

# Applied Machine Learning :Homework 2

Name: Priya Kumari

**Dataset:** The Alphabet dataset used below consists of images of alphabets represented by a total of 785 columns, where the first column represents the alphabet numbering from 0-25 as A-Z.

```
In [1]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn.pipeline import Pipeline
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.svm import SVC
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.ensemble import VotingClassifier
10 from sklearn.ensemble import BaggingClassifier
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.metrics import accuracy_score, f1_score
13 from sklearn.linear_model import LogisticRegression, SGDClassifier
14 import warnings
15 warnings.filterwarnings("ignore")
```

**Download and load the dataset. Split into train and test sets. Create a validation set . Display a few images.**

```
In [2]: 1 #Load the dataset
2 data=pd.read_csv("Alphabet_Data.csv")
3 data.head(10)
```

```
Out[2]:
```

	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	...	0.639	0.640	0.641	0.642	0.643	0.644	0.645	0.646	...
0	3	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
1	20	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
2	17	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
3	19	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
4	14	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
5	10	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
6	14	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
7	25	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0
9	14	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0

10 rows × 785 columns

The first column in the dataset, represented as '0', corresponds to the alphabetical characters from A-Z and is numbered from 0-25

```
In [3]: 1 data.shape
```

```
Out[3]: (93121, 785)
```

There are 93121 rows and 785 columns

```
In [4]: 1 #splitting data into training and testing sets.  
2 #The stratify parameter ensures that the label distribution is maintained between  
3 train,test = train_test_split(data, test_size=0.20, stratify=data['0'])  
4 train.shape,test.shape
```

```
Out[4]: ((74496, 785), (18625, 785))
```

```
In [5]: 1 #count the number of occurrences of each label in the train DataFrame  
2 label_counts = train['0'].value_counts()  
3 label_counts
```

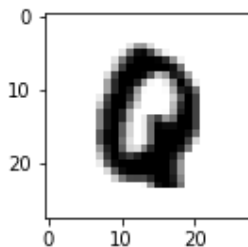
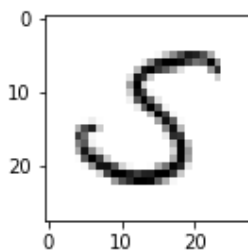
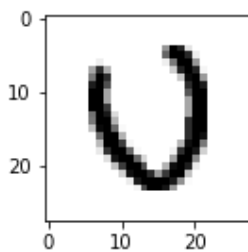
```
Out[5]: 14    11565  
18    9684  
20    5802  
2    4682  
19    4499  
15    3869  
13    3802  
0     2774  
12    2467  
11    2318  
17    2314  
4     2288  
24    2172  
22    2157  
3     2027  
1     1734  
9     1699  
7     1444  
23    1254  
25    1215  
16    1162  
6     1153  
10    1121  
21     837  
5      233  
8      224  
Name: 0, dtype: int64
```

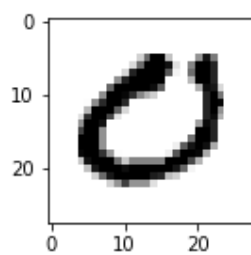
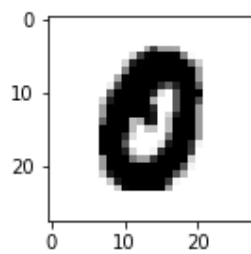
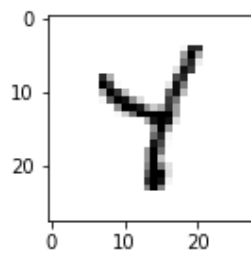
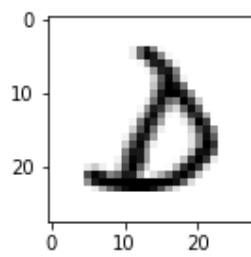
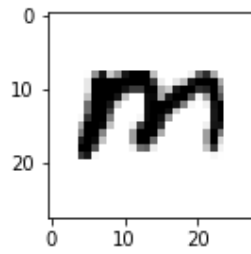
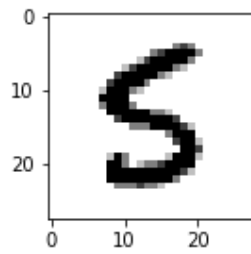
```
In [6]: 1 #dropping the first column ('0') from the train
        2 x = train.drop('0', axis=1)
        3 #extracting the label data from the train
        4 y = train['0']
        5
        6 #splitting the training data into a new training set and a validation set.
        7 x_train, x_val, y_train, y_val = train_test_split(x, y, test_size=0.16666, random_
        8
        9 x_test = test.drop('0', axis=1)
       10 y_test = test['0']
       11
```

