# Ensemble Techniques & it's types Assignment - 1

March 12, 2024

#

Question 1

## 0.1 ## Question 1: What is an ensemble technique in machine learning?

## 0.2 Answer:

### 0.2.1 An ensemble technique in machine learning is a method that combines the predictions of multiple individual models to improve the accuracy and robustness of the overall prediction. Ensemble methods are often used when a single model may not be sufficient to accurately capture the complexity of the data or when multiple models with different strengths and weaknesses are available.

### 0.2.2 There are several types of ensemble methods, including:

1. Bagging: This method involves training multiple instances of the same model on different subsets of the training data and then combining the predictions of each model. Bagging can improve model stability and reduce overfitting.

2. Boosting: This method involves iteratively training a sequence of weak models, with each model focusing on the examples that were misclassified by the previous model. The predictions of each model are combined to produce the final prediction.

3. Stacking: This method involves training multiple different models and using their predictions as input to a higher-level model that learns how to combine them optimally. Stacking can be used to leverage the strengths of different types of models.

### 0.2.3 Ensemble methods have been shown to be very effective in many applications of machine learning, particularly in areas such as computer vision, natural language processing, and speech recognition.

#

Question 2

## 0.3 ## Question 2: Why are ensemble techniques used in machine learning?

## 0.4 Answer :

### 0.4.1 Ensemble techniques are used in machine learning for several reasons:

1. Improved prediction accuracy: Ensemble methods can often achieve higher accuracy than individual models by combining the predictions of multiple models. This is because different

models may be better suited to different aspects of the data, and combining their predictions can lead to a more accurate overall prediction.

2. Robustness: Ensemble methods can help to reduce the impact of noisy or incorrect data by averaging out errors across multiple models. This can improve the robustness of the overall prediction.

3. Reduced overfitting: Ensemble methods can help to reduce overfitting by combining the predictions of multiple models trained on different subsets of the data. This can improve generalization performance on new data.

4. Improved model stability: Ensemble methods can help to improve the stability of models by reducing the impact of small changes in the training data or model parameters.

### 0.4.2 Overall, ensemble methods are a powerful tool in machine learning that can help to improve prediction accuracy, robustness, and stability. They are widely used in many applications, including computer vision, natural language processing, and speech recognition.

\#

Question 3

## 0.5 \#\# Question 3 : What is bagging?

## 0.6 Answer :

### 0.6.1 Bagging (short for Bootstrap Aggregating) is an ensemble technique in machine learning that involves training multiple instances of the same model on different subsets of the training data and then combining their predictions to make a final prediction.

### 0.6.2 The basic idea behind bagging is to reduce the variance of a single model by averaging out the predictions of multiple models trained on different subsets of the data. Each model in the ensemble is trained on a randomly sampled subset of the training data, with replacement. This means that some examples may appear multiple times in the same subset, while others may not appear at all. By sampling with replacement, we can create multiple different subsets of the data that are representative of the overall training set.

### 0.6.3 Once the models have been trained on their respective subsets of the training data, they are used to make predictions on the test data. The final prediction is then obtained by aggregating the predictions of all the models, usually by taking the average.

### 0.6.4 Bagging can be used with any type of model, including decision trees, neural networks, and support vector machines. It is particularly effective with models that are prone to overfitting, as it can help to reduce the variance of the model and improve its generalization performance.

\#

Question 4

## 0.7 ## Question 4 : What is boosting?

## 0.8 Answer :

**0.8.1** Boosting is an ensemble technique in machine learning that involves iteratively training a sequence of weak models, with each model focusing on the examples that were misclassified by the previous model. The predictions of each model are then combined to produce the final prediction.

**0.8.2** The basic idea behind boosting is to iteratively improve the overall prediction by focusing on the examples that are difficult to classify. The weak models in the ensemble are typically simple models, such as decision trees, that are trained on a subset of the training data. The models are then combined using a weighted sum, where the weights are adjusted at each iteration to give more weight to the models that perform well on the training data.

