embedding layer that maps each word to a dense vector representation.
- Example: ``Embedding(input_dim=vocab_size, output_dim=embedding_dim, input_length=max_sequence_length)``

2. **LSTM Layer**:

- One or more LSTM layers to capture temporal dependencies and contextual information from the sequences.
- Example: ``LSTM(units=128, return_sequences=True)`` (for stacked LSTM) or ``LSTM(units=128, return_sequences=False)`` (for single LSTM)

3. **Dropout Layer**:

- Add dropout to prevent overfitting.
- Example: ``Dropout(rate=0.5)``

4. **Fully Connected (Dense) Layer**:

- A dense layer to combine features learned by LSTM.
- Example: ``Dense(units=64, activation='relu')``

5. **Output Layer**:

- A dense layer with softmax activation for multi-class classification.
- Example: ``Dense(units=num_classes, activation='softmax')``

Training and Evaluation Process

1. **Train-Test Split**:

- Split the dataset into training, validation, and test sets (e.g., 70% training, 15% validation, 15% test).

2. **Compilation**:

- Choose an appropriate loss function and optimizer.
- Example: ``loss='categorical_crossentropy'`` for multi-class classification, ``optimizer='adam'``.

3. **Training**:

- Train the model using the training set and validate using the validation set.
- Track the loss and accuracy on both training and validation sets to monitor performance.
- Example: ``model.fit(X_train, y_train, validation_data=(X_val, y_val), epochs=20, batch_size=32)``

4. **Evaluation**:

- After training, evaluate the model on the test set to assess its generalization capability.
- Calculate metrics like accuracy, precision, recall, and F1-score.
- Example: `model.evaluate(X_test, y_test)`

Technique to Avoid Overfitting

****Dropout Regularization**:**

- Apply dropout layers within the network to randomly "drop out" neurons during training, which forces the network to learn more robust features.
- Example: `Dropout(rate=0.5)` applied after the LSTM or dense layers.

****Early Stopping**:**

- Monitor the validation loss during training, and if it stops improving for a predefined number of epochs (patience), stop the training early to prevent overfitting.
- Example: Using `EarlyStopping(monitor='val_loss', patience=3)`.

****Cross-Validation**:**

- Utilize k-fold cross-validation during training to ensure that the model is not overfitting to a particular subset of the data.

Summary

By implementing these steps, you can build an effective LSTM-based model for news classification, ensuring it generalizes well to unseen data. Key strategies like dropout, early stopping, and careful model tuning will help mitigate the risk of overfitting, leading to a robust and accurate model suitable for deployment in your fintech application.

ERROR TYPES FOR REG/CLASS

Single-Variable Regression: MSE, MAE, and RMSE are all applicable. MSE and RMSE are preferred if penalizing larger errors is important, while MAE is better if robustness to outliers is needed.

Multi-Variable Regression: MSE, MAE, and RMSE are suitable, with MSE and RMSE providing comprehensive error measures across multiple outputs. MAE offers robustness to outliers, and R-squared can be useful for evaluating the model's explanatory power.

When you're using Python to check your classification models, the best metrics to use will depend on your specific classification task. Here's a practical guide for choosing the right metrics:

Single-Variable (Single-Class) Classification

- **Accuracy:** `accuracy_score(y_true, y_pred)`
 - Best for balanced datasets.
- **Precision and Recall:** `precision_score(y_true, y_pred)`, `recall_score(y_true, y_pred)`
 - Use these when you need to minimize false positives or false negatives.
- **F1 Score:** `f1_score(y_true, y_pred)`
 - Ideal for imbalanced datasets where you need a balance between precision and recall.
- **AUC-ROC:** `roc_auc_score(y_true, y_prob)`
 - Great for evaluating the model's ability to distinguish between classes.
- **Log Loss:** `log_loss(y_true, y_prob)`
 - Useful for probabilistic models where you care about prediction confidence.

