

ML Model for Domian classifications based on profile summary and profile headline

In [2]:

```
1 import pandas as pd
2 import re
3 from sklearn.feature_selection import chi2
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import string
7 import nltk
8 from nltk.corpus import stopwords
9 from nltk.tokenize import word_tokenize
10 from sklearn.feature_extraction.text import TfidfVectorizer
11 from sklearn.model_selection import train_test_split
12 from sklearn.feature_extraction.text import CountVectorizer
13 from sklearn.feature_extraction.text import TfidfTransformer
14 from sklearn.naive_bayes import MultinomialNB
15 from sklearn.metrics import accuracy_score
16 from sklearn.svm import SVC
17 from sklearn.neighbors import KNeighborsClassifier
18 from sklearn.naive_bayes import GaussianNB
19 from sklearn.ensemble import RandomForestClassifier as rm
20 from sklearn.tree import DecisionTreeClassifier
21 from sklearn.linear_model import SGDClassifier
22 from sklearn.neural_network import MLPClassifier
23 from sklearn.gaussian_process import GaussianProcessClassifier
24 from sklearn.linear_model import LogisticRegression
25 from sklearn.ensemble import RandomForestClassifier
26 from sklearn.svm import LinearSVC
27 from sklearn.model_selection import cross_val_score
28 from sklearn.model_selection import StratifiedKFold
29 import warnings
30 warnings.filterwarnings("ignore")
```

In [5]:

```
1 df1 = pd.read_csv('Before_BERT(Preprocessed).csv')
```

In [6]:

```
1 df1['ProfileHeadline'] = df1['Profile Headline']
2 df1['ProfileSummary'] = df1['Profile Summary']
```

In [7]:

```
1 df1['ProfileSummary'].isnull().value_counts()
```

Out[7]:

```
False    3465
True      1445
Name: ProfileSummary, dtype: int64
```

In [8]:

```
1 df1.ProfileHeadline = df1.ProfileHeadline.fillna(' ')
2 df1.ProfileSummary = df1.ProfileSummary.fillna(' ')

```

In [9]:

```
1 #df1['total'] = df1['ProfileSummary']
2 df1['total'] = df1['ProfileHeadline'].str.cat(df1['ProfileSummary'],sep=" ")
3 df1.head()

```

Out[9]:

	Unnamed: 0	Job Title	Company Name	Industry	Company Location	New Job (90 Days)	Year Started	Profile Headline
0	0	Battery Designer	Rivian	Mechanical or Industrial Engineering	Dublin, Ohio, United States	False	2020.0	Mechanical Design Engineer, System Integration...
1	1	Digital DevOps Engineer	HSBC	Information Technology and Services	New York City Metropolitan Area	False	2018.0	Digital DevOps Engineer at HSBC
2	2	Product Designer	Two Point Conversions, Inc.	Information Technology and Services	Chicago, Illinois, United States	False	2018.0	Leading Product + UX at Remedy (Two Point Conv... hti
3	3	Product Designer	udaan.com	Information Technology and Services	Bangalore Urban, Karnataka, India	False	2018.0	Product Designer at udaan
4	4	Digital Technology Intern	GE	Information Technology and Services	Jaipur, Rajasthan, India	True	2021.0	Digital Technology Intern at General Electric ...

5 rows × 23 columns

In [10]:

```
1 df1['total'] = df1['total'].astype(str)

```

In [11]:

```
1 df1['job_category'] = df1['Domain'].str.lower()

```

In [12]:



```
1 # helper function to clean the text in the data
2 def preprocess_text(text):
3     text.lower()
4     # Remove urls
5     text = re.sub(r"http\S+|www\S+|https\S+", '', text, flags=re.MULTILINE)
6     # Remove user @ references and '#' from text
7     text = re.sub(r'\@\w+|\#','', text)
8     # Remove punctuations
9     text = text.translate(str.maketrans('', '', string.punctuation))
10    # Remove stopwords
11    text_tokens = word_tokenize(text)
12    filtered_words = [w for w in text_tokens if not w in stop_words]
13    return " ".join(filtered_words).lower()
```

