

Predicting Resale Value of Knives from a Texas Government Surplus Store

Using Machine Learning to Support an Ebay Store's Financial Success

Data Exploration and Modeling

Author: Dylan Dey

Model

```
In [1]: 1 from sklearn.model_selection import train_test_split
2 import os
3 from collections import Counter
4
5 import pandas as pd
6 import json
7 import requests
8 import numpy as np
9 import matplotlib.pyplot as plt
10 %matplotlib inline
11 import seaborn as sns
12 import ast
13 import re
14
15 from tensorflow.keras.preprocessing.text import Tokenizer
16 from tensorflow.keras.preprocessing.sequence import pad_sequences
17 from tensorflow.keras.layers import Dense, Input, GlobalMaxPooling1D
18 from tensorflow.keras.layers import LSTM, Embedding, Flatten, GRU
19 from tensorflow.keras.layers import Conv1D, MaxPooling1D, GlobalMaxPooling2D
20 from tensorflow.keras.layers import Conv2D, MaxPooling2D, Dropout, BatchNormalization
21 from tensorflow.keras.layers import SimpleRNN
22 from tensorflow.keras.models import Model
23 from keras import models
24 from keras import layers
25 import tensorflow as tf
26 from keras.utils import plot_model
27 from sklearn.metrics import mean_absolute_error
28 from keras_preprocessing.image import ImageDataGenerator

In [2]: 1 plt.style.use('dark_background')
```

Function Definition

```

In [3]: 1 def apply_iqr_filter(df):
2
3     price_Q1 = df['converted_price'].quantile(0.25)
4     price_Q3 = df['converted_price'].quantile(0.75)
5     price_iqr = price_Q3 - price_Q1
6
7     profit_Q1 = df['profit'].quantile(0.25)
8     profit_Q3 = df['profit'].quantile(0.75)
9     profit_iqr = profit_Q3 - profit_Q1
10
11     ROI_Q1 = df['ROI'].quantile(0.25)
12     ROI_Q3 = df['ROI'].quantile(0.75)
13     ROI_iqr = ROI_Q3 - ROI_Q1
14
15     price_upper_limit = price_Q3 + (1.5 * price_iqr)
16     price_lower_limit = price_Q1 - (1.5 * price_iqr)
17
18     profit_upper_limit = profit_Q3 + (1.5 * profit_iqr)
19     profit_lower_limit = profit_Q1 - (1.5 * profit_iqr)
20
21     ROI_upper_limit = ROI_Q3 + (1.5 * ROI_iqr)
22     ROI_lower_limit = ROI_Q1 - (1.5 * ROI_iqr)
23
24     # print(f'Brand: {df.brand[0]}')
25     # print(f'price upper limit: ${np.round(price_upper_limit,2)}')
26     # print(f'price lower limit: ${np.round(price_lower_limit,2)}')
27     # print('-----')
28     # print(f'profit upper limit: ${np.round(profit_upper_limit,2)}')
29     # print(f'profit lower limit: ${np.round(profit_lower_limit,2)}')
30     # print('-----')
31     # print(f'ROI upper limit: {np.round(ROI_upper_limit,2)}%')
32     # print(f'ROI lower limit: {np.round(ROI_lower_limit,2)}%')
33     # print('-----')
34
35     new_df = df[(df['converted_price'] < price_upper_limit) &
36                 (df['converted_price'] > price_lower_limit) &
37                 (df['profit'] < profit_upper_limit) &
38                 (df['profit'] > profit_lower_limit) &
39                 (df['ROI'] < ROI_upper_limit) &
40                 (df['ROI'] > ROI_lower_limit)]
41
42     return new_df
43
44 #download jpg urls from dataframe
45 def download(row):
46     filename = os.path.join(root_folder, str(row.name) + im_extension)
47
48     # create folder if it doesn't exist
49     os.makedirs(os.path.dirname(filename), exist_ok=True)
50
51     url = row.Image
52     # print(f"Downloading {url} to {filename}")
53
54     try:
55         r = requests.get(url, allow_redirects=True)
56         with open(filename, 'wb') as f:
57             f.write(r.content)
58     except:
59         print(f'{filename} error')
60
61
62 def cardinality_threshold(column,threshold=0.75,return_categories_list=True):
63     #calculate the threshold value using
64     #the frequency of instances in column
65     threshold_value=int(threshold*len(column))
66     #initialize a new list for lower cardinality column
67     categories_list=[]
68     #initialize a variable to calculate sum of frequencies
69     s=0
70     #Create a dictionary (unique_category: frequency)
71     counts=Counter(column)
72
73     #Iterate through category names and corresponding frequencies after sorting the categories
74     #by descending order of frequency
75     for i,j in counts.most_common():
76         #Add the frequency to the total sum
77         s += dict(counts)[i]
78         #append the category name to the categories list
79         categories_list.append(i)
80         #Check if the global sum has reached the threshold value, if so break the loop
81         if s >= threshold_value:
82             break
83         #append the new 'Other' category to list
84         categories_list.append('Other')
85
86     #Take all instances not in categories below threshold

```

```

87     #that were kept and lump them into the
88     #new 'Other' category.
89     new_column = column.apply(lambda x: x if x in categories_list else 'Other')
90
91     #Return the transformed column and
92     #unique categories if return_categories = True
93     if (return_categories_list):
94         return new_column, categories_list
95     #Return only the transformed column if return_categories=False
96     else:
97         return new_column
98
99 def fix(col):
100     dd = dict()
101     for d in col:
102         values = list(d.values())
103         if len(values) == 2:
104             dd[values[0]] = values[1]
105     return dd
106
107 #function for extracted item Specifics from Shopping API data
108 def transform_item_specifics(df, perc=90.0):
109
110     df.dropna(subset=['ItemSpecifics'], inplace=True)
111     df['ItemSpecifics'] = df['ItemSpecifics'].apply(lambda x: ast.literal_eval(x))
112     df['item_list'] = df['ItemSpecifics'].apply(lambda x: x['NameValueList'])
113
114     df['ItemSpecifics'] = df['ItemSpecifics'].apply(lambda x: [x['NameValueList']] if isinstance(x, 'NameValueList')
115
116     df['ItemSpecifics'] = df['ItemSpecifics'].apply(fix)
117
118     df = pd.json_normalize(df['ItemSpecifics'])
119
120     min_count = int(((100-perc)/100)*df.shape[0] + 1)
121     mod_df = df.dropna(axis=1,
122                       thresh=min_count)
123
124     return mod_df
125
126 # This function removes noisy data
127 #lots/sets/groups of knives can
128 #confuse the model from predicting
129 #the appropriate value of individual knives
130 def data_cleaner(df):
131     lot = re.compile('(?!-\S)lot(?:[^\s.,:?!])')
132     group = re.compile('(group)')
133     is_set = re.compile('(?!-\S)set(?:[^\s.,:?!])')
134     df['title'] = df['title'].str.lower()
135     trim_list = [lot, group, is_set]
136     for item in trim_list:
137         df.loc[df['title'].apply(lambda x: re.search(item, x)).notnull(), 'trim'] = 1
138     to_drop = df.loc[df['trim'] == 1].index
139     df.drop(to_drop, inplace=True)
140     df.drop('trim', axis=1, inplace=True)
141
142     return df
143
144
145
146 def prepare_listed(listed_data_df, Ids_df):
147     listed_data_df.drop('galleryPlusPictureURL', axis=1, inplace=True)
148
149     Ids_df.rename({'Title': 'title',
150                  'ItemID': 'itemId'},
151                  axis=1, inplace=True)
152
153     Ids_df.drop(['ConditionID', 'ConvertedCurrentPrice'],
154                axis=1, inplace=True)
155     Ids_df['title'] = Ids_df['title'].str.lower()
156
157     df_merged = listed_data_df.merge(Ids_df)
158
159     df_spec = transform_item_specifics(df_merged, perc=65.0)
160
161     df_spec.drop('Brand', axis=1, inplace=True)
162
163     tot_listed_df = df_merged.join(df_spec)
164
165     listed_knives = data_cleaner(tot_listed_df).copy()
166     listed_knives.drop(['sellingStatus', 'shippingInfo',
167                       'GalleryURL', 'ItemSpecifics',
168                       'item_list', 'listingInfo'],
169                       axis=1, inplace=True)
170     listed_used_knives = listed_knives.loc[list_knives['condition'] != 1000.0]
171     listed_used_knives.reset_index(drop=True, inplace=True)
172

```

```
173     return listed_used_knives
174
175
176 def prepare_tera_df(df, x, overhead_cost=3):
177     df['price_in_US'] = df['price_in_US'].str.replace("$", "")
178     df['price_in_US'] = df['price_in_US'].str.replace(",", "")
179     df['price_in_US'] = df['price_in_US'].apply(float)
180
181     df['shipping_cost'] = df['shipping_cost'].str.replace("$", "")
182     df['shipping_cost'] = df['shipping_cost'].str.replace(",", "")
183     df['shipping_cost'] = df['shipping_cost'].apply(float)
184
185     df['brand'] = list(bucket_dict.keys())[x]
186     df['converted_price'] = (df['price_in_US'] + df['shipping_cost'])
187     df['cost'] = list(bucket_dict.values())[x] + overhead_cost + 4.95
188     df['profit'] = ((df['converted_price']*.87) - df['cost'])
189     df['ROI'] = (df['profit'] / df['cost'])*100.0
190
191     return df
192
```

