## MODULE 03

1. **EXPLAIN WORKINNG PRINCIPLE OF SUPPOR VECTOR MACHINE(SVM)**

Support Vector Machine (SVM) is a powerful supervised machine learning algorithm used for both classification and regression tasks. It works by finding the optimal hyperplane that best separates data points of different classes in a high-dimensional feature space. The key principles behind the working of Support Vector Machine are:

1. *Linear Separability:* SVM is designed to work with data that can be separated into different classes using a straight line (in 2D), a plane (in 3D), or a hyperplane (in higher dimensions). If the data is not linearly separable, SVM uses a technique called the kernel trick to transform the data into a higher-dimensional space where it might become linearly separable.

2. *Margin Maximization:* SVM aims to find the hyperplane that maximizes the margin between the closest data points of different classes. The margin is the perpendicular distance between the hyperplane and the data points (support vectors) nearest to the hyperplane. The idea behind maximizing the margin is to create a larger "safety margin" between the classes, which helps in improving the generalization ability of the SVM.

3. *Support Vectors:* The data points that lie closest to the hyperplane are known as support vectors. These points play a crucial role in defining the hyperplane and the margin. They are the critical points that determine the optimal solution.

4. *Soft Margin:* In real-world scenarios, data is often not perfectly separable. To handle such cases, SVM allows for a soft margin, which means that some misclassification errors are tolerated to achieve a better overall solution. The balance between maximizing the margin and allowing some misclassifications is controlled by a parameter called "C." A large C value will prioritize a smaller margin with fewer misclassifications, while a smaller C value will allow a larger margin with more misclassifications.

5. *Kernel Trick:* The kernel trick is a powerful technique employed by SVM to handle non-linearly separable data. It involves transforming the original feature space into a higher-dimensional space, where the data might become linearly separable. The most commonly used kernels are the polynomial kernel, radial basis function (RBF) kernel, and the sigmoid kernel.

6. *Kernel Function:* The choice of the kernel function is crucial in SVM. The kernel function computes the dot product between the transformed data points in the higher-dimensional space without explicitly transforming them. This allows SVM to work efficiently in high-dimensional spaces.

7. *Training and Prediction:* During the training phase, SVM finds the optimal hyperplane and support vectors using optimization techniques such as Quadratic

Programming. Once the SVM model is trained, it can be used to predict the class of new, unseen data points by determining which side of the hyperplane they lie on.

Overall, SVM is a versatile algorithm that can handle both linearly and non-linearly separable data by finding an optimal hyperplane in the transformed feature space, maximizing the margin between classes, and utilizing support vectors for efficient classification.

## 2. DISCUSS VARIOUS USAGE OF KERNEL USED IN SVM

Support Vector Machines (SVM) are powerful machine learning algorithms used for classification and regression tasks. SVM uses a kernel function to map data points from their original feature space into a higher-dimensional space. This mapping allows SVM to find a hyperplane that separates the data points of different classes with a maximal margin. The kernel trick is a key aspect of SVM, as it enables SVM to efficiently work in high-dimensional spaces without explicitly calculating the coordinates of data points in that space.

There are various types of kernel functions that can be used with SVM, each suitable for different types of data and problems. Here are some commonly used kernels in SVM:

1. Linear Kernel:
The linear kernel is the simplest and most basic kernel function. It computes the dot product between the input data points in the original feature space. The linear kernel is appropriate when the data is linearly separable.

Kernel function: $K(x, y) = x^T * y$

2. Polynomial Kernel:
The polynomial kernel is used when the data is not linearly separable and has a polynomial boundary. It maps the data into a higher-dimensional space using a polynomial function.

Kernel function: $K(x, y) = (alpha * x^T * y + c)^d$
where alpha and c are constants and d is the degree of the polynomial.

3. Radial Basis Function (RBF) Kernel:
The RBF kernel, also known as the Gaussian kernel, is one of the most commonly used kernels. It maps the data into an infinite-dimensional space. It is suitable for data that may have non-linear boundaries and works well with a wide range of datasets.