Multi-Variable (Multi-Class) Classification

- **Accuracy:** `accuracy_score(y_true, y_pred)`
 - Still useful, but consider class balance.
- **Macro-Averaged Precision, Recall, F1 Score:** `precision_score(y_true, y_pred, average='macro')`, `recall_score(y_true, y_pred, average='macro')`, `f1_score(y_true, y_pred, average='macro')`
 - Treats all classes equally, good when you want a balanced view across classes.
- **Micro-Averaged Precision, Recall, F1 Score:** `precision_score(y_true, y_pred, average='micro')`, `recall_score(y_true, y_pred, average='micro')`, `f1_score(y_true, y_pred, average='micro')`
 - Good when class sizes are very different, as it considers the overall class balance.
- **Weighted-Averaged Precision, Recall, F1 Score:** `precision_score(y_true, y_pred, average='weighted')`, `recall_score(y_true, y_pred, average='weighted')`, `f1_score(y_true, y_pred, average='weighted')`
 - Useful when classes are imbalanced, giving more weight to larger classes.
- **AUC-ROC (One-vs-Rest):** `roc_auc_score(y_true, y_prob, multi_class='ovr')`

- Evaluates the ability to distinguish between classes.

Binary Classification

- **Accuracy:** `accuracy_score(y_true, y_pred)`
 - Standard for balanced datasets.
- **Precision, Recall, and F1 Score:** `precision_score(y_true, y_pred)`, `recall_score(y_true, y_pred)`, `f1_score(y_true, y_pred)`
 - Essential for imbalanced datasets.
- **AUC-ROC:** `roc_auc_score(y_true, y_prob)`
 - Important when dealing with probabilities and class separation.
- **Log Loss:** `log_loss(y_true, y_prob)`
 - Key for probabilistic models.
- **Confusion Matrix:** `confusion_matrix(y_true, y_pred)`
 - To get a detailed look at true positives, false positives, true negatives, and false negatives: `plot_confusion_matrix` from `sklearn`.

Implementation in Python:

```
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score,
roc_auc_score, log_loss, confusion_matrix, plot_confusion_matrix
```

These metrics cover most scenarios, and you can choose based on your model's specifics and the type of errors that matter most to your task.

RELU/SOFTMAX/SIGMOID

Summary

Hidden Layers:

ReLU is commonly used due to its effectiveness in deep networks and simplicity.

Tanh or Sigmoid might be used in specific scenarios, such as in RNNs or when outputs are normalized.

Output Layers:

Linear Activation: For regression tasks where the output is continuous.

Softmax: For multi-class classification where you need to predict probabilities for multiple classes.

Sigmoid: For binary classification where the output is a probability of the positive class.

```
model = Sequential()
model.add(LSTM(100, activation='relu', input_shape=(X_train.shape[1], X_train.shape[2]),
return_sequences=True))
model.add(LSTM(50, activation='relu'))
model.add(Dense(1, activation='linear')) # Suitable for regression tasks
```

In this example, ReLU is used in the hidden layers of the LSTM network to introduce non-linearity and help with learning complex patterns.

Linear activation in the output layer is chosen because the task is a regression problem, where the goal is to predict continuous values.

OPTIMISERS

For General Use

Adam:

Use case: General-purpose optimizer, often used for a wide range of problems due to its adaptive learning rates.

Example: Good for most neural network models, including those with large datasets.

SGD:

Use case: Basic optimizer, good for simple models or when fine control over learning rates is needed.

Example: Useful for traditional machine learning models or when combined with learning rate schedules.

For Specific Scenarios

Momentum:

Use case: Accelerates gradient descent and helps navigate local minima.

Example: Effective when dealing with complex loss surfaces.

NAG (Nesterov Accelerated Gradient):

Use case: Provides a lookahead step to improve convergence speed.

Example: Useful in cases where momentum alone is not enough.

RMSprop:

Use case: Adapts learning rates for each parameter, especially useful for RNNs.

Example: Ideal for problems with noisy gradients or varying feature scales.

Adagrad:

Use case: Adjusts learning rates based on gradient history, good for sparse data.