Load Data

```

In [4]: 1 #load Finding API data
2 df_bench = pd.read_csv("listed_data/df_bench.csv")
3 df_buck = pd.read_csv("listed_data/df_buck.csv")
4 df_case = pd.read_csv("listed_data/df_case.csv")
5 df_caseXX = pd.read_csv("listed_data/df_CaseXX.csv")
6 df_crkt = pd.read_csv("listed_data/df_crkt.csv")
7 df_kersh = pd.read_csv("listed_data/df_kershaw.csv")
8 df_sog = pd.read_csv("listed_data/df_sog.csv")
9 df_spyd = pd.read_csv("listed_data/df_spyderco.csv")
10 df_vict = pd.read_csv("listed_data/df_victorinox.csv")
11
12 #load Shopping API data
13 bench = pd.read_csv("listed_data/benchIds.csv")
14 buck = pd.read_csv("listed_data/buckIds.csv")
15 case = pd.read_csv("listed_data/caseIds.csv")
16 caseXX = pd.read_csv("listed_data/caseXXIds.csv")
17 crkt = pd.read_csv("listed_data/crktIds.csv")
18 kershaw = pd.read_csv("listed_data/kershawIds.csv")
19 sog = pd.read_csv("listed_data/sogIds.csv")
20 spyd = pd.read_csv("listed_data/spydIds.csv")
21 vict = pd.read_csv("listed_data/victIds.csv")
22
23 #Load scraped terapeak sold data
24 sold_bench = pd.read_csv("terapeak_data/bench_scraped2.csv")
25 sold_buck1 = pd.read_csv("terapeak_data/buck_scraped2.csv")
26 sold_buck2 = pd.read_csv("terapeak_data/buck_scraped2_reversed.csv")
27 sold_case = pd.read_csv("terapeak_data/case_scraped2.csv")
28 sold_caseXX1 = pd.read_csv("terapeak_data/caseXX_scraped2.csv")
29 sold_caseXX2 = pd.read_csv("terapeak_data/caseXX2_reversed.csv")
30 sold_crkt = pd.read_csv("terapeak_data/crkt_scraped.csv")
31 sold_kershaw1 = pd.read_csv("terapeak_data/kershaw_scraped2.csv")
32 sold_kershaw2 = pd.read_csv("terapeak_data/kershaw_scraped2_reversed.csv")
33 sold_sog = pd.read_csv("terapeak_data/SOG_scraped2.csv")
34 sold_spyd = pd.read_csv("terapeak_data/spyd_scraped2.csv")
35 sold_vict1 = pd.read_csv("terapeak_data/vict_scraped.csv")
36 sold_vict2 = pd.read_csv("terapeak_data/vict_reversed.csv")
37
38 sold_list = [sold_bench,sold_buck1,
39              sold_buck2,sold_case,
40              sold_caseXX1,sold_caseXX2,
41              sold_crkt,sold_kershaw1,
42              sold_kershaw2,sold_sog,
43              sold_spyd, sold_vict1,
44              sold_vict2]
45
46
47 listed_df = pd.concat([df_bench,df_buck,
48                        df_case,df_caseXX,
49                        df_crkt,df_kersh,
50                        df_sog,df_spyd,
51                        df_vict])
52
53 used_listed_df = listed_df.loc[listed_df['condition'] != 1000.0].copy()
54
55 cols = ['title','pictureURLLarge','converted_price','brand','profit','ROI']
56 used_listed = used_listed_df[cols].copy()
57 used_listed.dropna(subset=['pictureURLLarge'], inplace=True)
58
59 used_listed.reset_index(drop=True, inplace=True)
60
61
62 bucket_dict = {'benchmade': 45.0,
63                'buck': 20.0,
64                'case': 20.0,
65                'crkt': 15.0,
66                'kershaw': 15.0,
67                'sog': 15.0,
68                'spyderco': 30.0,
69                'victorinox': 20.0
70                }

```

Prepare Data

```
In [5]: 1 for dataframe in sold_list:
2         dataframe.rename({'Text': 'title',
3                           'shipping_': 'shipping_cost'},
4                           axis=1, inplace=True)
5
6         dataframe['date_sold'] = pd.to_datetime(dataframe['date_sold'])
7
8     sold_buck = pd.concat([sold_buck1,sold_buck2])
9     sold_caseXX = pd.concat([sold_caseXX1,sold_caseXX2])
10    sold_kershaw = pd.concat([sold_kershaw1,sold_kershaw2])
11    sold_vict = pd.concat([sold_vict1,sold_vict2])
12
13    sold_bench = prepare_tera_df(sold_bench, 0)
14    sold_buck = prepare_tera_df(sold_buck, 1)
15    sold_case = prepare_tera_df(sold_case, 2)
16    sold_caseXX = prepare_tera_df(sold_caseXX, 2)
17    sold_crkt = prepare_tera_df(sold_crkt, 3)
18    sold_kershaw = prepare_tera_df(sold_kershaw, 4)
19    sold_sog = prepare_tera_df(sold_sog, 5)
20    sold_spyd = prepare_tera_df(sold_spyd, 6)
21    sold_vict = prepare_tera_df(sold_vict, 7)
```

```
In [6]: 1 for dataframe in sold_list:
2         dataframe['title'] = dataframe['title'].str.lower()
3         dataframe['title'] = dataframe['title'].str.strip()
4         dataframe.drop_duplicates(
5             subset = ['date_sold','price_in_US',
6                      'shipping_cost'],
7             keep = 'last', inplace=True)
```

```
In [7]: 1 sold_df = pd.concat([sold_bench, sold_buck,
2                         sold_case, sold_caseXX,
3                         sold_crkt, sold_kershaw,
4                         sold_sog, sold_spyd,
5                         sold_vict])
6
7     sold_knives = data_cleaner(sold_df).copy()
8
9
10    df = pd.concat([sold_knives,used_listed]).copy()
11    df['Image'].fillna(df['pictureURLLarge'], inplace=True)
12
13    df = apply_iqr_filter(df).copy()
14    df.reset_index(drop=True, inplace=True)
```

```
In [8]: 1 def clean_text(x):
2         pattern = r'^a-zA-z0-9\s]'
3         text = re.sub(pattern, '', x)
4         return text
```

```
In [9]: 1 df['title'] = df['title'].apply(clean_text)
```

```
In [10]: 1 df['title'].sample(20).apply(print)
```

```
case xx pocket knife folding 2bladecase xx pocket knife folding 2blade
victorinox swiss army knife classic sd pink camo 58mm multitool nice
case xx 6318 ss red delrin stockman pocket knife used case xx 6318 ss red delrin stockman pocket knife used
crkt nirk tighe 2 folding knife 5240 excellentcrkt nirk tighe 2 folding knife 5240 excellent
case xx copperhead knife 6249 bone handles unused 1979case xx copperhead knife 6249 bone handles unused 1979
crkt crawford falcon 6242 folding pocket knifecrkt crawford falcon 6242 folding pocket knife
case xx baby swamp rat knife 19401964 era reddish bone handles 1 2 carry nr case xx baby swamp rat knife 19401964 era
reddish bone handles 1 2 carry nr
kershaw ao cryo 1555ti folding pocket knifekershaw ao cryo 1555ti folding pocket knife
benchmade 710 mchenry williams d2 folding pocket knifebenchmade 710 mchenry williams d2 folding pocket knife
2015 case xx saw cut amber bone medium stockman knife 6318cv chrome vanadium 2015 case xx saw cut amber bone medium s
tockman knife 6318cv chrome vanadium
case xx stockman pocket worn 25th anniversary knife 6347 sscase xx stockman pocket worn 25th anniversary knife 6347
ss
victorinox alox electrician swiss army knife 93 mmvictorinox alox electrician swiss army knife 93 mm
victorinox evolution grip s18 swiss army knife 73victorinox evolution grip s18 swiss army knife 73
2019 case xx trapper knife keep your hands sharp tour lasered bone serial 0032019 case xx trapper knife keep your ha
nds sharp tour lasered bone serial 003
original 1994 buck crosslock 180 knife serrated blade usa madeoriginal 1994 buck crosslock 180 knife serrated blade u
sa made
buck usa bantam mossy oak camo 284 folding pocket knife
kershaw steven seagal knife 1680 japan aus8a ken onion designkershaw steven seagal knife 1680 japan aus8a ken onion d
esign
vintage buck 110 381 1980s hunter knife w box sheath paperwork very nicevintage buck 110 381 1980s hunter knife w b
ox sheath paperwork very nice
kershaw 1660ol leek pocket knife olive green assisted opening liner lock usakershaw 1660ol leek pocket knife olive
green assisted opening liner lock usa
new kershaw 1555ti cryo hinderer designed assisted opening knife lnnew kershaw 1555ti cryo hinderer designed assisted
opening knife ln
```

```
Out[10]: 28429      None
75017      None
27880      None
35157      None
23059      None
33358      None
21677      None
43359      None
3243       None
29370      None
22637      None
61215      None
65114      None
29457      None
9123       None
67824      None
47152      None
4874       None
42321      None
38806      None
Name: title, dtype: object
```

```
In [11]: 1 df['title_len'] = df['title'].apply(lambda x: len(x))
2 df['word_count'] = df['title'].apply(lambda x: len(x.split()))
```

```
In [12]: 1 def avg_word_len(x):
2         words = x.split()
3         word_len = 0
4         for word in words:
5             word_len += len(word)
6
7         return word_len / len(words)
```

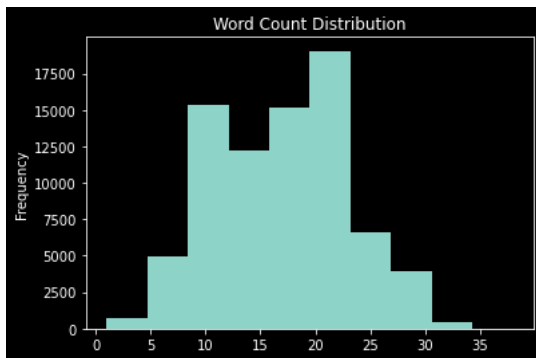
```
In [13]: 1 df['avg_word_len'] = df['title'].apply(lambda x: avg_word_len(x))
```

```
In [14]: 1 # pd.options.plotting.backend = "plotly"
```



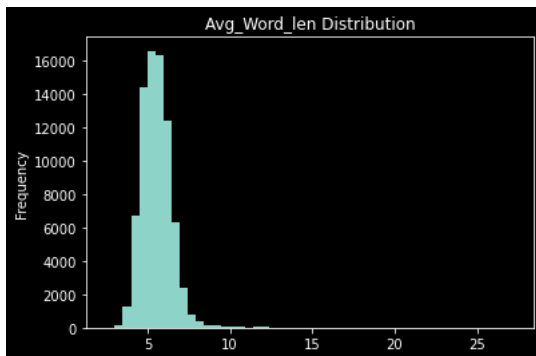
```
In [15]: 1 df['word_count'].plot(kind='hist', title='Word Count Distribution')
```

```
Out[15]: <AxesSubplot:title={'center':'Word Count Distribution'}, ylabel='Frequency'>
```

