# Support Vector Machines (SVMs)

# 1. Linear SVM for Classification

```python
import numpy as np
import matplotlib.pyplot as plt
from sklearn.svm import SVC
from sklearn.datasets import make_classification

# Create synthetic data (two classes)
X, y =
Make_classification
(
        n_samples=50,
        n_features=2,
        n_informative=2,
        n_redundant=0,
        random_state=1,
        class_sep=1.5
)

# Train a linear SVM classifier
clf = SVC(kernel='linear', C=1.0)  # C is the regularization parameter
clf.fit(X, y)

# Plot the data and decision boundary
plt.scatter(
X[:, 0],
X[:, 1],
c=y,
cmap='viridis')

# Get the separating hyperplane
w = clf.coef_[0]
b = clf.intercept_[0]
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min = (-b - w[0] * x_min) / w[1]
y_max = (-b - w[0] * x_max) / w[1]

plt.plot([x_min, x_max], [y_min, y_max], "k-")  # Plot the hyperplane

plt.xlabel("Feature 1")
```

```
plt.ylabel("Feature 2")
plt.title("Linear SVM Classification")
plt.show()


# Example of making predictions
new_data = np.array([[2, 1]])  # Example data point
prediction = clf.predict(new_data)
print(f"Prediction for {new_data}: {prediction}")
```

—--------------------------------------------------------------------------------------------------------------

**Explanation:**

```
# Create synthetic data (two classes)
X,   y   =   make_classification(n_samples=50,   n_features=2,   n_informative=2,
n_redundant=0, random_state=1, class_sep=1.5)
```

- **Data Generation:** `make_classification` creates a synthetic dataset with two classes. `class_sep` controls the separation between classes.

  **Here's an explanation of each parameter:**

- `n_samples=50`: This determines the total number of data points (samples) to generate. In this case, it will create 50 data points.

- `n_features=2`: This specifies the number of features (attributes) for each data point. Here, each data point will have 2 features. Think of these as the dimensions of your data. If you were visualizing this data, it would be 2-dimensional, making it easy to plot.

- `n_informative=2`: This is a crucial parameter. It indicates the number of features that are actually informative for the classification task. In this example, *both* of the 2 features are informative, meaning they contribute to distinguishing between the classes.

- `n_redundant=0`: This specifies the number of redundant features. Redundant features are features that don't add any new information and can be predicted from the other informative features. Here, there are 0 redundant features. Since `n_informative` is 2 and `n_redundant` is 0, all 2 features are used to create the classes.

- `random_state=1`: This sets the seed for the random number generator. Using a fixed `random_state` ensures that the data generated is reproducible. If you run the code multiple times with the same `random_state`, you'll get the same dataset. This is very important for testing and comparing different machine learning models.

- `class_sep=1.5`: This controls how well separated the classes are. A larger value for `class_sep` means the classes will be more easily distinguishable (well-separated), while a smaller value means they will be more intertwined or overlapping. 1.5 indicates a

reasonable separation between the classes.

**Return Values:**

The `make_classification` function returns two values:

- **X:** This is a NumPy array containing the generated data. It will have a shape of (50, 2) in this case (50 samples, 2 features).
- **y:** This is a NumPy array containing the class labels for each data point. It will have a shape of (50,) in this case. The labels will typically be 0 and 1 (or -1 and 1, depending on the specific function).

**In summary:** This line of code generates a synthetic dataset with 50 data points, each having 2 features. Both features are informative for classifying the data points into two classes, and the classes are reasonably well-separated. The `random_state` ensures reproducibility. This kind of synthetic data is very useful for testing and visualizing machine learning algorithms like SVMs.

—--------------------------------------------------------------------------------------------------------------

```
# Train a linear SVM classifier
clf = SVC(kernel='linear', C=1.0)  # C is the regularization parameter
clf.fit(X, y)
```

1. `clf = SVC(kernel='linear', C=1.0)`

- **SVC:** This stands for Support Vector Classifier. It's the class in scikit-learn used to create an SVM classifier.

- `kernel='linear':` This specifies that we want to use a linear kernel for the SVM. A linear kernel means the SVM will try to find a linear hyperplane to separate the data points of different classes. This is suitable when the classes are linearly separable, or at least approximately so.