In [13]:



```
1 df1['total'][2]
```

Out[13]:

'Leading Product + UX at Remedy (Two Point Conversions, Inc.) <http://aroonmathai.com>' (<http://aroonmathai.com>)'

In [14]:



```
1 stop_words = stopwords.words('english')
```

In [15]:



```
1 # import nltk
2 # nltk.download('stopwords')
```

In [16]:



```
1 df1['preprocessed'] = df1['total'].apply(preprocess_text)
2 df1['preprocessed'][2]
```

Out[16]:

'leading product ux remedy two point conversions inc'

In [17]:

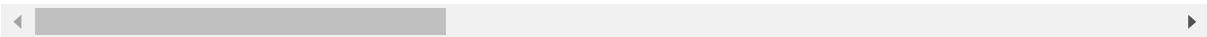


```
1 # creating encoded label of Domain in new column
2 df1['category_id'] = df1['job_category'].factorize()[0]
3 category_id_df = df1[['job_category', 'category_id']].drop_duplicates().sort_values('category_id')
4 category_to_id = dict(category_id_df.values)
5 id_to_category = dict(category_id_df[['category_id', 'job_category']].values)
6 df1.head()
```

Out[17]:

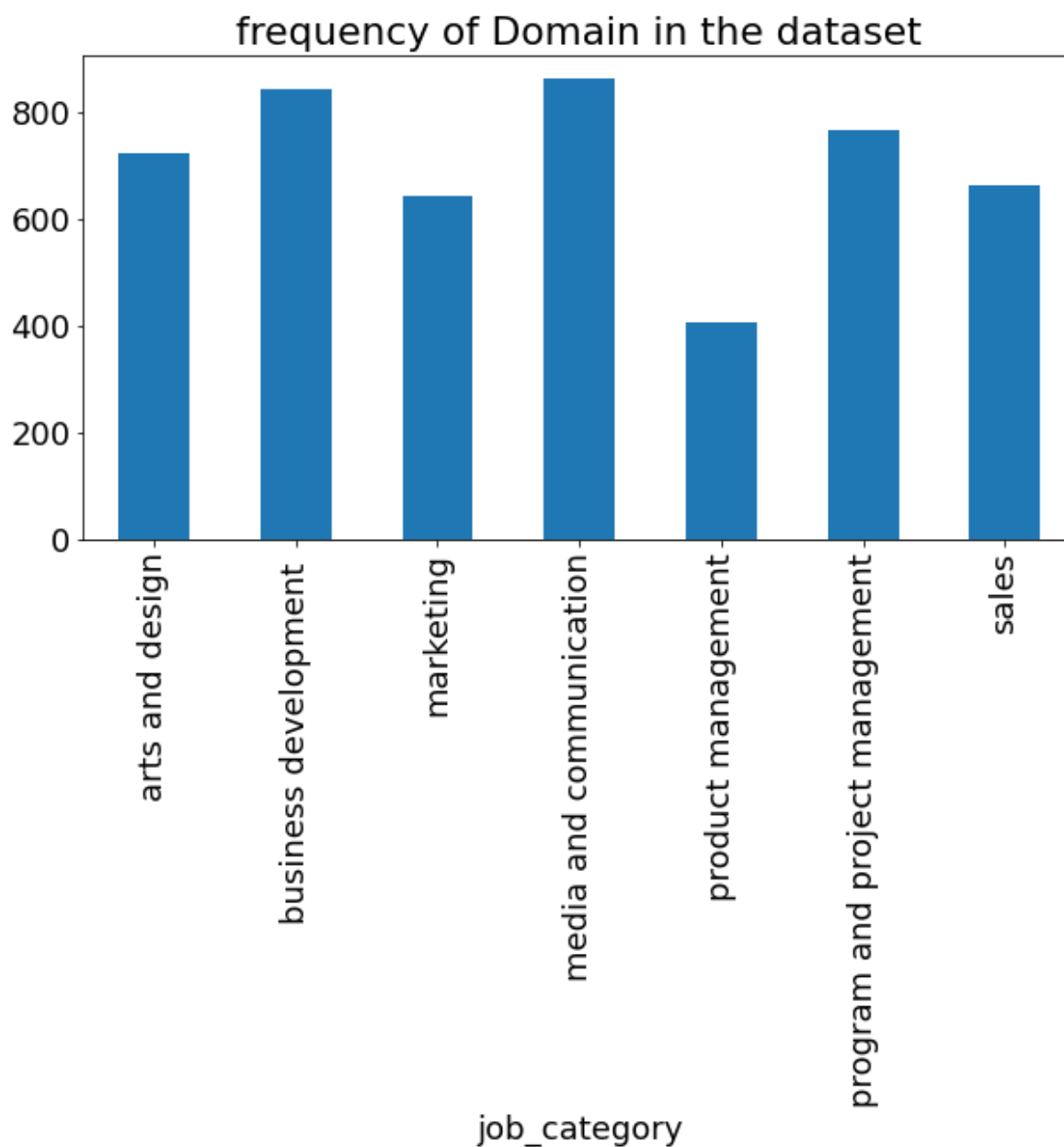
Unnamed: 0		Job Title	Company Name	Industry	Company Location	New Job (90 Days)	Year Started	Profile Headline
0	0	Battery Designer	Rivian	Mechanical or Industrial Engineering	Dublin, Ohio, United States	False	2020.0	Mechanical Design Engineer, System Integration...
1	1	Digital DevOps Engineer	HSBC	Information Technology and Services	New York City Metropolitan Area	False	2018.0	Digital DevOps Engineer at HSBC
2	2	Product Designer	Two Point Conversions, Inc.	Information Technology and Services	Chicago, Illinois, United States	False	2018.0	Leading Product + UX at Remedy (Two Point Conv... hti
3	3	Product Designer	udaan.com	Information Technology and Services	Bangalore Urban, Karnataka, India	False	2018.0	Product Designer at udaan
4	4	Digital Technology Intern	GE	Information Technology and Services	Jaipur, Rajasthan, India	True	2021.0	Digital Technology Intern at General Electric ...