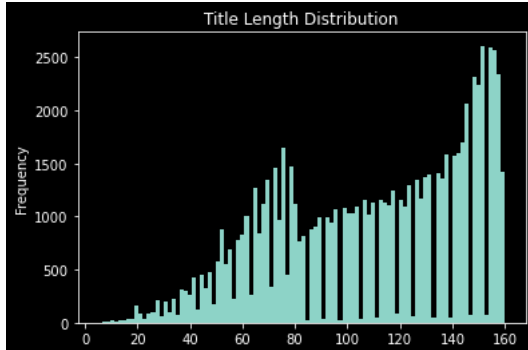


```
In [16]: 1 df['avg_word_len'].plot(kind='hist', bins = 50, title = 'Avg_Word_len Distribution')
```

```
Out[16]: <AxesSubplot:title={'center':'Avg_Word_len Distribution'}, ylabel='Frequency'>
```



```
In [17]: 1 df['title_len'].plot(kind='hist', bins= 100,title = 'Title Length Distribution');
```



```
In [18]: 1 df.describe()
```

```
Out[18]:
```

	price_in_US	shipping_cost	converted_price	cost	profit	ROI	title_len	word_count	avg_word_len
count	66768.000000	66768.000000	78330.000000	66768.000000	78330.000000	78330.000000	78330.000000	78330.000000	78330.000000
mean	45.969092	4.315289	50.044127	28.754727	15.728080	54.837503	109.446062	17.086340	5.534031
std	35.153294	3.598604	36.026730	7.272042	28.695811	103.741175	36.561769	6.116623	0.982619
min	0.010000	0.000000	6.940000	22.950000	-38.812500	-73.853309	5.000000	1.000000	2.500000
25%	19.990000	0.000000	22.182500	22.950000	-5.600000	-22.325490	78.000000	12.000000	4.888889
50%	35.000000	4.950000	38.740000	27.950000	7.065000	26.064401	114.000000	17.000000	5.466667
75%	62.000000	5.900000	67.490000	27.950000	29.204950	102.866646	144.000000	22.000000	6.000000
max	167.000000	80.000000	167.490000	52.950000	109.153000	605.666667	160.000000	38.000000	27.200000

Neural network with "title" column as input

```
In [19]: 1 df_title = df.loc[:, ['title', 'converted_price']]
          2
          3
          4 df_title.rename({'title': 'data',
          5                     'converted_price': 'labels'},
          6                     axis=1, inplace=True)
```

```
In [20]: 1 df_title.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 78330 entries, 0 to 78329
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  ---
0    data    78330 non-null     object
1   labels  78330 non-null     float64
dtypes: float64(1), object(1)
memory usage: 1.2+ MB
```

```
In [21]: 1 # df_title['labels'] = (df_title['labels']/mean_price)
          2 Y = df_title['labels'].values
```

```
In [22]: 1 df_title['data'].sample(10).apply(print)
```

2000 case xx usa small stockman pocket knife 6327 ss blue jig bone made in usa
 2000 case xx usa small stockman pocket knife 6327 ss blue jig bone made in usa
 kershaw 1670rd blur pocket knife speedsafe liner lock usa
 kershaw 1670rd blur pocket knife speedsafe liner lock usa
 kershaw folding pocket knife kai 2415
 kershaw folding pocket knife kai 2415
 victorinox signature lite with lgb memory card swiss army knife 54
 victorinox signature lite with lgb memory card swi
 ss army knife 54
 spyderco clipitool standard folding knife c208gp plain edge blade black g10
 spyderco clipitool standard folding knife c208gp plain edge blade black g10
 ntsa vntg195286 swiss army victorinox champion plus mfunction pkt knif
 ntsa vntg195286 swiss army victorinox champ
 ion plus mfunction pkt knife
 9 empty benchmade knife boxes
 9 empty benchmade knife boxes
 kershaw 1670olblk ken onion design folding pocket knife assisted open
 kershaw 1670olblk ken onion design folding pocke
 t knife assisted open
 case xx mini copperlock smooth abalone pocket knife usa 12020
 case xx mini copperlock smooth abalone pocket knife usa 12020
 used buck 284 pocket knife
 used buck 284 pocket knife

```
Out[22]: 27488      None
          49043      None
          38988      None
          60169      None
          54246      None
          58302      None
          2386       None
          38622      None
          24044      None
          15348      None
          Name: data, dtype: object
```

```
In [23]: 1 df_train, df_test, Ytrain, Ytest = train_test_split(df_title['data'],
          2                                                    Y,
          3                                                    test_size=0.3,
          4                                                    random_state=51)
```

```
In [24]: 1 X_val, X_test, Y_val, Y_test = train_test_split(df_test,
          2                                                    Ytest,
          3                                                    test_size=0.5,
          4                                                    random_state=51)
```

GRU

```
In [25]: 1 #Vectorize vocab
          2 voc_size = 25000
          3 max_len = 30
          4 embedding_features = 35
          5 tokenizer = Tokenizer(num_words=voc_size, oov_token = '<OOV>')
          6 tokenizer.fit_on_texts(df_train)
          7 sequences_train = tokenizer.texts_to_sequences(df_train)
          8 sequences_val = tokenizer.texts_to_sequences(X_val)
          9 sequences_test = tokenizer.texts_to_sequences(X_test)
```

```
In [26]: 1 #add padding to ensure all inputs are the same size
2 data_train = pad_sequences(sequences_train, maxlen=max_len, padding= 'post', truncating = 'post')
3 data_val = pad_sequences(sequences_val, maxlen=max_len, padding= 'post', truncating = 'post')
4 data_test = pad_sequences(sequences_test, maxlen=max_len, padding= 'post', truncating = 'post')
```

```
In [27]: 1 data_train.shape
```

Out[27]: (54831, 30)

```
In [28]: 1 model = models.Sequential()
2 model.add(Embedding(voc_size, embedding_features, input_length = max_len))
3 model.add(GRU(100))
4 model.add(Dense(64, activation = 'relu'))
5 # model.add(Dropout(0.3))
6 model.add(Dense(32, activation = 'relu'))
7 model.add(Dropout(0.3))
8 model.add(Dense(1, activation = 'relu'))
9 model.summary()
```

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, 30, 35)	875000
gru (GRU)	(None, 100)	41100
dense (Dense)	(None, 64)	6464
dense_1 (Dense)	(None, 32)	2080
dropout (Dropout)	(None, 32)	0
dense_2 (Dense)	(None, 1)	33
=====		
Total params: 924,677		
Trainable params: 924,677		
Non-trainable params: 0		

```

In [29]: 1 # Compile and fit
          2 model.compile(
          3     loss='MSE',
          4     optimizer='adam',
          5     metrics=['mae']
          6 )
          7
          8
          9 print('Training model...')
         10 r = model.fit(
         11     data_train,
         12     Ytrain,
         13     epochs=10,
         14     validation_data=(data_val, Y_val)
         15 )

Training model...
Epoch 1/10
1714/1714 [=====] - 28s 16ms/step - loss: 1406.4950 - mae: 29.1902 - val_loss: 1286.9244 - val_mae: 28.9204
Epoch 2/10
1714/1714 [=====] - 27s 16ms/step - loss: 1300.3254 - mae: 27.8673 - val_loss: 620.8804 - val_mae: 17.8693
Epoch 3/10
1714/1714 [=====] - 28s 16ms/step - loss: 580.3170 - mae: 16.7247 - val_loss: 463.5613 - val_mae: 14.9725
Epoch 4/10
1714/1714 [=====] - 29s 17ms/step - loss: 453.9511 - mae: 14.5761 - val_loss: 444.3378 - val_mae: 14.1846
Epoch 5/10
1714/1714 [=====] - 29s 17ms/step - loss: 393.9201 - mae: 13.4119 - val_loss: 438.7681 - val_mae: 14.0230
Epoch 6/10
1714/1714 [=====] - 27s 16ms/step - loss: 345.1503 - mae: 12.5439 - val_loss: 435.7736 - val_mae: 13.8277
Epoch 7/10
1714/1714 [=====] - 27s 16ms/step - loss: 319.4977 - mae: 12.0130 - val_loss: 440.4136 - val_mae: 13.9308
Epoch 8/10
1714/1714 [=====] - 27s 16ms/step - loss: 296.8629 - mae: 11.5484 - val_loss: 451.5638 - val_mae: 14.0264
Epoch 9/10
1714/1714 [=====] - 27s 16ms/step - loss: 272.7390 - mae: 11.0679 - val_loss: 453.2909 - val_mae: 13.9109
Epoch 10/10
1714/1714 [=====] - 28s 16ms/step - loss: 262.2486 - mae: 10.8240 - val_loss: 459.1286 - val_mae: 14.0011

```

```

In [30]: 1 pred=model.predict(data_test)

