

## **Lab 3**

Abhishek Sridhar

Connor Smith

Tom Li

### Q1: Summary

Information can be compressed - represented with less than the original amount of bits - without loss of information. There is a mathematical limit on the amount of data compression that can be done without losing information. It is  $H(x) = - \sum(p(x) * \log(p(x)))$  where  $p(x)$  is the probability of the random variable. The author mentions several approximation methods used in sentence information compression. The N-gram method is called. Words are grouped into N word groups and for each group a probability is produced to estimate the probability of the next word being a specific word given the previous N-1 words. For example,  $P(w_1, w_2, \dots, w_{n-1}, w_n)$  will give the probability that the next word is  $w_n$  given the previous n-1 words are  $w_1, w_2, \dots, w_{n-1}$ . This has a connection with a stochastic process where previous states will influence subsequent states. The stochastic process can be represented in a graphical format where each state is a dot and there are arrows connecting the states with probabilities of the next states. There are at most  $n^2$  in a trigram Markov stochastic diagram. Besides stochastic processes, there is also ergodic process. Ergodic process produces sequences that have the same statistical properties: letter frequencies, diagram frequencies. Entropy is given by the formula above and there are some interesting observations from that formula. When  $p=0$  or  $1$ ,  $H=0$ . Meaning when the value is completely certain the entropy of the system is 0. Conversely, entropy is maximized to 1 when  $p=1/2$  where there is no predictive advantage using probability. Relative entropy is the maximum value of a source while still staying with the same alphabet code. Relative entropy has one drawback - it has redundancy when compared to normal entropy because of conforming to the existing alphabet code which might not be optimal.

```

1 import math
2 import os
3 from collections import Counter
4 from random import choices
5 from urllib.parse import urljoin
6
7 import bs4
8 import requests
9 import texttract
10
11 # part 1
12 url = "http://proceedings.mlr.press/v70/"
13
14 folder_location = r'p2files'
15 if not os.path.exists(folder_location): os.mkdir(
    folder_location)
16
17 response = requests.get(url)
18 soup = bs4.BeautifulSoup(response.text, "html.
    parser")
19
20 for link in soup.select("a[href$='.pdf']"):
21     #Name the pdf files using the last portion of
    each link which are unique in this case
22     filename = os.path.join(folder_location, link['
    href'].split('/')[-1])
23     with open(filename, 'wb') as f:
24         print(f)
25         f.write(requests.get(urljoin(url, link['href
    '])).content)
26
27
28 text = ''
29 for filename in os.listdir('p2files'):
30     size = os.path.getsize('p2files/' + filename)
31     if size > 0:
32         text += texttract.process('p2files/' +
    filename).decode('utf-8')
33
34 text = text.replace('\x00', '')
35 text = text.split()
36
37 Counters_found = Counter(text)
38 most_occur = Counters_found.most_common(10)

```

```

39 print(most_occur)
40
41 # part 2
42 all_words = list(Counter(text).items())
43
44 total_word_count = sum(Counters_found.values())
45 entropy_sum = 0
46
47 # entropy
48 for i in all_words:
49     entropy_sum -= (i[1] / total_word_count) * math
        .log((i[1] / total_word_count), 2)
50
51 print(entropy_sum)
52
53 # part 3
54 words = []
55 weights = []
56
57 for i in Counters_found.keys():
58     words.append(i.replace('.', ''))
59
60 for i in Counters_found.values():
61     weights.append(i / total_word_count)
62
63 random_paragraph = ''
64 for i in range(100):
65     if i % 10 == 0 and (i > 0):
66         random_paragraph += '. '
67         random_paragraph += choices(words, weights
        )[0].capitalize()
68
69     elif i % 10 != 0:
70         random_paragraph += ' '
71         random_paragraph += choices(words, weights
        )[0]
72
73     elif (i % 10) == 0 and i == 0:
74         random_paragraph += choices(words, weights
        )[0].capitalize()
75
76 print(random_paragraph)
77

```

```

1 import matplotlib
2 import numpy as np # linear algebra
3 import pandas as pd # data processing, CSV file I/O
  (e.g. pd.read_csv)
4 import xgboost as xgb
5
6 from scipy.stats import skew
7 from collections import OrderedDict
8 # The error metric: RMSE on the log of the sale
  prices.
9 from sklearn.metrics import mean_squared_error
10
11 def rmse(y_true, y_pred):
12     return np.sqrt(mean_squared_error(y_true,
    y_pred))
13
14 def Q3_P1():
15
16     # =====
17     # read in the data
18
19     # =====
20
21     train_data = pd.read_csv('../input/train.csv',
    index_col=0)
22     test_data = pd.read_csv('../input/test.csv',
    index_col=0)
23
24     # =====
25
26     # here, for this simple demonstration we shall
    only use the numerical columns
27     # and ignore the categorical features
28
29     # =====
30
31     X_train = train_data.select_dtypes(include=['
    number']).copy()
32     X_train = X_train.drop(['SalePrice'], axis=1)
33     y_train = train_data["SalePrice"]
34     X_test = test_data.select_dtypes(include=['
    number']).copy()