5 rows × 26 columns



In [18]:

```
1 plt.rcParams.update({'font.size': 18})
2 fig = plt.figure(figsize=(10,5))
3 df1.groupby('job_category').preprocessed.count().plot.bar(ylim=0)
4 #plt.xticks(rotation=45)
5 plt.title('frequency of Domain in the dataset')
6 plt.show()
```



In [19]:



```
1 # making tfidf vector of the text in description column
2 tfidf = TfidfVectorizer(sublinear_tf=True, min_df=5, norm='l2', encoding='latin-1', ng
3 features = tfidf.fit_transform(df1.preprocessed).toarray()
4 labels = df1.category_id
5 features.shape
```

Out[19]:

(4910, 7245)

In [20]:



```

1 # finding out most common features of different Domain
2 from sklearn.feature_selection import chi2
3 import numpy as np
4 N = 2
5 for perprocessed, category_id in sorted(category_to_id.items()):
6     features_chi2 = chi2(features, labels == category_id)
7     indices = np.argsort(features_chi2[0])
8     feature_names = np.array(tfidf.get_feature_names())[indices]
9     unigrams = [v for v in feature_names if len(v.split(' ')) == 1]
10    bigrams = [v for v in feature_names if len(v.split(' ')) == 2]
11    print("# '{}':".format(perprocessed))
12    print("    . Most correlated unigrams:\n. {}".format('\n. '.join(unigrams[-N:])))
13    print("    . Most correlated bigrams:\n. {}".format('\n. '.join(bigrams[-N:])))
14    print(" ")

```

```

# 'arts and design':
. Most correlated unigrams:
. design
. designer
. Most correlated bigrams:
. design intern
. graphic designer

```

```

# 'business development ':
. Most correlated unigrams:
. deloitte
. analyst
. Most correlated bigrams:
. deloitte india
. analyst deloitte