```

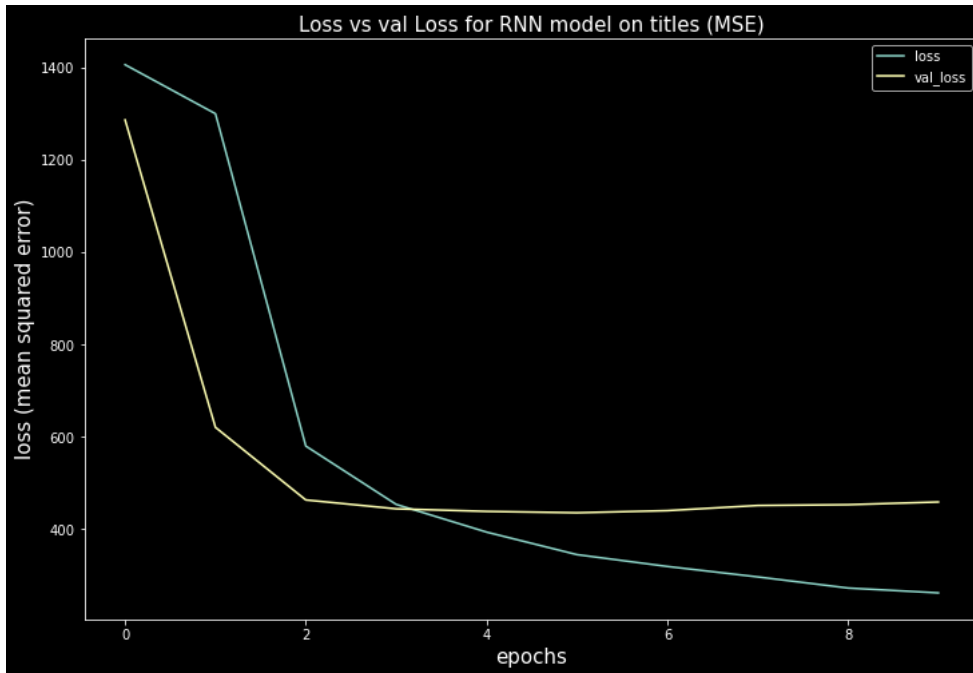
```

In [31]: 1 test_results = model.evaluate(data_test, Y_test)

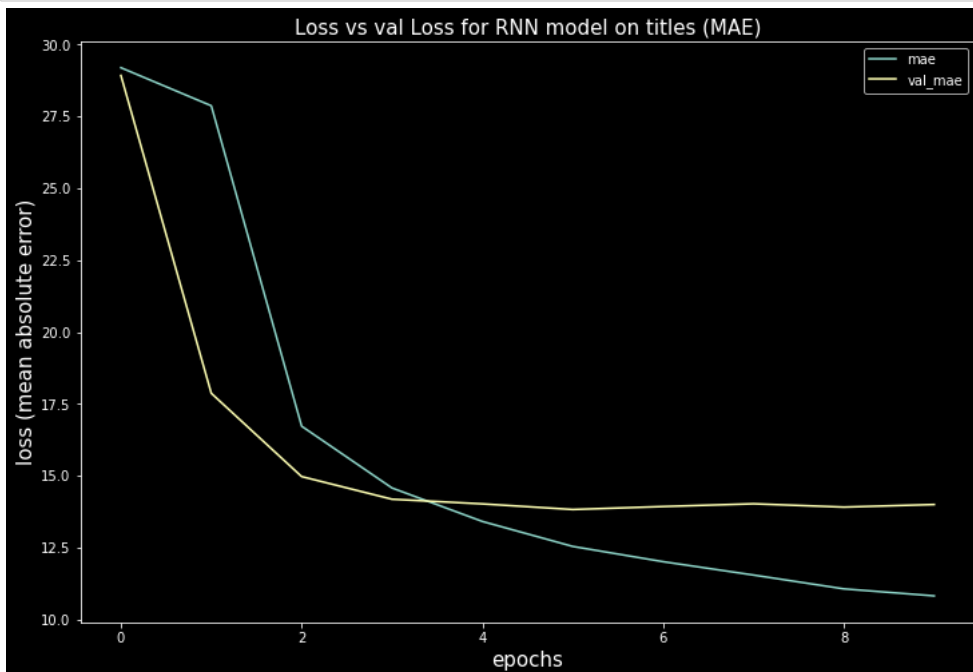
368/368 [=====] - 1s 3ms/step - loss: 465.1646 - mae: 13.9465

```

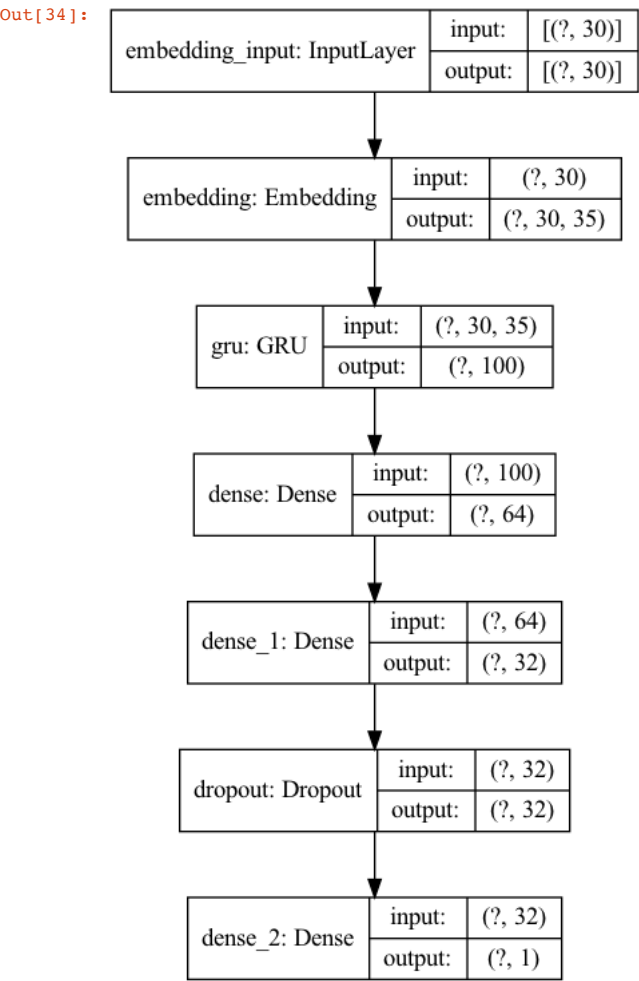
```
In [32]: 1 fig = plt.subplots(figsize=(12,8))
2 plt.plot(r.history['loss'], label='loss')
3 plt.plot(r.history['val_loss'], label='val_loss')
4 plt.title("Loss vs val Loss for RNN model on titles (MSE)", fontsize=15)
5 plt.xlabel("epochs", fontsize=15)
6 plt.ylabel("loss (mean squared error)", fontsize=15)
7 plt.legend();
8 plt.savefig('images/RNN_GRU_MSE1.png')
```



```
In [33]: 1 fig = plt.subplots(figsize=(12,8))
2 plt.plot(r.history['mae'], label='mae')
3 plt.plot(r.history['val_mae'], label='val_mae')
4 plt.title("Loss vs val Loss for RNN model on titles (MAE)", fontsize=15)
5 plt.xlabel("epochs", fontsize=15)
6 plt.ylabel("loss (mean absolute error)", fontsize=15)
7 plt.legend();
8 plt.savefig('images/RNN_GRU_MAE1.png')
```



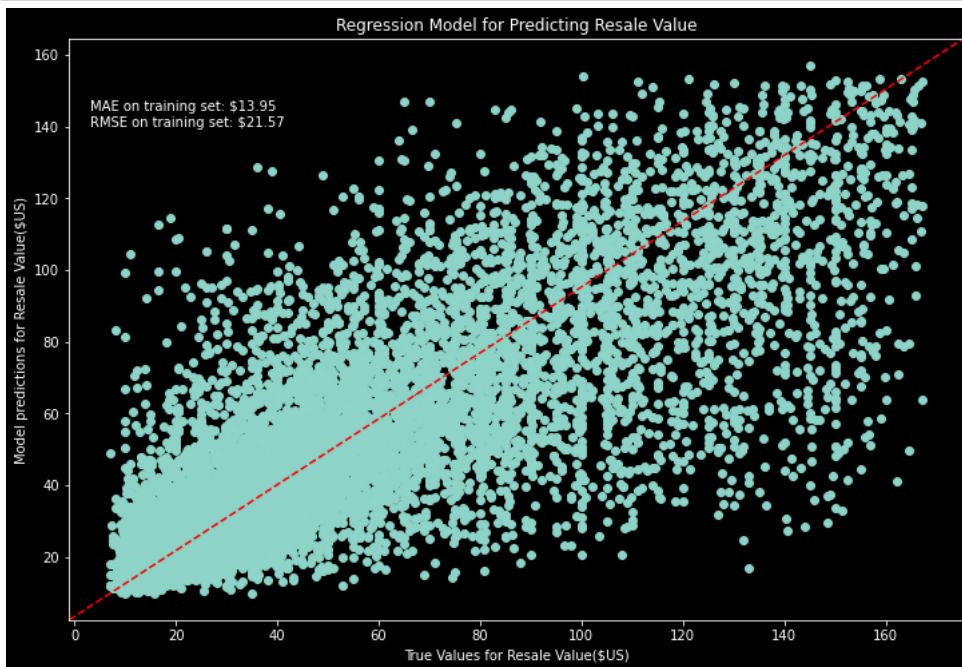
```
In [34]: 1 plot_model(model, show_shapes=True, to_file='images/RNN_GRU1_arc.png')
```



```
In [35]: 1 test_mae = mean_absolute_error(Y_test, pred)
```

```
In [36]: 1 RMSE = np.sqrt(test_results[0])
```

```
In [37]: 1 string_score = f'\nMAE on training set: ${test_mae:.2f}'
2 string_score += f'\nRMSE on training set: ${RMSE:.2f}'
3 fig, ax = plt.subplots(figsize=(12, 8))
4 plt.scatter(Y_test, pred)
5 ax.plot([0, 1], [0, 1], transform=ax.transAxes, ls="--", c="red")
6 plt.text(3, 140, string_score)
7 plt.title('Regression Model for Predicting Resale Value')
8 plt.ylabel('Model predictions for Resale Value($US)')
9 plt.xlabel('True Values for Resale Value($US)')
10 plt.savefig('images/regression_GRU_relu1.png');
```



```
In [38]: 1 df_title['labels'].describe()
```

```
Out[38]: count      78330.000000
mean         50.044127
std          36.026730
min           6.940000
25%          22.182500
50%          38.740000
75%          67.490000
max         167.490000
Name: labels, dtype: float64
```

```
In [39]: 1 df_title = df.loc[:, ['title', 'converted_price']]
2
3
4 df_title.rename({'title': 'data',
5                  'converted_price': 'labels'},
6                 axis=1, inplace=True)
```

```
In [40]: 1 # df_title['labels'] = (df_title['labels']/mean_price)
2 Y = df_title['labels'].values
```

```
In [41]: 1 df_train, df_test, Ytrain, Ytest = train_test_split(df_title['data'],
2                                                         Y,
3                                                         test_size=0.3,
4                                                         random_state=42)
```

```
In [42]: 1 X_val, X_test, Y_val, Y_test = train_test_split(df_test,
2                                                         Ytest,
3                                                         test_size=0.5,
4                                                         random_state=42)
```

LSTM

```
In [43]: 1 # Convert sentences to sequences
2 MAX_VOCAB_SIZE = 25000
3 tokenizer = Tokenizer(num_words=MAX_VOCAB_SIZE)
4 tokenizer.fit_on_texts(df_train)
5 sequences_train = tokenizer.texts_to_sequences(df_train)
6 sequences_val = tokenizer.texts_to_sequences(X_val)
7 sequences_test = tokenizer.texts_to_sequences(X_test)
```

```
In [44]: 1 # get word -> integer mapping
2 word2idx = tokenizer.word_index
3 V = len(word2idx)
4 print('Found %s unique tokens.' % V)
```

Found 32753 unique tokens.

```
In [45]: 1 # pad sequences so that we get a N x T matrix
2 data_train = pad_sequences(sequences_train)
3 print('Shape of data train tensor:', data_train.shape)
4
5 # get sequence length
6 T = data_train.shape[1]
```