Kernel function: $K(x, y) = \exp(-gamma * ||x - y||^2)$
where gamma is a hyperparameter that controls the width of the Gaussian.

4. Sigmoid Kernel:

The sigmoid kernel maps the data into a higher-dimensional space using a sigmoid function. It is primarily used for neural networks and can be used with SVM for special cases.

Kernel function: K(x, y) = tanh(alpha * x^T * y + c)

5. Laplacian Kernel:
The Laplacian kernel is similar to the RBF kernel but uses the L1 norm instead of the L2 norm to measure the distance between data points. It can be less sensitive to outliers compared to the RBF kernel.

Kernel function: K(x, y) = exp(-gamma * ||x - y||)

Choosing the appropriate kernel function depends on the specific problem and the nature of the data. In practice, it is common to try different kernels and tune their hyperparameters using cross-validation to find the best-performing combination for a given task. The kernel trick is what makes SVM such a versatile and powerful algorithm for various machine learning problems.

3. **USAGE OF GRADIENT DESCENT**

Gradient descent is a fundamental optimization algorithm used in various machine learning and deep learning techniques. Its main purpose is to minimize or maximize a function by iteratively adjusting the model's parameters in the direction of the steepest descent of the gradient.

Here's a discussion on the usage of gradient descent:

1. *Optimization in Machine Learning:* Gradient descent plays a crucial role in training machine learning models. The primary objective is to find the optimal set of model parameters that minimize the cost function (also known as the loss function). This cost function measures the discrepancy between the model's predictions and the actual target values. By minimizing this function, the model becomes more accurate in making predictions on new, unseen data.

2. *Gradient Computation:* To perform gradient descent, the algorithm requires computing the gradient (or derivative) of the cost function with respect to each model parameter. The gradient represents the direction and magnitude of the steepest increase of the function. It indicates how much and in which direction the parameters need to be adjusted to reduce the cost.

3. *Learning Rate:* The learning rate is a hyperparameter in gradient descent that determines the step size at each iteration. It controls the magnitude of parameter updates. A small learning rate may lead to slow convergence, while a large learning rate can cause overshooting and may prevent the algorithm from converging. Selecting an appropriate learning rate is crucial for effective training.

4. *Batch, Stochastic, and Mini-batch Gradient Descent:* There are different variations of gradient descent based on the number of training samples used in each iteration.

- Batch Gradient Descent: It computes the gradient using the entire training dataset in each iteration. This approach can be computationally expensive and may require a lot of memory for large datasets.

- Stochastic Gradient Descent (SGD): It computes the gradient using only one random training sample at a time. SGD can have higher variance but is computationally more efficient and can work well in large datasets.

- Mini-batch Gradient Descent: This method strikes a balance between batch and stochastic gradient descent by using a small subset of the training data (mini-batch) to compute the gradient. This is the most commonly used approach in practice, as it offers a compromise between efficiency and stability.

5. *Convex and Non-Convex Optimization:* Gradient descent is particularly effective for convex optimization problems, where there is a single global minimum. However, in non-convex optimization problems (common in neural networks), the algorithm can get stuck in local minima or saddle points. To overcome this, advanced optimization techniques like momentum, RMSprop, and Adam have been developed.

6. *Convergence and Early Stopping:* Gradient descent iterations continue until the cost function converges to a minimum or a predetermined number of epochs is reached. In practice, it's common to monitor the cost on a validation set during training and apply early stopping when the performance stops improving, to prevent overfitting.

7. *Limitations:* Gradient descent may converge slowly in certain situations, especially if the cost function has plateaus, narrow valleys, or sharp cliffs. To address this, adaptive learning rate algorithms and momentum-based methods can be used.

In summary, gradient descent is a versatile and essential optimization algorithm in the field of machine learning. It enables models to learn from data and find the optimal set of parameters for a given task, allowing them to make accurate predictions on new, unseen data. Researchers and practitioners continue to explore and develop advanced variants of gradient descent to improve training efficiency and convergence for complex models and large datasets.

