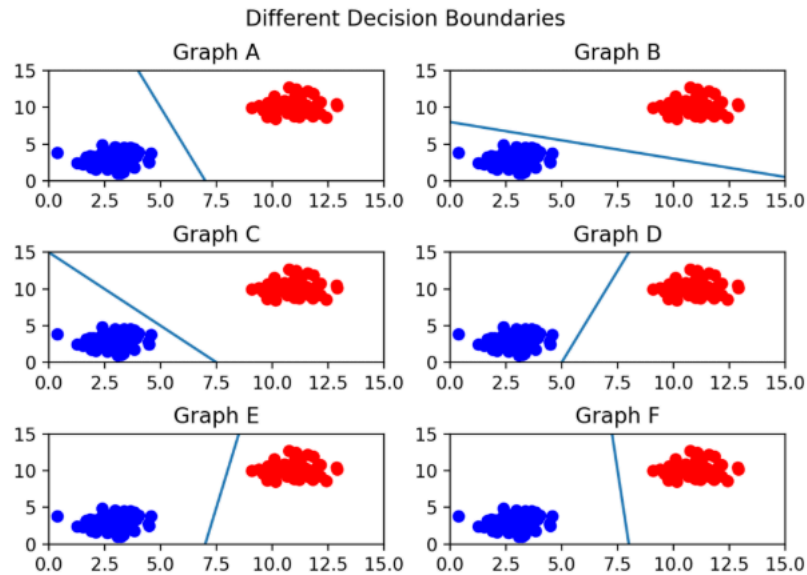


SVM

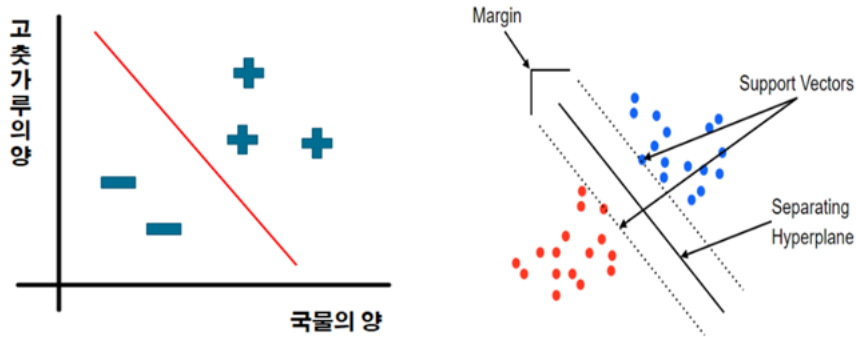
- [분류와 회귀 문제에 사용할 수 있는 강력한 머신러닝 지도학습 모델](#)
 - 결정 경계(Decision Boundary), 즉 분류를 위한 기준 선을 정의하는 모델
 - 분류되지 않은 새로운 점이 나타나면 어느 쪽에 속하는지 확인해서 분류 과제를 수행할 수 있게 됨
 - 결정 경계와 가까이 있는 데이터 포인트들을 의미
- [결국 이 결정 경계라는 걸 어떻게 정의하고 계산하는지 이해하는게 더 중요하다는 뜻임!!!!!!](#)



생각 문제

- [빨간색, 파란색 label을 분류하는 선형\(Linear\)을 찾는다고 가정해보자!!!](#)
 - 기본적으로 SVM은 두 개의 다른 클래스를 분류할 수 있는 여러가지의 선형식이 있을테지만 이 중에서 [가장 잘 분류할 수 있는 하나의\(Unique\) 선형식을 찾으려고 함](#)

- 그렇다면 가장 잘 분류할 수 있는 기준은 무엇일까?
- 그림에서 Separating Hyperplane(일종의 Decision boundary)가 두 개의 클래스를 가장 잘 분류할 수 있는 선형식을 찾는것
 - 이 Decision boundary를 찾기 위한 조건이 있음?
 - 여러가지의 선형식 후보들중에서 Support Vector들과의 Margin이 가장 Maximum이 되게해야 함
 - 이때 Support Vector는 반드시 최소 2개 이상의 데이터가 존재해야 함

