

[Open in app](#)

Medium



Search



Write



# Architecture : Combination of tabular, textual and image inputs to predict house prices.

Dave Cote, M.Sc. · [Follow](#)

10 min read · Nov 19, 2022

126

2



Modality refers to a certain type of information and/or the representation format in which information is stored.

According to wikipedia, « *Multimodal learning attempts to model the combination of different modalities of data, often rising in real-world applications. An example of joint data is combining text (typically represented as discrete word count vectors) with imaging data consisting of pixel intensities and annotation tags. As these modalities have fundamentally different statistical properties, combining them is non-trivial, which is why specialized modelling strategies and algorithms are required.*

Can we simultaneously train structured (tabular) and unstructured (text, image) data in the same neural network model while optimizing the same output target ?

We'll perform this experiment using a very simple example by implementing a « concat » combination method that will concatenate text features, tabular features, and image features all at once before our final regressor layer(s).

Specifically, we will build a multimodal neural network that contains 3 input heads :

- **Image** : Inception V3 input layers (Transfer Learning from Imagenet)
- **Text** : Bidirectionnal LSTM with Attention layers pretrained on GLOVe embeddings
- **Structured** : Multi-layers perceptron

And 1 output head :

- **House Price**

Of course, more advanced methods and strategies exist to solve multimodal problems. For example, the great HuggingFace offers a [multimodal-toolkit](#) that you can experiment with.

## Table of contents

- Dataset
- Tabular data engineering
- Text cleaning
- Image cleaning — Fine-tune VGG16 to predict good/bad images
- Architecture#1 — Structured Data fine-tuning with MLP (Python code)
- Architecture#2 — Text Data fine-tuning using GlOve BiLSTM with Attention (Python code)
- Architecture#3 — Image Data fine-tuning with Inception V3 (Python code)
- Final Architecture — Load each previously fine-tuned weights and create the final hybrid neural networks training architecture by using the concat method. (Python Code)
- Results and Conclusion
- Final words

## Dataset

First, let's take a quick look at the dataset we're going to use. I searched a very long time to find a dataset containing images, text and tabular data and

finally found one! The link to download the dataset is here:

<https://www.kaggle.com/datasets/ericpierce/austinhousingprices> (You must have a Kaggle account to download it).

As it is written on the Kaggle page:

*“Context — Austin Housing*

*This data was originally sourced in support of my capstone project at Northwestern University. The Austin Housing market is one of the hottest markets in 2021, and these listings show how that market has changed over the past couple of years.”*

The data contains a csv file with all the tabular features, a text description feature and a column that specify the image filename associated with the dataset row. The csv dataset comes with a huge zipped file containing the house images. Unfortunatly, there is only one image available for each house but we will work with that.

The goal here is to predict house **price** using all the data we have.

## **Tabular data engineering**

Although it is said that the data has been cleaned, there are still many data quality issues left. For example, there are impossible feature values like houses with 0 bathrooms or 0 bedrooms or negative square foot of living area.

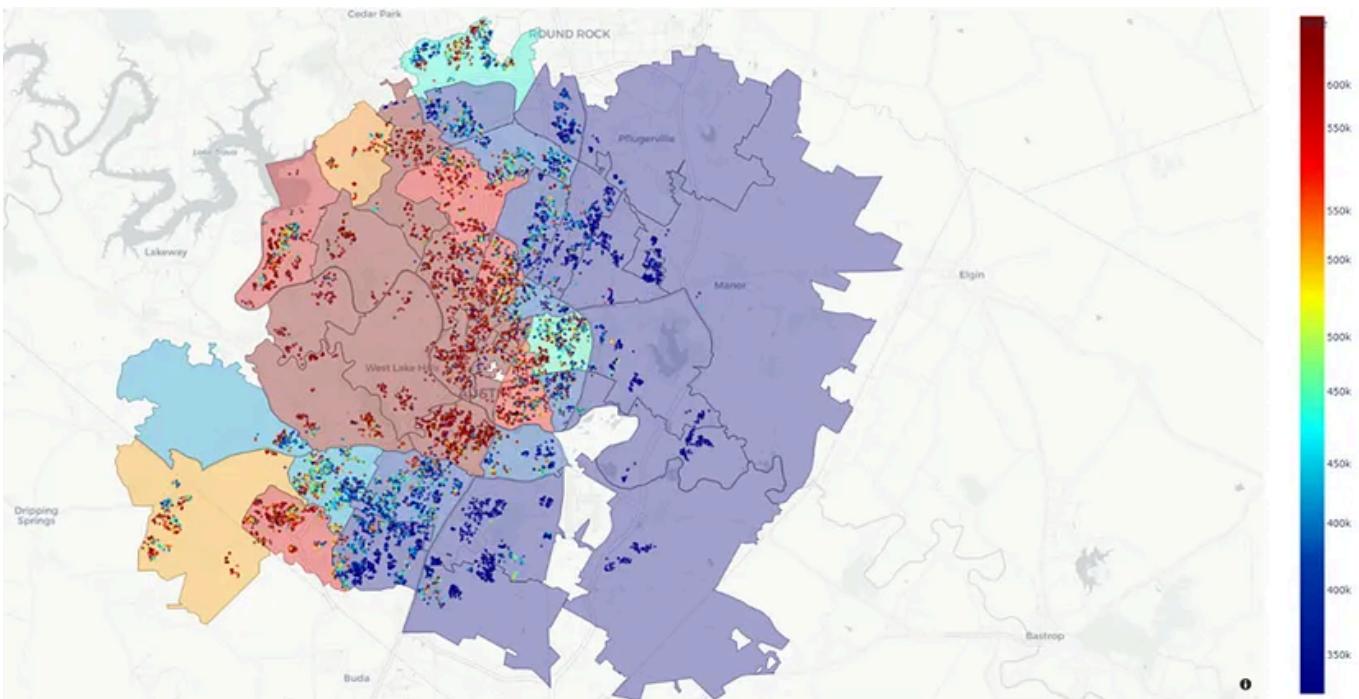
Here is a table that resume some of the major data transformations and cleaning steps performed:

Feature name	Quality problem / Opportunity	Transformation
numOfBathrooms	Impossible to have 0 bathrooms	> 1989= 2, <= 1989 = 1
numOfBedrooms	Impossible to have 0 bedroms	> 1989= 2.5, <= 1989 = 2
garageSpaces	More than 3 garageSpaces and price < 1M\$ and homeType = Single Family	max = 3
parkingSpaces	More than 3 parkingSpaces and price < 1M\$ and homeType = Single Family	max = 3
lotSizeSqFt	lotSizeSqFt <= 0, outliers > IQR*1.6, outliers < IQR*1.6	delete
livingAreaSqFt	livingAreaSqFt <= 0, outliers > IQR*1.6, outliers < IQR*1.6	delete
latestPrice (target)	Concept Drift (Price Inflation)	Adjust past price with inflation
latitude	Outliers < 30.12	delete
zipcode	Outliers outside Austin	delete
numOfAccessibilityFeatures	Few when > 0	convert to binary, > 0
numOfPatioAndPorchFeatures	Few when > 0	convert to binary, > 0
numOfSecurityFeatures	Few when > 0	convert to binary, > 0
numOfWaterfrontFeatures	Few when > 0	convert to binary, > 0
numOfWindowFeatures	Few when > 0	convert to binary, > 0
numOfCommunityFeatures	Few when > 0	convert to binary, > 0
streetAddress	Unexploitable text format	Extract type of address (drive, lane, street, cove, road, avenue, trail, circle, court, way, other)
description	Dirty text description	Clean text (normalization, stop words, lemmatization...)
longitude, latitude	Opportunity: Geospatial clusters	Create clusters with Kmeans
zipcode	Opportunity: Census data	Map census data (crime_reports, median_household_income(\$), population_below_poverty_level(%), unemployment(%)...)

An important step is to clean up our target value, **house prices**, using inflation rates (concept drift problem). Past selling prices do not match actual selling prices due to inflation; the definition of 1\$ unit is not the same over time. For this purpose, I used a table of inflation rates to correct the price amount.

Once cleaned, our dataset contains **56 features** and **10537 rows**.

Here is a geospatial plot (plotly) of all Austin Houses we have:



We can see that there are some zipcode clusters where the price is higher (red) than other regions (blue). I've created a new feature that represent the zipcode clusters called **regionCluster** using the kmeans algorithm:

```

1  from sklearn.cluster import KMeans
2
3  kmeans = KMeans(n_clusters=25, random_state=0).fit(data[['longitude', 'latitude']])
4  data['regionCluster'] = kmeans.predict(data[['longitude', 'latitude']])

```

creating\_region\_cluster.py hosted with ❤ by GitHub

[view raw](#)

## Text cleaning

To clean up house text descriptions, we apply normalization, remove stop words and find the root of each word using lemmatization.

Here is an example of the result:

```
data_text_raw_vs_clean[data_text_raw_vs_clean.index == 0]['description'].values[0]
```

"Absolutely GORGEOUS 4 Bedroom home with 2 full baths located in the tucked away Pflugerville community The Lakes at Northtown. This home was recently updated with new carpets, new tile kitchen back splash, and fresh paint throughout the entire house. This charming abode is a MUST SEE if you're looking for an affordable home in the Austin area. Something very hard to come by in this market! Only 14 miles from downtown, this home will accommodate downtown commuting as well as a work from home environment. This community is tucked away and yet only minutes from major retailers and shopping. Do not miss this opportunity!"

```
data_text_raw_vs_clean[data_text_raw_vs_clean.index == 0]['description_clean'].values[0]
```

' absolutely gorgeous locate tuck away community lake northtown recently update new carpet new tile kitchen splash fresh paint entire house charming abode yoube look affordable area hard come market mile downtown accommodate downtown commuting as as work environment community tuck away minute major retailer shopping miss opportunity'

## Image cleaning — Fine-tune VGG16 to predict good/bad images

This one is a bit tricky. After randomly exploring some images, there seem to be very different / heterogenous types of house images. Here are a few:



Nice picture!



Nice river but... which house are we talking about?

Sorry, we have no imagery here.

Not very useful...



Interior picture...



Would have a fresh beer on that patio but not very useful...



**What is this!!? A cemetery? ...**

As you can see, there are several heterogeneous types of house images. Another challenge is that we only have one image available per house, so the image must be useful.

Since there are about 10,000 images, we can't look at all of them to decide which data to keep and which data to delete. To help us, we will fine-tune a **VGG16** to predict good/bad images for us.