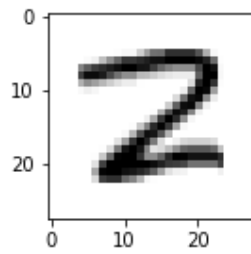
```
In [7]: 1 x_train.shape,x_val.shape,y_train.shape,y_val.shape
```

```
Out[7]: ((62080, 784), (12416, 784), (62080,), (12416,))
```

```
In [33]: 1 '''
2 Displaying a few more images. We are padding the pixel data into a 28x28 array and
3 plot the image. We are displaying 10 random images from the train dataset.
4 '''
5
6 disp_index = [1,141,100,811,512,3003,1111,807,3166,2294]
7
8 for m in disp_index:
9     img_data = []
10    k = []
11    cnt = 1
12    for z in x_train.loc[m]:
13        if(int(cnt)%28 == 0):
14            k.append(z)
15            img_data.append(k)
16            cnt+=1
17            k=[]
18        else:
19            k.append(z)
20            cnt+=1
21    plt.figure(figsize=(2, 2))
22    plt.imshow(img_data, cmap=plt.cm.binary)
23    plt.show()
24
```







## Multinomial Logistic Regression (softmax regression)

In [122]:

```
1  '''
2  We will be changing the values of 3 hyper-parameters:- penalty, max_iter and C, to
3  The solver parameter specifies the algorithm to use in the optimization problem.
4  The penalty parameter specifies the type of regularization to be used.
5  The max_iter parameter specifies the maximum number of iterations for the solver to
6  The C parameter specifies the inverse of the regularization strength.
7  '''
8
9  penalty = ['none', 'l2']
10 max_iter = [100, 500, 1000]
11 C = [0.001, 1.0, 100.0]
12
13 for a in penalty:
14     for b in max_iter:
15         for c in C:
16             model_1 = LogisticRegression(multi_class='multinomial', solver='lbfgs')
17             model_1.fit(x_train, y_train)
18             y_pred_train = model_1.predict(x_train)
19             accuracy_train = accuracy_score(y_train, y_pred_train)
20             f1_train=f1_score(y_train, y_pred_train, average='weighted')
21             y_pred_val = model_1.predict(x_val)
22             accuracy_val = accuracy_score(y_val, y_pred_val)
23             f1_val = f1_score(y_val, y_pred_val, average='weighted')
24             y_pred_test = model_1.predict(x_test)
25             accuracy_test = accuracy_score(y_test, y_pred_test)
26             f1_test = f1_score(y_test, y_pred_test, average='weighted')
27             print("Model with penalty =",a,"max_iter =",b,"C =",c)
28             print("Train Accuracy:", accuracy_train, "Validation Accuracy:", accuracy_val)
29             print("Train F1 score:", f1_train, "Validation F1 score:", f1_val, "Test F1 score:", f1_test)
30             print("\n")
```

Model with penalty = none max\_iter = 100 C = 0.001  
Train Accuracy: 0.892284149484536 Validation Accuracy: 0.8663820876288659 Test Accuracy: 0.8686174496644296  
Train F1 score: 0.891647500222681 Validation F1 score: 0.8651989850299983 Test f1 score: 0.8677401775515268

Model with penalty = none max\_iter = 100 C = 1.0  
Train Accuracy: 0.892284149484536 Validation Accuracy: 0.8663820876288659 Test Accuracy: 0.8686174496644296  
Train F1 score: 0.891647500222681 Validation F1 score: 0.8651989850299983 Test f1 score: 0.8677401775515268

Model with penalty = none max\_iter = 100 C = 100.0  
Train Accuracy: 0.892284149484536 Validation Accuracy: 0.8663820876288659 Test Accuracy: 0.8686174496644296  
Train F1 score: 0.891647500222681 Validation F1 score: 0.8651989850299983 Test f1 score: 0.8677401775515268

Model with penalty = none max\_iter = 500 C = 0.001  
Train Accuracy: 0.9029961340206185 Validation Accuracy: 0.8563949742268041 Test Accuracy: 0.8589530201342281  
Train F1 score: 0.9025757203615234 Validation F1 score: 0.8552185074120079 Test f1 score: 0.8581371423107979

Model with penalty = none max\_iter = 500 C = 1.0  
Train Accuracy: 0.9029961340206185 Validation Accuracy: 0.8563949742268041 Test Accuracy: 0.8589530201342281  
Train F1 score: 0.9025757203615234 Validation F1 score: 0.8552185074120079 Test f1 score: 0.8581371423107979