Shape of data train tensor: (54831, 41)

```
In [46]: 1 data_val = pad_sequences(sequences_val, maxlen=T)
2 print('Shape of data test tensor:', X_val.shape)
```

Shape of data test tensor: (11749,)

```
In [47]: 1 data_test = pad_sequences(sequences_test, maxlen=T)
2 print('Shape of data test tensor:', X_test.shape)
```

Shape of data test tensor: (11750,)

```
In [48]: 1 # Create the RNN model
2 # We get to choose embedding dimensionality
3 D = 30
4 # Hidden state dimensionality
5 M = 35
6 i = Input(shape=(T,))
7 x = Embedding(V + 1, D)(i)
8 x = LSTM(M, return_sequences=True)(x)
9 x = GlobalMaxPooling1D()(x)
10 x = Dense(62, activation='relu')(x)
11 x = Dense(32, activation='relu')(x)
12 x = Dropout(0.3)(x)
13 x = Dense(1, activation='relu')(x)
14 model = Model(i, x)
```



```
In [49]: 1 # Compile and fit
2 model.compile(
3     loss='MSE',
4     optimizer='adam',
5     metrics=['mae']
6 )
7
8
9 print('Training model...')
10 r = model.fit(
11     data_train,
12     Ytrain,
13     epochs=5,
14     validation_data=(data_val, Y_val)
15 )
```

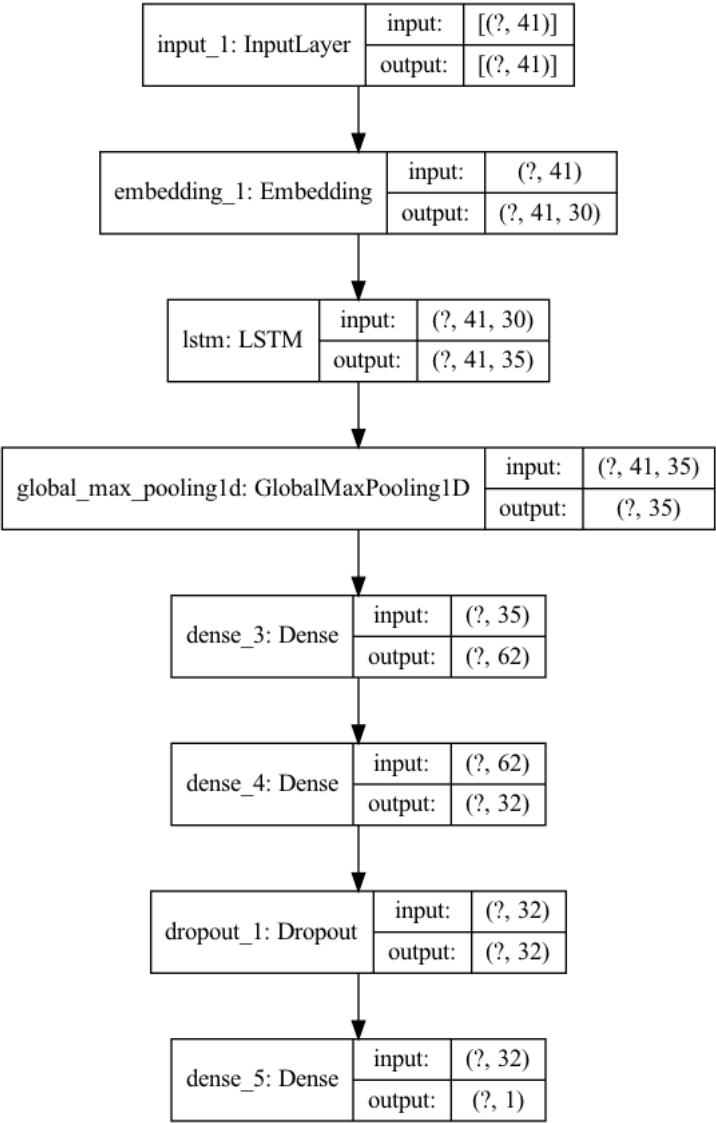
Training model...
Epoch 1/5
1714/1714 [=====] - 25s 15ms/step - loss: 1026.4459 - mae: 23.0135 - val_loss: 539.7323 - val_mae: 16.0943
Epoch 2/5
1714/1714 [=====] - 24s 14ms/step - loss: 581.2591 - mae: 16.6714 - val_loss: 482.6658 - val_mae: 14.6322
Epoch 3/5
1714/1714 [=====] - 24s 14ms/step - loss: 490.1191 - mae: 15.1249 - val_loss: 455.9258 - val_mae: 14.2108
Epoch 4/5
1714/1714 [=====] - 24s 14ms/step - loss: 439.0449 - mae: 14.2941 - val_loss: 439.0598 - val_mae: 13.8696
Epoch 5/5
1714/1714 [=====] - 24s 14ms/step - loss: 404.6346 - mae: 13.6278 - val_loss: 442.9078 - val_mae: 13.9808

```
In [50]: 1 model.summary()
```

input_1 (InputLayer)	(None, 11)	0
embedding_1 (Embedding)	(None, 41, 30)	982620
lstm (LSTM)	(None, 41, 35)	9240
global_max_pooling1d (Global	(None, 35)	0
dense_3 (Dense)	(None, 62)	2232
dense_4 (Dense)	(None, 32)	2016
dropout_1 (Dropout)	(None, 32)	0
dense_5 (Dense)	(None, 1)	33
=====		
Total params: 996,141		
Trainable params: 996,141		
Non-trainable params: 0		

```
In [51]: 1 plot_model(model, show_shapes=True, to_file='images/RNN_LSTM_arc.png')
```

Out[51]:



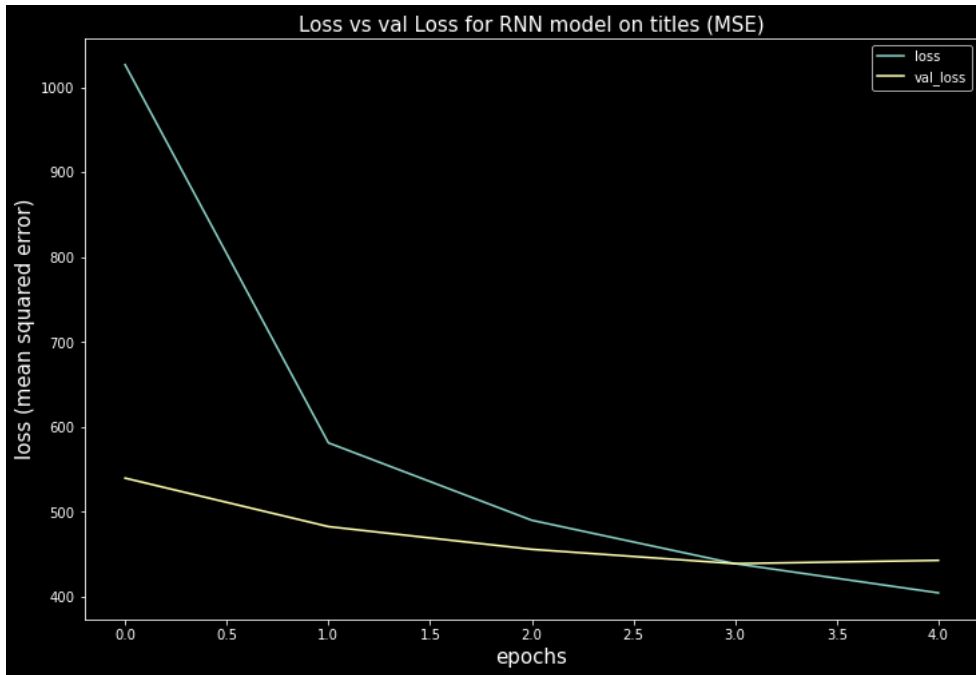
```
In [52]: 1 pred=model.predict(data_test)
```

```
In [53]: 1 test_results = model.evaluate(data_test, Y_test)

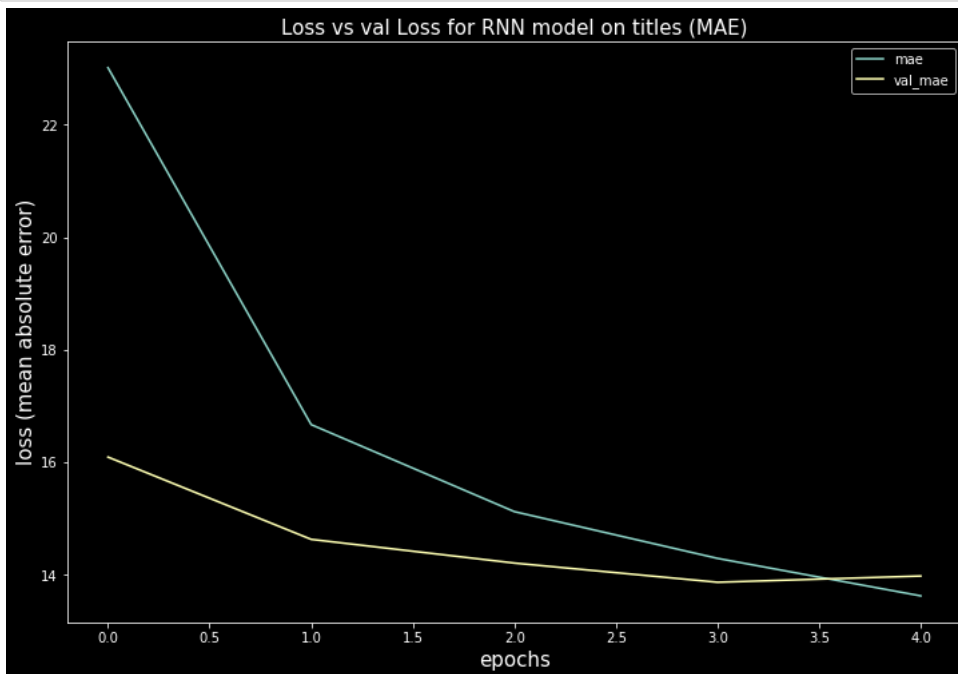
368/368 [=====] - 1s 3ms/step - loss: 436.1495 - mae: 13.9282
```

```
In [54]: 1 RMSE = np.sqrt(test_results[0])
```

```
In [55]: 1 fig = plt.subplots(figsize=(12,8))
2 plt.plot(r.history['loss'], label='loss')
3 plt.plot(r.history['val_loss'], label='val_loss')
4 plt.title("Loss vs val Loss for RNN model on titles (MSE)", fontsize=15)
5 plt.xlabel("epochs", fontsize=15)
6 plt.ylabel("loss (mean squared error)", fontsize=15)
7 plt.legend()
8 plt.savefig('images/MSE_LSTM_relu.png');
```



```
In [56]: 1 fig = plt.subplots(figsize=(12,8))
2 plt.plot(r.history['mae'], label='mae')
3 plt.plot(r.history['val_mae'], label='val_mae')
4 plt.title("Loss vs val Loss for RNN model on titles (MAE)", fontsize=15)
5 plt.xlabel("epochs", fontsize=15)
6 plt.ylabel("loss (mean absolute error)", fontsize=15)
7 plt.legend()
8 plt.savefig('images/MAE_LSTM_relu.png');
```



```
In [57]: 1 test_mae = mean_absolute_error(Y_test, pred)
```

```
In [58]: 1 string_score = f'\nMAE on training set: ${test_mae:.2f}'
2 string_score += f'\nMAE on training set: ${RMSE:.2f}'
3 fig, ax = plt.subplots(figsize=(12, 8))
4 plt.scatter(Y_test, pred)
5 ax.plot([0, 1], [0, 1], transform=ax.transAxes, ls="--", c="red")
6 plt.text(5, 135, string_score)
7 plt.title('Regression Model for Predicting Resale Value')
8 plt.ylabel('Model predictions for Resale Value($US)')
9 plt.xlabel('True Values for Resale Value($US)')
10 plt.savefig("images/regression_LSTM_relu.png")
```