Example: Suitable for tasks like text or image classification with sparse features.

Adadelat:

Use case: Reduces aggressive learning rate decay of Adagrad.

Example: Helpful when Adagrad's learning rate becomes too small.

FTRL (Follow The Regularized Leader):

Use case: Designed for online learning and large-scale problems.

Example: Effective for models that need to update frequently with new data.

Nadam:

Use case: Combines Adam and NAG benefits, improving convergence.

Example: Useful when both adaptive learning and momentum benefits are needed.

These optimizers are chosen based on the nature of the problem, the model architecture, and the characteristics of the dataset.

Random Forest Algorithm Overview

Random Forest is an ensemble learning algorithm that combines multiple decision trees to improve the overall predictive performance and robustness compared to individual decision trees. It is widely used for both classification and regression tasks.

Key Components of Random Forest

1. Decision Trees

- A decision tree is the fundamental building block of a random forest. It is a tree-like structure where internal nodes represent decisions based on feature values, branches represent outcomes of these decisions, and leaf nodes represent final predictions.
- In a random forest, multiple decision trees are trained, each on a different subset of the data.

2. Bootstrap Aggregating (Bagging)

- Bagging is a technique used to create multiple subsets of the original dataset by sampling with replacement. This means that some data points might appear multiple times in a single subset, while others might not appear at all.
- Each decision tree in the random forest is trained on a different bootstrap sample. This introduces variability among the trees, reducing the likelihood that the model will overfit to a particular subset of the data.

3. Random Feature Selection

- When splitting nodes in each decision tree, instead of considering all available features, random forest considers only a random subset of features. This further decorrelates the trees, making the ensemble less prone to overfitting.
- The random selection of features ensures that the trees are diverse and reduces the chance of a single dominant feature influencing the model too much.

Advantages of Random Forest Compared to Individual Decision Trees

1. Improved Accuracy

- By averaging the predictions of multiple trees, random forest reduces the variance of the model. This leads to better generalization to unseen data compared to a single decision tree, which might overfit the training data.

2. Robustness

- Random forest is less sensitive to outliers and noisy data because the decision trees in the ensemble are built on different subsets of the data. This reduces the impact of any single outlier or noise in the training data on the overall model.

3. Reduced Overfitting

- Individual decision trees are prone to overfitting, especially when they are deep and complex. Random forest mitigates this by averaging across multiple trees, making the overall model more robust and less likely to fit noise in the data.

4. Handling High Dimensionality

- Random forest can handle large datasets with many features effectively. The random feature selection process ensures that not all features are used for every tree, making the algorithm computationally efficient and better at handling datasets with high dimensionality.

Feature Selection and Generalization in Random Forest

1. Feature Importance

- Random forest provides a measure of feature importance by averaging the decrease in impurity (like Gini impurity or entropy) across all trees in the forest. Features that consistently contribute to making good splits in the decision trees are deemed more important.

- This feature importance metric can be used for feature selection, helping to identify and focus on the most relevant features in the dataset.

**Summary

- **Yes**, Random Forest algorithms **reduce Entropy and Gini Impurity** to create effective splits in decision trees.
- **Impurity reduction** is fundamental to how decision trees and Random Forests **learn patterns from data**.
- **Feature importance** is derived from **aggregating impurity reductions** contributed by each feature across all trees, helping in **model interpretation and feature selection**.
- Choosing between **Entropy and Gini** depends on specific use-cases and sometimes computational efficiency, but both aim to **maximize the purity of child nodes** after each split.

2. Generalization

- The use of multiple trees, each built on a different subset of data and with a random subset of features, allows the random forest to generalize better to unseen data. This means that random forest models tend to perform well not just on the training data, but also on new, unseen data, reducing the risk of overfitting.

- The aggregation of predictions (either by voting in classification or averaging in regression) ensures that the random forest produces more stable and accurate predictions than individual decision trees.