Model with penalty = none max\_iter = 500 C = 100.0  
Train Accuracy: 0.9029961340206185 Validation Accuracy: 0.8563949742268041 Test Accuracy: 0.8589530201342281  
Train F1 score: 0.9025757203615234 Validation F1 score: 0.8552185074120079 Test f1 score: 0.8581371423107979

Model with penalty = none max\_iter = 1000 C = 0.001  
Train Accuracy: 0.9040753865979382 Validation Accuracy: 0.8564755154639175 Test Accuracy: 0.8577181208053691  
Train F1 score: 0.9036727198096002 Validation F1 score: 0.8553038059290389 Test f1 score: 0.8568789033267602

Model with penalty = none max\_iter = 1000 C = 1.0  
Train Accuracy: 0.9040753865979382 Validation Accuracy: 0.8564755154639175 Test Accuracy: 0.8577181208053691  
Train F1 score: 0.9036727198096002 Validation F1 score: 0.8553038059290389 Test f1 score: 0.8568789033267602

Model with penalty = none max\_iter = 1000 C = 100.0  
Train Accuracy: 0.9040753865979382 Validation Accuracy: 0.8564755154639175 Test Accuracy: 0.8577181208053691  
Train F1 score: 0.9036727198096002 Validation F1 score: 0.8553038059290389 Test f1 score: 0.8568789033267602



ore: 0.8568789033267602

Model with penalty = 12 max\_iter = 100 C = 0.001

Train Accuracy: 0.8929284793814433 Validation Accuracy: 0.8682345360824743 Test Accuracy: 0.8694228187919463

Train F1 score: 0.8924003626502173 Validation F1 score: 0.8671584528559252 Test f1 score: 0.8686607434404423

Model with penalty = 12 max\_iter = 100 C = 1.0

Train Accuracy: 0.8917525773195877 Validation Accuracy: 0.8675096649484536 Test Accuracy: 0.8686711409395973

Train F1 score: 0.8911845070194099 Validation F1 score: 0.8663975061525153 Test f1 score: 0.8679463086828001

Model with penalty = 12 max\_iter = 100 C = 100.0

Train Accuracy: 0.8919942010309279 Validation Accuracy: 0.8663820876288659 Test Accuracy: 0.8679731543624161

Train F1 score: 0.8913908258754442 Validation F1 score: 0.8652234984330588 Test f1 score: 0.8671435139756593

Model with penalty = 12 max\_iter = 500 C = 0.001

Train Accuracy: 0.9032538659793814 Validation Accuracy: 0.8571198453608248 Test Accuracy: 0.8610469798657718

Train F1 score: 0.9028353032399865 Validation F1 score: 0.8560013192151189 Test f1 score: 0.8602397770722425

Model with penalty = 12 max\_iter = 500 C = 1.0

Train Accuracy: 0.9029317010309278 Validation Accuracy: 0.8572003865979382 Test Accuracy: 0.859489932885906

Train F1 score: 0.9025047919685264 Validation F1 score: 0.8560037496028221 Test f1 score: 0.8586540573011576

Model with penalty = 12 max\_iter = 500 C = 100.0

Train Accuracy: 0.9032216494845361 Validation Accuracy: 0.8569587628865979 Test Accuracy: 0.8597046979865772

Train F1 score: 0.9027724217156703 Validation F1 score: 0.8557489475135232 Test f1 score: 0.8588944089479061

Model with penalty = 12 max\_iter = 1000 C = 0.001

Train Accuracy: 0.90402706185567 Validation Accuracy: 0.8574420103092784 Test Accuracy: 0.8588993288590604

Train F1 score: 0.9036204765152476 Validation F1 score: 0.8562521326806717 Test f1 score: 0.8580541457801041

Model with penalty = 12 max\_iter = 1000 C = 1.0

Train Accuracy: 0.9041398195876289 Validation Accuracy: 0.8556701030927835 Test Accuracy: 0.8581476510067114

Train F1 score: 0.9037383504921748 Validation F1 score: 0.8544409275325945 Test f1 score: 0.8573709508617379

Model with penalty = 12 max\_iter = 1000 C = 100.0

Train Accuracy: 0.9043009020618556 Validation Accuracy: 0.8556701030927835 Test Accuracy:

```
acy: 0.8572348993288591
```

```
Train F1 score: 0.9038944094755548 Validation F1 score: 0.8544561984338221 Test f1 score: 0.8564640963192844
```

The logistic regression model with  $\text{penalty}=\text{l2}$ ,  $\text{max\_iter}=100$ , and  $C=0.001$  was found to be the best-performing model based on the results obtained. It achieved higher test accuracy (0.869) and test F1 score (0.8686) than all other models. The use of regularization through  $\text{penalty}=\text{l2}$  generally led to better performance across the models. Although logistic regression achieved an accuracy of about 85-86% on the test set, there may be other models that could potentially outperform it and produce better results.