CNN Titles

```
In [59]: 1 # Create the CNN model
2
3 # We get to choose embedding dimensionality
4 D = 256
5
6
7
8 i = Input(shape=(T,))
9 x = Embedding(V + 1, D)(i)
10 x = Conv1D(32, 3, activation='relu')(x)
11 x = MaxPooling1D(3)(x)
12 x = Conv1D(64, 3, activation='relu')(x)
13 x = MaxPooling1D(3)(x)
14 x = Conv1D(128, 3, activation='relu')(x)
15 x = GlobalMaxPooling1D()(x)
16 x = Dense(1, activation='relu')(x)
17
18 model = Model(i, x)
```

```
In [60]: 1 # Compile and fit
2 model.compile(
3     loss='MSE',
4     optimizer='adam',
5     metrics=['mae']
6 )
7
8
9 print('Training model...')
10 r = model.fit(
11     data_train,
12     Ytrain,
13     epochs=5,
14     validation_data=(data_val, Y_val)
15 )
16
```

Training model...

Epoch 1/5

1714/1714 [=====] - 79s 46ms/step - loss: 702.6185 - mae: 18.8140 - val_loss: 554.2573 - val_mae: 16.6334

Epoch 2/5

1714/1714 [=====] - 83s 48ms/step - loss: 450.1896 - mae: 14.6215 - val_loss: 525.1744 - val_mae: 15.6066

Epoch 3/5

1714/1714 [=====] - 79s 46ms/step - loss: 352.3939 - mae: 12.7084 - val_loss: 517.3931 - val_mae: 15.5277

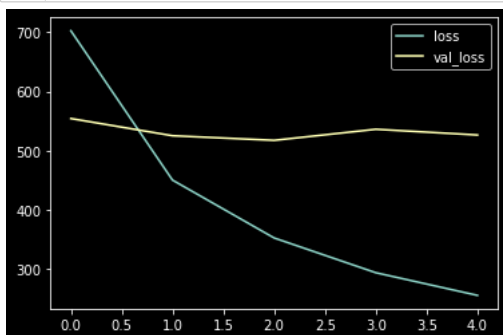
Epoch 4/5

1714/1714 [=====] - 79s 46ms/step - loss: 293.6541 - mae: 11.4138 - val_loss: 536.0966 - val_mae: 15.3805

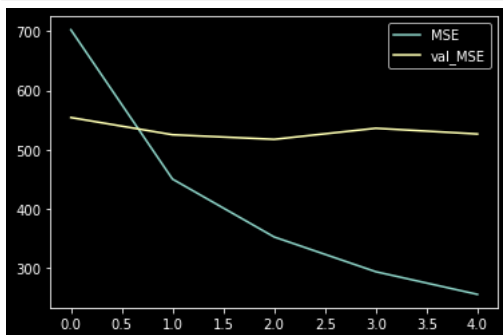
Epoch 5/5

1714/1714 [=====] - 80s 47ms/step - loss: 255.3643 - mae: 10.5055 - val_loss: 526.4434 - val_mae: 15.5761

```
In [61]: 1 # Plot loss per iteration
2 import matplotlib.pyplot as plt
3 plt.plot(r.history['loss'], label='loss')
4 plt.plot(r.history['val_loss'], label='val_loss')
5 plt.legend();
```



```
In [64]: 1 # Plot accuracy per iteration
2 plt.plot(r.history['loss'], label='MSE')
3 plt.plot(r.history['val_loss'], label='val_MSE')
4 plt.legend();
```



CNN using images as input

```
In [65]: 1 df_imgs = df.drop(['title', 'url',
2                        'date_sold', 'profit',
3                        'ROI', 'brand', 'cost',
4                        'pictureURLLarge'],
5                        axis=1).copy()
```

```
In [66]: 1 df_imgs.dropna(subset=['Image'], inplace=True)
```

```
In [67]: 1 df_imgs.reset_index(drop=True, inplace=True)
```

```
In [68]: 1 df_imgs['file_index'] = df_imgs.index.values
2 df_imgs['file_index'] = df_imgs['file_index'].astype(str)
```

```
In [69]: 1 df_imgs['filename'] = df_imgs['file_index'] + '.jpg'
```

```
In [70]: 1 def download(row):
2         filename = row.filepath
3
4         # create folder if it doesn't exist
5         # os.makedirs(os.path.dirname(filename), exist_ok=True)
6
7         url = row.Image
8         # print(f"Downloading {url} to {filename}")
9
10        try:
11            r = requests.get(url, allow_redirects=True)
12            with open(filename, 'wb') as f:
13                f.write(r.content)
14        except:
15            print(f'{filename} error')
```

```
In [71]: 1 root_folder = 'C:/Users/12108/Documents/GitHub/Neural_Network_Predicting_Reseller_Success_Ebay/nn_images/'
2 df_imgs['filepath'] = root_folder + df_imgs['filename']
```

```
In [72]: 1 df_imgs['filepath'].sample(2).apply(print)
```

```
C:/Users/12108/Documents/GitHub/Neural_Network_Predicting_Reseller_Success_Ebay/nn_images/74992.jpg
C:/Users/12108/Documents/GitHub/Neural_Network_Predicting_Reseller_Success_Ebay/nn_images/12140.jpg
```

```
Out[72]: 74992    None
12140    None
Name: filepath, dtype: object
```

```
In [73]: 1 # df_imgs.apply(download, axis=1)
```

All image files are stored locally for this project. The below markdown code is for reference.

```
img_list = os.listdir('C:/Users/12108/Documents/GitHub/Neural_Network_Predicting_Reseller_Success_Ebay/nn_images/')

img_df = df_imgs.loc[df_imgs['filename'].isin(img_list)].copy()

img_df.reset_index(drop=True, inplace=True)

img_df.rename({'Image': 'data',
              'converted_price': 'labels'},
              axis=1, inplace=True)
```

```
df_train, df_test, Ytrain, Ytest = train_test_split(img_df, Y, test_size=0.20)
datagen=ImageDataGenerator(rescale=1./255.,validation_split=0.20)

train_generator=datagen.flow_from_dataframe(
dataframe=df_train,
directory= None,
x_col="filepath",
y_col="labels",
subset="training",
batch_size=100,
seed=55,
shuffle=True,
class_mode="raw")

valid_generator=datagen.flow_from_dataframe(
dataframe=df_train,
directory=None,
x_col="filepath",
y_col="labels",
subset="validation",
batch_size=100,
seed=55,
shuffle=True,
class_mode="raw")

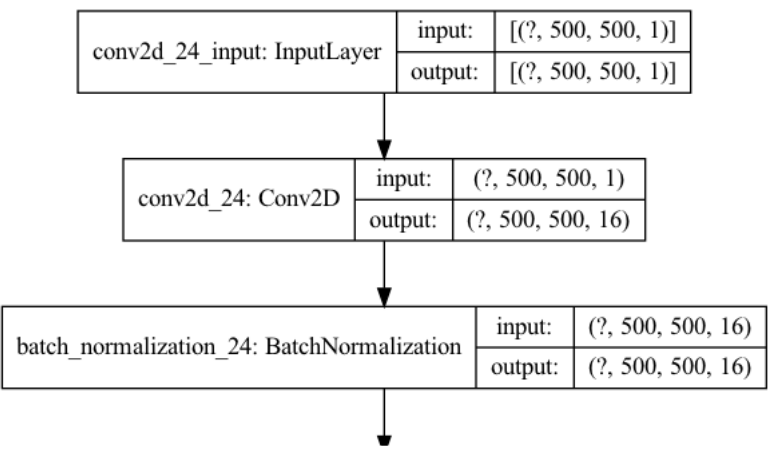
test_datagen=ImageDataGenerator(rescale=1./255.)
test_generator=test_datagen.flow_from_dataframe(
dataframe=df_test,
directory=None,
x_col="filepath",
y_col="labels",
batch_size=100,
seed=55,
shuffle=False,
class_mode="raw")
```

```
In [74]: 1 # model = models.Sequential()
2
3 # model.add(layers.Conv2D(16, (3, 3), padding='same', activation='relu',
4 #                       input_shape=(256, 256, 3)))
5 # model.add(layers.BatchNormalization())
6 # model.add(layers.Conv2D(16, (3, 3), activation='relu', padding='same'))
7 # model.add(layers.BatchNormalization())
8 # model.add(layers.MaxPooling2D((2, 2)))
9
10 # model.add(layers.Conv2D(32, (3, 3), padding='same', activation='relu',
11 #                       input_shape=(256, 256, 3)))
12 # model.add(layers.BatchNormalization())
13 # model.add(layers.Conv2D(32, (3, 3), activation='relu', padding='same'))
14 # model.add(layers.BatchNormalization())
15 # model.add(layers.MaxPooling2D((2, 2)))
16
17 # model.add(layers.Conv2D(64, (3, 3), activation='relu', padding='same'))
18 # model.add(layers.BatchNormalization())
19 # model.add(layers.Conv2D(64, (3, 3), activation='relu', padding='same'))
20 # model.add(layers.BatchNormalization())
21 # model.add(layers.MaxPooling2D((2, 2)))
22
23 # model.add(layers.Conv2D(128, (3, 3), activation='relu', padding='same'))
24 # model.add(layers.BatchNormalization())
25 # model.add(layers.Conv2D(128, (3, 3), activation='relu', padding='same'))
26 # model.add(layers.BatchNormalization())
27 # model.add(layers.MaxPooling2D((2, 2)))
28
29 # model.add(layers.Flatten())
30
31 # model.add(Dense(512, activation='relu'))
32 # model.add(Dropout(0.1))
33 # model.add(Dense(256, activation='relu'))
34 # model.add(Dropout(0.1))
35 # model.add(Dense(128, activation='relu'))
36 # model.add(Dense(1, activation='linear'))
37
38 # model.compile(loss='MSE',
39 #               optimizer='Adam',
40 #               metrics=['mae', 'mse'])
41
42 # summary = model.fit(train_generator, epochs=3, validation_data=valid_generator)
```

```
In [75]: 1 model = tf.keras.models.load_model('cnn_grayscale_relu1.h5', compile=False)
```

```
In [76]: 1 plot_model(model, show_shapes=True, to_file="images/CNN_architecture.png")
```