Conclusion

Random forest is a powerful and versatile algorithm that leverages the strengths of multiple decision trees to improve predictive accuracy and robustness. By incorporating bootstrap aggregating (bagging) and random feature selection, it reduces the risk of overfitting, handles high-dimensional data effectively, and provides insights into feature importance. These advantages make random forest a popular choice for a wide range of machine learning tasks.

Similarities Between R2 Score and F1 Score:

- **Evaluation Metrics:** Both are metrics used to evaluate the performance of machine learning models.
- **Scale:** Both metrics range from 0 to 1, where higher values indicate better model performance.
- **Model Assessment:** Both provide a way to assess how well a model is performing on a given task, whether regression or classification.

Differences Between R2 Score and F1 Score:

- **Type of Task:**
 - **R2 Score:** Used for regression tasks.
 - **F1 Score:** Used for classification tasks.

- **Nature of Prediction:**
 - **R2 Score:** Measures how well the model's predictions match continuous, numeric data.
 - **F1 Score:** Measures the balance between precision and recall for binary or multiclass classification.
- **Interpretation:**
 - **R2 Score:** Indicates the proportion of variance in the dependent variable explained by the model.
 - **F1 Score:** Represents the harmonic mean of precision and recall, focusing on the balance between them.
- **Sensitivity to Data:**
 - **R2 Score:** Sensitive to how well the model captures the variance in continuous outcomes.
 - **F1 Score:** Sensitive to imbalanced data, where precision and recall are critical.
- **Calculation:**
 - **R2 Score:** Based on the sum of squared differences between actual and predicted values.
 - **F1 Score:** Derived from precision and recall, focusing on true positives, false positives, and false negatives.

Use Cases:

- **R2 Score:**
 - **Regression Analysis:** Predicting continuous outcomes, such as house prices, sales, or temperatures.
 - **Model Fit:** Assessing how well a regression model explains the variance in the target variable.
- **F1 Score:**
 - **Classification Tasks:** Especially useful in binary classification, such as spam detection or fraud detection.
 - **Imbalanced Data:** Situations where positive and negative classes are not equally represented, like rare disease detection.

Key Characteristics:

- **R2 Score:**
 - **Proportion of Variance:** Indicates how much of the total variation in the dependent variable is captured by the model.
 - **Values:** Ranges from 0 to 1, where 1 means perfect prediction.
- **F1 Score:**
 - **Precision and Recall Balance:** Balances the need for accurate positive predictions (precision) with the ability to capture all actual positives (recall).
 - **Harmonic Mean:** More sensitive to lower values of precision and recall compared to the arithmetic mean.

1. Accuracy

- **Definition:** Accuracy is the ratio of correctly predicted instances (both true positives and true negatives) to the total number of instances.

- **Formula:**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

-
- **When to Use:** Best for balanced datasets where the cost of false positives and false negatives is similar.
- **Limitations:** Can be misleading in imbalanced datasets, as it might not reflect the true performance of the model.

2. Precision

- **Definition:** Precision, also known as Positive Predictive Value (PPV), is the ratio of true positive predictions to the total number of positive predictions (i.e., how many selected items are relevant).

- **Formula:**

$$\text{Precision} = \frac{TP}{TP + FP}$$

-
- **When to Use:** Useful when the cost of false positives is high (e.g., spam detection, where marking a legitimate email as spam is costly).
- **Limitations:** Ignores false negatives, so it doesn't account for how many actual positives were missed.

3. Recall (Sensitivity, True Positive Rate)

- **Definition:** Recall is the ratio of true positive predictions to the total number of actual positives (i.e., how many relevant items are selected).

- **Formula:**

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **When to Use:** Important when the cost of false negatives is high (e.g., medical diagnosis, where missing a positive case can have serious consequences).
- **Limitations:** Ignores false positives, so it doesn't account for how many incorrect positive predictions were made.

4. F1 Score

- **Definition:** The F1 score is the harmonic mean of precision and recall, providing a single metric that balances both concerns.
- **Formula**

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- **When to Use:** Best used when you need to balance precision and recall, particularly in situations with imbalanced datasets.
- **Limitations:** Doesn't capture the true negatives, so it might not fully represent the performance of the model.