- `C=1.0:` This is the regularization parameter, often called the "penalty parameter." It controls the trade-off between maximizing the margin (the distance between the decision boundary and the nearest data points) and minimizing the classification error [1] (the number of misclassified data points).
  - A *smaller* value of `C` (e.g., `C=0.1`) means *stronger* regularization. The SVM will prioritize a wider margin, even if it means misclassifying some data points. This can help prevent overfitting, especially if you have noisy data.
  - A *larger* value of `C` (e.g., `C=10` or `C=100`) means *weaker* regularization. The SVM will try to correctly classify as many training points as possible, even if it means a narrower margin. This can lead to overfitting if the data is noisy or if you don't have enough training examples.

- Choosing the right value for `C` is important and often done using techniques like cross-validation.

- **`clf = ...`:** This creates an instance of the `SVC` class and assigns it to the variable `clf`. `clf` is now our classifier object.

## 2. `clf.fit(X, y)`

- **`clf.fit()`:** This is the method used to train the SVM classifier.

- **X:** This is the training data. It's a NumPy array (or a similar data structure) where each row represents a data point, and each column represents a feature. In the previous example using `make_classification`, X was the array of shape (50, 2) containing the feature values for the 50 data points.

- **y:** These are the corresponding class labels for the training data. It's a 1D array where each element represents the class label (e.g., 0 or 1, -1 or 1) for the corresponding data point in X. In the previous example, y was the array of shape (50,) containing the class labels.

**In summary:**

The first line creates an SVM classifier object (`clf`) with a linear kernel and a regularization parameter `C` set to 1.0. The second line trains the classifier using the training data `X` and the corresponding class labels `y`. After `clf.fit(X, y)` is executed, the `clf` object is "fitted" or "trained," meaning it has learned the parameters of the separating hyperplane. You can then use `clf` to make predictions on new, unseen data using `clf.predict()`.

—--------------------------------------------------------------------------------------------------------------------

# Plot the data and decision boundary
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')

This line of code uses the `matplotlib.pyplot` library (usually imported as `plt`) to create a scatter plot of your data. Let's break down what each part does:

- **`plt.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')`:** This is the core function call that generates the scatter plot.

  - **`plt.scatter()`:** This function is specifically for creating scatter plots.

  - **`X[:, 0]`:** This extracts all rows (`:`) and the first column (`0`) from your data X. This assumes X is a 2D array-like structure (e.g., NumPy array). This provides the x-coordinates for the scatter plot.

- ○ **`X[:, 1]`**: This extracts all rows and the second column (`1`) from `X`. This provides the y-coordinates for the scatter plot. So, you're plotting the first feature against the second feature.

- ○ **`c=y`**: This sets the color of each point in the scatter plot based on the corresponding value in the `y` array. `y` presumably contains the class labels or target variable. Points belonging to different classes will be colored differently.

- ○ **`cmap='viridis'`**: This specifies the colormap to use for the colors. `viridis` is a colormap that is perceptually uniform (meaning changes in color intensity correspond directly to changes in the data value) and colorblind-friendly. There are many other colormaps available in Matplotlib (e.g., 'jet', 'plasma', 'magma', 'inferno', 'gray'). You can explore them in the matplotlib documentation.

**In summary:** This code takes your data `X` (assuming it has at least two features) and plots it as a scatter plot. The x-coordinates are taken from the first feature, the y-coordinates from the second feature, and the color of each point is determined by its class label in `y`. The `viridis` colormap is used to visualize the class labels.

—-------------------------------------------------------------------------------------------------------------------

# Get the separating hyperplane
w = clf.coef_[0]
b = clf.intercept_[0]
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min = (-b - w[0] * x_min) / w[1]
y_max = (-b - w[0] * x_max) / w[1]

This code snippet calculates the equation of the separating hyperplane for a linear SVM classifier (like the one you trained earlier with `kernel='linear'`). Let's break down what each part does:

- ● **`w = clf.coef_[0]`**: `clf.coef_` contains the coefficients of the hyperplane. Since we're dealing with a linear SVM, `clf.coef_` will be a 2D array if you have multiple classes (one row per class), and a 1D array if you have only two classes. `[0]` accesses the coefficients for the first class (or the only class, in the two-class case). This gives you the weights `w` of the features in the hyperplane equation. For a 2D dataset, `w` will have two elements: `w[0]` and `w[1]`.