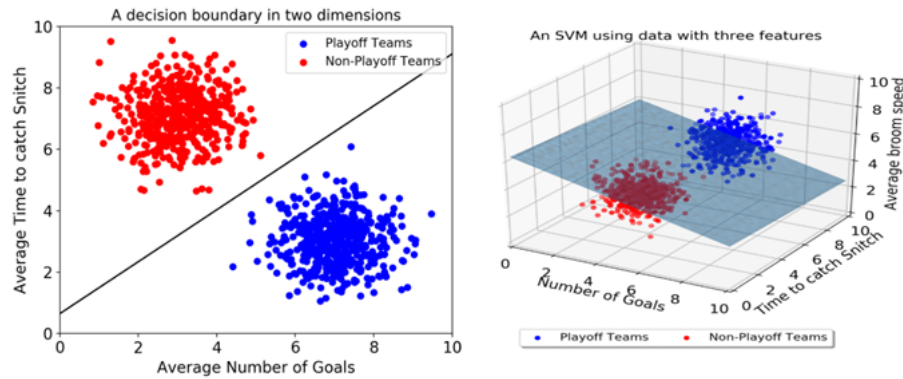


▽ Support Vector

- 두 가지 카테고리에 각각 해당되는 data set들('-'샘플이 모인곳/'+'샘플이 모인곳)의 최외각에 있는 샘플들을 의미
 - 이 가장 최외각에 있는 벡터들을 토대로 margin을 구할 수 있기 때문에 중요한 벡터들임

▽ 하이퍼플레인(hyperplane)

- N차원 공간을 두 부분으로 나누는 N-1차원의 부분공간
 - 2차원 공간(평면)에서 하이퍼플레인은 1차원의 선
 - 3차원 공간에서는 2차원의 평면
- 머신러닝과 특히 서포트 벡터 머신(SVM)에서 중요한 개념으로, 하이퍼플레인은 데이터 클래스를 구분하는 결정 경계로 사용됨



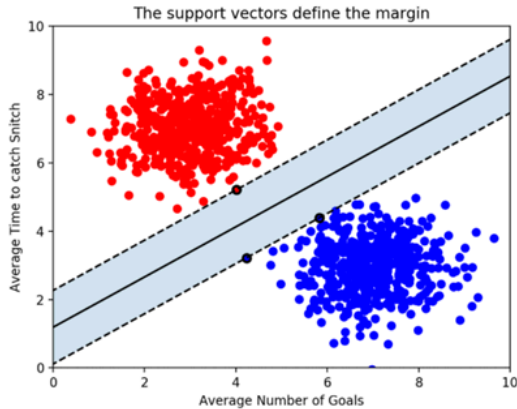
- [하이퍼플레인의 정의는 주로 선형 방정식을 통해 이루어짐](#)
 - 다음과 같은 방정식으로 표현될 수 있음
 - w 는 가중치 벡터, x 는 입력 벡터, b 는 편향

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

- 하이퍼플레인은 고차원 데이터를 분류할 때 중요한 역할을 함
 - 특히, 서포트 벡터 머신에서는 가장 가까운 데이터 포인트(서포트 벡터) 사이의 [마진을 최대화하는 하이퍼플레인을 찾음으로써, 두 클래스를 가장 잘 구분할 수 있는 결정 경계를 찾음](#)
 - 이러한 성질 때문에 하이퍼플레인은 패턴 인식, 이미지 분석, 바이오인포매틱스 등 다양한 머신러닝 응용 분야에서 활용됨

▽ [마진\(Margin\)](#)

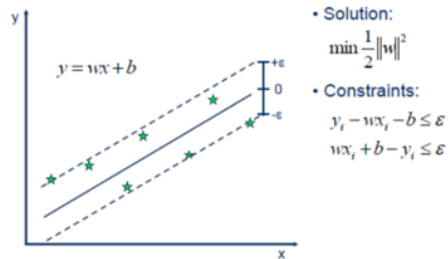
- [결정 경계와 서포트 벡터 사이의 거리를 의미](#)