4. **DIFF. OF POLYNIMIAL REGRESSION FROM LINEAR REGRESSION**
Polynomial regression and linear regression are both techniques used in statistical modeling and machine learning to establish a relationship between an independent variable (or variables) and a dependent variable. However, they differ in how they model this relationship.

1. Linear Regression:
Linear regression models the relationship between the independent variable(s) and the dependent variable as a straight line. It assumes a linear relationship, which means that the change in the dependent variable is proportional to the change in the independent variable. The equation for a simple linear regression is of the form:

y = mx + b

Where:
- y is the dependent variable.
- x is the independent variable.
- m is the slope of the line, representing the rate of change of y with respect to x.
- b is the y-intercept, representing the value of y when x is 0.

Linear regression is suitable when the relationship between the variables is approximately linear. If the data follows a straight line pattern, linear regression can provide a good fit.

2. Polynomial Regression:
Polynomial regression, on the other hand, models the relationship between the independent variable(s) and the dependent variable as an nth-degree polynomial. It assumes a nonlinear relationship, which means that the change in the dependent variable is not proportional to the change in the independent variable. The equation for a polynomial regression is of the form:

y = b0 + b1 * x + b2 * x^2 + b3 * x^3 + ... + bn * x^n

Where:
- y is the dependent variable.
- x is the independent variable.
- b0, b1, b2, ..., bn are the coefficients of the polynomial equation.
- n is the degree of the polynomial, determining the order of the curve.

Polynomial regression can capture more complex patterns in the data compared to linear regression. It allows for a more flexible curve, which can lead to a better fit when the relationship between the variables is nonlinear.

It's important to note that while polynomial regression can provide a better fit to the training data, it may be more prone to overfitting (capturing noise in the data rather than the underlying pattern). Therefore, the choice between linear regression and polynomial regression depends on the nature of the data and the complexity of the relationship between the variables. In some cases, a combination of the two approaches, such as using higher-order polynomials or adding polynomial features to a linear regression model, may be used to strike the right balance between flexibility and avoiding overfitting.

5. **LOSS FUNC. OF LOGISTIC REGRESSION**
In logistic regression, the loss function is a fundamental component used to measure the model's performance and guide the process of learning the optimal parameters. The goal of logistic regression is to predict the probability that a given input belongs to a certain class (binary classification) based on the input features.

The logistic regression model uses a logistic function (also known as the sigmoid function) to map the input features to probabilities. The sigmoid function takes any real-valued number and "squashes" it into the range (0, 1). It is defined as follows:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

where $z$ is a linear combination of the input features and their corresponding weights. Mathematically, it can be written as:

$$z = w_0 + w_1x_1 + w_2x_2 + \ldots + w_nx_n$$

where:
- $w_0, w_1, \ldots, w_n$ are the model's weights (parameters).
- $x_1, x_2, \ldots, x_n$ are the input features.
- $n$ is the number of features.

The output of the sigmoid function, $\sigma(z)$, represents the probability that the input $x$ belongs to the positive class (class 1 in binary classification).

Now, in order to train the logistic regression model, we need a loss function that quantifies the difference between the predicted probabilities and the actual class labels in the training data.

One common loss function used in logistic regression is the *binary cross-entropy loss, also known as **log loss*. For a single training example with true class label $y$ (where $y = 0$ or $y = 1$) and predicted probability $\hat{y}$, the binary cross-entropy loss is defined as:

$$\text{Loss}(y, \hat{y}) = - \left( y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}) \right)$$

The negative log likelihood is used here to ensure that maximizing the likelihood (or minimizing the negative log-likelihood) is equivalent to minimizing the cross-entropy loss.

To train the logistic regression model, the goal is to minimize the overall average loss across all the training examples. The training process typically involves using optimization algorithms, such as gradient descent, to update the model's weights iteratively until convergence or a certain number of epochs.

Minimizing the cross-entropy loss encourages the model to produce high probabilities (close to 1) for positive class instances when they are observed in the training data and low probabilities (close to 0) for negative class instances. This way, the model learns to make accurate predictions and separate the two classes effectively.