5. Negative Predictive Value (NPV)

- **Definition:** NPV is the ratio of true negative predictions to the total number of negative predictions (i.e., how many identified negative items are relevant).
- **Formula:**

$$\text{NPV} = \frac{TN}{TN + FN}$$

- **When to Use:** Useful when the cost of false negatives is particularly high, and you want to ensure that negative predictions are reliable.
- **Limitations:** Ignores true positives and false positives, so it doesn't give a complete picture of the model's performance.

6. Specificity (True Negative Rate)

- **Definition:** Specificity measures the proportion of actual negatives that are correctly identified (i.e., the ratio of true negatives to all actual negatives).
- **Formula:**

$$\text{Specificity} = \frac{TN}{TN + FP}$$

- **When to Use:** Important in scenarios where correctly identifying the negative class is critical (e.g., screening tests where false positives need to be minimized).

- **Limitations:** Focuses only on the negative class and ignores the positive class performance.

7. Balanced Accuracy

- **Definition:** Balanced accuracy is the average of sensitivity and specificity. It is particularly useful when dealing with imbalanced datasets.
- **Formula:**

$$\text{Balanced Accuracy} = \frac{\text{Sensitivity} + \text{Specificity}}{2}$$

- **When to Use:** Useful when class imbalance is present, providing a more balanced view than simple accuracy.
- **Limitations:** May still be influenced by class distribution, though less so than raw accuracy.

Similarities and Differences

- **Similarities:**
 - All metrics rely on the confusion matrix (TP, TN, FP, FN) to calculate their values.
 - Precision, Recall, F1 Score, and Specificity are all related to the model's ability to correctly identify the positive and negative classes.
 - Metrics like Accuracy, F1 Score, and ROC-AUC aim to provide a single number that summarizes the model's performance, though they do so in different ways.
- **Differences:**
 - **Focus:** Different metrics emphasize different types of errors. Precision focuses on false positives, Recall focuses on false negatives, Specificity focuses on true negatives, etc.
 - **Applicability:** Some metrics (like ROC-AUC) are better suited for comparing models across different thresholds, while others (like F1 Score) are better for summarizing performance in specific situations.
 - **Sensitivity to Imbalance:** Accuracy can be misleading in imbalanced datasets, whereas metrics like Precision, Recall, and MCC provide a clearer picture in such cases.

Summary

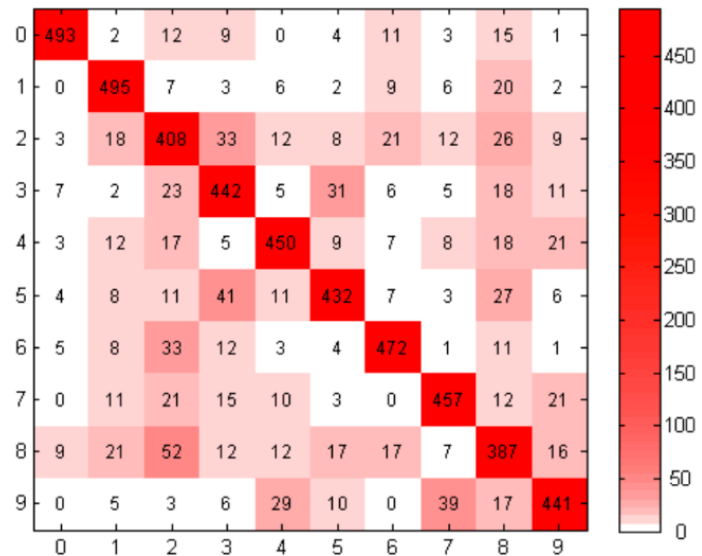
- **Use Accuracy** when the dataset is balanced and the cost of different errors is similar.
- **Use Precision** when false positives are costly.
- **Use Recall** when false negatives are costly.
- **Use F1 Score** when you need to balance Precision and Recall.
- **Use Specificity** when identifying negatives correctly is critical.