```

```

29
30     regressor = xgb.XGBRegressor()
31
32     # =====
33     # exhaustively search for the optimal
    hyperparameters
34
35     # =====
36     from sklearn.model_selection import
    GridSearchCV
37     # set up our search grid
38     param_grid = {"max_depth": [3, 4],
39                   "n_estimators": [600, 700],
40                   "learning_rate": [0.015, 0.020, 0.
41                                     025]}
42
43     # try out every combination of the above values
44     search = GridSearchCV(regressor, param_grid, cv
45                           =5).fit(X_train, y_train)
46
47     print("The best hyperparameters are ", search.
48           best_params_)
49
50     regressor = xgb.XGBRegressor(learning_rate=
51     search.best_params_["learning_rate"],
52     n_estimators=
53     search.best_params_["n_estimators"],
54     max_depth=search.
55     best_params_["max_depth"])
56
57     regressor.fit(X_train, y_train)
58
59     # =====
60
61     # use the model to predict the prices for the
62     test data
63
64     # =====

```

```

56 =====
57     predictions = regressor.predict(X_test)
58     # read in the ground truth file
59     solution = pd.read_csv('../input/solution.csv')
60     y_true = solution["SalePrice"]
61
62     from sklearn.metrics import
mean_squared_log_error
63     RMSLE = np.sqrt(mean_squared_log_error(y_true,
predictions))
64     print("The score of forum post is %.5f" % RMSLE
)
65
66     # =====
67     # write out CSV submission file
68
69     # =====
70
71     output = pd.DataFrame({"Id": test_data.index, "
SalePrice": predictions})
72     output.to_csv('forum_submission.csv', index=
False)
73
74     """What I Learned
75     I learned about using CV to optimize parameters
. More specifically I learned about the GridSearch
CV
76     which can optimize hyperparameters for any
model. Additionally, I gained a greater understand
of XGBoost and how to use it"""
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

```

84 6.0)
85     prices = pd.DataFrame({"price": train["
SalePrice"], "log(price + 1)": np.log1p(train["
SalePrice"])}))
86     # prices.hist()
87
88     # log transform the target:
89     train["SalePrice"] = np.log1p(train["SalePrice
"])
90
91     # log transform skewed numeric features:
92     numeric_feats = all_data.dtypes[all_data.
dtypes != "object"].index
93
94     skewed_feats = train[numeric_feats].apply(
lambda x: skew(x.dropna())) # compute skewness
95     skewed_feats = skewed_feats[skewed_feats > 0.
75]
96     skewed_feats = skewed_feats.index
97
98     all_data[skewed_feats] = np.log1p(all_data[
skewed_feats])
99
100    all_data = pd.get_dummies(all_data)
101
102    # filling NA's with the mean of the column:
103    all_data = all_data.fillna(all_data.mean())
104
105    # creating matrices for sklearn:
106    X_train = all_data[:train.shape[0]]
107    X_test = all_data[train.shape[0]:]
108    y = train.SalePrice
109
110    regressor = xgb.XGBRegressor()
111
112    # =====
=====
113    # exhaustively search for the optimal
hyperparameters
114
115    # =====
=====
115    from sklearn.model_selection import

```

```

115 GridSearchCV
116     # set up our search grid
117     param_grid = {"max_depth": [3, 4],
118                  "n_estimators": [700, 1000, 5000
119 ],
120                  "learning_rate": [0.01, 0.04, .
121 05]}
122
123     # try out every combination of the above
124     values
125     #search = GridSearchCV(regressor, param_grid,
126 cv=5).fit(X_train, y)
127
128     #print("The best hyperparameters are ", search
129 .best_params_)
130
131     #regressor = xgb.XGBRegressor(learning_rate=
132 search.best_params_["learning_rate"],
133 #                                n_estimators=
134 search.best_params_["n_estimators"],
135 #                                max_depth=search
136 .best_params_["max_depth"])
137     regressor = xgb.XGBRegressor(learning_rate=.01
138 ,
139                                n_estimators=5000
140 ,
141                                max_depth=4)
142
143     regressor.fit(X_train, y)
144
145     # =====
146     # use the model to predict the prices for the
147     test data
148     # =====
149
150     predictions = regressor.predict(X_test)
151     # read in the ground truth file
152     solution = pd.read_csv('../input/solution.csv'
153 )
154     y_true = solution["SalePrice"]
155

```



