

```
In [110]: import pandas
import numpy as np
```

## Task 1 Bag of Words and simple Features

### Linear model with non-review based features and L2 regularisation

```
In [111]: df = pd.read_csv('data/train.csv')

In [112]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 89970 entries, 0 to 89969
Data columns (total 14 columns):
country      89970 non-null object
description  89970 non-null object
designation   89970 non-null object
points       89970 non-null int64
price        89970 non-null float64
province     89970 non-null object
region_1     89970 non-null object
region_2     89970 non-null object
taster_name  71843 non-null object
taster_twitter_handle  68373 non-null object
title        89970 non-null object
variety      89970 non-null object
winery       89970 non-null object
dtypes: float64(1), int64(2), object(11)
memory usage: 9.64 MB

In [5]: df.head()

Out[5]:
```

	id	country	description	designation	points	price	province	region_1	region_2	taster_name	taster_twitter_handle	tit	
	0	60712	France	A buoyant nose, this has delicate acidity and...	La Riviera	89	17.0	Provence	Côtes de Provence	NaN	Roger Voss	Domai de Sangliu 2015 Rue	
	1	1171	Germany	Deeping in the wine, this seems...	Narstener Bergknecht	86	20.0	Rheinessen	NaN	Joe Czerwinski	@JoeCz	Lo. Guts 20X Newtun Berghof Kabi	
	2	129161	US	Sweet, simple and light-bodied wine...	Estate	84	23.0	California	Sierra Foothills	Sierra Foothills	Jim Gordon	@gordone_cellars	Naggi Esta Musz Cab (Bar
	3	116814	US	Tastes like a sweet anise sauce mixed into...	NaN	84	19.0	California	California	California Other	NaN	NaN	One Ho Sauvign Rose (Californ
	4	116519	Spain	As a one-liter party you get one-third more...	NaN	86	9.0	Central Spain	Vino de la Dena de Castilla	NaN	Michael Schachner	@wineschach	Bodega Equin 20X Tempranillo de Ti

```
In [113]: def dummy(df, col):
    drop = np.union1d(df[col][1:-1]
    dummy = pd.get_dummies(df[col])
    new_col = []
    for i in dummy.columns:
        new_col.append(col+'_'+str(i))
    dummy.columns = new_col
    df = pd.concat([df, dummy], axis=1)
    df.drop([col, new_col[1:-1]], inplace=True, axis=1)
    return df

In [114]: df.columns

Out[114]: Index(['id', 'country', 'description', 'designation', 'points', 'price',
        'province', 'region_1', 'region_2', 'taster_name',
        'taster_twitter_handle', 'title', 'variety', 'winery'],
        dtype='object')

In [115]: df['country'].fillna(df['country'].mode()[0], inplace=True)
df['designation'].fillna(df['designation'].mode()[0], inplace=True)
df['province'].fillna(df['province'].mode()[0], inplace=True)
df['region_1'].fillna(df['region_1'].mode()[0], inplace=True)
df['region_2'].fillna(df['region_2'].mode()[0], inplace=True)
df['taster_name'].fillna(df['taster_name'].mode()[0], inplace=True)
df['taster_twitter_handle'].fillna(df['taster_twitter_handle'].mode()[0], inplace=True)

In [116]: df['price'].fillna(df['price'].median(), inplace=True)

In [117]: cols = ['country', 'taster_name', 'variety']

In [118]: empty = []
    in cols:
        try:
            df = dummy(df, col)
        except:
            empty.append(col)

In [119]: empty

Out[119]: []

In [120]: y = df['points'].values

In [121]: df.drop(['id', 'country', 'description', 'designation',
        'points',
        'price',
        'region_1',
        'region_2',
        'taster_twitter_handle',
        'title',
        'winery'], axis=1, inplace=True)

In [184]: from sklearn.preprocessing importMinMaxScaler
scaler =MinMaxScaler()
df_scaled = scaler.fit_transform(df)

In [185]: df_sub = df.iloc[:, df_scaled.var() > 0.02].values

In [:]: (df_scaled.var() > 0.02).values

In [186]: df_sub['price'] = df['price']

In [187]: df_sub.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 89970 entries, 0 to 89969
Data columns (total 31 columns):
country      89970 non-null int64
country_Austria      89970 non-null int64
country_Chile      89970 non-null int64
country_France      89970 non-null int64
country_Italy      89970 non-null int64
country_Portugal      89970 non-null int64
country_Spain      89970 non-null int64
country_US      89970 non-null int64
taster_name      89970 non-null int64
taster_name_Anna Lee C. Iijima      89970 non-null int64
taster_name_Anne Krebshel MW      89970 non-null int64
taster_name_Jim Gordon      89970 non-null int64
taster_name_Joe Czerwinski      89970 non-null int64
taster_name_Kerin O'Keefe      89970 non-null int64
taster_name_Matt Kettmann      89970 non-null int64