## Support Vector Machines

## SGD as LinearSVC

In [11]:

```

1  C = [0.001, 1.0, 100.0]
2  for c in C:
3      sgd = SGDClassifier(loss='hinge', alpha=1/(len(x_train)*c), fit_intercept=True)
4      sgd.fit(x_train, y_train)
5      y_pred_train = sgd.predict(x_train)
6      accuracy_train = accuracy_score(y_train, y_pred_train)
7      f1_train=f1_score(y_train, y_pred_train, average='weighted')
8      y_pred_val = sgd.predict(x_val)
9      accuracy_val = accuracy_score(y_val, y_pred_val)
10     f1_val = f1_score(y_val, y_pred_val, average='weighted')
11     y_pred_test = sgd.predict(x_test)
12     accuracy_test = accuracy_score(y_test, y_pred_test)
13     f1_test = f1_score(y_test, y_pred_test, average='weighted')
14     print("Model with C =",c)
15     print("Train Accuracy:", accuracy_train, "Validation Accuracy:", accuracy_val,
16     print("Train F1 score:", f1_train, "Validation F1 score:", f1_val, "Test f1 score:", f1_test)
17     print("\n")

```

Model with C = 0.001

Train Accuracy: 0.7882087628865979 Validation Accuracy: 0.7629671391752577 Test Accuracy: 0.7706308724832215

Train F1 score: 0.7892802458571851 Validation F1 score: 0.7616655877082997 Test f1 score: 0.7711497492937867

Model with C = 1.0

Train Accuracy: 0.7020779639175257 Validation Accuracy: 0.6805734536082474 Test Accuracy: 0.6904161073825503

Train F1 score: 0.6966198646842083 Validation F1 score: 0.6721897550830878 Test f1 score: 0.6833812781515608

Model with C = 100.0

Train Accuracy: 0.7655605670103093 Validation Accuracy: 0.7417042525773195 Test Accuracy: 0.7481342281879194

Train F1 score: 0.7668890017104963 Validation F1 score: 0.7420549745762961 Test f1 score: 0.7485420978720004

The value of the regularization parameter C determines the trade-off between the model's ability to fit the training data closely and its ability to generalize well to new data. When C is set to a very small value like 0.001, the model is highly regularized and may not fit the training data well, but it may generalize better to new data. However, if C is set too high, the model may overfit to the training data and not generalize well. Therefore, finding the optimal value of C requires careful experimentation and evaluation on multiple datasets. In the above scenario, model with c=0.001 performs better with test accuracy as 0.77

## Using different Kernels

In [13]:

```
1 kernel = ['linear', 'rbf', 'poly']
2 C = [0.001, 1.0, 100.0]
3 gamma = ['scale', 'auto']
4
5
6 batch_size = 10000
7 num_batches = int(x_train.shape[0] / batch_size)
8
9 for a in kernel:
10     for b in C:
11         for c in gamma:
12             model_2 = SVC(kernel=a, gamma=c, C=b)
13             for i in range(num_batches):
14                 batch_x = x_train[i*batch_size:(i+1)*batch_size]
15                 batch_y = y_train[i*batch_size:(i+1)*batch_size]
16                 model_2.fit(batch_x, batch_y)
17             y_pred_train = model_2.predict(x_train)
18             accuracy_train = accuracy_score(y_train, y_pred_train)
19             f1_train = f1_score(y_train, y_pred_train, average='weighted')
20             y_pred_val = model_2.predict(x_val)
21             accuracy_val = accuracy_score(y_val, y_pred_val)
22             f1_val = f1_score(y_val, y_pred_val, average='weighted')
23             y_pred_test = model_2.predict(x_test)
24             accuracy_test = accuracy_score(y_test, y_pred_test)
25             f1_test = f1_score(y_test, y_pred_test, average='weighted')
26             print("Model with kernel =", a, " C=", b, " gamma=", c)
27             print("Train Accuracy:", accuracy_train, "Validation Accuracy:", accuracy_val)
28             print("Train f1 score:", f1_train, "Validation f1 score:", f1_val, "Test f1 score:", f1_test)
29             print("\n")
```