First, we need to determine what is a good image for us. I decided to keep only images where you can see the complete architecture of the house from the outside:



Bad images would be any interior image or any image where we cannot see the full architecture of the house from the outside.

As we are using **transfer learning**, fewer samples from each class should do the job.

We need to create two directories, one for the good images and one for the bad ones, then manually dispatch some of them in the right place. After an hour, I sorted around **800 images** for the train set (400 good and 400 bad) and around **250 images** for the test set.

We are now ready to fine-tune our VGG16 model!

Python Code:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 import os
6 import glob
7 import datetime
8 import pickle
9 # keras imports
10 from keras.applications.vgg16 import VGG16, preprocess_input
11 from keras.preprocessing import image
12 from keras.models import Model
13 from keras.layers import Layer
14 from keras import backend as K
15 # other imports
16 from sklearn.preprocessing import LabelEncoder
17 from sklearn.model_selection import train_test_split
18 from sklearn.linear_model import LogisticRegression
19 from sklearn.metrics import confusion_matrix
20 from sklearn.metrics import accuracy_score
21 import numpy as np
22 import h5py
23 import json
24 import time
25
26 # filter warnings
27 import warnings
28 warnings.simplefilter(action="ignore", category=FutureWarning)
29
30 def make_confusion_matrix(cf,
31                         group_names=None,
32                         categories='auto',
33                         count=True,
34                         percent=True,
35                         cbar=True,
36                         xyticks=True,
37                         xyplotlabels=True,
38                         sum_stats=True,
39                         figsize=None,
40                         cmap='Blues',
41                         title=None):
42     ...
43     This function will make a pretty plot of an sklearn Confusion Matrix cm using a Seaborn heat
44     Arguments
45     -----
```

```

46     cf:           confusion matrix to be passed in
47     group_names: List of strings that represent the labels row by row to be shown in each square
48     categories: List of strings containing the categories to be displayed on the x,y axis. Default is None.
49     count:        If True, show the raw number in the confusion matrix. Default is False.
50     normalize:   If True, show the proportions for each category. Default is False.
51     cbar:         If True, show the color bar. The cbar values are based off the values in the confusion matrix. Default is False.
52
53     xyticks:     If True, show x and y ticks. Default is True.
54     xyplotlabels: If True, show 'True Label' and 'Predicted Label' on the figure. Default is False.
55     sum_stats:   If True, display summary statistics below the figure. Default is False.
56     figsize:    Tuple representing the figure size. Default will be the matplotlib rcParams values.
57     cmap:        Colormap of the values displayed from matplotlib.pyplot.cm. Default is 'Blues'
58
59
60     title:       Title for the heatmap. Default is None.
61     ...
62
63
64 # CODE TO GENERATE TEXT INSIDE EACH SQUARE
65 blanks = ['' for i in range(cf.size)]
66
67 if group_names and len(group_names)==cf.size:
68     group_labels = ["{}\n".format(value) for value in group_names]
69 else:
70     group_labels = blanks
71
72 if count:
73     group_counts = ["{0:0.0f}\n".format(value) for value in cf.flatten()]
74 else:
75     group_counts = blanks
76
77 if percent:
78     group_percentages = ["{0:.2%}".format(value) for value in cf.flatten()/np.sum(cf)]
79 else:
80     group_percentages = blanks
81
82 box_labels = [f"{v1}{v2}{v3}".strip() for v1, v2, v3 in zip(group_labels,group_counts,group_percentages)]
83 box_labels = np.asarray(box_labels).reshape(cf.shape[0],cf.shape[1])
84
85
86 # CODE TO GENERATE SUMMARY STATISTICS & TEXT FOR SUMMARY STATS
87 if sum_stats:
88     #Accuracy is sum of diagonal divided by total observations
89     accuracy = np.trace(cf) / float(np.sum(cf))

```

```
90
91     #if it is a binary confusion matrix, show some more stats
92     if len(cf)==2:
93         #Metrics for Binary Confusion Matrices
94         precision = cf[1,1] / sum(cf[:,1])
95         recall    = cf[1,1] / sum(cf[1,:])
96         f1_score  = 2*precision*recall / (precision + recall)
97         stats_text = "\n\nAccuracy={:0.3f}\nPrecision={:0.3f}\nRecall={:0.3f}\nF1 Score={:0.
98             accuracy,precision,recall,f1_score)
99     else:
100        stats_text = "\n\nAccuracy={:0.3f}".format(accuracy)
101    else:
102        stats_text = ""
103
104
105    # SET FIGURE PARAMETERS ACCORDING TO OTHER ARGUMENTS
106    if figsize==None:
107        #Get default figure size if not set
108        figsize = plt.rcParams.get('figure.figsize')
109
110    if xyticks==False:
111        #Do not show categories if xyticks is False
112        categories=False
113
114
115    # MAKE THE HEATMAP VISUALIZATION
116    plt.figure(figsize=figsize)
117    sns.heatmap(cf,annot=box_labels,fmt="",cmap=cmap,cbar=cbar,xticklabels=categories,yticklabe]
118
119    if xyplotlabels:
120        plt.ylabel('True label')
121        plt.xlabel('Predicted label' + stats_text)
122    else:
123        plt.xlabel(stats_text)
124
125    if title:
126        plt.title(title)
127
128
129
130
131 config={
132     "model"          : "vgg16",
133     "weights"        : "imagenet",
134     "include_top"    : False,
```

```
135
136     "train_path"      : "/data/classifier_cleaner/train_house/",
137     "test_path"        : "/data/classifier_cleaner/test_house/",
138     "features_path"   : "/data/classifier_cleaner/features.h5",
139     "labels_path"      : "/data/classifier_cleaner/labels.h5",
140     "results"          : "/data/classifier_cleaner/results.txt",
141     "classifier_path" : "/data/classifier_cleaner/classifier.pickle",
142     "model_path"       : "/data/classifier_cleaner/model",
143     "label_encoder_path": "/data/classifier_cleaner/label_encoder.pickle",
144
145     "test_size"        : 0.20,
146     "seed"              : 9,
147     "num_classes"      : 2
148 }
149
150
151
152 # config variables
153 model_name    = config["model"]
154 weights       = config["weights"]
155 include_top   = config["include_top"]
156 train_path    = config["train_path"]
157 features_path = config["features_path"]
158 labels_path   = config["labels_path"]
159 test_size     = config["test_size"]
160 results       = config["results"]
161 model_path    = config["model_path"]
162 label_encoder_path = config["label_encoder_path"]
163
164 # Here we are loading the base VGG16 model with weights and then excluding the top dense layer
165 if model_name == "vgg16":
166     base_model = VGG16(weights=weights), include_top=True)
167     model = Model(base_model.input, base_model.get_layer('fc1').output)
168     image_size = (224, 224)
169 else:
170     base_model = None
171
172 print ("Successfully loaded base model and model...")
173
174 train_labels = os.listdir(train_path)
175 train_labels
176
177 # path to training dataset
178 #train_labels = os.listdir(train_path)
```

```
179     train_labels = ['good', 'bad']
180
181     # encode the labels
182     print ("Encoding labels...")
183     le = LabelEncoder()
184     le.fit([tl for tl in train_labels])
185
186     # variables to hold features and labels
187     features = []
188     labels    = []
189
190
191     #print(file_names)
192     # loop over all the labels in the folder
193     count = 1
194     for i, label in enumerate(train_labels):
195         cur_path = train_path + label + '/'
196         relevant_path = cur_path
197         included_extensions = ['jpg']
198         file_names = [fn for fn in os.listdir(relevant_path)
199                         if any(fn.endswith(ext) for ext in included_extensions)]
200         #print(cur_path)
201         #print(cur_path)
202         #print(len(labels))
203         count = 1
204         for image_path in file_names:
205
206             img = image.load_img(train_path + label + '/' + image_path, target_size=image_size)
207             x = image.img_to_array(img)
208             #print(x)
209             x = np.expand_dims(x, axis=0)
210             x = preprocess_input(x)
211             feature = model.predict(x)
212             flat = feature.flatten()
213             features.append(flat)
214             labels.append(label)
215             print ("Processed - " + str(count))
216             count += 1
217             print ("Completed label - " + label)
218
219     # encode the labels using LabelEncoder
220     le = LabelEncoder()
221     le_labels = le.fit_transform(labels)
222
223     # get the shape of training labels
```

```
224     print ("Training labels: {}".format(le_labels))
225     print ("Training labels shape: {}".format(le_labels.shape))
226
227     model.summary()
228
229     # save features and labels
230     h5f_data = h5py.File(features_path, 'w')
231     h5f_data.create_dataset('dataset_1', data=np.array(features))
232
233     h5f_label = h5py.File(labels_path, 'w')
234     h5f_label.create_dataset('dataset_1', data=np.array(le_labels))
235
236     h5f_data.close()
237     h5f_label.close()
238
239     # save model and weights
240     model_json = model.to_json()
241     with open(model_path + str(test_size) + ".json", "w") as json_file:
242         json_file.write(model_json)
243
244     pickle.dump(le, open(label_encoder_path, 'wb'))
245
246     # save weights
247     model.save_weights(model_path + str(test_size) + ".h5")
248     print("Saved model and weights to disk..")
249
250     print ("Features and labels saved..")
251
252     # end time
253     end = time.time()
254     print ("End time - {}".format(datetime.datetime.now().strftime("%Y-%m-%d %H:%M")))
255
256
257
258
259     ##### PREDICT
260     config={
261         "model" : "vgg16",
262         "weights" : "imagenet",
263         "features_path" : "/data/classifier_cleaner/features.h5",
264         "labels_path" : "/data/classifier_cleaner/labels.h5",
265         "classifier_path" : "/data/classifier_cleaner/classifier.pickle",
266         "model_path" : "/data/classifier_cleaner/model",
267         "label_encoder_path" : "/data/classifier_cleaner/label_encoder.pickle",
268         "batch_size" : 32
```