- ● **`b = clf.intercept_[0]`**: `clf.intercept_` contains the intercept term of the hyperplane. Similar to `clf.coef_`, it's a 1D array with one element per class. `[0]` accesses the intercept for the first class (or the only class). This is the bias `b` in the hyperplane equation.

- ● **`x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1`**: This calculates the minimum and maximum values of the first feature (column 0) in your data `X`. It then adds and subtracts 1 to create a small margin around the data points for visualization

purposes. These `x_min` and `x_max` values will be used to define the range of the x-axis when plotting the hyperplane.

- **`y_min = (-b - w[0] * x_min) / w[1]`**: This calculates the y-coordinate of the hyperplane at `x_min`. It uses the equation of the hyperplane: `w[0]*x + w[1]*y + b = 0`. Rearranging for `y`, we get: `y = (-b - w[0]*x) / w[1]`. Plugging in `x_min` gives us the corresponding `y_min`.

- **`y_max = (-b - w[0] * x_max) / w[1]`**: This similarly calculates the y-coordinate of the hyperplane at `x_max`.

**In summary:** This code calculates the equation of the separating hyperplane (`w[0]*x + w[1]*y + b = 0`) learned by the linear SVM. It then determines two points on the hyperplane (`(x_min, y_min)` and `(x_max, y_max)`) which can be used to plot the hyperplane as a line on your scatter plot.

—-------------------------------------------------------------------------------------------------------------------

plt.plot([x_min, x_max], [y_min, y_max], "k-")  # Plot the hyperplane

This line of code plots the separating hyperplane on your scatter plot. Let's break it down:

- **`plt.plot([x_min, x_max], [y_min, y_max], "k-")`**: This uses the `plt.plot()` function to draw a line.

  - **`plt.plot()`**: This function is used for plotting lines and/or markers.

  - **`[x_min, x_max]`**: This provides the x-coordinates of the two points that define the line. Remember, `x_min` and `x_max` were calculated earlier to span the range of your data (with a small margin).

  - **`[y_min, y_max]`**: This provides the y-coordinates of the two points. `y_min` and `y_max` were calculated using the equation of the hyperplane, corresponding to `x_min` and `x_max` respectively. So, these two points (`(x_min, y_min)` and `(x_max, y_max)`) lie on the hyperplane.

  - **`"k-"`**: This is a string that specifies the style of the line.

    - `"k"`: Represents the color black.
    - `"-"`: Represents a solid line style.

**In summary:** This code draws a black, solid line on your plot. The line connects the points `(x_min, y_min)` and `(x_max, y_max)`, which both lie on the separating hyperplane. Because these points span the x-range of your data, the line visually represents the hyperplane across the plot.

—-------------------------------------------------------------------------------------------------------------------

```
plt.xlabel("Feature 1")
plt.ylabel("Feature 2")
plt.title("Linear SVM Classification")
plt.show()
```

These lines of code add labels to the axes and a title to your plot, and then display the plot. Let's break down each line:

- **`plt.xlabel("Feature 1")`**: This sets the label for the x-axis. `"Feature 1"` is the text that will be displayed next to the x-axis. You should replace `"Feature 1"` with a more descriptive name relevant to your data (e.g., "Sepal Length," "Income," etc.).

- **`plt.ylabel("Feature 2")`**: This sets the label for the y-axis, similar to `plt.xlabel()`. Replace `"Feature 2"` with a descriptive name for your y-axis.

- **`plt.title("Linear SVM Classification")`**: This sets the title of the plot. `"Linear SVM Classification"` is the text that will be displayed at the top of the plot. Again, it's good practice to use a more specific title that reflects the content of your plot.

- **`plt.show()`**: This command is essential. It tells Matplotlib to display the plot you've created. Without `plt.show()`, the plot is generated in memory but not shown to the user.