Model with kernel = linear C= 0.001 gamma= scale  
Train Accuracy: 0.8890302835051547 Validation Accuracy: 0.8658182989690721 Test Accuracy: 0.8653959731543625  
Train f1 score: 0.889267013426636 Validation f1 score: 0.865771021155259 Test f1 score: 0.8655351371142919

Model with kernel = linear C= 0.001 gamma= auto  
Train Accuracy: 0.8890302835051547 Validation Accuracy: 0.8658182989690721 Test Accuracy: 0.8653959731543625  
Train f1 score: 0.889267013426636 Validation f1 score: 0.865771021155259 Test f1 score: 0.8655351371142919

Model with kernel = linear C= 1.0 gamma= scale  
Train Accuracy: 0.8890302835051547 Validation Accuracy: 0.8658182989690721 Test Accuracy: 0.8653959731543625  
Train f1 score: 0.889267013426636 Validation f1 score: 0.865771021155259 Test f1 score: 0.8655351371142919

Model with kernel = linear C= 1.0 gamma= auto  
Train Accuracy: 0.8890302835051547 Validation Accuracy: 0.8658182989690721 Test Accuracy: 0.8653959731543625  
Train f1 score: 0.889267013426636 Validation f1 score: 0.865771021155259 Test f1 score: 0.8655351371142919

Model with kernel = linear C= 100.0 gamma= scale  
Train Accuracy: 0.8890302835051547 Validation Accuracy: 0.8658182989690721 Test Accuracy: 0.8653959731543625  
Train f1 score: 0.889267013426636 Validation f1 score: 0.865771021155259 Test f1 score: 0.8655351371142919

Model with kernel = linear C= 100.0 gamma= auto  
Train Accuracy: 0.8890302835051547 Validation Accuracy: 0.8658182989690721 Test Accuracy: 0.8653959731543625  
Train f1 score: 0.889267013426636 Validation f1 score: 0.865771021155259 Test f1 score: 0.8655351371142919

Model with kernel = rbf C= 0.001 gamma= scale  
Train Accuracy: 0.1552673969072165 Validation Accuracy: 0.15512242268041238 Test Accuracy: 0.1552751677852349  
Train f1 score: 0.04173573080463085 Validation f1 score: 0.04166305760458433 Test f1 score: 0.041739627758042415

Model with kernel = rbf C= 0.001 gamma= auto  
Train Accuracy: 0.1552673969072165 Validation Accuracy: 0.15512242268041238 Test Accuracy: 0.1552751677852349  
Train f1 score: 0.04173573080463085 Validation f1 score: 0.04166305760458433 Test f1 score: 0.041739627758042415

Model with kernel = rbf C= 1.0 gamma= scale  
Train Accuracy: 0.9306056701030928 Validation Accuracy: 0.9186533505154639 Test Accuracy: 0.9254765100671141  
Train f1 score: 0.9299107157880034 Validation f1 score: 0.917618645244788 Test f1 score: 0.9254765100671141

re: 0.9246227986391252

Model with kernel = rbf C= 1.0 gamma= auto

Train Accuracy: 0.30908505154639176 Validation Accuracy: 0.17694909793814434 Test Accuracy: 0.17540939597315436

Train f1 score: 0.3082936584989006 Validation f1 score: 0.08488258897942177 Test f1 score: 0.08174153251994004

Model with kernel = rbf C= 100.0 gamma= scale

Train Accuracy: 0.9477931701030928 Validation Accuracy: 0.9358086340206185 Test Accuracy: 0.9385234899328859

Train f1 score: 0.9475220171112791 Validation f1 score: 0.9353216091484924 Test f1 score: 0.9381063875868415

Model with kernel = rbf C= 100.0 gamma= auto

Train Accuracy: 0.30908505154639176 Validation Accuracy: 0.17694909793814434 Test Accuracy: 0.17540939597315436

Train f1 score: 0.3082936584989006 Validation f1 score: 0.08488258897942177 Test f1 score: 0.08174153251994004

Model with kernel = poly C= 0.001 gamma= scale

Train Accuracy: 0.24370167525773195 Validation Accuracy: 0.2506443298969072 Test Accuracy: 0.24515436241610739

Train f1 score: 0.15088529091655406 Validation f1 score: 0.1542103195833361 Test f1 score: 0.15112744449369034

Model with kernel = poly C= 0.001 gamma= auto

Train Accuracy: 0.9256926546391753 Validation Accuracy: 0.9088273195876289 Test Accuracy: 0.9128053691275168

Train f1 score: 0.9254142864550521 Validation f1 score: 0.9084467673485823 Test f1 score: 0.9125670014327856