6. **what is gradient descent algorithm and discuss its various types.**

Gradient descent is an optimization algorithm used to minimize a function by iteratively adjusting the parameters of a model or system. It is widely used in machine learning

and deep learning to find the optimal values of the model's parameters that minimize the cost or loss function. The fundamental idea behind gradient descent is to follow the negative gradient (or steepest descent) of the function in the parameter space, which leads to the direction of the steepest decrease in the function's value.

The steps involved in the gradient descent algorithm are as follows:

1. Initialize parameters: Start with some initial values for the model's parameters.

2. Compute the gradient: Calculate the gradient of the cost function with respect to each parameter. The gradient indicates the direction in which the function's value increases the most.

3. Update parameters: Adjust the parameters in the opposite direction of the gradient to minimize the cost function. The learning rate determines the step size in this direction.

4. Repeat steps 2 and 3: Keep updating the parameters until convergence or a predefined number of iterations is reached.

There are various types of gradient descent algorithms, which differ in how they update the parameters and handle the learning rate:

1. Batch Gradient Descent:
   - In batch gradient descent, the entire training dataset is used to compute the gradient in each iteration.
   - It provides a more accurate estimate of the gradient but can be computationally expensive for large datasets.

2. Stochastic Gradient Descent (SGD):
   - In SGD, only one random training sample is used to compute the gradient in each iteration.
   - It is computationally efficient and can work well for large datasets.
   - However, the updates can be very noisy, leading to a more erratic convergence path.

3. Mini-batch Gradient Descent:
   - Mini-batch gradient descent is a compromise between batch gradient descent and SGD.
   - It uses a small random subset (mini-batch) of the training data to compute the gradient.
   - It strikes a balance between efficiency and convergence smoothness.

4. Momentum Gradient Descent:
   - Momentum is a modification of gradient descent that helps accelerate convergence, especially in shallow ravines and noisy gradients.
   - It adds a fraction of the previous update to the current update, which helps the algorithm maintain momentum in the relevant directions.

5. Nesterov Accelerated Gradient (NAG):
   - NAG is a variation of momentum gradient descent that reduces the oscillations experienced during convergence.
   - It computes the gradient at the position slightly ahead of the current parameter values, based on the momentum direction.

6. Adagrad (Adaptive Gradient Algorithm):
   - Adagrad adapts the learning rate for each parameter based on the historical gradients.
   - It provides larger updates for infrequent parameters and smaller updates for frequent parameters.

7. RMSprop (Root Mean Square Propagation):
   - RMSprop is an improvement over Adagrad, addressing its tendency to decrease the learning rate excessively.
   - It uses an exponentially decaying average of squared gradients to adjust the learning rate.

8. Adam (Adaptive Moment Estimation):
   - Adam combines the concepts of momentum and RMSprop.
   - It maintains running averages of both the gradient and the squared gradient and uses bias correction terms.

These variations of gradient descent help improve convergence speed and handling of different types of optimization landscapes. The choice of the gradient descent algorithm depends on the specific problem and dataset at hand. Adam, in particular, has gained popularity due to its robustness and good performance in many scenarios.

7. **In Regularized Linear Models illustrate the three different methods to constrain the weights.**

In regularized linear models, the aim is to add a penalty term to the cost function, which helps in constraining the weights and preventing overfitting. The three different methods to constrain the weights are:

1. Ridge Regression (L2 Regularization):

Ridge Regression adds the L2 norm of the weight vector to the cost function. The L2 norm is the sum of squares of all the weights. By adding this regularization term, the model is penalized for having large weights. The cost function for Ridge Regression is:

Cost = Least Squares Cost + α * Σ(w_i^2)

Where:
- Least Squares Cost is the traditional linear regression cost function.
- α is the regularization parameter (also known as lambda) that controls the strength of the penalty. Higher values of α lead to more constrained (smaller) weights.

The regularization term tends to push the weights towards zero but does not make them exactly zero, allowing all features to be considered in the model.