taster_name_Michael Schachner      89970 non-null int64
taster_name_Paul Gregutt      89970 non-null int64
taster_name_Roger Voss      89970 non-null int64
taster_name_Sean P. Sullivan      89970 non-null int64
variety_Bordeaux-style Red Blend      89970 non-null int64
variety_Cabernet Sauvignon      89970 non-null int64
variety_Chardonnay      89970 non-null int64
variety_Merlot      89970 non-null int64
variety_Nebbiolo      89970 non-null int64
variety_Pinet Noir      89970 non-null int64
variety_Red Blend      89970 non-null int64
variety_Riesling      89970 non-null int64
variety_Rose      89970 non-null int64
variety_Sangiovese      89970 non-null int64
variety_Sauvignon Blanc      89970 non-null int64
variety_Syrah      89970 non-null int64
price      89970 non-null float64
dtypes: float64(1), int64(30)
memory usage: 3.3 MB

In [126]: X = df_sub.values

In [127]: from sklearn.linear_model import Ridge
clf = Ridge(alpha = 0.01, normalize = False)
clf.fit(X, y)

Out[127]: Ridge(alpha=0.01, copy_X=True, fit_intercept=True, max_iter=None,
        normalize=False, random_state=None, solver='auto', tol=0.001)

In [128]: clf.score(X, y)

Out[128]: 0.237723646398423183

In [129]: alpha = 10**np.linspace(10,-2,100)*0.5

Out[129]: array([5.00000000e+09, 1.28353816e+09, 2.86118386e+08, 1.16438064e+08,
        6.37274588e+07, 3.29308711e+07, 1.68737081e+07, 8.70873708e+06,
        4.56133611e+06, 2.38593664e+06, 1.23079442e+06, 6.36133611e+05,
        3.29308711e+05, 1.68737081e+05, 8.70873708e+04, 4.56133611e+04,
        2.38593664e+04, 1.23079442e+04, 6.36133611e+03, 3.29308711e+03,
        1.68737081e+03, 8.70873708e+02, 4.56133611e+02, 2.38593664e+02,
        1.23079442e+02, 6.36133611e+01, 3.29308711e+01, 1.68737081e+01,
        8.70873708e+00, 4.56133611e+00, 2.38593664e+00, 1.23079442e+00,
        6.36133611e-01, 3.29308711e-01, 1.68737081e-01, 8.70873708e-02,
        4.56133611e-02, 2.38593664e-02, 1.23079442e-02, 6.36133611e-03,
        3.29308711e-03, 1.68737081e-03, 8.70873708e-04, 4.56133611e-04,
        2.38593664e-04, 1.23079442e-04, 6.36133611e-05, 3.29308711e-05,
        1.68737081e-05, 8.70873708e-06, 4.56133611e-06, 2.38593664e-06,
        1.23079442e-06, 6.36133611e-07, 3.29308711e-07, 1.68737081e-07,
        8.70873708e-08, 4.56133611e-08, 2.38593664e-08, 1.23079442e-08,
        6.36133611e-09, 3.29308711e-09, 1.68737081e-09, 8.70873708e-10,
        4.56133611e-10, 2.38593664e-10, 1.23079442e-10, 6.36133611e-11,
        3.29308711e-11, 1.68737081e-11, 8.70873708e-12, 4.56133611e-12,
        2.38593664e-12, 1.23079442e-12, 6.36133611e-13, 3.29308711e-13,
        1.68737081e-13, 8.70873708e-14, 4.56133611e-14, 2.38593664e-14,
        1.23079442e-14, 6.36133611e-15, 3.29308711e-15, 1.68737081e-15,
        8.70873708e-16, 4.56133611e-16, 2.38593664e-16, 1.23079442e-16,
        6.36133611e-17, 3.29308711e-17, 1.68737081e-17, 8.70873708e-18,
        4.56133611e-18, 2.38593664e-18, 1.23079442e-18, 6.36133611e-19,
        3.29308711e-19, 1.68737081e-19, 8.70873708e-20, 4.56133611e-20,
        2.38593664e-20, 1.23079442e-20, 6.36133611e-21, 3.29308711e-21,
        1.68737081e-21, 8.70873708e-22, 4.56133611e-22, 2.38593664e-22,
        1.23079442e-22, 6.36133611e-23, 3.29308711e-23, 1.68737081e-23,
        8.70873708e-24, 4.56133611e-24, 2.38593664e-24, 1.23079442e-24,
        6.36133611e-25, 3.29308711e-25, 1.68737081e-25, 8.70873708e-26,
        4.56133611e-26, 2.38593664e-26, 1.23079442e-26, 6.36133611e-27,
        3.29308711e-27, 1.68737081e-27, 8.70873708e-28, 4.56133611e-28,
        2.38593664e-28, 1.23079442e-28, 6.36133611e-29, 3.29308711e-29,
        1.68737081e-29, 8.70873708e-30, 4.56133611e-30, 2.38593664e-30,
        1.23079442e-30, 6.36133611e-31, 3.29308711e-31, 1.68737081e-31,
        8.70873708e-32, 4.56133611e-32, 2.38593664e-32, 1.23079442e-32,
        6.36133611e-33, 3.29308711e-33, 1.68737081e-33, 8.70873708e-34,
        4.56133611e-34, 2.38593664e-34, 1.23079442e-34, 6.36133611e-35,
        3.29308711e-35, 1.68737081e-35, 8.70873708e-36, 4.56133611e-36,
        2.38593664e-36, 1.23079442e-36, 6.36133611e-37, 3.29308711e-37,
        1.68737081e-37, 8.70873708e-38, 4.56133611e-38, 2.38593664e-38,
        1.23079442e-38, 6.36133611e-39, 3.29308711e-39, 1.68737081e-39,
        8.70873708e-40, 4.56133611e-40, 2.38593664e-40, 1.23079442e-40,
        6.36133611e-41, 3.29308711e-41, 1.68737081e-41, 8.70873708e-42,
        4.56133611e-42, 2.38593664e-42, 1.23079442e-42, 6.36133611e-43,