**In summary:** These lines enhance the readability and understanding of your plot by adding informative labels and a title. `plt.show()` makes the plot visible. It's always recommended to include these elements in your plotting code for clear communication of your results.

—-----------------------------------------------------------------------------------------------------------------

```
# Example of making predictions
new_data = np.array([[2, 1]])  # Example data point
prediction = clf.predict(new_data)
print(f"Prediction for {new_data}: {prediction}")
```

This code snippet demonstrates how to use the trained SVM classifier (`clf`) to make predictions on new, unseen data. Let's break it down:

- **`new_data = np.array([[2, 1]])`**: This creates a NumPy array representing the new data point you want to classify. It's crucial that the shape of `new_data` is consistent with the data the classifier was trained on. Since your training data `X` had two features (two columns), `new_data` must also have two columns. The double square brackets `[[ ]]` create a 2D array, even if you only have one data point. This is important because `clf.predict()` expects a 2D array-like input, even for a single prediction.

- **`prediction = clf.predict(new_data)`**: This is the core line that performs the prediction.

  - `clf`: The trained SVM classifier object.
  - `predict(new_data)`: This method uses the trained model to predict the class label for the `new_data` point. It returns an array containing the predicted class labels. Even if you only predict for one data point, the result will be in an array.
- **`print(f"Prediction for {new_data}: {prediction}")`**: This line prints the prediction result in a user-friendly format. The f-string formatting (using `f"..."`) allows you to embed the values of variables directly within the string.

**In summary:** This code takes a new data point `new_data`, uses the trained SVM classifier `clf` to predict its class label, and then prints the result.

# 2. Non-Linear SVM with Kernel Trick

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

# Create non-linearly separable data
X, y = make_circles(n_samples=100, factor=0.5, noise=0.05, random_state=1)

# Train an SVM with a radial basis function (RBF) kernel
clf = SVC(kernel='rbf', C=1.0, gamma='scale') # gamma controls the width of the RBF kernel
clf.fit(X, y)

# Plot the data and decision boundary
plt.scatter(X[:, 0], X[:, 1], c=y, cmap='viridis')

# Create a meshgrid to plot the decision boundary
xx, yy = np.meshgrid(np.linspace(X[:, 0].min() - 1, X[:, 0].max() + 1, 100),
                np.linspace(X[:, 1].min() - 1, X[:, 1].max() + 1, 100))
Z = clf.decision_function(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)
plt.contourf(xx, yy, Z, levels=np.linspace(Z.min(), Z.max(), 7), alpha=0.3, cmap='coolwarm')
#Filled contour

plt.xlabel("Feature 1")
plt.ylabel("Feature 2")
plt.title("Non-Linear SVM (RBF Kernel)")
plt.show()
```

**Explanation:**

- **Non-Linearly Separable Data:** `make_circles` creates a dataset where the classes are not linearly separable.
- **RBF Kernel:** `SVC(kernel='rbf', C=1.0, gamma='scale')` uses the Radial Basis Function (RBF) kernel. The RBF kernel maps the data to a higher-dimensional space where it becomes linearly separable. `gamma` controls the width of the RBF. `gamma='scale'` is a good default.
- **Decision Boundary Plotting:** The code creates a meshgrid and uses `clf.decision_function` to get the decision values for each point on the grid. The `contourf` function plots the decision boundary.

**Key Concepts and Parameters:**

- **Kernel Trick:** The kernel trick allows SVMs to work in high-dimensional spaces without explicitly computing the coordinates of the data in that space. Common kernels include linear, polynomial, and RBF.
- **Regularization Parameter (C):** Controls the trade-off between maximizing the margin and minimizing classification errors.
- **Gamma (for RBF kernel):** Controls the width of the RBF. A small gamma means a wider kernel, capturing more of the overall data structure. A large gamma means a narrower kernel, focusing more on individual points.

**Choosing the right kernel and parameters is crucial for SVM performance. Cross-validation is often used to find the optimal values.** These examples provide a foundation for understanding and applying SVMs to various machine learning problems. Remember to explore the scikit-learn documentation for more advanced features and options.