2. Lasso Regression (L1 Regularization):
Lasso Regression adds the L1 norm of the weight vector to the cost function. The L1 norm is the sum of the absolute values of all the weights. Similar to Ridge Regression, this penalty term discourages large weights but has a different effect on the weight selection. The cost function for Lasso Regression is:

Cost = Least Squares Cost + $\alpha * \Sigma(|w\_i|)$

Where:
- Least Squares Cost is the traditional linear regression cost function.
- $\alpha$ is the regularization parameter (lambda) controlling the strength of the penalty. Higher values of $\alpha$ increase the regularization effect.

Lasso Regression can drive some weights to exactly zero, effectively performing feature selection by ignoring irrelevant features. This makes Lasso Regression useful when we want to create sparse models.

3. Elastic Net Regression:
Elastic Net Regression combines both L1 and L2 regularization terms in the cost function. This helps to overcome some limitations of Ridge and Lasso Regression. The cost function for Elastic Net Regression is:

Cost = Least Squares Cost + $\alpha * \Sigma(w\_i^2) + \beta * \Sigma(|w\_i|)$

Where:
- Least Squares Cost is the traditional linear regression cost function.
- $\alpha$ and $\beta$ are the two regularization parameters, controlling the strength of the L2 and L1 penalties, respectively.

Elastic Net combines the strengths of Ridge and Lasso regularization while mitigating their weaknesses. It can handle highly correlated features better than Lasso, and when the number of features is greater than the number of samples, Elastic Net can select at most the number of samples as non-zero coefficients.

Each of these regularization techniques provides different ways to control the model's complexity and avoid overfitting, allowing for better generalization to new data. The choice of which method to use depends on the specific problem and the characteristics of the dataset.

8. **With respect to Nonlinear SVM Classification, explain Polynomial Kernel Gaussian and RBF Kernel along with code snippet.**
Sure! Both the polynomial kernel and the radial basis function (RBF) kernel are used in Support Vector Machine (SVM) classification for handling non-linearly separable data.

1. Polynomial Kernel:
The polynomial kernel is used to map the data into a higher-dimensional space, making it easier to separate the data that is not linearly separable in the original feature space. However, using a high degree polynomial can lead to overfitting, and using a low degree may not provide enough complexity to capture the underlying patterns. Therefore, it is crucial to choose an appropriate degree for the polynomial kernel to achieve the best results.

2. RBF Kernel:
The radial basis function (RBF) kernel uses a similarity measure to compute the distance between data points in a higher-dimensional space. It is commonly used because it does not suffer from the "curse of dimensionality" and can effectively handle non-linearly separable data. The RBF kernel has a parameter called gamma ($\gamma$), which controls the smoothness of the decision boundary. Higher values of gamma can lead to overfitting.

Now, let's demonstrate how to use the polynomial and RBF kernels with a simple example in Python using the popular machine learning library, Scikit-learn.

```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import make_circles
from sklearn.model_selection import train_test_split
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score

# Generate some non-linearly separable data (circles)
X, y = make_circles(n_samples=500, noise=0.1, random_state=42)

# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# SVM classifier with polynomial kernel
poly_svc = SVC(kernel='poly', degree=3, gamma='auto', random_state=42)
poly_svc.fit(X_train, y_train)

# SVM classifier with RBF kernel
rbf_svc = SVC(kernel='rbf', gamma=0.1, random_state=42)
rbf_svc.fit(X_train, y_train)

# Make predictions
y_pred_poly = poly_svc.predict(X_test)
y_pred_rbf = rbf_svc.predict(X_test)
```

```
# Calculate accuracy
accuracy_poly = accuracy_score(y_test, y_pred_poly)
accuracy_rbf = accuracy_score(y_test, y_pred_rbf)

print("Accuracy with Polynomial Kernel:", accuracy_poly)
print("Accuracy with RBF Kernel:", accuracy_rbf)

# Plot the decision boundaries
def plot_decision_boundary(clf, X, y, title):
    h = 0.01
    x_min, x_max = X[:, 0].min() - 0.1, X[:, 0].max() + 0.1
    y_min, y_max = X[:, 1].min() - 0.1, X[:, 1].max() + 0.1
    xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
    Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
    Z = Z.reshape(xx.shape)
    plt.contourf(xx, yy, Z, cmap=plt.cm.coolwarm, alpha=0.8)
    plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.coolwarm, edgecolors='k')
    plt.xlabel('Feature 1')
    plt.ylabel('Feature 2')
    plt.title(title)
    plt.show()

# Plot decision boundaries for polynomial kernel
plot_decision_boundary(poly_svc, X, y, title="SVM with Polynomial Kernel")

# Plot decision boundaries for RBF kernel
plot_decision_boundary(rbf_svc, X, y, title="SVM with RBF Kernel")
```