**0.8.3** Boosting can be thought of as a form of adaptive weighting, where the weights of the examples in the training data are adjusted based on their difficulty to classify. Examples that are difficult to classify are given higher weights, while examples that are easy to classify are given lower weights. This allows the weak models to focus on the examples that are most important for improving the overall prediction.

**0.8.4** Boosting has been shown to be very effective in many applications of machine learning, particularly in areas such as computer vision, natural language processing, and speech recognition. It is often used with decision trees, but can also be used with other types of models, such as neural networks and support vector machines.

\#

Question 5

## 0.9 ## Question 5 : What are the benefits of using ensemble techniques?

## 0.10 Answer :

### 0.10.1 Ensemble techniques offer several benefits in machine learning:

1. Improved accuracy: Ensemble techniques can improve the accuracy of predictions by combining the predictions of multiple models. This is because different models may be better suited to different aspects of the data, and combining their predictions can lead to a more accurate overall prediction.

2. Reduced overfitting: Ensemble techniques can help to reduce overfitting by combining the predictions of multiple models trained on different subsets of the data. This can improve generalization performance on new data.

3. Improved robustness: Ensemble techniques can help to reduce the impact of noisy or incorrect data by averaging out errors across multiple models. This can improve the robustness of the overall prediction.

4. Increased stability: Ensemble techniques can help to improve the stability of models by reducing the impact of small changes in the training data or model parameters.

5. Can leverage different strengths of models: Ensemble techniques can combine the strengths of different models that are good at different aspects of the data and combine their predictions to get a more accurate overall prediction.

**0.10.2 Overall, ensemble techniques are a powerful tool in machine learning that can improve prediction accuracy, robustness, stability, and generalization performance. They are widely used in many applications, including computer vision, natural language processing, and speech recognition.**

#

Question 6

## 0.11 ## Question 6 : Are ensemble techniques always better than individual models?

## 0.12 Answer :

**0.12.1 Ensemble techniques are not always better than individual models. While ensemble techniques can improve the accuracy, robustness, and generalization performance of models, there are cases where they may not be effective.**

**0.12.2 For example, if the individual models in the ensemble are very similar to each other, then the ensemble may not provide much benefit over a single model. Additionally, if the individual models are very complex and overfit the training data, then the ensemble may also overfit the data and not generalize well to new data.**

**0.12.3 Furthermore, ensemble techniques can be computationally expensive and may require more resources to train and deploy compared to individual models.**

**0.12.4 Therefore, it is important to carefully evaluate the effectiveness of ensemble techniques in a given application and compare their performance against individual models before deciding to use them. It is also important to consider the computational cost and feasibility of implementing ensemble techniques in the given application.**

#

Question 7

## 0.13 ## Question 7 : How is the confidence interval calculated using bootstrap?

## 0.14 Answer :

**0.14.1 The confidence interval can be calculated using the bootstrap method in the following way:**

1. Randomly sample the original dataset with replacement to create a new dataset of the same size as the original dataset. This is called a bootstrap sample.

2. Compute the statistic of interest, such as the mean or median, on the bootstrap sample.

3. Repeat steps 1 and 2 many times, typically several thousand times, to generate a large number of bootstrap samples and their corresponding statistics.

4. Calculate the standard error of the statistic using the bootstrap samples. This is an estimate of the variability of the statistic.

5. Use the standard error to calculate the confidence interval using the desired level of confidence and a normal distribution or t-distribution, depending on the sample size and assumptions.

### 0.14.2 For example, suppose we want to calculate the 95% confidence interval for the mean of a dataset. We can use the bootstrap method to estimate the standard error of the mean and then calculate the confidence interval using a normal distribution. The steps would be as follows:

1. Randomly sample the original dataset with replacement to create a new dataset of the same size as the original dataset. Repeat this process many times, such as 10,000 times, to create a large number of bootstrap samples.