        3.29308711e-43, 1.68737081e-43, 8.70873708e-44, 4.56133611e-44,
        2.38593664e-44, 1.23079442e-44, 6.36133611e-45, 3.29308711e-45,
        1.68737081e-45, 8.70873708e-46, 4.56133611e-46, 2.38593664e-46,
        1.23079442e-46, 6.36133611e-47, 3.29308711e-47, 1.68737081e-47,
        8.70873708e-48, 4.56133611e-48, 2.38593664e-48, 1.23079442e-48,
        6.36133611e-49, 3.29308711e-49, 1.68737081e-49, 8.70873708e-50,
        4.56133611e-50, 2.38593664e-50, 1.23079442e-50, 6.36133611e-51,
        3.29308711e-51, 1.68737081e-51, 8.70873708e-52, 4.56133611e-52,
        2.38593664e-52, 1.23079442e-52, 6.36133611e-53, 3.29308711e-53,
        1.68737081e-53, 8.70873708e-54, 4.56133611e-54, 2.38593664e-54,
        1.23079442e-54, 6.36133611e-55, 3.29308711e-55, 1.68737081e-55,
        8.70873708e-56, 4.56133611e-56, 2.38593664e-56, 1.23079442e-56,
        6.36133611e-57, 3.29308711e-57, 1.68737081e-57, 8.70873708e-58,
        4.56133611e-58, 2.38593664e-58, 1.23079442e-58, 6.36133611e-59,
        3.29308711e-59, 1.68737081e-59, 8.70873708e-60, 4.56133611e-60,
        2.38593664e-60, 1.23079442e-60, 6.36133611e-61, 3.29308711e-61,
        1.68737081e-61, 8.70873708e-62, 4.56133611e-62, 2.38593664e-62,
        1.23079442e-62, 6.36133611e-63, 3.29308711e-63, 1.68737081e-63,
        8.70873708e-64, 4.56133611e-64, 2.38593664e-64, 1.23079442e-64,
        6.36133611e-65, 3.29308711e-65, 1.68737081e-65, 8.70873708e-66,
        4.56133611e-66, 2.38593664e-66, 1.23079442e-66, 6.36133611e-67,
        3.29308711e-67, 1.68737081e-67, 8.70873708e-68, 4.56133611e-68,
        2.38593664e-68, 1.23079442e-68, 6.36133611e-69, 3.29308711e-69,
        1.68737081e-69, 8.70873708e-70, 4.56133611e-70, 2.38593664e-70,
        1.23079442e-70, 6.36133611e-71, 3.29308711e-71, 1.68737081e-71,
        8.70873708e-72, 4.56133611e-72, 2.38593664e-72, 1.23079442e-72,
        6.36133611e-73, 3.29308711e-73, 1.68737081e-73, 8.70873708e-74,
        4.56133611e-74, 2.38593664e-74, 1.23079442e-74, 6.36133611e-75,
        3.29308711e-75, 1.68737081e-75, 8.70873708e-76, 4.56133611e-76,
        2.38593664e-76, 1.23079442e-76, 6.36133611e-77, 3.29308711e-77,
        1.68737081e-77, 8.70873708e-78, 4.56133611e-78, 2.38593664e-78,
        1.23079442e-78, 6.36133611e-79, 3.29308711e-79, 1.68737081e-79,
        8.70873708e-80, 4.56133611e-80, 2.38593664e-80, 1.23079442e-80,
        6.36133611e-81, 3.29308711e-81, 1.68737081e-81, 8.70873708e-82,
        4.56133611e-82, 2.38593664e-82, 1.23079442e-82, 6.36133611e-83,
        3.29308711e-83, 1.68737081e-83, 8.70873708e-84, 4.56133611e-84,
        2.38593664e-84, 1.23079442e-84, 6.36133611e-85, 3.29308711e-85,
        1.68737081e-85, 8.70873708e-86, 4.56133611e-86, 2.38593664e-86,
        1.23079442e-86, 6.36133611e-87, 3.29308711e-87, 1.68737081e-87,
        8.70873708e-88, 4.56133611e-88, 2.38593664e-88, 1.23079442e-88,
        6.36133611e-89, 3.29308711e-89, 1.68737081e-89, 8.70873708e-90,
        4.56133611e-90, 2.38593664e-90, 1.23079442e-90, 6.36133611e-91,
        3.29308711e-91, 1.68737081e-91, 8.70873708e-92, 4.56133611e-92,
        2.38593664e-92, 1.23079442e-92, 6.36133611e-93, 3.29308711e-93,
        1.68737081e-93, 8.70873708e-94, 4.56133611e-94, 2.38593664e-94,
        1.23079442e-94, 6.36133611e-95, 3.29308711e-95, 1.68737081e-95,
        8.70873708e-96, 4.56133611e-96, 2.38593664e-96, 1.23079442e-96,
        6.36133611e-97, 3.29308711e-97, 1.68737081e-97, 8.70873708e-98,
        4.56133611e-98, 2.38593664e-98, 1.23079442e-98, 6.36133611e-99,
        3.29308711e-99, 1.68737081e-99, 8.70873708e-100, 4.56133611e-100,
        2.38593664e-100, 1.23079442e-100, 6.36133611e-101, 3.29308711e-101,
        1.68737081e-101, 8.70873708e-102, 4.56133611e-102, 2.38593664e-102,
        1.23079442e-102, 6.36133611e-103, 3.29308711e-103, 1.68737081e-103,
        8.70873708e-104, 4.56133611e-104, 2.38593664e-104, 1.23079442e-104,
        6.36133611e-105, 3.29308711e-105, 1.68737081e-105, 8.70873708e-106,
        4.56133611e-106, 2.38593664e-106, 1.23079442e-106, 6.36133611e-107,
        3.29308711e-107, 1.68737081e-107, 8.70873708e-108, 4.56133611e-108,
        2.38593664e-108, 1.23079442e-108, 6.36133611e-109, 3.29308711e-109,
        1.68737081e-109, 8.70873708e-110, 4.56133611e-110, 2.38593664e-110,
        1.23079442e-110, 6.36133611e-111, 3.29308711e-111, 1.68737081e-111,
        8.70873708e-112, 4.56133611e-112, 2.38593664e-112, 1.23079442e-112,
        6.36133611e-113, 3.29308711e-113, 1.68737081e-113, 8.70873708e-114,