- ✓ x축과 y축 2개의 속성을 가진 데이터로 결정 경계를 그었는데, 총 3개의 데이터 포인트(서포트 벡터)가 필요함
- ✓ 즉, n개의 속성을 가진 데이터에는 최소 n+1개의 서포트 벡터가 존재한다는 걸 알 수 있음
- ✓ 대부분 머신러닝 지도 학습 알고리즘은 학습 데이터 모두를 사용하여 모델을 학습하는데 SVM에서는 결정 경계를 정의하는 게 결국 서포트 벡터이기 때문에 데이터 포인트 중에서 서포트 벡터만 잘 골라내면 나머지 쓸 데 없는 수많은 데이터 포인트들을 무시할 수 있기 때문에 매우 빠름

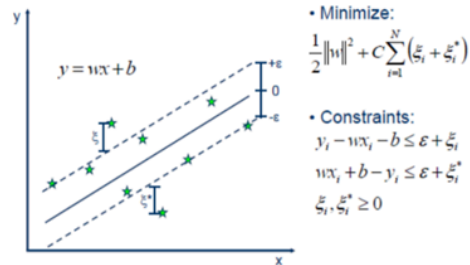
● 하드 마진 SVR 회귀모델

- 분류와는 다르게 예측에서 오차가 0이 되는 직선은 이론상으로도 존재하기 어려움
- 따라서 서포트 벡터 회귀에서는 마진을 최대화하면서 동시에 오차 범위 내($\pm\epsilon$)로 샘플이 들어오도록 예측 직선을 구축함



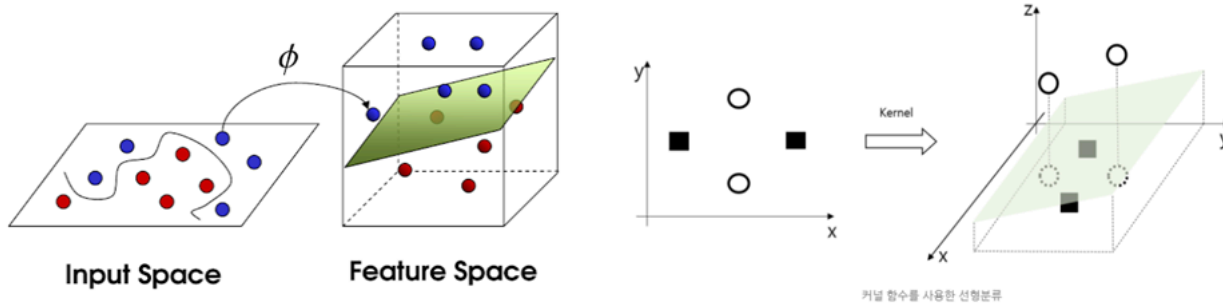
● 소프트 마진 SVR 회귀모델

- 그러나 마진을 최대화하면서 동시에 오차 범위 내($\pm\epsilon$)로 샘플이 들어오도록 예측 직선도 거의 존재하지 않아, 추가적인 예측 오차(양의 오차 ξ_i 와 음의 오차 ξ_i^*)를 허용함



✓ 커널 함수(kernel function)

- Support Vector Machine(SVM)과 같은 머신러닝 모델에서 사용되는 함수
- SVM(Support Vector Machine)은 두 개의 클래스를 가진 학습 데이터들을 구분하기 위해 두 범주 사이의 거리를 최대화 시키는 최적의 초평면을 찾는 이진분류기법이라 할 수 있음
- 일반적으로 분류하고 싶은 데이터는 비선형일 경우가 많기 때문에 SVM은 기본적으로 초평면을 사용하는 선형분류기법이지만 선형분류가 되지 않는 경우, 이를 해결하기 위해 그림과 같이 커널 함수를 통해 저차원의 데이터를 고차원으로 매핑하여 선형분류 가능하게 함