```

143     from sklearn.metrics import
mean_squared_log_error
144     RMSLE = np.sqrt(mean_squared_log_error(y_true
, np.expm1(predictions)))
145     print("The score of improved forum post is %.
5f" % RMSLE)
146
# =====
=====
147     # write out CSV submission file
148
# =====
=====
149     output = pd.DataFrame({"Id": test.Id, "
SalePrice": np.expm1(predictions)})
150     output.to_csv('improved_forum_submission.csv'
, index=False)
151     """Our Approach:
152     Our approach was to improve the test score was
to add data preprocessing steps, and to tune the
hyperparameters.
153     We tried many different value combinations
when using GridSearchCV to find the best
hyperparameters.
154     For preprocessing, we used a log transform on
the numerical data while also replacing NA values
with the mean."""
155
156
157
158 def Q3_P3():
159     # PREPROCESSING
160
161     train = pd.read_csv("../input/train.csv")
162     test = pd.read_csv("../input/test.csv")
163     train.head()
164     all_data = pd.concat((train.loc[:, 'MSSubClass
': 'SaleCondition'],
165                             test.loc[:, 'MSSubClass'
: 'SaleCondition']))
166     matplotlib.rcParams['figure.figsize'] = (12.0
, 6.0)
167     prices = pd.DataFrame({"price": train["
SalePrice"], "log(price + 1)": np.log1p(train["

```

```

167 SalePrice"]]))
168     # prices.hist()
169
170     # log transform the target:
171     train["SalePrice"] = np.log1p(train["SalePrice
172     "])
173     # log transform skewed numeric features:
174     numeric_feats = all_data.dtypes[all_data.
175     dtypes != "object"].index
176     skewed_feats = train[numeric_feats].apply(
177     lambda x: skew(x.dropna())) # compute skewness
178     skewed_feats = skewed_feats[skewed_feats > 0.
179     75]
180     skewed_feats = skewed_feats.index
181
182     all_data[skewed_feats] = np.log1p(all_data[
183     skewed_feats])
184
185     all_data = pd.get_dummies(all_data)
186
187     # filling NA's with the mean of the column:
188     all_data = all_data.fillna(all_data.mean())
189
190     # creating matrices for sklearn:
191     X_train = all_data[:train.shape[0]]
192     X_test = all_data[train.shape[0]:]
193     y = train.SalePrice
194
195     # Models
196     from sklearn.linear_model import Ridge,
197     RidgeCV, ElasticNet, Lasso, LassoLarsCV
198     from sklearn.model_selection import
199     cross_val_score
200     from sklearn.metrics import mean_squared_error
201
202     #UNDERFITTING MODEL
203     model_lasso = Lasso(1000)
204     model_lasso.fit(X_train, y)
205     l1_pred = model_lasso.predict(X_test)
206     l1_pred_train = model_lasso.predict(X_train)
207     rms = np.sqrt(mean_squared_error(y,
208     l1_pred_train))

```

```

203     print("Lasso Underfit Training Score:", rms)
204     solution = pd.DataFrame({"id": test.Id, "
SalePrice": np.expm1(l1_pred)})
205     solution.to_csv("Lasso_Regression_Underfit.csv
", index=False)
206
207     #OVERFITTING MODEL
208
209     dtrain = xgb.DMatrix(X_train, label=y)
210     dtest = xgb.DMatrix(X_test)
211
212     params = {"max_depth": 100, "eta": 0.1}
213     model_xgb = xgb.XGBRegressor(n_estimators=360
, max_depth=100, learning_rate=0.1) # the params
were tuned using xgb.cv
214     model_xgb.fit(X_train, y)
215     xgb_preds = np.expm1(model_xgb.predict(X_test
))
216     xgb_preds_train = model_xgb.predict(X_train)
217
218     solution = pd.DataFrame({"id": test.Id, "
SalePrice": xgb_preds})
219     solution.to_csv("XGB_Regression.csv", index=
False)
220
221     print("XGBoost Overfit Training Score:", rmse(
y, xgb_preds_train))
222
223
224
225
226
227 # Press the green button in the gutter to run the
script.
228 if __name__ == '__main__':
229     Q3_P1()
230     Q3_P2()
231     Q3_P3()

```