        4.56133611e-114, 2.38593664e-114, 1.23079442e-114, 6.36133611e-115,
        3.29308711e-115, 1.68737081e-115, 8.70873708e-116, 4.56133611e-116,
        2.38593664e-116, 1.23079442e-116, 6.36133611e-117, 3.29308711e-117,
        1.68737081e-117, 8.70873708e-118, 4.56133611e-118, 2.38593664e-118,
        1.23079442e-118, 6.36133611e-119, 3.29308711e-119, 1.68737081e-119,
        8.70873708e-120, 4.56133611e-120, 2.38593664e-120, 1.23079442e-120,
        6.36133611e-121, 3.29308711e-121, 1.68737081e-121, 8.70873708e-122,
        4.56133611e-122, 2.38593664e-122, 1.23079442e-122, 6.36133611e-123,
        3.29308711e-123, 1.68737081e-123, 8.70873708e-124, 4.56133611e-124,
        2.38593664e-124, 1.23079442e-124, 6.36133611e-125, 3.29308711e-125,
        1.68737081e-125, 8.70873708e-126, 4.56133611e-126, 2.38593664e-126,
        1.23079442e-126, 6.36133611e-127, 3.29308711e-127, 1.68737081e-127,
        8.70873708e-128, 4.56133611e-128, 2.38593664e-128, 1.23079442e-128,
        6.36133611e-129, 3.29308711e-129, 1.68737081e-129, 8.70873708e-130,
        4.56133611e-130, 2.38593664e-130, 1.23079442e-130, 6.36133611e-131,
        3.29308711e-131, 1.68737081e-131, 8.70873708e-132, 4.56133611e-132,
        2.38593664e-132, 1.23079442e-132, 6.36133611e-133, 3.29308711e-133,
        1.68737081e-133, 8.70873708e-134, 4.56133611e-134, 2.38593664e-134,
        1.23079442e-134, 6.36133611e-135, 3.29308711e-135, 1.68737081e-135,
        8.70873708e-136, 4.56133611e-136, 2.38593664e-136, 1.23079442e-136,
        6.36133611e-137, 3.29308711e-137, 1.68737081e-137, 8.70873708e-138,
        4.56133611e-138, 2.38593664e-138, 1.23079442e-138, 6.36133611e-139,
        3.29308711e-139, 1.68737081e-139, 8.70873708e-140, 4.56133611e-140,
        2.38593664e-140, 1.23079442e-140, 6.36133611e-141, 3.29308711e-141,
        1.68737081e-141, 8.70873708e-142, 4.56133611e-142, 2.38593664e-142,
        1.23079442e-142, 6.36133611e-143, 3.29308711e-143, 1.68737081e-143,
        8.70873708e-144, 4.56133611e-144, 2.38593664e-144, 1.23079442e-144,
        6.36133611e-145, 3.29308711e-145, 1.68737081e-145, 8.70873708e-146,
        4.56133611e-146, 2.38593664e-146, 1.23079442e-146, 6.36133611e-147,
        3.29308711e-147, 1.68737081e-147, 8.70873708e-148, 4.56133611e-148,
        2.38593664e-148, 1.23079442e-148, 6.36133611e-149, 3.29308711e-149,
        1.68737081e-149, 8.70873708e-150, 4.56133611e-150, 2.38593664e-150,
        1.23079442e-150, 6.36133611e-151, 3.29308711e-151, 1.68737081e-151,
        8.70873708e-152, 4.56133611e-152, 2.38593664e-152, 1.23079442e-152,
        6.36133611e-153, 3.29308711e-153, 1.68737081e-153, 8.70873708e-154,
        4.56133611e-154, 2.38593664e-154, 1.23079442e-154, 6.36133611e-155,
        3.29308711e-155, 1.68737081e-155, 8.70873708e-156, 4.56133611e-156,
        2.38593664e-156, 1.23079442e-156, 6.36133611e-157, 3.29308711e-157,
        1.68737081e-1
```



```
In [218]: import re # for regular expressions
import pandas as pd
pd.set_option("display.max_colwidth", 200)
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import string
import nltk # for text manipulation
import warnings
warnings.filterwarnings("ignore", category=DeprecationWarning)

import matplotlib inline

def remove_pattern(input_txt, pattern):
    r = re.findall(pattern, input_txt)
    for i in r:
        input_txt = re.sub(i, '', input_txt)
    return input_txt

public['description'] = np.vectorize(remove_pattern)(public['description'], "@{w}**")

public['description'] = public['description'].str.replace("@(a-zA-Z)", " ")

public['description'] = public['description'].str.replace('#', '')

public['description'] = public['description'].apply(lambda x: ' '.join([w for w in x.split() if len(w) > 2]))

tokenized_tweet = public['description'].apply(lambda x: x.split()) # tokenizing
tokenized_tweet.head()

from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
tokenized_tweet.apply(lambda x: [lemmatizer.lemmatize(i) for i in x]) # stemming
#lemmatized_output = ' '.join([lemmatizer.lemmatize(token) for token in tokens])

for i in range(len(tokenized_tweet)):
    tokenized_tweet[i] = ' '.join(tokenized_tweet[i])

public['description'] = tokenized_tweet

all_words = ' '.join([text for text in public['description']])

bow_public = bow_vectorizer.transform(public['description'])
bow_public.shape

bow_public = np.array(bow_public.todense())
y = public['points']

mean_squared_error(y, ridge.predict(bow_public))

ridge.score(bow_public, y)

Out[218]: 0.41728530345424897
```