✓ [scikit-learn 라이브러리를 이용한 SVM 구현 - SVC\(1\)](#)

```
1 #SVC는 SVM에서 M을 C(Classification)으로 바꾼 것, 회귀는 SVR(Regression)임
2 from sklearn.svm import SVC
3 classifier = SVC(kernel = 'linear')
4 training_points=[[1,2],[1,6],[2, 2],[7,5],[9,4],[8,2]]
5 labels=[1,1,1,0,0,0] #레이블은 빨간 1, 파란 0
6 classifier.fit(training_points, labels)
```



SVC ⓘ ?
SVC(kernel='linear')

```
1 print(classifier.predict([[3, 2]])) # [3,2]는 빨간 점 1로 분류
```



[1]

```
1 #서포트 벡터(결정 경계)를 정의하는 서포트 벡터를 확인하려면
2 #classifier.support_vectors_를 point해보면 됨
3 classifier.support_vectors_
```



```
array([[7., 5.],
       [1., 6.],
       [2., 2.]])
```

✓ [scikit-learn 라이브러리를 이용한 SVM 구현 - SVC\(2\)](#)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
```

```

3 from sklearn import svm
4 from sklearn.datasets import make_blobs #from sklearn.datasets.samples_generator import make_blobs를 제공하지 않음
5 X, y = make_blobs(n_samples=40, centers=2, random_state=20)

```

```

1 clf = svm.SVC(kernel='linear')
2 clf.fit(X, y)

```

↗ SVC

SVC(kernel='linear')

```

1 newData = [[3,4]]
2 print(clf.predict(newData))

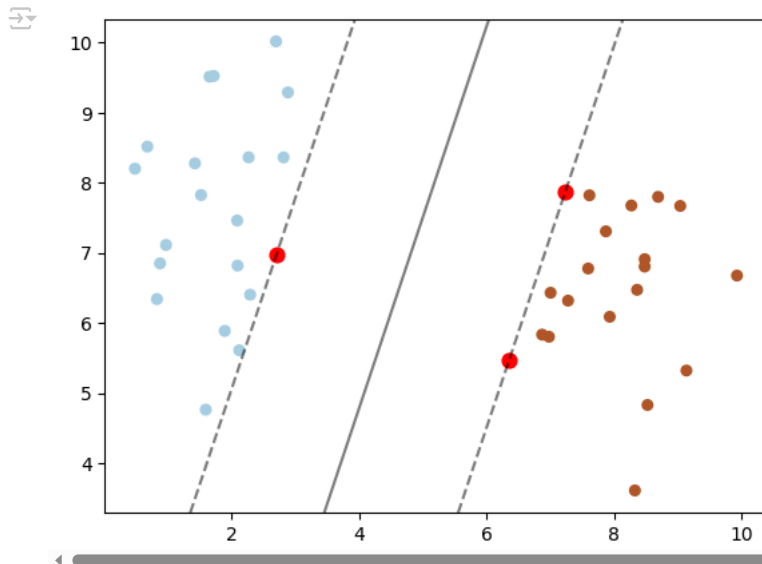
```

↗ [0]

```

1 # 샘플 데이터 표현
2 plt.scatter(X[:,0], X[:,1], c=y, s=30, cmap=plt.cm.Paired)
3 # 초평면(Hyper-Plane) 표현
4 ax = plt.gca()
5 xlim = ax.get_xlim()
6 ylim = ax.get_ylim()
7 xx = np.linspace(xlim[0], xlim[1], 30)
8 yy = np.linspace(ylim[0], ylim[1], 30)
9 YY, XX = np.meshgrid(yy, xx)
10 xy = np.vstack([XX.ravel(), YY.ravel()]).T
11 Z = clf.decision_function(xy).reshape(XX.shape)
12 ax.contour(XX, YY, Z, colors='k', levels=[-1,0,1], alpha=0.5, linestyles=['--', '-', '--'])
13 # 지지벡터(Support Vector) 표현
14 ax.scatter(clf.support_vectors_[:,0], clf.support_vectors_[:,1], s=60, facecolors='r')
15 plt.show()

```



✓ scikit-learn 라이브러리를 이용한 SVM 구현 - SVR