2. Calculate the mean of each bootstrap sample.

3. Compute the standard error of the mean using the bootstrap samples. This is typically calculated as the standard deviation of the bootstrap sample means.

4. Calculate the lower and upper bounds of the confidence interval using the standard error and a normal distribution. For example, if the standard error is 0.1 and we want a 95% confidence interval, we would use a z-score of 1.96 (the z-score corresponding to the 97.5th percentile of a standard normal distribution) and calculate the confidence interval as: mean +/- zstandard error = mean +/- 1.96*std_error

5. The resulting confidence interval would be the range of values from the lower bound to the upper bound.

#

Question 8

## 0.15 ## Question 8 : How does bootstrap work and What are the steps involved in bootstrap?

## 0.16 Answer :

### 0.16.1 Bootstrap is a statistical method for estimating the uncertainty of a statistic or model parameter by resampling from the original data. The steps involved in the bootstrap method are as follows:

1. Collect the original dataset: The first step is to collect the original dataset of size N.

2. Sample from the original dataset with replacement: From the original dataset, we randomly sample N observations with replacement to create a new dataset of the same size N as the original dataset. This is called a bootstrap sample.

3. Compute the statistic of interest: We compute the statistic of interest, such as the mean, median, variance, or any other parameter we want to estimate, on the bootstrap sample.

4. Repeat step 2 and 3 many times: We repeat steps 2 and 3 many times, typically 1,000 or more, to create many bootstrap samples and their corresponding statistics.

5. Compute the standard error of the statistic: We compute the standard error of the statistic using the bootstrap samples. This is an estimate of the variability of the statistic.

6. Construct the confidence interval: We use the standard error to construct a confidence interval around the estimate of the statistic. For example, we can construct a 95% confidence interval by taking the middle 95% of the bootstrap distribution of the statistic.

**0.16.2  The bootstrap method allows us to estimate the variability of a statistic or parameter without making assumptions about the underlying distribution of the data or the model. It can be used for a wide range of statistical applications, such as hypothesis testing, regression analysis, and machine learning. The bootstrap method is particularly useful when the sample size is small, or the underlying distribution is unknown or nonparametric.**

#

Question 9

**0.17  ## Question 9 : A researcher wants to estimate the mean height of a population of trees. They measure the height of a sample of 50 trees and obtain a mean height of 15 meters and a standard deviation of 2 meters. Use bootstrap to estimate the 95% confidence interval for the population mean height.**

**0.18  Answer :**

**0.18.1  Below is code in python to calculate 95% interval for above bootstrap**

```python
# Given Data
samples = 50
sample_mean = 15
sample_std = 2
confidence_level = 0.95

# Calculate the t value for desired level of confidence
import scipy.stats as stats
alpha = 1 - confidence_level
dof = samples-1
t_value = stats.t.ppf(1 - alpha/2, dof)

# calculate the standard error and margin of error
import math
std_error = sample_std / math.sqrt(samples)
margin_of_error = t_value * std_error

# calculate the confidence interval bounds
lower_bound = sample_mean - margin_of_error
```

```python
upper_bound = sample_mean + margin_of_error

# print 95% confidence interval
print(f'Sample mean height for {samples} Trees is {sample_mean} and Sample
  ↪Standard Deviation is {sample_std}')
print('\n==========================================================================\n')
print(f'T-Statistic with {confidence_level*100}% condifence interval for dof
  ↪{dof} : {t_value:.4f}')
print(f'Standard Error : {std_error:.4f}')
print(f'Margin of error : {margin_of_error:.4f}')
print(f'\nEstimated Population mean with 95% confidence interval is
  ↪({lower_bound:.2f} , {upper_bound:.2f})')
```

Sample mean height for 50 Trees is 15 and Sample Standard Deviation is 2

========================================================================

T-Statistic with 95.0% condifence interval for dof 49 : 2.0096
Standard Error : 0.2828
Margin of error : 0.5684

Estimated Population mean with 95% confidence interval is (14.43 , 15.57)