## Checking accuracy of linear model using review based features, Count vectoriser and L2 regularisation on private dataset

```
In [219]: def remove_pattern(input_txt, pattern):
    r = re.findall(pattern, input_txt)
    for i in r:
        input_txt = re.sub(i, '', input_txt)
    return input_txt

private['description'] = np.vectorize(remove_pattern)(private['description'], "@{w}**")

private['description'] = private['description'].str.replace("@(a-zA-Z)", " ")

private['description'] = private['description'].str.replace('#', '')

private['description'] = private['description'].apply(lambda x: ' '.join([w for w in x.split() if len(w) > 2]))

tokenized_tweet = private['description'].apply(lambda x: x.split()) # tokenizing
tokenized_tweet.head()

from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
tokenized_tweet.apply(lambda x: [lemmatizer.lemmatize(i) for i in x]) # stemming
#lemmatized_output = ' '.join([lemmatizer.lemmatize(token) for token in tokens])

for i in range(len(tokenized_tweet)):
    tokenized_tweet[i] = ' '.join(tokenized_tweet[i])

private['description'] = tokenized_tweet

all_words = ' '.join([text for text in private['description']])

bow_private = bow_vectorizer.transform(private['description'])

bow_private = np.array(bow_private.todense())
y = private['points']

mean_squared_error(y, ridge.predict(bow_private))

ridge.score(bow_private, y)

Out[219]: 0.42575339762328936
```