```
1 #https://medium.com/@niousha.rf/support-vector-regressor-theory-and-coding-exercise-in-python-ca6a7dfa927
2 #Dataset
3 import pandas as pd
4
5 df = pd.read_csv('/content/drive/MyDrive/머신러닝&딥러닝/머신러닝-7-SVM과 결정트리/Student_Marks.csv')
6 df.head()
```



	number_courses	time_study	Marks
0	3	4.508	19.202
1	4	0.096	7.734
2	4	3.133	13.811
3	6	7.909	53.018
4	8	7.811	55.299

```
1 from sklearn.model_selection import train_test_split
2
3 train, test = train_test_split(df, test_size=0.2, random_state=42)
4
5 # train and test datasets are sorted for plotting purpose
6 train = train.sort_values('time_study')
7 test = test.sort_values('time_study')
8
9 X_train, X_test = train[['time_study']], test[['time_study']]
10 y_train, y_test = train['Marks'], test['Marks']

1 #Feature scaling
2 from sklearn.preprocessing import StandardScaler
3
4 ### When using StandardScaler(), fit() method expects a 2D array-like input
5 scaler = StandardScaler().fit(X_train)
6 X_train_scaled = scaler.transform(X_train)
7 X_test_scaled = scaler.transform(X_test)

1 #Fitting SVR model
2 from sklearn.svm import SVR
3
4 svr_lin = SVR(kernel = 'linear')
5 svr_rbf = SVR(kernel = 'rbf')
6 svr_poly = SVR(kernel = 'poly')
7
8 svr_lin.fit(X_train_scaled, y_train)
```

```

9 svr_rbf.fit(X_train_scaled, y_train)
10 svr_poly.fit(X_train_scaled, y_train)

```



SVR ⓘ ?

SVR(kernel='poly')

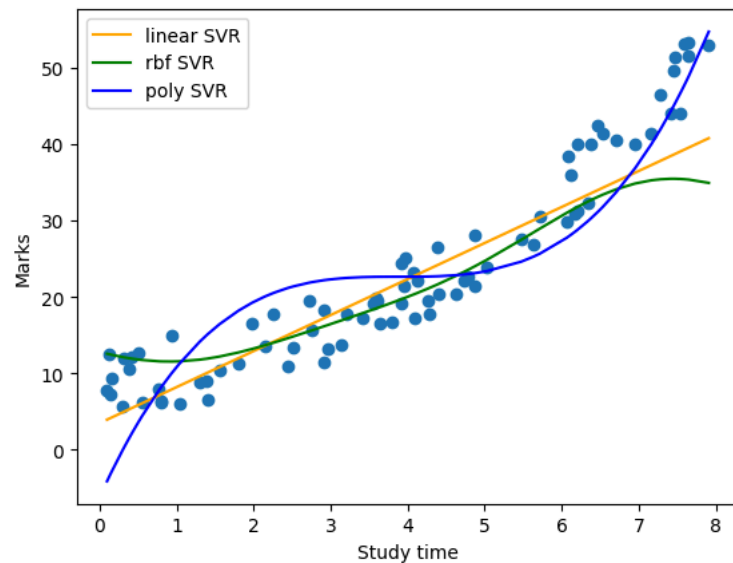
```

1 #Evaluating model performance
2 from matplotlib import pyplot as plt
3
4 ##### Model prediction for train dataset #####
5 train['linear_svr_pred'] = svr_lin.predict(X_train_scaled)
6 train['rbf_svr_pred'] = svr_rbf.predict(X_train_scaled)
7 train['poly_svr_pred'] = svr_poly.predict(X_train_scaled)
8
9 ##### Visualization #####
10 plt.scatter(train['time_study'], train['Marks'])
11 plt.plot(train['time_study'], train['linear_svr_pred'], color = 'orange', label = 'linear SVR')
12 plt.plot(train['time_study'], train['rbf_svr_pred'], color = 'green', label = 'rbf SVR')
13 plt.plot(train['time_study'], train['poly_svr_pred'], color = 'blue', label = 'poly SVR')
14 plt.legend()
15 plt.xlabel('Study time')
16 plt.ylabel('Marks')

```



Text(0, 0.5, 'Marks')