Out[76]:




```
In [77]: 1 model.summary()
```

Model: "sequential_4"

Layer (type)	Output Shape	Param #
conv2d_24 (Conv2D)	(None, 500, 500, 16)	160
batch_normalization_24 (Batch Normalization)	(None, 500, 500, 16)	64
max_pooling2d_16 (MaxPooling2D)	(None, 250, 250, 16)	0
conv2d_25 (Conv2D)	(None, 250, 250, 32)	4640
batch_normalization_25 (Batch Normalization)	(None, 250, 250, 32)	128
max_pooling2d_17 (MaxPooling2D)	(None, 125, 125, 32)	0
conv2d_26 (Conv2D)	(None, 125, 125, 64)	18496
batch_normalization_26 (Batch Normalization)	(None, 125, 125, 64)	256
max_pooling2d_18 (MaxPooling2D)	(None, 62, 62, 64)	0
conv2d_27 (Conv2D)	(None, 62, 62, 128)	73856
batch_normalization_27 (Batch Normalization)	(None, 62, 62, 128)	512
max_pooling2d_19 (MaxPooling2D)	(None, 31, 31, 128)	0
flatten_4 (Flatten)	(None, 123008)	0
dense_12 (Dense)	(None, 512)	62980608
dense_13 (Dense)	(None, 128)	65664
dense_14 (Dense)	(None, 1)	129
Total params: 63,144,513		
Trainable params: 63,144,033		
Non-trainable params: 480		

```
In [ ]: 1 model.evaluate(valid_generator)
```

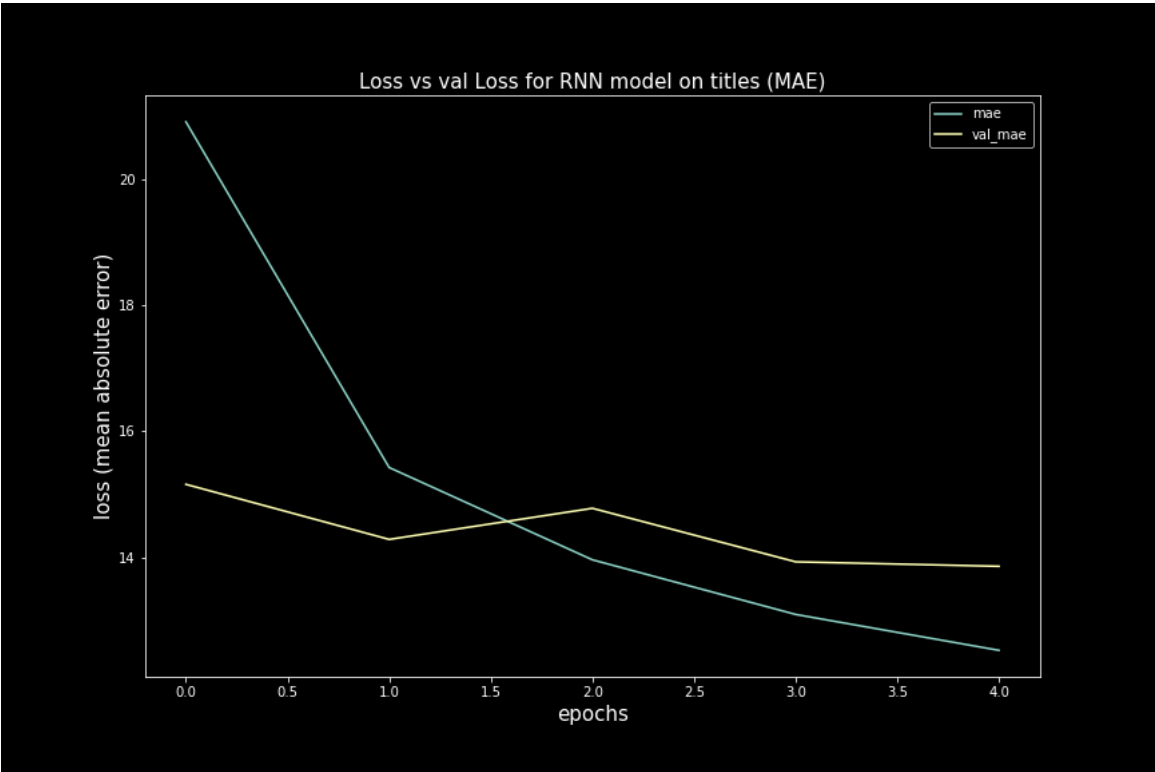
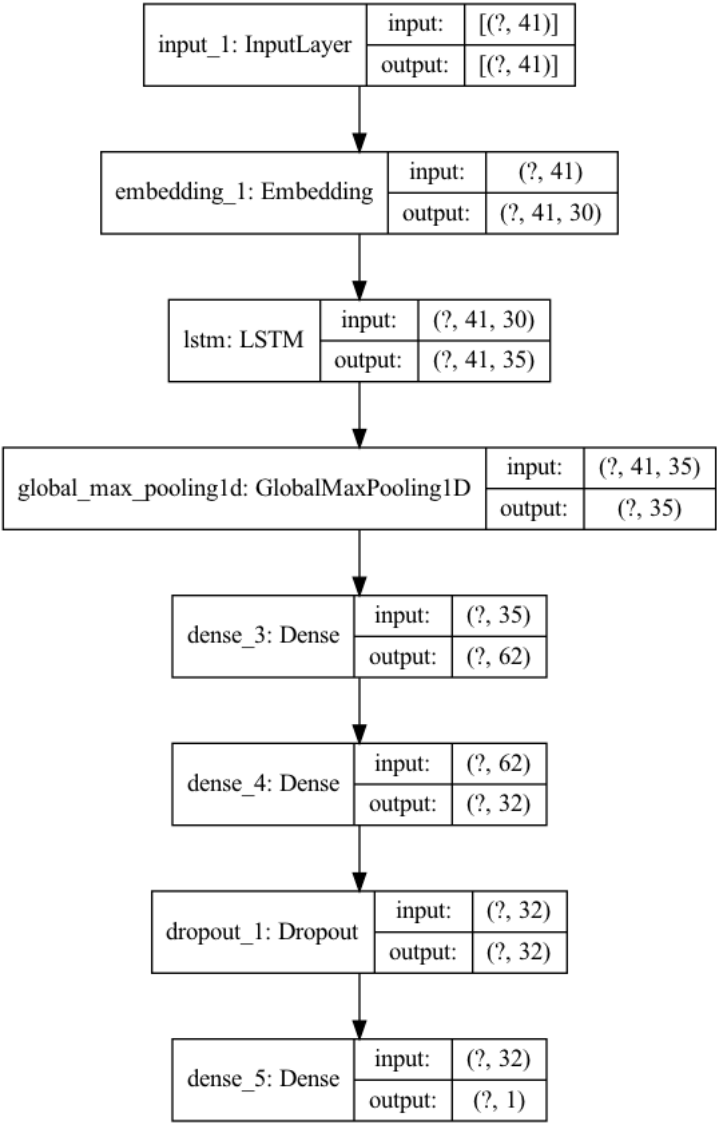
```
In [ ]: 1 test_generator.reset()
2 pred=model.predict(test_generator,verbose=1)
```

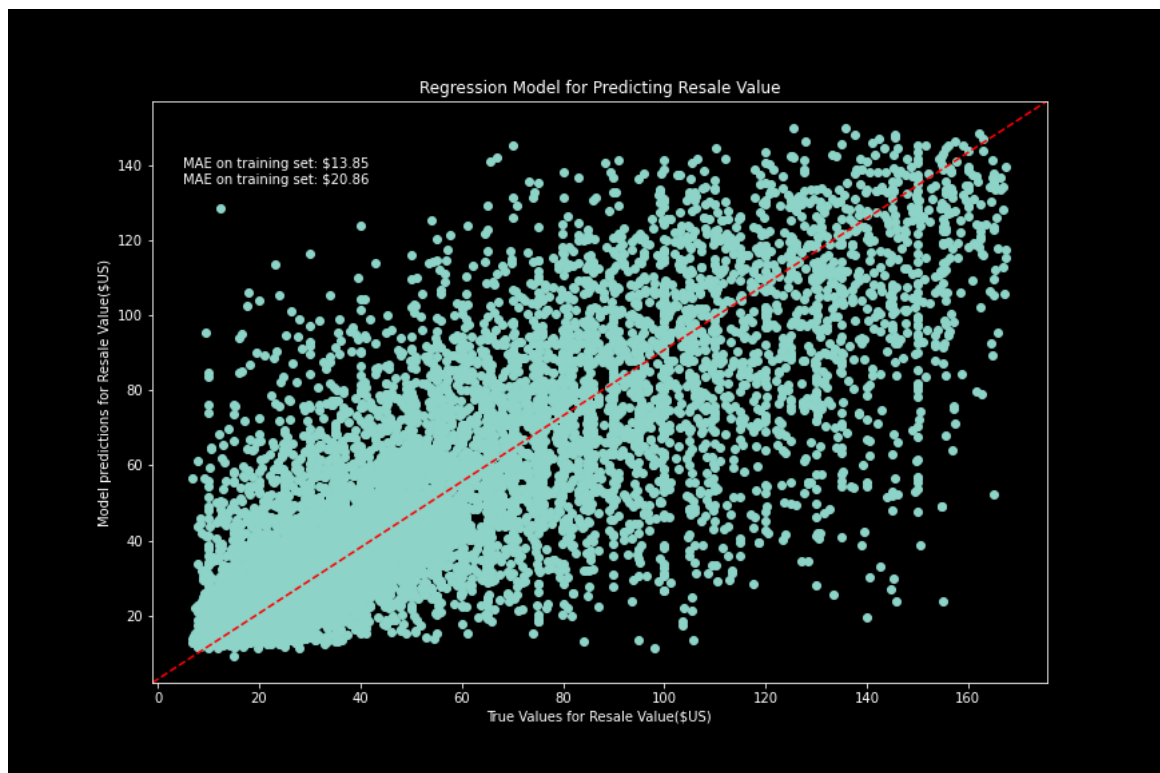
```
In [ ]: 1 test_results = model.evaluate(test_generator)
```

```
In [ ]: 1 fig = plt.figure(figsize=(12,8))
2 plt.plot(summary.history['loss'])
3 plt.plot(summary.history['val_loss'])
4 plt.plot
5 plt.title('model loss')
6 plt.ylabel('loss(mean absolute error)')
7 plt.xlabel('epoch')
8 plt.legend(['train_loss', 'val_loss'], loc='upper right')
9 plt.show();
```