```

```

# 'marketing':
. Most correlated unigrams:
. marketer
. marketing
. Most correlated bigrams:
. digital marketing
. marketing intern

```

```

# 'media and communication':
. Most correlated unigrams:
. content
. writer
. Most correlated bigrams:
. medical writer
. content writer

```

```

# 'product management':
. Most correlated unigrams:
. manager
. product
. Most correlated bigrams:
. product management
. product manager

```

```

# 'program and project management':
. Most correlated unigrams:
. manager

```

```
. project
. Most correlated bigrams:
. program manager
. project manager

# 'sales':
. Most correlated unigrams:
. account
. sales
. Most correlated bigrams:
. account manager
. sales manager
```

Linear SVC

In [21]:

```
1 # creating word vector and model for prediction
2 X_train, X_test, y_train, y_test = train_test_split(df1['preprocessed'], df1['job_category'],
3 count_vect = CountVectorizer()
4 X_train_counts = count_vect.fit_transform(X_train)
5 tfidf_transformer = TfidfTransformer()
6 X_train_tfidf = tfidf_transformer.fit_transform(X_train_counts)
7 clf = LinearSVC().fit(X_train_tfidf, y_train)
```

In [22]:

```
1 # checking model prediction in data
2 print(clf.predict(count_vect.transform(["product designer udaan"])))
```

['arts and design']

In [23]:

```
1 df1[df1['preprocessed'] == "product designer udaan"]
```

Out[23]:

	Unnamed: 0	Job Title	Company Name	Industry	Company Location	New Job (90 Days)	Year Started	Profile Headline	Profile Summary
3	3	Product Designer	udaan.com	Information Technology and Services	Bangalore Urban, Karnataka, India	False	2018.0	Product Designer at udaan	NaN
25	25	Product Designer	udaan.com	Information Technology and Services	Bangalore Urban, Karnataka, India	False	2018.0	Product Designer at udaan	NaN

2 rows × 26 columns

In [24]:



```
1 # training accuracy of the model
2 print(accuracy_score(y_train,clf.predict(X_train_tfidf)))
```

0.9098316132536665

In [25]:



```
1 # testing accuracy of the model
2 print(accuracy_score(y_test,clf.predict(count_vect.transform(X_test))))
```

0.5488599348534202

Multinomial Naive Bayes

In [26]:



```
1 clf = MultinomialNB().fit(X_train_tfidf, y_train)
```

In [27]:



```
1 # checking model prediction in data
2 print(clf.predict(count_vect.transform(["product designer udaan"])))
```

['arts and design']

In [28]:



```
1 # training accuracy of the model
2 print(accuracy_score(y_train,clf.predict(X_train_tfidf)))
```

0.7045084193373167

In [29]:



```
1 # testing accuracy of the model
2 print(accuracy_score(y_test,clf.predict(count_vect.transform(X_test))))
```

0.501628664495114

In [30]:



```

1 # testing some more models
2 CV = StratifiedKFold(n_splits=5)
3 models = [
4     RandomForestClassifier(n_estimators=200, max_depth=3, random_state=0),
5     LinearSVC(),
6     MultinomialNB(),
7 ]
8 #CV = 5
9 cv_df = pd.DataFrame(index=range(5 * len(models)))
10 entries = []
11 for model in models:
12     model_name = model.__class__.__name__
13     accuracies = cross_val_score(model, features, labels, scoring='accuracy', cv=CV,)
14     for fold_idx, accuracy in enumerate(accuracies):
15         entries.append((model_name, fold_idx, accuracy))
16 cv_df = pd.DataFrame(entries, columns=['model_name', 'fold_idx', 'accuracy'])
17 import seaborn as sns
18 sns.boxplot(x='model_name', y='accuracy', data=cv_df)
19 sns.stripplot(x='model_name', y='accuracy', data=cv_df,
20              size=8, jitter=True, edgecolor="gray", linewidth=2)
21 plt.xticks(rotation=90)
22 plt.show()

```

