
Cross validation methods in Deep Learning

Bhupen Sinha

25-07-2024

GROK KERS
AI FOR EVERYONE

TABLE OF CONTENTS

common cross-validation techniques.....	3
1. K-Fold Cross-Validation	3
2. Stratified K-Fold Cross-Validation	4
3. Leave-One-Out Cross-Validation (LOO)	4
4. Leave-P-Out Cross-Validation (LPO).....	5
5. Time Series Split.....	5
Best practices on cross validation	6
Useful Links and references	8

GROKWKERS
AI FOR EVERYONE

COMMON CROSS-VALIDATION TECHNIQUES

1. K-FOLD CROSS-VALIDATION

Description:

- K-Fold Cross-Validation splits the dataset into K equally sized folds. The model is trained K times, each time using K-1 folds for training and the remaining fold for validation. The final performance metric is the average of the metrics from each fold.

Purpose:

- To ensure that every data point gets to be in the validation set at least once and in the training set K-1 times.

Pros:

- Provides a better estimate of model performance compared to a single train-test split.
- Reduces variance in the performance estimate.

Cons:

- Computationally expensive, especially with large datasets and complex models.

When to Apply:

- When you have a moderate-sized dataset and you want a reliable estimate of model performance.

2. STRATIFIED K-FOLD CROSS-VALIDATION

Description:

- Similar to K-Fold but ensures that each fold has the same proportion of classes as the original dataset.

Purpose:

- To maintain the class distribution in each fold, especially important for imbalanced datasets.

Pros:

- Provides more accurate and stable performance estimates for imbalanced datasets.

Cons:

- Slightly more complex to implement than regular K-Fold.

When to Apply:

- When dealing with imbalanced datasets to ensure each fold is representative of the whole dataset.

3. LEAVE-ONE-OUT CROSS-VALIDATION (LOO)

Description:

- Each data point is used as a single validation sample while the remaining data points form the training set. This process is repeated for all data points.

Purpose:

- To provide an almost unbiased estimate of model performance.

Pros:

- Best use of data since it uses almost all the data points for training.

Cons:

- Extremely computationally expensive for large datasets.

When to Apply:

- For small datasets where it's feasible to train the model as many times as there are data points.

4. LEAVE-P-OUT CROSS-VALIDATION (LPO)

Description:

- Similar to LOO, but P data points are left out for validation each time. This process is repeated for all possible combinations.

Purpose:

- To provide a very thorough performance estimate, though more feasible for small values of P.

Pros:

- Uses nearly all data for training each time, giving a detailed performance estimate.

Cons:

- Computationally infeasible for large P or large datasets.

When to Apply:

- When P is small and the dataset is not too large.

5. TIME SERIES SPLIT

Description:

- Specifically designed for time series data where the order of data points matters. The data is split into train and validation sets while preserving the time order.

Purpose:

- To ensure that the validation set is always ahead in time compared to the training set, mimicking real-world scenarios.

Pros:

- Ensures the model is validated on future data, providing a realistic performance estimate.

Cons:

- Can be less effective with small datasets as the number of validation sets is limited.

When to Apply:

- When dealing with time series data or any data where the order of observations is important.

BEST PRACTICES ON CROSS VALIDATION

- **Use Stratified Sampling for Imbalanced Data**

When dealing with classification problems where class distribution is imbalanced, use stratified sampling methods such as Stratified K-Fold Cross-Validation.

Reason: Ensures each fold has a representative ratio of classes, preventing misleading performance metrics.

- **Maintain Temporal Order for Time Series Data**

Use time series-specific cross-validation methods like [TimeSeriesSplit](#).

- **Use Multiple Metrics**

Evaluate model performance using multiple metrics.

Reason: Different metrics can provide different insights into model performance, especially in classification tasks where accuracy might not tell the whole story.

- **Consistent Data Preprocessing**

Ensure that data preprocessing steps (e.g., scaling, encoding) are consistently applied within the cross-validation loop.

Inconsistent preprocessing can lead to data leakage and inaccurate performance estimates.

- **Use Sufficient Number of Folds**

Use an appropriate number of folds, typically 5 or 10.

Reason: Provides a good balance between bias and variance, offering a reliable estimate of model performance.

- **Avoid Overlapping Data Splits**

Ensure data splits do not overlap in a way that could lead to data leakage (e.g., ensuring training and validation sets are completely separate).

Reason: Prevents the model from learning patterns from the validation set, leading to overly optimistic performance estimates.

- **Account for Computational Efficiency**

Consider the computational cost of cross-validation, especially with complex models and large datasets.

Reason: Some cross-validation techniques (like Leave-One-Out) can be computationally expensive and impractical for large datasets.

- **Report the Mean and Standard Deviation of Metrics**

Report both the mean and standard deviation of performance metrics across all folds.

Reason: Provides a more comprehensive view of model performance and its stability.

- **Visualize Performance Across Folds**

Visualize the performance metrics across different folds.

Reason: Helps to understand the variance and stability of the model performance.

- **Consider the Impact of Random Seed**

Set a random seed for reproducibility.

Reason: Ensures that cross-validation results are consistent and reproducible.

USEFUL LINKS AND REFERENCES

-
-

GROKWKERS
AI FOR EVERYONE

INDEX

No index entries found.

GROKWKERS
AI FOR EVERYONE