## Combine non-review and review based count vectoriser features

```
In [54]: X.shape, bow.shape, type(bow), type(X)

Out[54]: ((89970, 31), (89970, 200), numpy.ndarray, numpy.ndarray)

In [55]: X_bow = np.hstack((X, bow))

In [56]: X_train, X_test, y_train, y_test = train_test_split(X_bow, y, test_size=0.2, random_state=1)

In [57]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[57]: 0.006609705742330144

In [58]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

Out[58]: 4.3330388080809418

In [60]: ridge.score(X_test, y_test)

Out[60]: 0.5268639740985578
```

## Trying review based tf-idf vectorizer

```
In [220]: tfidf_vectorizer = TfidfVectorizer(max_df=0.90, min_df=2, max_features=200, stop_words='english')
tfidf = tfidf_vectorizer.fit_transform(combi['description'])
tfidf.shape

Out[220]: (89970, 200)

In [223]: X_train, X_test, y_train, y_test = train_test_split(tfidf, combi['points'], test_size=0.2, random_state=1)

In [224]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[224]: 0.005

In [225]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

Out[225]: 5.138117016528831

In [226]: ridge.score(X_test, y_test)

Out[226]: 0.438955353808066313
```

## Checking accuracy of tfidf vectoriser based model on public dataset

```
In [227]: tfidf_public = tfidf_vectorizer.transform(public['description'])

mean_squared_error(public['points'], ridge.predict(tfidf_public))

ridge.score(tfidf_public, public['points'])

Out[227]: 0.42952393763233454
```

## Checking accuracy of tfidf vectoriser based model on private dataset

```
In [228]: tfidf_private = tfidf_vectorizer.transform(private['description'])

mean_squared_error(private['points'], ridge.predict(tfidf_private))

ridge.score(tfidf_private, private['points'])

Out[228]: 0.43694473915367305
```

## Combining non review and review based tfidf features

```
In [68]: tfidf = np.array(tfidf.todense())
X.shape, tfidf.shape, type(tfidf), type(X)

Out[68]: ((89970, 31), (89970, 200), numpy.ndarray, numpy.ndarray)

In [69]: X_tfidf = np.hstack((X, tfidf))

X_train, X_test, y_train, y_test = train_test_split(X_tfidf, y, test_size=0.2, random_state=1)

ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

ridge.score(X_test, y_test)

Out[69]: 0.5395942473832831
```

## Task 3: Custom Word Vectors

### Custom word to vec embeddings using skip-gram model

```
In [70]: from sklearn.feature_extraction.text import TfidfVectorizer, CountVectorizer
import gensim

In [229]: tokenized_tweet = combi['description'].apply(lambda x: x.split()) # tokenizing

model_w2v = gensim.models.Word2Vec(
    tokenized_tweet,
    size=50, # desired no. of features/independent variables
    window=5, # context window size
    min_count=2,
    sg = 1, # 1 for skip-gram model
    hs = 0,
    negative = 10, # for negative sampling
    workers= 2, # no. of cores
    seed = 30)

In [230]: model_w2v.train(tokenized_tweet, total_examples= len(combi['description']), epochs=20)

Out[230]: (49396884, 61990780)

model_w2v.shape

In [231]: def word_vector(tokens, size):
    vec = np.zeros(size).reshape((1, size))
    count = 0
    for word in tokens:
        try:
            vec += model_w2v[word].reshape((1, size))
            count += 1
        except KeyError: # handling the case where the token is not in vocabulary
            continue
    if count != 0:
        vec /= count
    return vec

In [232]: wordvec_arrays = np.zeros((len(tokenized_tweet), 50))

for i in range(len(tokenized_tweet)):
    wordvec_arrays[i,:] = word_vector(tokenized_tweet[i], 50)

wordvec_df = pd.DataFrame(wordvec_arrays)
wordvec_df.shape

Out[232]: (89970, 50)

In [236]: X_train, X_test, y_train, y_test = train_test_split(wordvec_df, combi['points'], test_size=0.2, random_state=1)

In [237]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[237]: 0.005

In [238]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

Out[238]: 4.208923794471053

In [239]: ridge.score(X_test, y_test)

Out[239]: 0.5404164300815004
```