Model with kernel = poly C= 1.0 gamma= scale

Train Accuracy: 0.8939432989690722 Validation Accuracy: 0.8805573453608248 Test Accuracy: 0.8838120805369127

Train f1 score: 0.8947350879540427 Validation f1 score: 0.8807599846558871 Test f1 score: 0.8845162296370496

Model with kernel = poly C= 1.0 gamma= auto

Train Accuracy: 0.9256926546391753 Validation Accuracy: 0.9088273195876289 Test Accuracy: 0.9128053691275168

Train f1 score: 0.9254142864550521 Validation f1 score: 0.9084467673485823 Test f1 score: 0.9125670014327856

Model with kernel = poly C= 100.0 gamma= scale

Train Accuracy: 0.9257731958762887 Validation Accuracy: 0.909068943298969 Test Accuracy: 0.913020134228188

Train f1 score: 0.925484797501365 Validation f1 score: 0.9086807176284423 Test f1 score: 0.9127686168427475

Model with kernel = poly C= 100.0 gamma= auto

Train Accuracy: 0.9256926546391753 Validation Accuracy: 0.9088273195876289 Test Accuracy:

```
acy: 0.9128053691275168
```

```
Train f1 score: 0.9254142864550521 Validation f1 score: 0.9084467673485823 Test f1 score: 0.9125670014327856
```

From the results above, it appears that the linear kernel with different combinations of C and gamma values produce similar accuracy and f1 score results for all three datasets, whereas the rbf kernel yielded better results, particularly with higher C values. The most optimal model was identified as the rbf kernel with C=100 and gamma=scale. It is worth noting that the rbf kernel with gamma=auto performed poorly, potentially due to the default calculation of  $1 / (n\_features * X.var())$  not being the most appropriate for this dataset. Overall, the rbf kernel with C=100 and gamma=scale emerged as the most effective model out of all the models examined with an accuracy of 0.9385.

## Polynomial Kernel

In [9]:

```
1 degrees = [1, 2, 3, 5, 7, 9]
2
3 batch_size = 10000
4 num_batches = int(x_train.shape[0] / batch_size)
5 for degree_value in degrees:
6     model_2 = SVC(kernel='poly', degree=degree_value)
7     for i in range(num_batches):
8         batch_x = x_train[i*batch_size:(i+1)*batch_size]
9         batch_y = y_train[i*batch_size:(i+1)*batch_size]
10        model_2.fit(batch_x, batch_y)
11    y_pred_train = model_2.predict(x_train)
12    accuracy_train = accuracy_score(y_train, y_pred_train)
13    f1_train = f1_score(y_train, y_pred_train, average='weighted')
14    y_pred_val = model_2.predict(x_val)
15    accuracy_val = accuracy_score(y_val, y_pred_val)
16    f1_val = f1_score(y_val, y_pred_val, average='weighted')
17    y_pred_test = model_2.predict(x_test)
18    accuracy_test = accuracy_score(y_test, y_pred_test)
19    f1_test = f1_score(y_test, y_pred_test, average='weighted')
20    print("Model with degree =", degree_value)
21    print("Train Accuracy:", accuracy_train, "Validation Accuracy:", accuracy_val)
22    print("Train f1 score:", f1_train, "Validation f1 score:", f1_val, "Test f1 score:", f1_test)
23    print("\n")
```



Model with degree = 1

Train Accuracy: 0.8852770618556701 Validation Accuracy: 0.8770135309278351 Test Accuracy: 0.88048322147651

Train f1 score: 0.8841154669943685 Validation f1 score: 0.8753495745400449 Test f1 score: 0.8789631799001123

Model with degree = 2

Train Accuracy: 0.9165431701030928 Validation Accuracy: 0.9060083762886598 Test Accuracy: 0.9111946308724832

Train f1 score: 0.9157512542693804 Validation f1 score: 0.9048787243771336 Test f1 score: 0.9102608282787787

Model with degree = 3

Train Accuracy: 0.8939432989690722 Validation Accuracy: 0.8805573453608248 Test Accuracy: 0.8838120805369127

Train f1 score: 0.8947350879540427 Validation f1 score: 0.8807599846558871 Test f1 score: 0.8845162296370496

Model with degree = 5

Train Accuracy: 0.7978898195876288 Validation Accuracy: 0.7771423969072165 Test Accuracy: 0.7822281879194631

Train f1 score: 0.8181555529923229 Validation f1 score: 0.7967383096408899 Test f1 score: 0.8017207081000921

Model with degree = 7

Train Accuracy: 0.7019813144329897 Validation Accuracy: 0.6777545103092784 Test Accuracy: 0.6812885906040268