```
269     "seed"      : 9,
270 }
271
272 # config variables
273 test_size = config["test_size"]
274 seed = config["seed"]
275 features_path = config["features_path"]
276 labels_path = config["labels_path"]
277 classifier_path = config["classifier_path"]
278
279 # import features and labels
280 h5f_data = h5py.File(features_path, 'r')
281 h5f_label = h5py.File(labels_path, 'r')
282
283 features_string = h5f_data['dataset_1']
284 labels_string = h5f_label['dataset_1']
285
286 features = np.array(features_string)
287 labels = np.array(labels_string)
288
289 h5f_data.close()
290 h5f_label.close()
291
292 # verify the shape of features and labels
293 print ("[INFO] features shape: {}".format(features.shape))
294 print ("[INFO] labels shape: {}".format(labels.shape))
295
296 print ("[INFO] training started...")
297 # split the training and testing data
298 (trainData, testData, trainLabels, testLabels) = train_test_split(np.array(features),
299                                         np.array(labels),
300                                         test_size=test_size,
301                                         random_state=seed)
302
303 print ("[INFO] splitted train and test data...")
304 print ("[INFO] train data  : {}".format(trainData.shape))
305 print ("[INFO] test data   : {}".format(testData.shape))
306 print ("[INFO] train labels: {}".format(trainLabels.shape))
307 print ("[INFO] test labels : {}".format(testLabels.shape))
308
309 # Use logistic regression as the model
310 print ("[INFO] creating model...")
311 model = LogisticRegression(max_iter=5500)
312
```

```

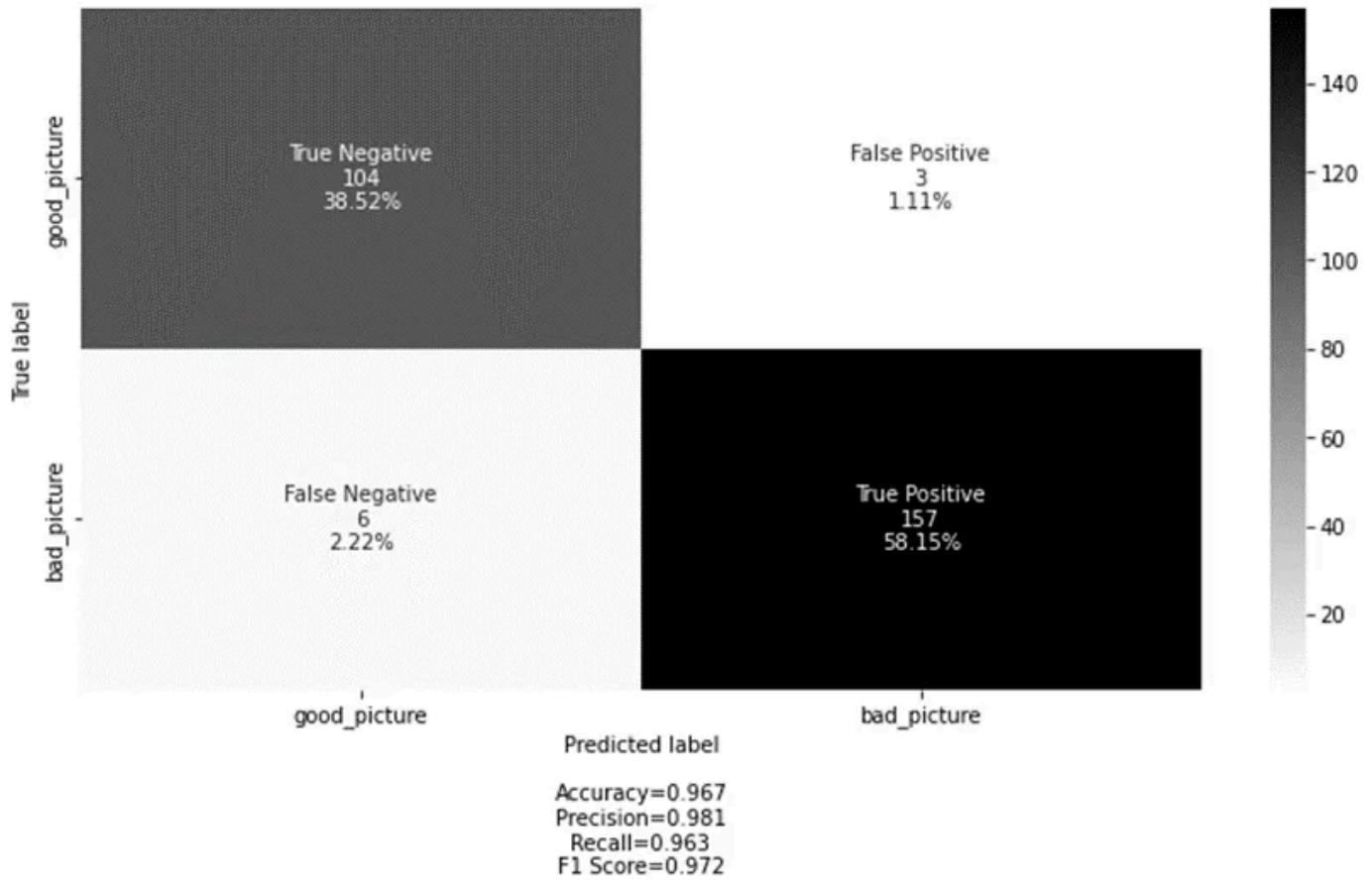
313 # evaluate the model of test data
314 preds = model.predict(testData)
315
316 # dump classifier to file
317 print ("[INFO] saving model...")
318 pickle.dump(model, open(classifier_path, 'wb'))
319
320 # plot the confusion matrix
321 cm = confusion_matrix(testLabels, preds)
322
323 categories = ['good_picture', 'bad_picture']
324 labels = ['True Negative', 'False Positive', 'False Negative', 'True Positive']
325 make_confusion_matrix(cm, figsize=(12,6),
326                         group_names=labels,
327                         categories=categories,
328                         cmap='binary')

```

vgg16\_image\_cleaner.py hosted with ❤ by GitHub

[view raw](#)

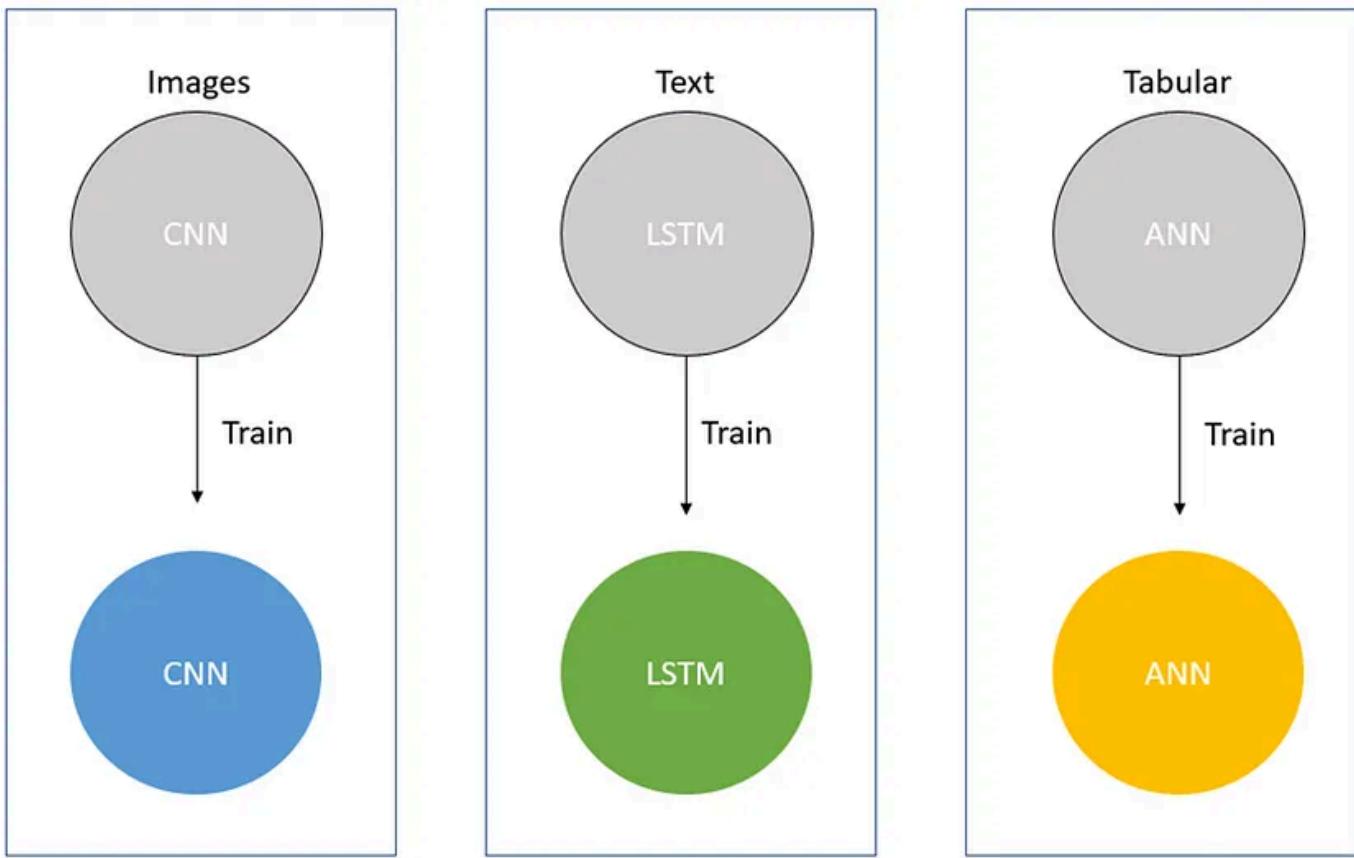
Here is the final **confusion matrix** from the test data sample:



Confusion Matrix from the VGG16 image cleaning test set

We can now predict every image we have using our “**VGG16 cleaner**” and remove all rows of our dataset where the available image has been classified as “bad”.

For the next steps, the strategy we will use is to individually train each neural network, save the weights and then load the trained weights into our final hybrid model and retrain with all the parts merged together.



Fine-tuning each neural network architecture individually first

## Architecture#1 — Structured Data fine-tuning with MLP:

Python code to create the neural network architecture, fine-tune, and save the weights:

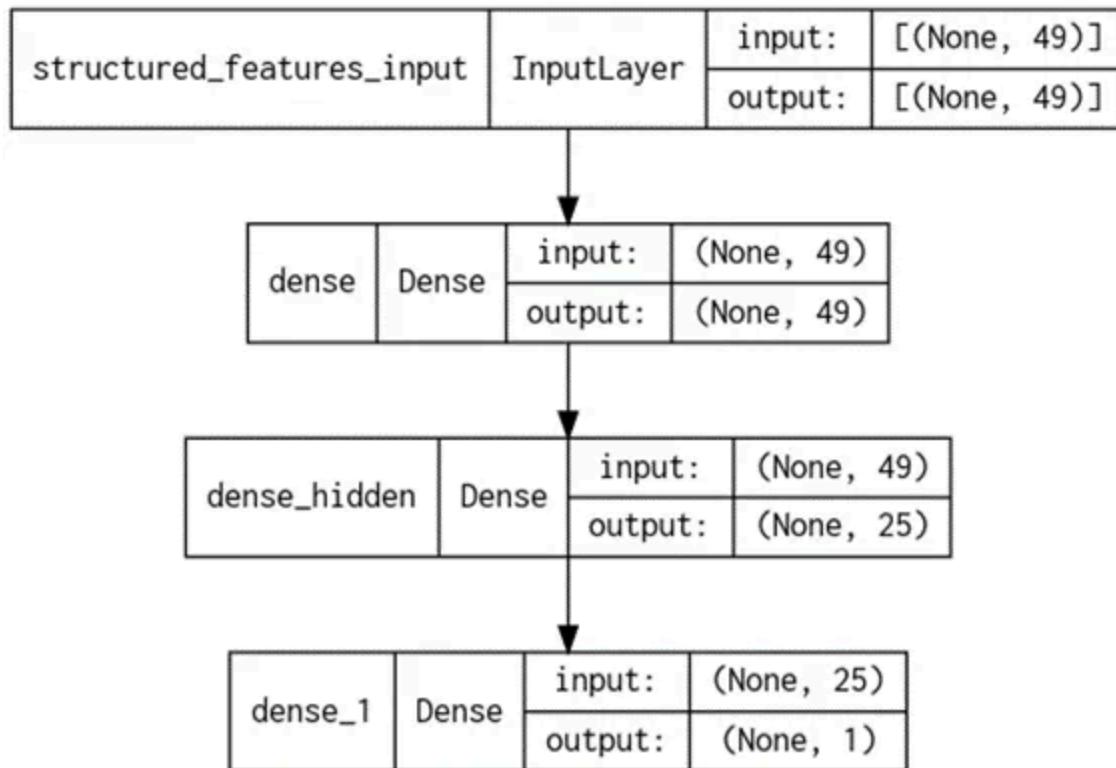
```
1  from tensorflow.keras import layers
2  from keras.models import Model
3  from keras.layers import Dense
4  from keras.callbacks import EarlyStopping, ModelCheckpoint
5  from tensorflow.keras.optimizers import Adam
6  from keras.losses import MeanSquaredError
7  from tensorflow.keras.optimizers.schedules import ExponentialDecay
8  import pandas as pd
9  from keras.utils.vis_utils import plot_model
10 import pickle as pkl
11 import warnings
12 warnings.simplefilter(action="ignore", category=FutureWarning)
13
14 X_train_structured_std = pd.read_csv('/data/X_train_structured.csv', sep = ';')
15 X_test_structured_std = pd.read_csv('/data/X_test_structured.csv', sep = ';')
16
17 y_train_std = pd.read_csv('/data/y_train_std.csv', sep = ';')
18 y_test_std = pd.read_csv('/data/y_test_std.csv', sep = ';')
19
20
21 file = open("/data/model/y_numerical_scaler.pkl",'rb')
22 y_numerical_scaler = pkl.load(file)
23 file.close()
24
25 ##### MODEL ARCHITECTURE
26 structured_input = layers.Input(shape=(X_train_structured_std.shape[1]),
27                                 , name='structured_features_input'
28 )
29
30 x = layers.Dense(X_train_structured_std.shape[1],
31                   activation='relu',
32                   kernel_initializer='he_normal')(structured_input) #80
33 dense_hidden = Dense(25, activation='relu', name='dense_hidden')(x)
34 last_layer = Dense(1, activation='linear')(dense_hidden)
35
36
37 # declare the final model inputs and outputs
38 final_model = Model(inputs=structured_input, outputs=last_layer)
39
40 # print a summary of the model
41 print(final_model.summary())
42
43 # set up learning rate decay schedule
44 initial_learning_rate = 0.1
45 lr_schedule = ExponentialDecay(
```