## Checking accuracy of custom word to vec embeddings using skip-gram model on private dataset

```
In [241]: tokenized_tweet_private = private['description'].apply(lambda x: x.split()) # tokenizing

wordvec_arrays_private = np.zeros((len(tokenized_tweet_private), 50))

for i in range(len(tokenized_tweet_private)):
    wordvec_arrays_private[i,:] = word_vector(tokenized_tweet_private[i], 50)

wordvec_df_private = pd.DataFrame(wordvec_arrays_private)
wordvec_df_private.shape

mean_squared_error(private['points'], ridge.predict(wordvec_df_private))

ridge.score(wordvec_df_private, private['points'])

Out[241]: 0.5018984673156571
```

## Checking accuracy of custom word to vec embeddings using skip-gram model on public dataset

```
In [242]: tokenized_tweet_public = public['description'].apply(lambda x: x.split()) # tokenizing

wordvec_arrays_public = np.zeros((len(tokenized_tweet_public), 50))

for i in range(len(tokenized_tweet_public)):
    wordvec_arrays_public[i,:] = word_vector(tokenized_tweet_public[i], 50)

wordvec_df_public = pd.DataFrame(wordvec_arrays_public)
wordvec_df_public.shape

mean_squared_error(public['points'], ridge.predict(wordvec_df_public))

ridge.score(wordvec_df_public, public['points'])

Out[242]: 0.4933246314798058
```

## Combine custom skip gram word2vec and bow features

```
In [82]: X_bow_sg_w2v = np.hstack((wordvec_df, bow))

In [83]: X_bow_sg_w2v

Out[83]: array([[[-0.00262551, -0.23499813, 0.2317796, ..., 0.
, 0.
, 0.],
[-0.27572221, -0.10347272, 0.18937457, ..., 0.
, 0.
, 0.],
[-0.03521368, -0.15286143, 0.28467456, ..., 0.
, 0.
, 0.],
...,
[-0.17225228, -0.19234901, 0.18830725, ..., 0.
, 0.
, 0.],
[-0.16598068, -0.23316762, 0.32500752, ..., 0.
, 0.
, 0.],
[-0.32364449, -0.09797542, 0.10864892, ..., 0.
, 0.
, 0.]]])

In [84]: X_train, X_test, y_train, y_test = train_test_split(X_bow_sg_w2v, y, test_size=0.2, random_state=1)

In [85]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[85]: 0.005

In [86]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

Out[86]: 3.545981161827917

In [87]: ridge.score(X_test, y_test)

Out[87]: 0.6128048972144841
```

## Custom word to vec embeddings using cbow model

```
In [243]: tokenized_tweet = combi['description'].apply(lambda x: x.split()) # tokenizing

model_w2v = gensim.models.Word2Vec(
    tokenized_tweet,
    size=50, # desired no. of features/independent variables
    window=5, # context window size
    min_count=2,
    sg = 0, # 1 for skip-gram model
    hs = 0,
    negative = 10, # for negative sampling
    workers= 2, # no. of cores
    seed = 34)

In [244]: model_w2v.train(tokenized_tweet, total_examples= len(combi['description']), epochs=20)

Out[244]: (49396884, 61990780)

In [245]: wordvec_arrays = np.zeros((len(tokenized_tweet), 50))

for i in range(len(tokenized_tweet)):
    wordvec_arrays[i,:] = word_vector(tokenized_tweet[i], 50)

wordvec_df = pd.DataFrame(wordvec_arrays)
wordvec_df.shape

Out[245]: (89970, 50)

In [247]: X_train, X_test, y_train, y_test = train_test_split(wordvec_df, combi['points'], test_size=0.2, random_state=1)

In [248]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[248]: 0.005

In [249]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

Out[249]: 4.650935144232271

In [250]: ridge.score(X_test, y_test)

Out[250]: 0.49215203661954066
```