- Use **ROC-AUC** when comparing models across different thresholds.
- Use **MCC** when you need a balanced metric, especially for imbalanced datasets.
- Use **Balanced Accuracy** when dealing with class imbalance.

Each metric serves a specific purpose and should be chosen based on the context of the problem you're addressing.

Confusion matrix

- Is an error matrix, that allows visualization of the performance of an algorithm. Each row of the matrix represents the instances in a predicted class while each column represents the instances in an actual class (or vice versa).



Користиме confusion matrix кога не можеме да ги користиме стандардните recall/accuracy/precision..., овој случај најчесто се јавува кога имаме повеќе класи и ваквите стандардни метрики не можат да се користат. Кога имаме голем overlap, значи дека моделот ја толкува едната класа како друга класа, т.е. имаме големи грешки. Треба да имаме overlap ако предвидената и реалната класа е тоа квадратче, а ако не е, значи предвидената и вистинската класа не се поклопуваат во моделот и можеби се поклопува каде што не треба. Овој пример е пример со бројки, можеби рачно напишаи моделот да ги предвиди, и тука можеме да видиме како на пр 7 многу се поклопува со 9 што е грешка.

Breaking Down the Differences Between merge, aggregate, and concatenate in Pandas

Understanding the differences between merge, aggregate, and concatenate is crucial for effective data manipulation in Pandas. Here's how each function works:

1. merge

- **Purpose:** The merge function is used to combine two DataFrames based on one or more keys (columns). It's similar to SQL joins and supports different types of joins like inner, outer, left, and right.
- **Use Case:** Use merge when you need to combine datasets based on a common column(s) and align rows accordingly.
- **Example:**

```
merged_data = data1.merge(data2, how='inner', on='Datetime')
```

This code merges data1 and data2 where the 'Datetime' column matches in both DataFrames.

2. aggregate

- **Purpose:** The aggregate function is used to perform aggregation operations on the columns of a DataFrame, such as calculating the sum, mean, or count.
- **Use Case:** Use aggregate when you want to apply one or more aggregation functions to each group in a DataFrame, often after a groupby operation.
- **Example:**

```
grouped_data = data.groupby('Category').aggregate({'Value': 'sum'})
```

This code sums the values in the 'Value' column for each unique 'Category'.

3. concatenate

- **Purpose:** The concatenate function, typically accessed via pd.concat, is used to combine two or more DataFrames along a particular axis (either rows or columns).
- **Use Case:** Use concatenate when you need to stack DataFrames either vertically (along rows) or horizontally (along columns).
- **Example:**

```
concatenated_data = pd.concat([data1, data2], axis=0)
```

This code stacks data1 and data2 vertically, with the rows of data2 added below data1.

Summary of Differences:

- **merge:** Use for joining DataFrames based on common columns.
- **aggregate:** Use for applying functions like sum or mean to groups within a DataFrame.
- **concatenate:** Use for stacking DataFrames either vertically or horizontally.

Each function serves a specific purpose in data manipulation, and choosing the right one depends on the operation you need to perform on your data.

THE PROF EXAMPLES ARE OLD, have to add tensorflow. (before all keras imports)

```
from tensorflow.keras.layers import Input, Dense, LSTM, Embedding, concatenate
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Sequential, Model
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
```

```
-----
ModuleNotFoundError                                Traceback (most recent call last)
<ipython-input-1-96d23b25bb57> in <cell line: 7>()
      5 from sklearn.preprocessing import MinMaxScaler
      6 from sklearn.model_selection import train_test_split
----> 7 from keras.preprocessing.text import Tokenizer
      8 from keras.preprocessing.sequence import pad_sequences
      9 from keras.models import Sequential, Model

ModuleNotFoundError: No module named 'keras.preprocessing.text'
```