```
46     initial_learning_rate,
47         decay_steps=100000,
48         decay_rate=0.96,
49         staircase=True)
50
51 stop = EarlyStopping(monitor="val_loss", patience=50, restore_best_weights=True, mode='min', verbose=1)
52 best = ModelCheckpoint(filepath='/data/model/best_structured_model.hdf5', save_best_only=True, mode='min')
53
54 # compile the model
55 final_model.compile(optimizer=Adam(learning_rate=lr_schedule, epsilon=1),
56                     loss="mean_squared_error",
57                     metrics=[MeanSquaredError()])
58
59 plot_model(final_model, show_shapes=True, to_file='model1_structured.png')
60
61
62 results = final_model.fit(
63     X_train_structured_std, y_train_std,
64     epochs=500,
65     callbacks=[stop, best],
66     validation_data=(X_test_structured_std, y_test_std))
67 )
68
69 final_model.save('/data/model/best_structured_model.h5')
```

[fine\\_tune\\_structured\\_ann.py](#) hosted with ❤ by GitHub

[view raw](#)

## Multi-layer perceptron (MLP) architecture:



Structured Data Neural Network Architecture (MultiLayers Perceptron)

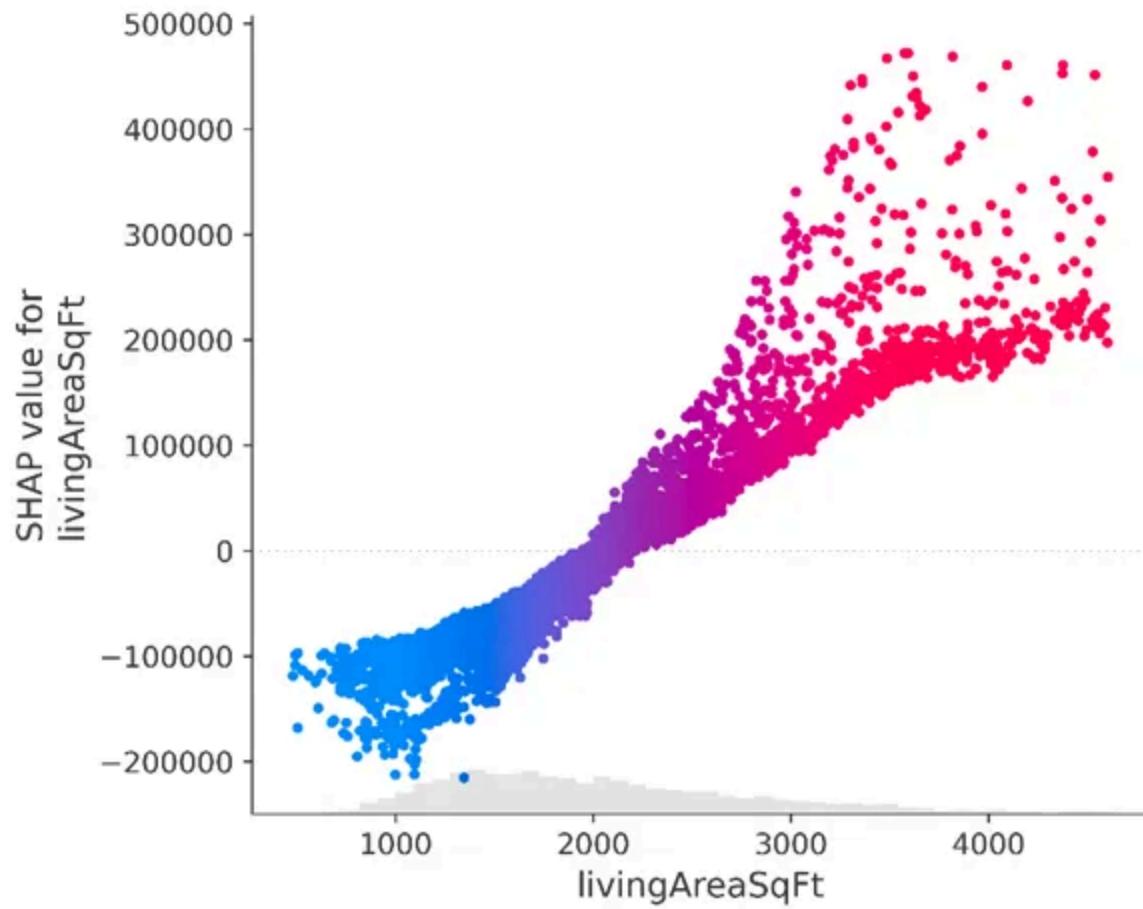
Here are the most important shapley values for our tabular features (Using the DeepLift algorithm (DeepExplainer) from the awesome [SHAP package](#)) :



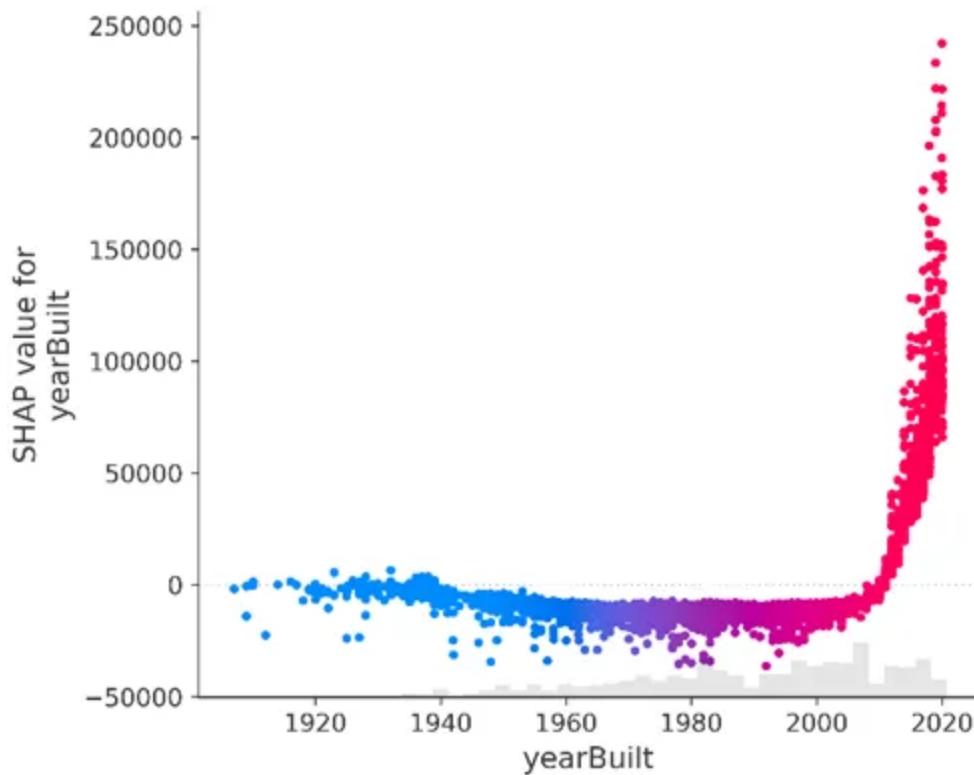
This barplot shows us that the 3 most important features to predict price of our structured data neural network are: **livingAreaSqFt**, **zipcode** and **regionCluster**.

**numPriceChanges**, **yearBuilt** and **lotSizeSqFt** follow next.

If we look at some of these features:



The biggest livingAreaSqFt, the biggest shapley value, the biggest price (seems logical!)



The more recent yearBuilt, the biggest shapley value, the biggest price (seems logical!)

## Architecture#2 — Text Data fine-tuning GLOve BiLSTM with Attention:

Python code to create the neural network architecture, fine-tune, and save the weights:

```
1 #Import Libraries
2
3 import pandas as pd
4 import numpy as np
5 import pickle as pkl
6
7 import tensorflow as tf
8 from tensorflow import keras
9 from tensorflow.keras import layers
10
11 # keras imports
12 from keras.models import Model
13
14
15 # filter warnings
16 import warnings
17 warnings.simplefilter(action="ignore", category=FutureWarning)
18
19 from keras.layers import Input, Attention, Embedding, Dropout, Dense , Bidirectional, LSTM
20 from keras.initializers import Constant
21 from keras.callbacks import EarlyStopping, ModelCheckpoint
22 from keras.preprocessing.text import Tokenizer
23 from tensorflow.keras.optimizers import Adam
24 from keras.losses import MeanSquaredError
25 #from keras.utils import plot_model, Sequence
26 from tensorflow.keras.optimizers.schedules import ExponentialDecay
27
28
29 from keras.utils.vis_utils import plot_model
30
31
32 from keras.preprocessing.sequence import pad_sequences
33
34
35 from category_encoders import *
36
37
38 class TransformerBlock(layers.Layer):
39     def __init__(self, embed_dim = 32, num_heads = 2, ff_dim = 32, rate=0.1):
40         super(TransformerBlock, self).__init__()
41         self.embed_dim = embed_dim
42         self.num_heads = num_heads
43         self.ff_dim = ff_dim
44         self.rate = rate
```

```
46         self.att = layers.MultiHeadAttention(num_heads=num_heads, key_dim=embed_dim)
47         self.ffn = keras.Sequential(
48             [layers.Dense(ff_dim, activation="relu"), layers.Dense(embed_dim),]
49         )
50         self.layernorm1 = layers.LayerNormalization(epsilon=1e-6)
51         self.layernorm2 = layers.LayerNormalization(epsilon=1e-6)
52         self.dropout1 = layers.Dropout(rate)
53         self.dropout2 = layers.Dropout(rate)
54
55     def get_config(self):
56
57         config = super().get_config().copy()
58         config.update({
59             'embed_dim': self.embed_dim,
60             'num_heads': self.num_heads,
61             'ff_dim': self.ff_dim,
62             'rate': self.rate,
63             'att': self.att,
64             'ffn': self.ffn,
65             'layernorm1': self.layernorm1,
66             'layernorm2': self.layernorm2,
67             'dropout1': self.dropout1,
68             'dropout2': self.dropout2
69         })
70         return config
71
72     def call(self, inputs, training):
73         attn_output = self.att(inputs, inputs)
74         attn_output = self.dropout1(attn_output, training=training)
75         out1 = self.layernorm1(inputs + attn_output)
76         ffn_output = self.ffn(out1)
77         ffn_output = self.dropout2(ffn_output, training=training)
78         return self.layernorm2(out1 + ffn_output)
79
80     class TokenAndPositionEmbedding(layers.Layer):
81         def __init__(self, maxlen = 100, vocab_size = 1000, embed_dim = 125):
82             super(TokenAndPositionEmbedding, self).__init__()
83             self maxlen = maxlen
84             self.vocab_size = vocab_size
85             self.embed_dim = embed_dim
86             self.token_emb = layers.Embedding(input_dim=vocab_size, output_dim=embed_dim)
87             self.pos_emb = layers.Embedding(input_dim=maxlen, output_dim=embed_dim)
88
89         def get_config(self):
```

```
90
91     config = super().get_config().copy()
92     config.update({
93         'maxlen': self maxlen,
94         'vocab_size': self vocab_size,
95         'embed_dim': self embed_dim,
96         'token_emb': self token_emb,
97         'pos_emb': self pos_emb
98     })
99     return config
100
101 def call(self, x):
102     maxlen = tf.shape(x)[-1]
103     positions = tf.range(start=0, limit=maxlen, delta=1)
104     positions = self.pos_emb(positions)
105     x = self.token_emb(x)
106     return x + positions
107
108 def prepare_text_data_train(dataf):
109     for col in dataf.columns:
110         dataf[col] = dataf[col].str.lower()
111
112     tokenizer = Tokenizer(num_words=2096) #5000
113     tokenizer.fit_on_texts(dataf[col])
114     pk1.dump(tokenizer, open('/data/model/text_tokenizer.pkl','wb'))
115
116     vocab_size = len(tokenizer.word_index) + 1
117     print(vocab_size)
118
119     pk1.dump(vocab_size, open('/data/model/text_vocab_size.pkl','wb'))
120
121     maxlen = int(np.mean(dataf[col].str.len()))
122     pk1.dump(maxlen, open('/data/model/text_maxlen.pkl','wb'))
123
124 def prepare_text_data_inference(dataf):
125     for col in dataf.columns:
126         dataf[col] = dataf[col].str.lower()
127
128         file = open("/data/model/text_tokenizer.pkl",'rb')
129         tokenizer = pk1.load(file)
130         file.close()
131
132         data_tokens = tokenizer.texts_to_sequences(dataf[col])
133
134         file = open("/data/model/text_vocab_size.pkl",'rb')
```

```
135     vocab_size = pkl.load(file)
136     file.close()
137
138     file = open("/data/model/text_maxlen.pkl",'rb')
139     maxlen = pkl.load(file)
140     file.close()
141
142     data_tokens_pad = pad_sequences(data_tokens, padding='post', maxlen=max(maxlen, 300), tr
143
144     return data_tokens_pad, data_tokens, vocab_size, max(maxlen, 300)
145
146
147 X_train_text = pd.read_csv('/data/X_train_text.csv', sep = ';')
148 X_test_text = pd.read_csv('/data/X_test_text.csv', sep = ';')
149
150 y_train_std = pd.read_csv('/data/y_train_std.csv', sep = ';')
151 y_test_std = pd.read_csv('/data/y_test_std.csv', sep = ';')
152
153 del X_train_text['Unnamed: 0']
154 del X_test_text['Unnamed: 0']
155
156 print('TRAIN', len(X_train_text), len(y_train_std))
157 print('TEST', len(X_test_text), len(y_test_std))
158
159
160 prepare_text_data_train(X_train_text)
161 train_tokens_pad, train_tokens, vocab_size, maxlen = prepare_text_data_inference(X_train_text)
162 test_tokens_pad, test_tokens, vocab_size, maxlen = prepare_text_data_inference(X_test_text)
163
164 file = open("/data/model/text_tokenizer.pkl",'rb')
165 tokenizer = pkl.load(file)
166 file.close()
167
168
169
170 ##### LOAD GLOVE EMBEDDINGS
171 print('glove load...')
172 # identify the embedding filename; we are using the Glove 42B 300d embeddings
173 glove_file = "/data/model/glove.840B.300d.txt"
174 print('glove loaded!')
175 # create the embeddings index dictionary
176 embeddings_index = {} # create a lookup dictionary to store words and their vectors
177 f = open(glove_file, errors='ignore')# open our embedding file
178 for line in f: # for each line in the file
```

12/24, 8:44 PM Hybrid (multimodal) neural network architecture : Combination of tabular, textual and image inputs to predict house prices. | by Dav...

```

179     values = line.split(' ') #split the line on spaces between the word and its vectors
180     word = values[0] # the word is the first entry
181     if word in tokenizer.word_index.keys(): # we check if the word is in our tokenizer word index
182         coefs = np.asarray(values[1:], dtype='float32') # if so, get the word's vectors
183         embeddings_index[word] = coefs # add the word and its vectors to the embeddings_index dictionary
184 f.close()
185 print('Found %s word vectors.' % len(embeddings_index)) # report how many words in our corpus we found
186
187 # amount of vocabulary to use, will pick the top 10000 words seen in the corpus
188 features = 10000
189 # max text sequence length, must match tokens in transfer file, we are using glove 300d so it is 300
190 max_words = 300
191
192 num_tokens = (len(tokenizer.word_index) + 1) # for num tokens we always do the length of our word index + 1
193 hits = 0
194 misses = 0
195 embedding_dim = 300
196
197 embedding_matrix = np.zeros((num_tokens, max_words)) # setting up an array for our tokens with a size of num_tokens by max_words
198 for word, i in tokenizer.word_index.items(): # for each word in the tokenizer word index
199     embedding_vector = embeddings_index.get(word) #get the vector from the embeddings index dictionary
200     if embedding_vector is not None: # if the vector isn't None,
201         # words not found in embedding index will be all-zeros.
202         embedding_matrix[i] = embedding_vector # store the embedding vector in the matrix at the index i
203         hits += 1
204     else:
205         misses += 1
206 print("Converted %d words (%d misses)" % (hits, misses))
207
208
209 ##### CREATING TEXT NEURAL NETWORK ARCHITECTURE WITH GLOVE EMBEDDINGS
210 # for text embeddings the input shape can be none
211 nlp_input = Input(shape=(maxlen,), name='nlp_input')
212
213 # uses the embedding matrix dictionary to create word embeddings for the inputs
214 embedded_sequences = Embedding(
215     input_dim = num_tokens, # number of unique tokens
216     output_dim = embedding_dim, #number of features
217     embeddings_initializer=Constant(embedding_matrix), # initialize with Glove embeddings
218     input_length=maxlen,
219     trainable=False)(nlp_input)
220 # uses two bi-directional LSTM layers
221 lstm = Bidirectional(LSTM(300, return_sequences=True))(embedded_sequences)
222 att_lstm = Attention(units=300)(lstm)
223 nlp_out = Dropout(0.5)(att_lstm)