## Custom word to vec embeddings using cbow model on private dataset

```
In [251]: tokenized_tweet_private = private['description'].apply(lambda x: x.split()) # tokenizing

wordvec_arrays_private = np.zeros((len(tokenized_tweet_private), 50))

for i in range(len(tokenized_tweet_private)):
    wordvec_arrays_private[i,:] = word_vector(tokenized_tweet_private[i], 50)

wordvec_df_private = pd.DataFrame(wordvec_arrays_private)
wordvec_df_private.shape

mean_squared_error(private['points'], ridge.predict(wordvec_df_private))

ridge.score(wordvec_df_private, private['points'])

Out[251]: 0.47207511102727523
```

## Custom word to vec embeddings using cbow model on public dataset

```
In [252]: tokenized_tweet_public = public['description'].apply(lambda x: x.split()) # tokenizing

wordvec_arrays_public = np.zeros((len(tokenized_tweet_public), 50))

for i in range(len(tokenized_tweet_public)):
    wordvec_arrays_public[i,:] = word_vector(tokenized_tweet_public[i], 50)

wordvec_df_public = pd.DataFrame(wordvec_arrays_public)
wordvec_df_public.shape

mean_squared_error(public['points'], ridge.predict(wordvec_df_public))

ridge.score(wordvec_df_public, public['points'])

Out[252]: 0.4628463252114597
```

## Task 2 Pre-trained Word Vectors

### Using pre-trained word2vec embeddings created by Google

```
In [253]: model = gensim.models.KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)

In [254]: def word_vector_pretrained(tokens, size):
    vec = np.zeros(size).reshape((1, size))
    count = 0
    for word in tokens:
        try:
            vec += model[word].reshape((1, size))
            count += 1
        except KeyError: # handling the case where the token is not in vocabulary
            continue
    if count != 0:
        vec /= count
    return vec

In [255]: wordvec_arrays = np.zeros((len(tokenized_tweet), 300))

for i in range(len(tokenized_tweet)):
    wordvec_arrays[i,:] = word_vector_pretrained(tokenized_tweet[i], 300)

wordvec_df = pd.DataFrame(wordvec_arrays)
wordvec_df.shape

Out[255]: (89970, 300)

In [257]: X_train, X_test, y_train, y_test = train_test_split(wordvec_df, combi['points'], test_size=0.2, random_state=1)

In [258]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[258]: 0.008737642000038414

In [259]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

Out[259]: 5.239015673425837

In [261]: ridge.score(X_test, y_test)

Out[261]: 0.4279379614297319
```

### Using pre-trained word2vec embeddings on private dataset

```
In [262]: tokenized_tweet_private = private['description'].apply(lambda x: x.split()) # tokenizing

wordvec_arrays_private = np.zeros((len(tokenized_tweet_private), 300))

for i in range(len(tokenized_tweet_private)):
    wordvec_arrays_private[i,:] = word_vector_pretrained(tokenized_tweet_private[i], 300)

wordvec_df_private = pd.DataFrame(wordvec_arrays_private)
wordvec_df_private.shape

mean_squared_error(private['points'], ridge.predict(wordvec_df_private))

ridge.score(wordvec_df_private, private['points'])

Out[262]: 0.4846125449904594
```

### Using pre-trained word2vec embeddings on public dataset

```
In [264]: tokenized_tweet_public = public['description'].apply(lambda x: x.split()) # tokenizing

wordvec_arrays_public = np.zeros((len(tokenized_tweet_public), 300))

for i in range(len(tokenized_tweet_public)):
    wordvec_arrays_public[i,:] = word_vector_pretrained(tokenized_tweet_public[i], 300)

wordvec_df_public = pd.DataFrame(wordvec_arrays_public)
wordvec_df_public.shape

mean_squared_error(public['points'], ridge.predict(wordvec_df_public))

ridge.score(wordvec_df_public, public['points'])

Out[264]: 0.476505668795022
```

## Combine pretrained word embedding with bag of words features

```
In [103]: wordvec_bow = np.hstack((wordvec_df, bow))

In [107]: X_train, X_test, y_train, y_test = train_test_split(wordvec_bow, y, test_size=0.2, random_state=1)

In [108]: ridgecv = RidgeCV(alphas = alphas, scoring = 'neg_mean_squared_error', normalize = True)
ridgecv.fit(X_train, y_train)
ridgecv.alpha_

Out[108]: 0.01155064850041579

In [109]: ridge = Ridge(alpha = ridgecv.alpha_, normalize = True)
ridge.fit(X_train, y_train)
mean_squared_error(y_test, ridge.predict(X_test))

ridge.score(X_test, y_test)

Out[109]: 0.559809729277211

In [ ]:
```