```

1 import numpy as np
2 import pandas as pd
3 from sklearn.preprocessing import LabelEncoder

```



```

4 from sklearn.model_selection import train_test_split
5 from matplotlib import pyplot as plt
6 from sklearn.svm import LinearSVC
7 from sklearn.metrics import confusion_matrix
8 from mpl_toolkits.mplot3d import Axes3D

1 df_iris=pd.read_csv('/content/drive/MyDrive/머신러닝&딥러닝/머신러닝-7-SVM과 결정트리/Iris.csv')
2 df_iris.head()
3 #https://www.kaggle.com/code/xopxesalmon/svm-2dimension-hyperplane-visualization

```



	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	1	5.1	3.5	1.4	0.2	Iris-setosa
1	2	4.9	3.0	1.4	0.2	Iris-setosa
2	3	4.7	3.2	1.3	0.2	Iris-setosa
3	4	4.6	3.1	1.5	0.2	Iris-setosa
4	5	5.0	3.6	1.4	0.2	Iris-setosa

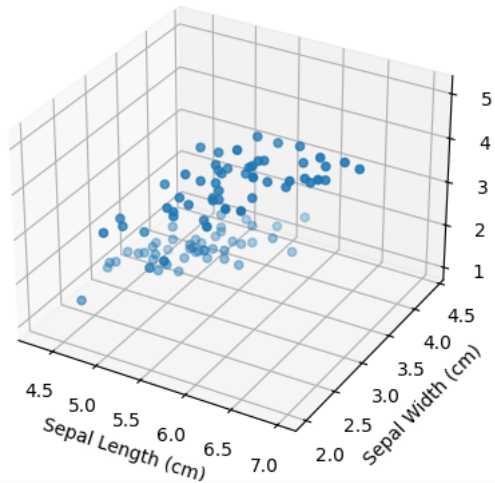
```

1 #discriminate between Iris Setosa and Iris Versicolor using 3 features: Sepal Length, Sepal width and petal length.
2 df=df_iris.iloc[:100,1:4]
3
4 ax = plt.axes(projection='3d')
5 ax.scatter(df['SepalLengthCm'], df['SepalWidthCm'], df['PetalLengthCm'])
6 ax.set_xlabel('Sepal Length (cm)')
7 ax.set_ylabel('Sepal Width (cm)')
8 ax.set_zlabel('Petal Length (cm)')

```



```
Text(0.5, 0, 'Petal Length (cm)')
```



```

1 #Features and labels
2 # The support vector machines in scikit-learn support both dense (numpy.ndarray and convertible to that by numpy.asarray)
3 X=df.to_numpy()
4
5 # Converting string value to int type for labels: Setosa = 0, Versicolor = 1
6 y=df_iris.iloc[:100,-1]
7 y = LabelEncoder().fit_transform(y)

```

```

1 #Splitting of the dataset and raining the model
2 X_train, X_test, y_train, y_test = train_test_split(X,y, random_state=0)
3
4 svc = LinearSVC()
5 svc.fit(X_train, y_train)

```



▼ LinearSVC ⓘ ?

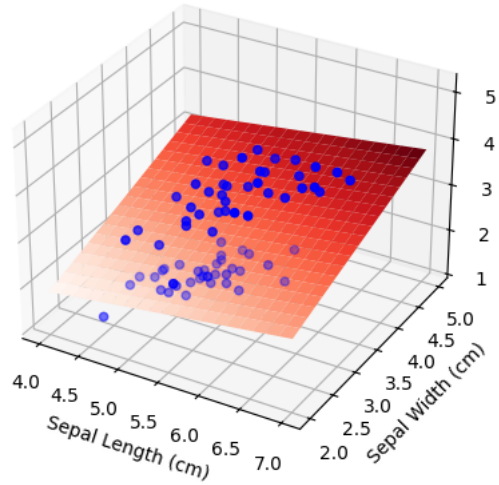
LinearSVC()