Results

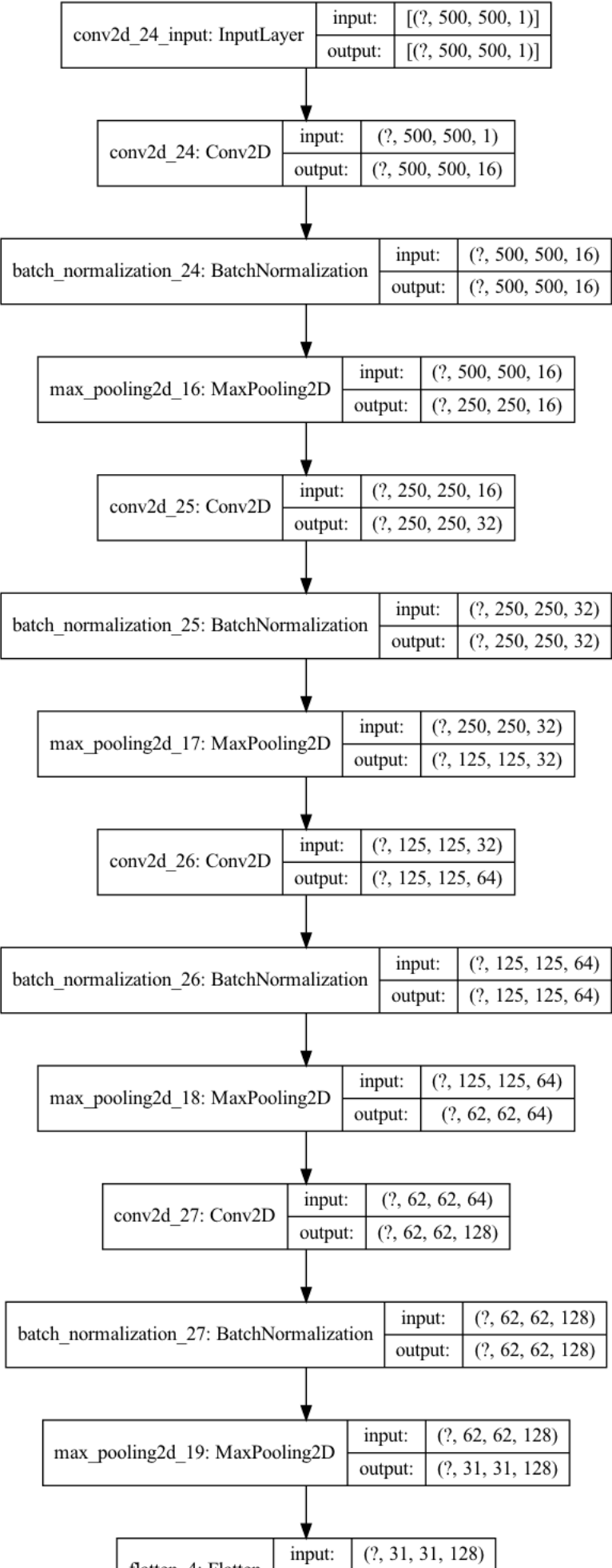
Recurrent Neural Network (Long Short Term Memory)

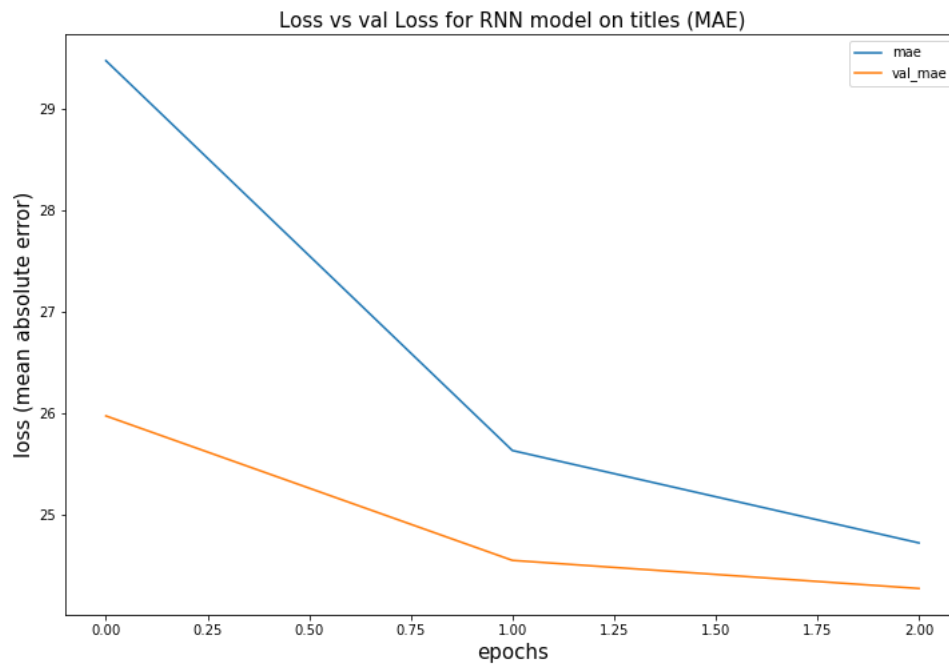
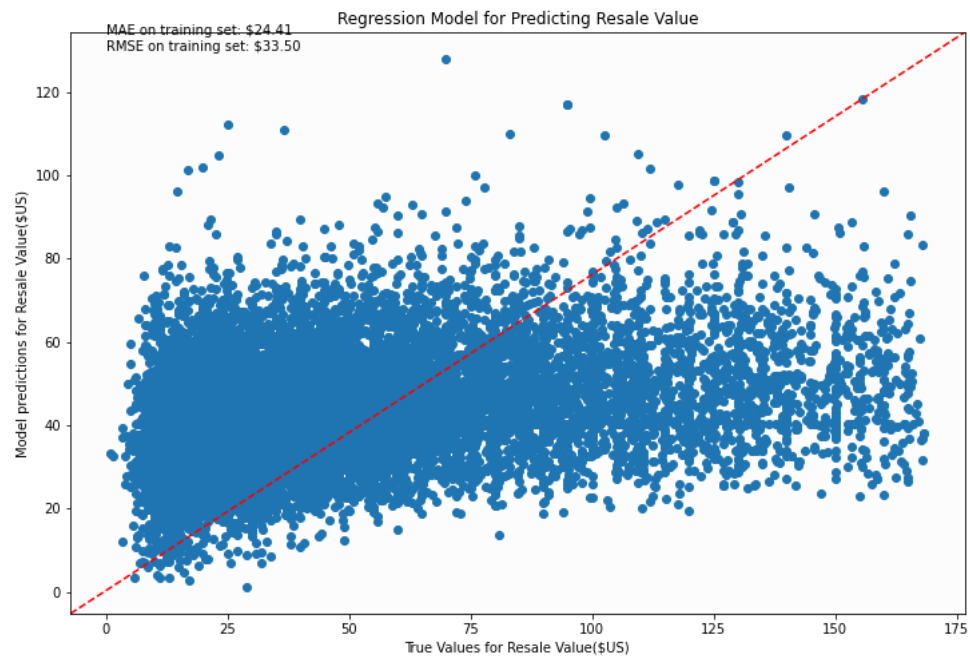




- The mean price of the 8 brands of knives sold on ebay is around 50.00. *A mean absolute error of about plus or minus 13.80 is acceptable.*

Convolutud Neural Network on Grayscale Images





- The MAE when testing the CNN was roughly \$25.00. That is an error of plus or minus about 50% of the mean price of knives sold. Not acceptable yet as compared to the RNN with titles. Will address in future work.

Future Work

- Expand data to include other products readily purchasable at the Surplus Store.
- Attempt data augmentation on the CNN image network
- Attempt to obtain more aspect data for sold knives. Some important aspect data is limited access to sellers who average a certain amount of money per month.

```
In [ ]: 1 ## define two sets of inputs
2 # inputA = Input(shape=(32,))
3 # inputB = Input(shape=(128,))
4 ## the first branch operates on the first input
5 # x = Dense(8, activation="relu")(inputA)
6 # x = Dense(4, activation="relu")(x)
7 # x = Model(inputs=inputA, outputs=x)
8 ## the second branch operates on the second input
9 # y = Dense(64, activation="relu")(inputB)
10 # y = Dense(32, activation="relu")(y)
11 # y = Dense(4, activation="relu")(y)
12 # y = Model(inputs=inputB, outputs=y)
13 ## combine the output of the two branches
14 # combined = concatenate([x.output, y.output])
15 ## apply a FC layer and then a regression prediction on the
16 ## combined outputs
17 # z = Dense(2, activation="relu")(combined)
18 # z = Dense(1, activation="linear")(z)
19 ## our model will accept the inputs of the two branches and
20 ## then output a single value
21 # model = Model(inputs=[x.input, y.input], outputs=z)
```

```
In [ ]: 1
```