Train f1 score: 0.743294641947625 Validation f1 score: 0.7184499351498191 Test f1 score: 0.7215492091921462

Model with degree = 9

Train Accuracy: 0.6267557989690722 Validation Accuracy: 0.5932667525773195 Test Accuracy: 0.593503355704698

Train f1 score: 0.6810959867662455 Validation f1 score: 0.6454973531595835 Test f1 score: 0.6466184503337715

The results indicate that the polynomial model with degree=2 performs the best on both the validation and test sets, with the highest accuracy and F1 score. The models with degree=1 and 3 also perform well, but not as well as the degree=2 model. On the other hand, the models with degree=5, 7 and 9 perform worse than the models with degree=1, 2, and 3. Therefore, it appears that a polynomial model with degree=2 is the most suitable choice for this particular classification task with an accuracy of 0.91

## Random Forest Classifier

```

In [12]: 1 '''
2 We will be changing the values of 3 hyper-parameters:- n_estimators, max_depth and
3 '''
4
5 n_estimators = [10, 50, 100]
6 max_depth = [5, 10, 20]
7 min_samples_split = [2, 5, 10]
8
9 for a in n_estimators:
10     for b in max_depth:
11         for c in min_samples_split:
12             model_3 = RandomForestClassifier(n_estimators=a, max_depth=b, min_samp
13             model_3.fit(x_train, y_train)
14             y_pred_train = model_3.predict(x_train)
15             accuracy_train = accuracy_score(y_train, y_pred_train)
16             y_pred_val = model_3.predict(x_val)
17             accuracy_val = accuracy_score(y_val, y_pred_val)
18             y_pred_test = model_3.predict(x_test)
19             accuracy_test = accuracy_score(y_test, y_pred_test)
20             print("Model with n_estimators =",a,"max_depth =",b,"min_samples_split
21             print("Train Accuracy:", accuracy_train, "Validation Accuracy:", accur
22             print("\n")

```

Train Accuracy: 0.6402867268041237 Validation Accuracy: 0.6315238402061856 Test Acc  
 uracy: 0.6374228187919463

Model with n\_estimators = 50 max\_depth = 10 min\_samples\_split = 2  
 Train Accuracy: 0.921520618556701 Validation Accuracy: 0.8817654639175257 Test Accu  
 racy: 0.8878389261744967

Model with n\_estimators = 50 max\_depth = 10 min\_samples\_split = 5  
 Train Accuracy: 0.9209729381443299 Validation Accuracy: 0.8848260309278351 Test Acc  
 uracy: 0.8897718120805369

Model with n\_estimators = 50 max\_depth = 10 min\_samples\_split = 10  
 Train Accuracy: 0.918347293814433 Validation Accuracy: 0.8784632731958762 Test Accu  
 racy: 0.8896107382550336

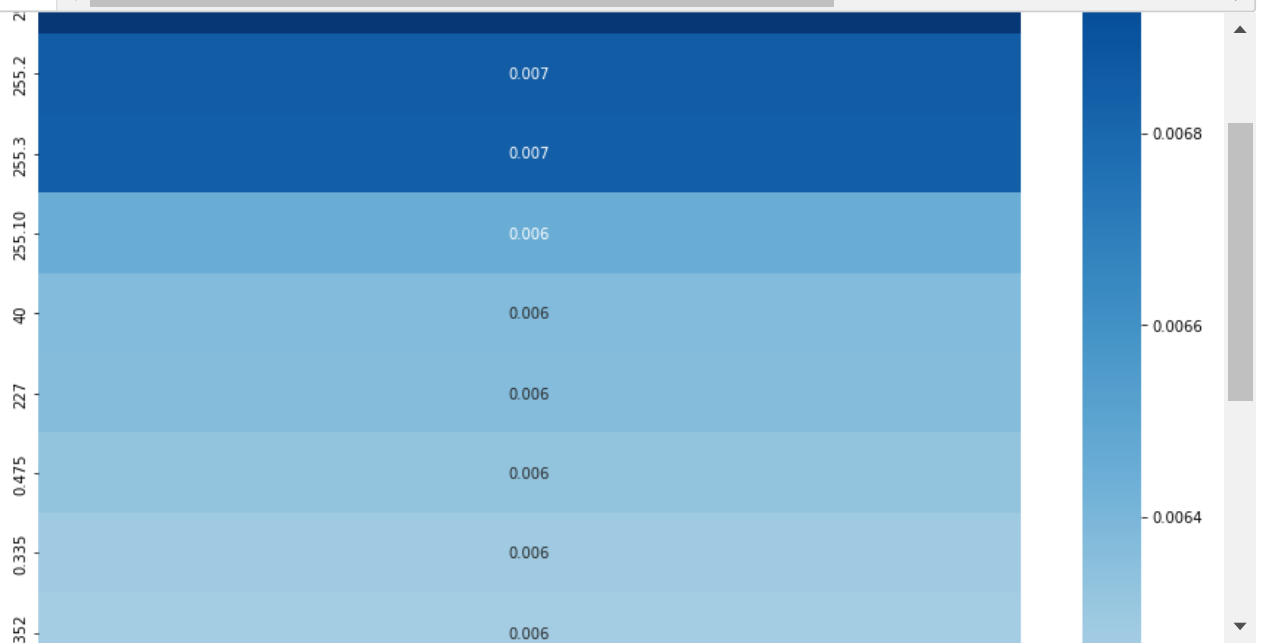
Model with n\_estimators = 50 max\_depth = 20 min\_samples\_split = 2