Remember to scale the data before you try to use KNN!

```
[98] from sklearn.preprocessing import MinMaxScaler
      scaler = MinMaxScaler()
      scaled_data = scaler.fit_transform(encoded_df[['workclass', 'occupation']])
```

```
[99] from sklearn.impute import SimpleImputer, KNNImputer

      knn_imputer = KNNImputer(n_neighbors=5)
      imputed_data = knn_imputer.fit_transform(scaled_data)
      encoded_df[['workclass', 'occupation']] = scaler.inverse_transform(imputed_data)
```

Correlation statistics, very useful

```
correlation_matrix = encoded_df.corr()
print(correlation_matrix['income'].sort_values(ascending=True))
```

1. Regression

- **Scatter Plot:** Used to visualize the relationship between two continuous variables. Ideal for showing the trend and distribution of data points.
- **Line Plot:** Shows the trend over time or continuous values. Useful for time series data or continuous regression lines.
- **Residual Plot:** Displays residuals on the y-axis versus fitted values or another variable on the x-axis. Helps in diagnosing the fit of the regression model and identifying patterns in residuals.
- **Regression Line Plot:** Combines scatter and line plots to show the relationship between variables and the fitted regression line.

2. Classification

- **Confusion Matrix:** Shows the performance of a classification model by displaying true vs. predicted classifications. Helps to understand the errors and accuracy of the model.
- **ROC Curve (Receiver Operating Characteristic Curve):** Plots true positive rate vs. false positive rate. Useful for evaluating the performance of binary classifiers.
- **Precision-Recall Curve:** Shows the trade-off between precision and recall for different thresholds. Useful for imbalanced datasets.
- **Pair Plot:** Visualizes relationships between pairs of features in a dataset. Useful for exploring how different features interact and how classes are separated.

3. Binary Outcomes

- **Bar Plot:** Compares the count or proportion of binary outcomes (e.g., success vs. failure).
- **Histogram:** Shows the distribution of binary outcomes. Useful for understanding the frequency of each outcome.
- **Box Plot:** Compares distributions of continuous variables for binary outcomes. Useful for understanding the differences in distributions between the two classes.

4. Correlations

- **Heatmap:** Displays the correlation matrix of features, with colors representing correlation strength. Useful for identifying patterns and relationships between variables.
- **Pair Plot:** Also useful for visualizing pairwise correlations and distributions between features.

5. Comparisons

- **Box Plot:** Compares distributions across categories. Useful for understanding the spread, central tendency, and outliers in different groups.
- **Violin Plot:** Similar to a box plot but shows the density of the distribution. Useful for comparing distributions and understanding data spread.
- **Bar Plot:** Compares categorical data or aggregated values across groups. Useful for visualizing counts, means, or other summary statistics.
- **Line Plot:** Compares trends over time or categories. Useful for visualizing changes and patterns across different groups or time points.

6. General Visualization Tips

- **Histograms:** Good for visualizing the distribution of a single variable.

- **Pie Charts:** Used to show proportions of a whole, but often less preferred due to difficulty in comparing angles.
- **Scatter Plots with Regression Lines:** Useful for visualizing the relationship and fit of a regression model.

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
```

```
# Filter data for the two age groups
```

```
combined_df = pd.concat([
    encoded_df[(encoded_df['age'] >= 41) & (encoded_df['age'] <= 50)].assign(age_group='41-50'),
    encoded_df[(encoded_df['age'] >= 51) & (encoded_df['age'] <= 60)].assign(age_group='51-60')
])
```

```
# Calculate proportions for each income category
```

```
proportions = combined_df.groupby(['age_group', 'income']).size().reset_index(name='count')
proportions['proportion'] = proportions['count'] /
proportions.groupby('age_group')['count'].transform('sum')
```

```
# Visualization
```

```
plt.figure(figsize=(10, 6))
sns.barplot(x='age_group', y='proportion', hue='income', data=proportions)
plt.title('Proportion of Income Categories by Age Group')
plt.xlabel('Age Group')
plt.ylabel('Proportion')
plt.legend(title='Income')
plt.show()
```