In this code snippet, we first create a non-linearly separable dataset using `make_circles` from Scikit-learn. We then split the data into training and testing sets. Next, we create two SVM classifiers: one with a polynomial kernel (`SVC(kernel='poly')`) and another with an RBF kernel (`SVC(kernel='rbf')`). We fit the classifiers on the training data and evaluate their performance on the testing data by calculating the accuracy. Finally, we plot the decision boundaries for both classifiers to visualize how they separate the data.

9. **Show that how SVMs make predictions using Quadratic Programming and Kernelized SVM.**
Support Vector Machines (SVMs) are a class of supervised machine learning algorithms used for classification and regression tasks. SVMs make predictions by finding the decision boundary that maximizes the margin between different classes. When dealing with linearly separable data, SVM uses a method called "quadratic programming" to find the optimal hyperplane. For non-linearly separable data, SVM employs a technique called "kernel trick" to extend the feature space and make it linearly separable.

1. SVM using Quadratic Programming (Linear SVM):

Linear SVM aims to find the optimal hyperplane that separates the data into different classes. The decision boundary is represented by the equation:

$$w \cdot x + b = 0$$

where 'w' is the weight vector perpendicular to the hyperplane, 'x' is the input data vector, and 'b' is the bias term. The objective of the SVM is to maximize the margin (distance) between the two classes while minimizing classification errors.

To find 'w' and 'b', we formulate the SVM optimization problem as follows:

Minimize $(1/2) \|w\|^2$ subject to $y_i(w \cdot x_i + b) \geq 1$ for all i

where '$y_i$' is the class label (+1 or -1) of the data point '$x_i$'. This is a convex quadratic programming problem, and various optimization algorithms can be used to solve it efficiently.

2. Kernelized SVM (Non-linear SVM):

The kernel trick is used to extend the feature space in order to make the data linearly separable in a higher-dimensional space. Instead of directly mapping the data to that higher-dimensional space, SVM introduces a kernel function that calculates the dot product between the data points in the higher-dimensional space.

The kernel function $K(x_i, x_j)$ takes two data points '$x_i$' and '$x_j$' and calculates their dot product in the transformed space without explicitly calculating the transformation. Common kernel functions include the polynomial kernel ($K(x_i, x_j) = (x_i \cdot x_j + c)^d$) and the radial basis function (RBF) kernel ($K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$).

With the kernel trick, the SVM optimization problem becomes:

Minimize $(1/2) \Sigma \alpha_i \alpha_j y_i y_j K(x_i, x_j) - \Sigma \alpha_i$

subject to $\Sigma \alpha_i y_i = 0$

and $0 \leq \alpha_i \leq C$ for all i

where $\alpha_i$ is the Lagrange multiplier associated with each data point '$x_i$', and 'C' is the regularization parameter that controls the trade-off between maximizing the margin and allowing some misclassifications.

To make predictions with kernelized SVM, we use the learned support vectors' coefficients ($\alpha_i$) and the kernel function to calculate the decision function:

$$f(x) = \Sigma \alpha_i y_i K(x_i, x) + b$$

where 'b' is the bias term. The class label of the new data point 'x' is determined based on the sign of the decision function value. If f(x) is positive, the data point belongs to one class, and if it is negative, it belongs to the other class.

In practice, libraries like Scikit-learn in Python encapsulate these complexities, allowing users to easily use SVMs with different kernels without having to explicitly deal with quadratic programming or kernel transformations. The underlying optimization and kernel calculations are efficiently implemented in these libraries.