Based on the above results, it was observed that the model with n\_estimators = 100, max\_depth = 20, and min\_samples\_split = 2 had the best accuracy with value of 0.944 on both the validation and test sets. While increasing the number and depth of decision trees tends to increase accuracy, excessive depth may result in overfitting. Additionally, an increase in the minimum number of samples required to split an internal node can prevent overfitting.

```

In [16]: 1 '''
2 Feature importance analysis using in-built function for random forest. From this, v
3 impacting and features which have the least influence on the model.
4 '''
5 import seaborn as sns
6 clf_rf = RandomForestClassifier(n_estimators=100, max_depth=20, min_samples_split=
7 clf_rf.fit(x_train, y_train)
8 feature_importance = dict(zip(x_train.columns, clf_rf.feature_importances_))
9 indices = np.argsort(feature_importance)[::-1]
10
11 feature_importance = dict(zip(x_train.columns, clf_rf.feature_importances_))
12
13 df_feature_importance = pd.DataFrame.from_dict(feature_importance, orient='index',
14 df_feature_importance = df_feature_importance.sort_values('Importance', ascending=
15
16 # Selecting top 15 features
17 top_features = df_feature_importance[:15]
18
19
20 fig, ax = plt.subplots(figsize=(15,15))
21 sns.heatmap(top_features, cmap='Blues', annot=True, fmt='.3f', ax=ax)
22 plt.title('Feature Importance (Top 15)')
23 plt.show()

```



The above plot shows top 15 features which are selected based on their importance scores, and a heatmap is plotted. The heatmap exhibits the top 15 features with their respective importance scores, where a higher score implies a greater importance in predicting the target variable. The column '0.372' has the highest feature importance with a value of 0.007.

## Ensemble

```
In [17]: 1 '''
2 We will now create an ensemble model out of the 3 best performing models which are
3 Multinomial Logistic Regression (Softmax Regression) with penalty = 'l2' | max_iter
4 SVM with kernel as 'rbf'
5 Random Forest Classifier with n_estimators = 100 | max_depth = 20 | min_samples_sp
6 '''
7
8 clf_lr = LogisticRegression(multi_class='multinomial', solver='lbfgs', penalty='l2
9 clf_svm = SVC(kernel='rbf', probability=True)
10 clf_rf = RandomForestClassifier(n_estimators=100, max_depth=20, min_samples_split=
```

```
In [18]: 1 voting_clf = VotingClassifier(
2     estimators=[('lr', clf_lr), ('svm', clf_svm), ('rf', clf_rf)],
3     voting='soft')
4 voting_clf.fit(x_train, y_train)
```

```
Out[18]: VotingClassifier(estimators=[('lr',
                                         LogisticRegression(C=0.001,
                                                                multi_class='multinomial')),
                                       ('svm', SVC(probability=True)),
                                       ('rf',
                                         RandomForestClassifier(max_depth=20,
                                                                min_samples_split=10))],
                           voting='soft')
```

```
In [19]: 1 y_pred = voting_clf.predict(x_test)
2 accuracy = accuracy_score(y_test, y_pred)
3 print("Accuracy:", accuracy)
```

Accuracy: 0.9495838926174497

The above voting classifier utilizes three models, namely logistic regression, support vector machine, and random forest, to combine their predictions using soft voting. In soft voting, the predicted class is selected based on the highest probability of a class across all models. The accuracy of the voting classifier is reported to be almost 0.95. This accuracy is better than all the individual accuracies achieved by the three models used in the ensemble, suggesting that the combination of models and the voting method helped to improve the overall accuracy of the predictions.