```

1 #Visualization of the hyperplane¶
2 #generates a plot of the vectors from X_train, and the SVM hyperplane.
3 plt.figure()
4 ax = plt.axes(projection='3d')
5 ax.scatter(X_train[:,0], X_train[:,1], X_train[:,2], c='b')
6 ax.set_xlabel('Sepal Length (cm)')
7 ax.set_ylabel('Sepal Width (cm)')
8 ax.set_zlabel('Petal Length (cm)')
9
10 zz = lambda xx,yy: (-svc.intercept_[0]-svc.coef_[0][0]*xx-svc.coef_[0][1]*yy) / svc.coef_[0][2]
11 tmpx = np.linspace(4, 7, 20)
12 tmpy = np.linspace(2, 5, 20)
13 xx,yy = np.meshgrid(tmpx,tmpy)
14 ax.plot_surface(xx, yy, zz(xx,yy), cmap='Reds')

```

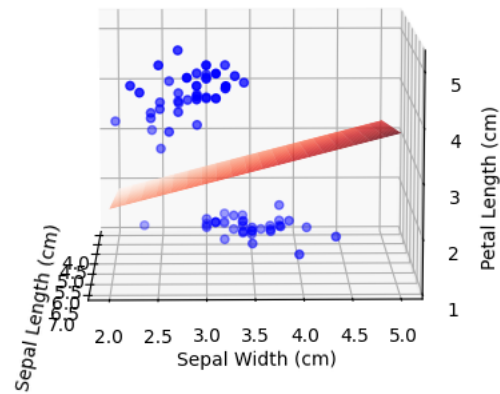
 <mpl_toolkits.mplot3d.art3d.Poly3DCollection at 0x78595f6958d0>



```

1 #The hyperplane shown from another angle:
2 plt.figure()
3 ax = plt.axes(projection='3d')
4 ax.scatter(X_train[:,0], X_train[:,1], X_train[:,2], c='b')
5 ax.set_xlabel('Sepal Length (cm)')
6 ax.set_ylabel('Sepal Width (cm)')
7 ax.set_zlabel('Petal Length (cm)')
8
9 zz = lambda xx,yy: (-svc.intercept_[0]-svc.coef_[0][0]*xx-svc.coef_[0][1]*yy) / svc.coef_[0][2]
10 tmpx = np.linspace(4, 7, 20)
11 tmpy = np.linspace(2, 5, 20)
12 xx,yy = np.meshgrid(tmpx,tmpy)
13 ax.plot_surface(xx, yy, zz(xx,yy), cmap='Reds')
14 for ii in range(0,360,1):
15     ax.view_init(elev=10., azim=ii)

```



```

1 #And from another angle
2 plt.figure()
3 ax = plt.axes(projection='3d')
4 ax.scatter(X_train[:,0], X_train[:,1], X_train[:,2], c='b')
5 ax.set_xlabel('Sepal Length (cm)')
6 ax.set_ylabel('Sepal Width (cm)')
7 ax.set_zlabel('Petal Length (cm)')
8
9 zz = lambda xx,yy: (-svc.intercept_[0]-svc.coef_[0][0]*xx-svc.coef_[0][1]*yy) / svc.coef_[0][2]
10 tmpx = np.linspace(4, 7, 20)
11 tmpy = np.linspace(2, 5, 20)
12 xx,yy = np.meshgrid(tmpx,tmpy)
13 ax.plot_surface(xx, yy, zz(xx,yy), cmap='Reds')
14 for ii in range(0,30,1):
15     ax.view_init(elev=25, azim=ii)

```



```
1 #Confusion Matrix:  
2 y_pred = svc.predict(X_test)  
3 confusion_matrix(y_test, y_pred)
```