# 16 March Assignment

May 21, 2023

[ ]: Q1: Define overfitting **and** underfitting **in** machine learning. What are the␣
    ↪consequences of each, **and** how
    can they be mitigated?

ANS -

[ ]:    Overfitting **and** underfitting are two common problems **in** machine learning␣
    ↪that can lead to poor model performance. Overfitting occurs
    when a model **is** too **complex** **and** fits the training data too well, resulting **in**␣
    ↪poor generalization to new data. Underfitting occurs when a
    model **is** too simple **and** cannot capture the underlying trend of the data,␣
    ↪resulting **in** poor performance on both training **and** test data.
       The consequences of overfitting are that the model will perform well on the␣
    ↪training data but poorly on new data. This **is** because the
    model has learned the noise **in** the training data instead of the underlying␣
    ↪pattern. The consequences of underfitting are that the model will
    perform poorly on both training **and** test data because it cannot capture the␣
    ↪underlying trend of the data.
       To mitigate overfitting, one can use regularization techniques such **as** L1␣
    ↪**or** L2 regularization, dropout, **or** early stopping.
    To mitigate underfitting, one can use more **complex** models **or** increase the␣
    ↪number of features used by the model.

[ ]:

[ ]:

[ ]: Q2: How can we reduce overfitting? Explain **in** brief.

ANS -

[ ]:       Overfitting **is** a common problem **in** machine learning where a model **is**␣
    ↪trained too well on the training data **and** performs poorly on the
    test data. There are several ways to reduce overfitting **in** machine learning.␣
    ↪Some of them are:

    1.Cross-validation: Cross-validation **is** a powerful preventative measure against␣
    ↪overfitting. The idea **is** clever: Use your training data to

```
generate multiple mini train-test splits. Then, use these splits to tune your␣
 ↪model.

2.Train with more data: It won't work every time, but training with more data␣
 ↪can help algorithms detect the signal better.

3.Remove features: Remove irrelevant features from the dataset.

4.Early stopping: Stop training once the performance on a validation set starts␣
 ↪decreasing.

5.Regularization: Regularization is a technique used to reduce the complexity␣
 ↪of the model by adding a penalty term to the loss function.
 This penalty term discourages the model from fitting the training data too␣
 ↪closely.
```

[ ]:

[ ]:

[ ]:
```
Q3: Explain underfitting. List scenarios where underfitting can occur in ML.
```

ANS -

[ ]:
```
        Underfitting occurs when a model is too simple to capture the␣
 ↪underlying patterns in the data. In other words, the model is not
complex enough to fit the data well, resulting in poor performance on both the␣
 ↪training and test data. Underfitting can occur when the model
is not trained for long enough or when the model architecture is too simple.

Reasons for Underfitting:

-High bias and low variance
-The size of the training dataset used is not enough.
-The model is too simple.
-Training data is not cleaned and also contains noise in it.
```

[ ]:

[ ]:

[ ]:
```
Q4: Explain the bias-variance tradeoff in machine learning. What is the␣
 ↪relationship between bias and
variance, and how do they affect model performance?
```

ANS -

```
The bias-variance tradeoff is a fundamental concept in machine learning
that affects a supervised model's predictive performance and
accuracy. It is the tradeoff between a model's ability to fit the training data
set well (low bias) and its ability to generalize well to
new data (low variance). A model with high bias will underfit the training data
and have poor performance on both the training and test sets.
A model with high variance will overfit the training data and have good
performance on the training set but poor performance on the test set.
The goal is to find an optimal balance between bias and variance that minimizes
the total error of the model .
    In summary, bias refers to how well a model fits the training data, while
variance refers to how much the model's predictions vary for
different training sets. High bias means that the model is too simple and
cannot capture the complexity of the data. High variance means that
the model is too complex and captures noise in the training data. The optimal
balance between bias and variance depends on the specific
problem being solved.
```

```
Q5: Discuss some common methods for detecting overfitting and underfitting in
machine learning models.
How can you determine whether your model is overfitting or underfitting?
```

ANS -

```
    There are several ways to detect over- or under-fitting in a machine
learning model:

1.Plot the learning curves: Learning curves show the model's performance on
training and validation data over time as the model is being
trained. If the model is overfitting, the learning curve will show a low error
on the training data and a high error on the validation data.

2.Evaluate the model on a holdout set: A holdout set is a subset of the data
that is not used during training but is used to evaluate the
model after training. If the model is overfitting, it will perform well on the
training data but poorly on the holdout set.

3.Use cross-validation: Cross-validation is a technique that involves dividing
the data into multiple subsets and training the model on each
subset while using the remaining subsets for validation. If the model is
overfitting, it will perform well on the training subsets but poorly
on the validation subsets.
```

```
    To determine whether your model is overfitting or underfitting, you can
  ↪use these methods to evaluate its performance. If your model has
high accuracy on both training and validation sets, it is likely not
  ↪underfitting or overfitting. If your model has high accuracy on the
training set but low accuracy on the validation set, it is likely overfitting.
  ↪If your model has low accuracy on both sets, it is likely
underfitting.
```

[ ]:

[ ]:

[ ]: ```
Q6: Compare and contrast bias and variance in machine learning. What are some
  ↪examples of high bias
and high variance models, and how do they differ in terms of their performance?
```

ANS -

[ ]: ```
   In machine learning, bias refers to the difference between a model's
  ↪predictions and the actual distribution of the value it tries to
predict. Models with high bias oversimplify the data distribution rule/
  ↪function, resulting in high errors in both the training outcomes and
test data analysis results. Variance refers to the amount by which the model's
  ↪prediction would change if we trained it on a different data
set. Models with high variance are too sensitive to changes in the training
  ↪data and may overfit the training data, leading to poor
generalization performance on new data.
   A model that exhibits small variance and high bias will underfit the target,
  ↪while a model with high variance and little bias will overfit
the target. A model with high variance may represent the data set accurately
  ↪but could lead to overfitting to noisy or otherwise
unrepresentative training data.
   Some examples of machine learning algorithms with low bias are Decision
  ↪Trees, k-Nearest Neighbours and Support Vector Machines.
At the same time, an algorithm with high bias is Linear Regression, Linear
  ↪Discriminant Analysis and Logistic Regression.
```

[ ]:

[ ]:

[ ]: ```
Q7: What is regularization in machine learning, and how can it be used to
  ↪prevent overfitting? Describe
```

```
some common regularization techniques and how they work.
```

ANS -

```
[ ]:    Regularization is a technique used in machine learning to prevent␣
     ↪overfitting by adding extra information to the dataset.
     In the context of machine learning, regularization is the process which␣
     ↪regularizes or shrinks the coefficients towards zero.
     It discourages learning a more complex or flexible model to prevent overfitting.

     There are three commonly used regularization techniques to control the␣
     ↪complexity of machine learning models: L2 regularization,
     L1 regularization and Elastic Net.

     L2 regularization is also known as Ridge regression. It is used to reduce the␣
     ↪complexity of the model by introducing a small amount of bias
     so that we can get better long-term predictions.

     L1 regularization is also known as Lasso regularization. It is used to reduce␣
     ↪the complexity of the model by introducing sparsity in the model.
     It helps in feature selection by shrinking some coefficients to zero.

     Elastic Net regularization is a combination of L1 and L2 regularization. It␣
     ↪helps in feature selection and reduces the complexity of the
     model by shrinking some coefficients to zero while keeping others non-zero.
```

```
[ ]:
```