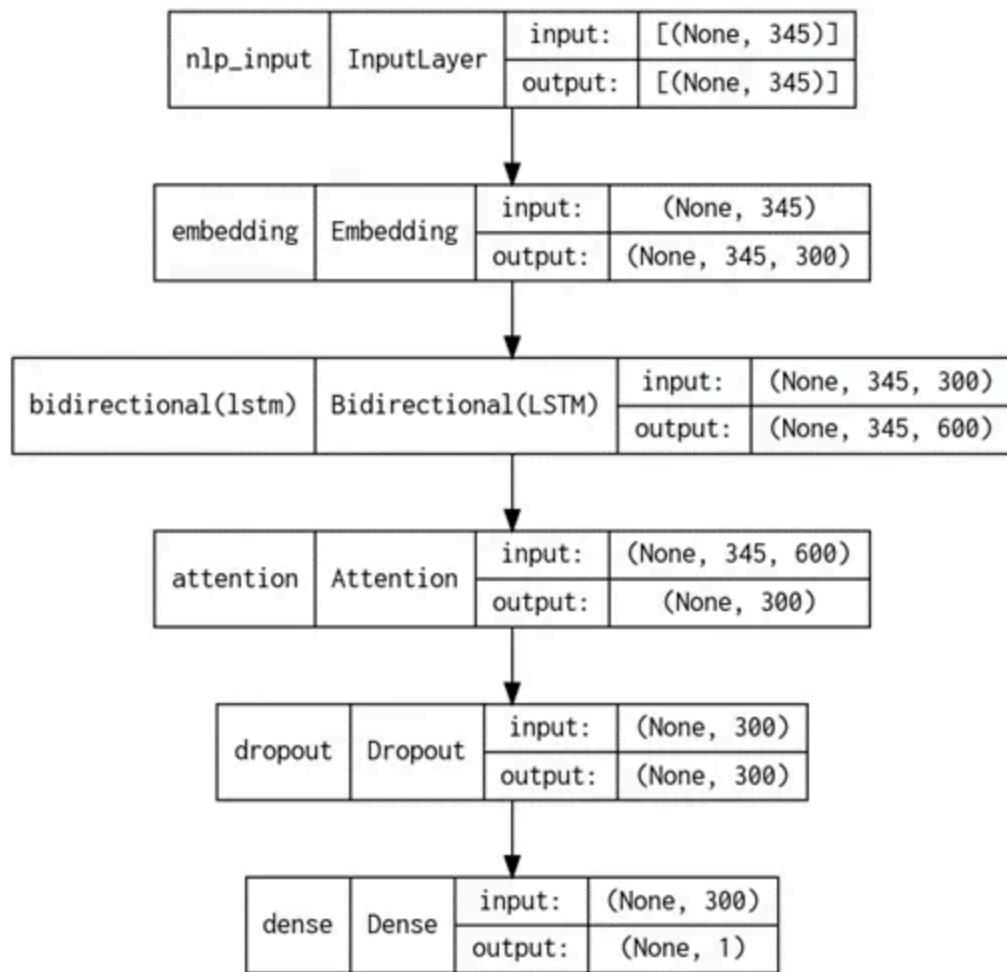
```

```
224 # adds a dense layer
225 #output_embed = Dense(350, activation="relu", kernel_initializer='he_normal')(x)
226 last_layer = Dense(1, activation='linear')(nlp_out)
227
228 # declare the final model inputs and outputs
229 final_model = Model(inputs=nlp_input, outputs=last_layer)
230
231 # print a summary of the model
232 print(final_model.summary())
233
234 # set up learning rate decay schedule
235 initial_learning_rate = 0.1
236 lr_schedule = ExponentialDecay(
237     initial_learning_rate,
238     decay_steps=100000,
239     decay_rate=0.96,
240     staircase=True)
241
242 stop = EarlyStopping(monitor="val_loss", patience=25, restore_best_weights=True, mode='min', ver
243 best = ModelCheckpoint(filepath='/data/model/best_nlp_model.hdf5', save_best_only=True, monitor=
244
245 # compile the model
246 final_model.compile(optimizer=Adam(learning_rate=lr_schedule, epsilon=1),
247                     loss="mean_squared_error",
248                     metrics=[MeanSquaredError()])
249
250 plot_model(final_model, show_shapes=True, to_file='model2_nlp_attention.png')
251
252 results = final_model.fit(
253     train_tokens_pad, y_train_std,
254     epochs=500,
255     callbacks=[stop, best],
256     validation_data=(test_tokens_pad, y_test_std))
257 )
258
259 final_model.save('/data/model/pretrained_nlp_model.h5')
260
261
262
263
```

fine\_tune\_bilstm\_attention\_glove.py hosted with ❤ by GitHub

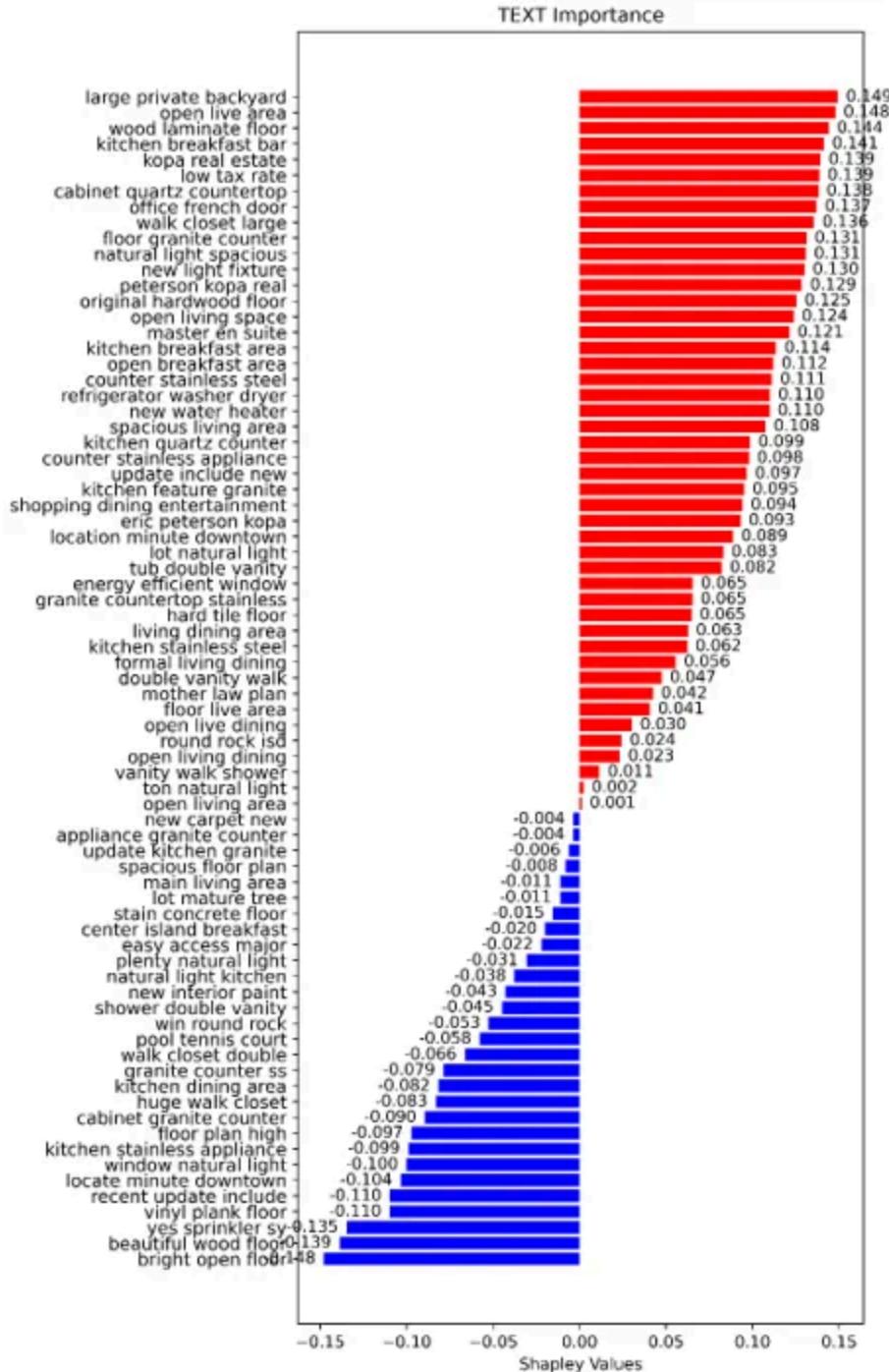
[view raw](#)

## BiLSTM Attention with GloVe embeddings architecture:



Text Neural Network Architecture (BiLSTM Attention with GLOVe Embeddings)

Again using shapley values, here are some useful **ngrams** of the house text description to predict price:

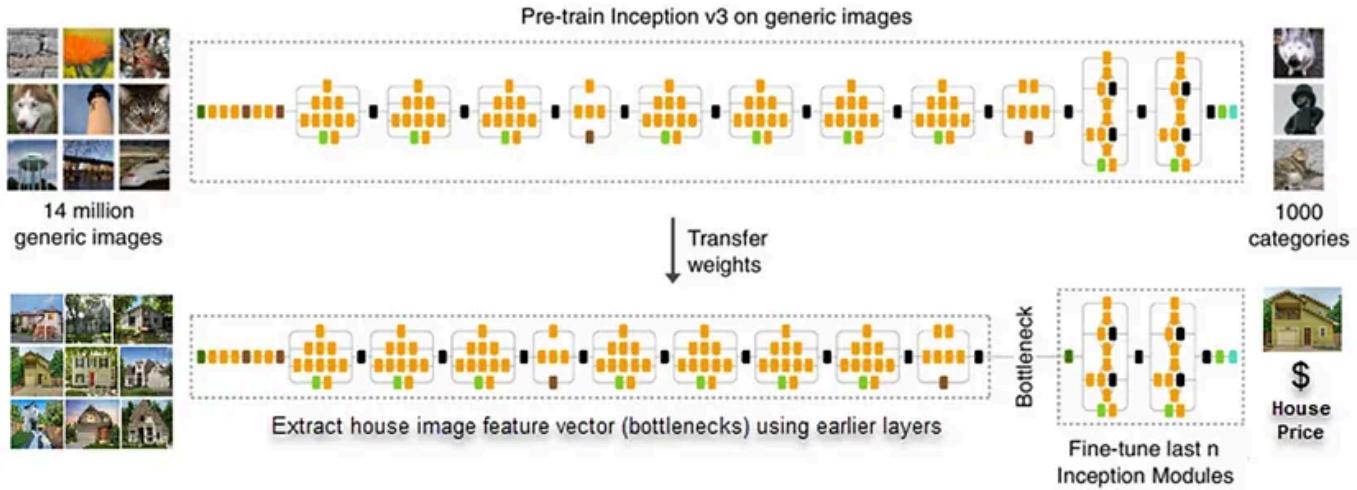


The “large private backyard”, “open live area” and “wood laminate floor” grams appear to have the most positive influence on price according to the shapley values.

In contrast, “bright open floor”, “beautiful wood floor”, “yes sprinkler sy” and “vinyl plank floor” seem to have the most negative influence on the price.

## Architecture#3 — Image Data fine-tuning with Inception V3:

We use the pre-trained InceptionV3 model on generic image (imagenet weights) and fine-tune only the head layer using our houses prices.



Python code to create the neural network architecture, fine-tune, and save the weights:

```
1 import pandas as pd
2 from keras.applications.inception_v3 import InceptionV3
3 from keras.applications.vgg16 import VGG16
4
5
6 from keras.layers import BatchNormalization, GlobalAveragePooling2D, Dense, Input
7 from keras.models import Model
8 from tensorflow.keras.optimizers import Adam
9 from keras.callbacks import EarlyStopping, ModelCheckpoint
10 from keras.losses import MeanSquaredError
11 from tensorflow.keras.optimizers.schedules import ExponentialDecay
12 from keras.utils.vis_utils import plot_model
13
14 from keras.layers import Dense, Flatten, Conv2D, Dropout, MaxPooling2D, MaxPooling1D
15 from keras.layers import Activation, Dropout, Flatten, BatchNormalization
16 from keras.models import Sequential
17 from keras.utils import np_utils
18
19 X_train_image = pd.read_csv('/data/X_train_image.csv', sep = ';')
20 X_test_image = pd.read_csv('/data/X_test_image.csv', sep = ';')
21
22 y_train_std = pd.read_csv('/data/y_train_std.csv', sep = ';')
23 y_test_std = pd.read_csv('/data/y_test_std.csv', sep = ';')
24
25
26 X_train_image['homeImage'] = X_train_image['homeImage'].apply(lambda x: '/data/img/' + str(x))
27 X_train_image['price'] = y_train_std.values
28 X_train_image
29
30 X_test_image['homeImage'] = X_test_image['homeImage'].apply(lambda x: '/data/img/' + str(x))
31 X_test_image['price'] = y_test_std.values
32 X_test_image
33
34 img_size_shape = 128
35
36 from keras.preprocessing.image import ImageDataGenerator
37 image_train_generator = ImageDataGenerator(
38     brightness_range=(0.75, 1),
39     shear_range=0.1,
40     zoom_range=[0.75, 1],
41     featurewise_center=True,
42     featurewise_std_normalization=True,
43     rotation_range=20,
44     width_shift_range=0.2,
45     height_shift_range=0.2)
```

```
46     horizontal_flip=True
47 )
48
49 # test/val have only the pixel data normalization
50 image_test_generator = ImageDataGenerator(rescale = 1./1)
51
52 # visualize an image augmentation sample
53 #visualize_augmentations(image_train_generator, images_train.iloc[1])
54
55 # specify where the train generator pulls batches
56 image_train_generator = image_train_generator.flow_from_dataframe(
57     dataframe=X_train_image,
58     x_col="homeImage", # this is where your image data is stored
59     y_col="price", # this is your target feature
60     class_mode="raw", # use "raw" for regressions
61     color_mode='rgb',
62     target_size=(img_size_shape, img_size_shape),
63     batch_size=32, # increase or decrease to fit your GPU,
64 )
65
66 # specify where the test generator pulls batches
67 image_test_generator = image_test_generator.flow_from_dataframe(
68     dataframe=X_test_image,
69     x_col="homeImage",
70     y_col="price",
71     class_mode="raw",
72     color_mode='rgb',
73     target_size=(img_size_shape, img_size_shape),
74     batch_size=32,
75 )
76
77
78
79 model = InceptionV3(weights="imagenet", include_top=False, input_tensor=Input(shape=(img_size_s
80
81 for layer in model.layers:
82     layer.trainable=False
83
84 img_input = Input(shape=(img_size_shape, img_size_shape, 3))
85 x = model(img_input, training=False)
86 x = GlobalAveragePooling2D(name="avg_pool")(x)
87 x = BatchNormalization()(x)
88 #x = Dropout(0.5)(x)
89 output_cnn = Dense(64, activation='relu', kernel_initializer='he_normal', name='dense_hidden')()
```

```
90     last_layer = Dense(1, activation='linear')(output_cnn)
91
92     model = Model(inputs=img_input, outputs=last_layer)
93
94     batch_size = 64
95     epochs = 100
96     verbose = 1
97
98
99     initial_learning_rate = 0.1
100
101 """
102     lr_schedule = ExponentialDecay(
103         initial_learning_rate,
104         decay_steps=100000,
105         decay_rate=0.96,
106         staircase=True)
107 """
108
109     model.compile(
110         optimizer=Adam(learning_rate=initial_learning_rate, epsilon=1),
111         metrics=['mse'],
112         loss='mse'
113     )
114
115     model.summary()
116
117     plot_model(model, show_shapes=True, to_file='inceptionV3_image.png')
118
119     stop = EarlyStopping(monitor="val_loss", patience=15, restore_best_weights=True, mode='min', ver
120     best = ModelCheckpoint(filepath='/data/model/best_iv3_model_simple.hdf5', save_best_only=True, n
121
122     history = model.fit(
123         image_train_generator,
124         validation_data = image_test_generator,
125         batch_size=batch_size,
126         verbose=verbose,
127         epochs=epochs,
128         callbacks=[stop, best]
129     )
130
131
132     model.save('/data/model/best_iv3_model_simple.h5')
```

The final architecture (InceptionV3 is presented here as a single layer but represents a more complex underlying architecture):

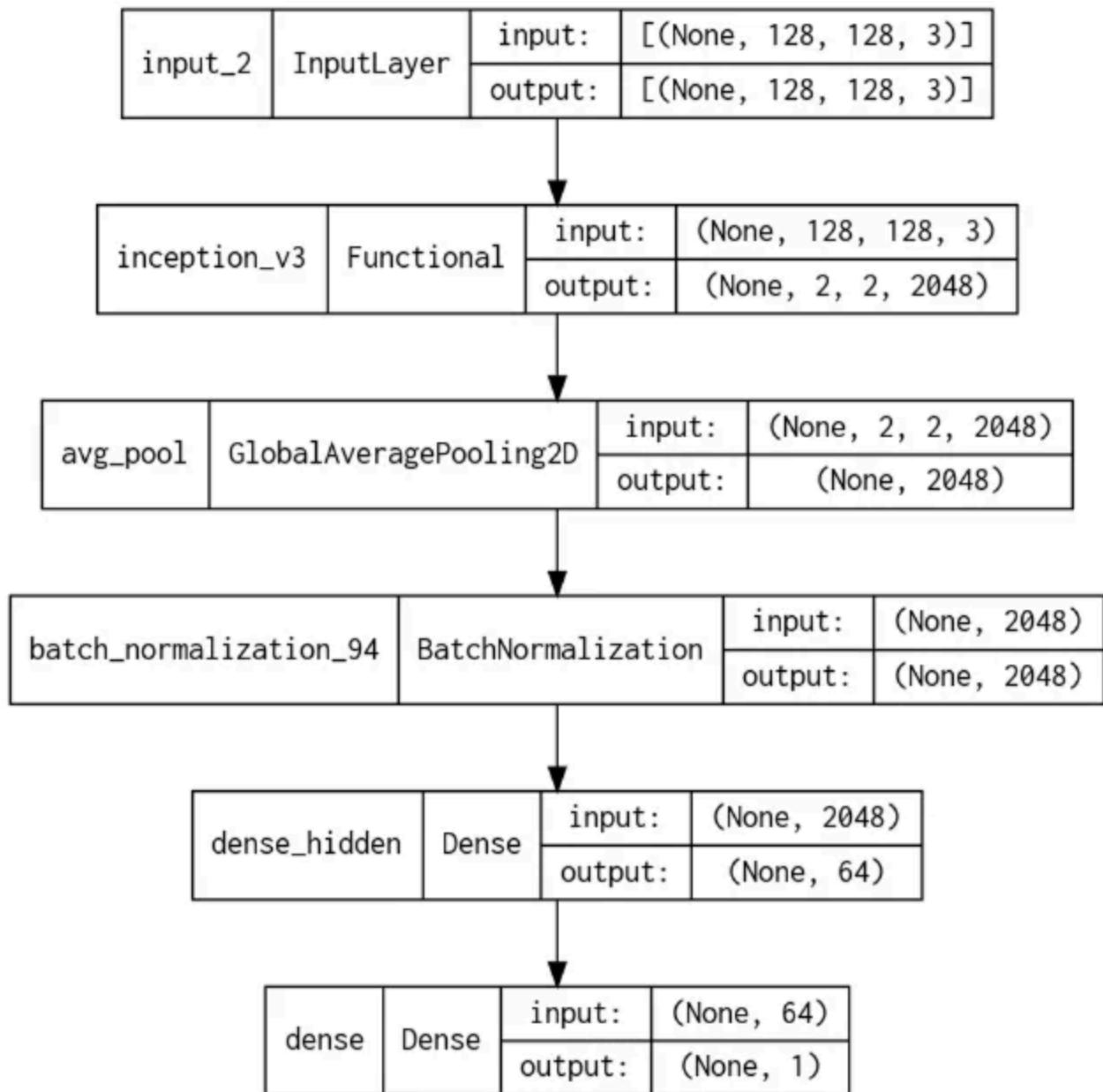
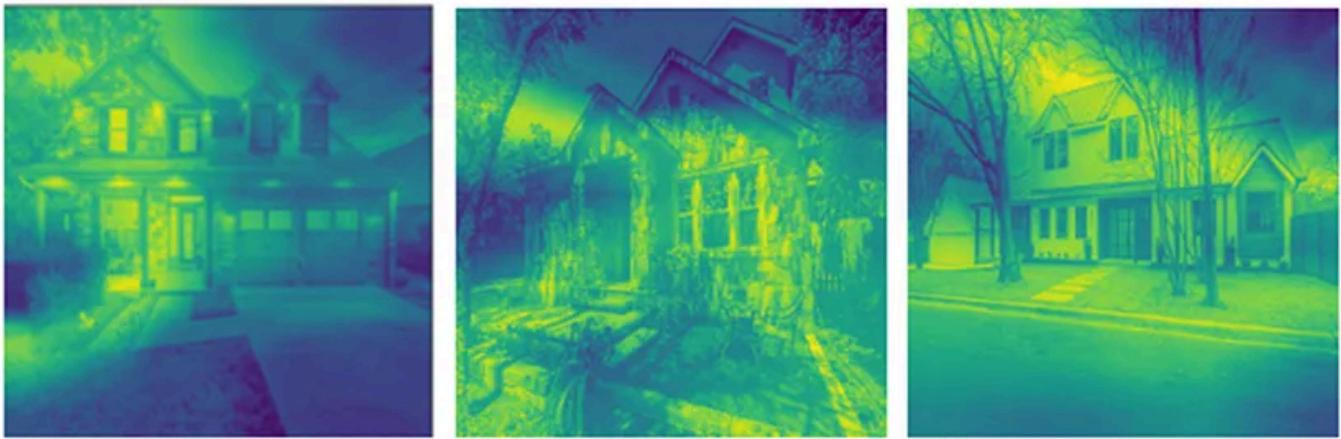


Image Neural Netowrk Architecture (InceptionV3)

Using the **gradcam** method, we can try to understand what our CNN is looking at:



The neural network seems to focus mostly on the house part of the images, which is fine, but there is no evidence which parts of the houses are more important than others... maybe our model is not trained enough !

Grad-CAM, *Gradient-weighted Class Activation Mapping*, uses the feature maps / class-specific gradient information flowing into the final convolutional layer of a CNN. It produces a visual explanation heat-maps / localization map that will help us understand how the model make decisions / the important regions in the image.

Python code for the gradcam:

```
1  from keras.models import Model
2  import keras
3  import cv2
4  import tensorflow as tf
5  import matplotlib.pyplot as plt
6  from keras.preprocessing.image import load_img
7  from keras.preprocessing.image import img_to_array
8  import numpy as np
9
10 def find_target_layer(model_pretrained):
11     # attempt to find the final convolutional layer in the network
12     # by looping over the layers of the network in reverse order
13     for layer in reversed(model_pretrained.layers):
14         # check to see if the layer has a 4D output
15         if len(layer.output_shape) == 4:
16             return layer.name
17     # otherwise, we could not find a 4D layer so the GradCAM
18     # algorithm cannot be applied
19     raise ValueError("Could not find 4D layer. Cannot apply GradCAM.")
20
21 def compute_heatmap_inceptionV3(model_pretrained, image, eps=1e-8):
22     # construct our gradient model by supplying (1) the inputs
23     # to our pre-trained model, (2) the output of the (presumably)
24     # final 4D layer in the network, and (3) the output of the
25     # softmax activations from the model
26
27     gradModel = Model(
28         inputs=[model_pretrained.get_layer('inception_v3').inputs],
29         outputs=[model_pretrained.get_layer('inception_v3').get_layer(find_target_layer(model_pretrained)).output])
30
31
32     # record operations for automatic differentiation
33     with tf.GradientTape() as tape:
34         # cast the image tensor to a float-32 data type, pass the
35         # image through the gradient model, and grab the loss
36         # associated with the specific class index
37         inputs = tf.cast(image, tf.float32)
38         (convOutputs, predictions) = gradModel(inputs)
39         loss = predictions[:, 0]
40
41     # use automatic differentiation to compute the gradients
42     grads = tape.gradient(loss, convOutputs)
43
44     # compute the guided gradients
45     castConvOutputs = tf.cast(convOutputs > 0, "float32")
46     rectGrads = tf.cast(grads < 0, "float32")\
```

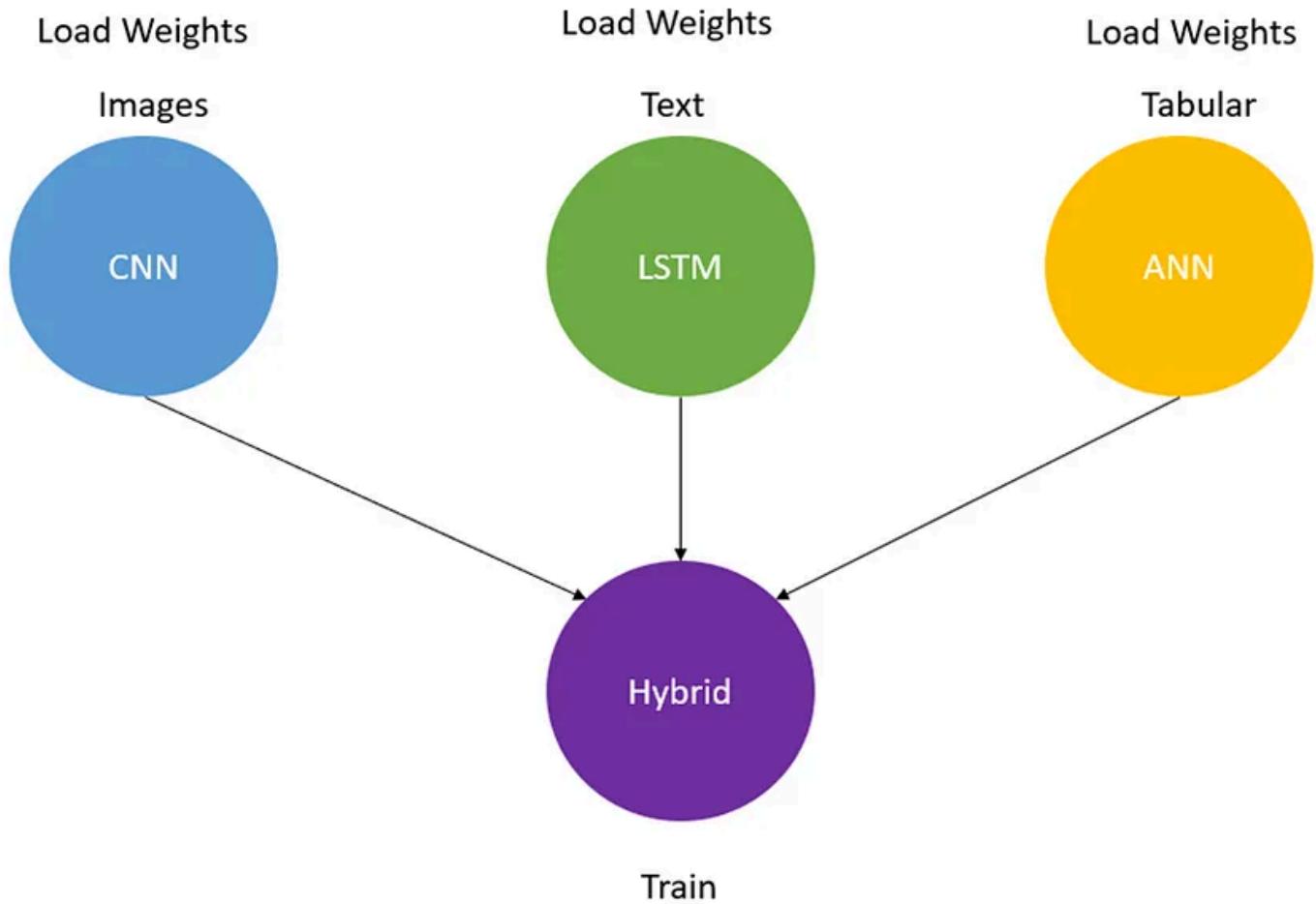
```
46 guidedGrads = castConvOutputs * castGrads * grads
47 # the convolution and guided gradients have a batch dimension
48 # (which we don't need) so let's grab the volume itself and
49 # discard the batch
50 convOutputs = convOutputs[0]
51 guidedGrads = guidedGrads[0]
52 # compute the average of the gradient values, and using them
53 # as weights, compute the ponderation of the filters with
54 # respect to the weights
55 weights = tf.reduce_mean(guidedGrads, axis=(0, 1))
56 cam = tf.reduce_sum(tf.multiply(weights, convOutputs), axis=-1)
57 # grab the spatial dimensions of the input image and resize
58 # the output class activation map to match the input image
59 # dimensions
60 (w, h) = (image.shape[2], image.shape[1])
61 heatmap = cv2.resize(cam.numpy(), (w, h))
62 # normalize the heatmap such that all values lie in the range
63 # [0, 1], scale the resulting values to the range [0, 255],
64 # and then convert to an unsigned 8-bit integer
65 numer = heatmap - np.min(heatmap)
66 denom = (heatmap.max() - heatmap.min()) + eps
67 heatmap = numer / denom
68 heatmap = (heatmap * 255).astype("uint8")
69 # return the resulting heatmap to the calling function
70 return heatmap
71
72 fine_tuned_inceptionv3 = keras.models.load_model('/data/model/best_iv3_model_simple.hdf5', compil
73
74 for img_file in ["/data/img/110537181_11732947ec76f3d1a8830ec70de79bac-p_f.jpg",
75             "/data/img/2129847996_fbfebce6f34e641b58aafbc76d66ddea-p_f.jpg",
76             "/data/img/29333742_612b27ee951dabb4b583c8b4ce87f445-p_f.jpg"]:
77     img = load_img(img_file, target_size=(299, 299))
78     img = img_to_array(img)
79     img = img.reshape((1, img.shape[0], img.shape[1], img.shape[2]))
80     img_cv2 = cv2.imread(img_file, 0)
81     img_cv2 = cv2.resize(img_cv2, (299, 299))
82     img_cv2.shape
83
84
85     img_heatmap = compute_heatmap_inceptionV3(fine_tuned_inceptionv3, img, eps=1e-8)
86     img_heatmap.shape
87
88     colormap=cv2.COLORMAP_TURBO
89     alpha = 0.43
```

12/24, 8:44 PM Hybrid (multimodal) neural network architecture : Combination of tabular, textual and image inputs to predict house prices. | by Dav...  
90 heatmap = cv2.applyColorMap(img\_heatmap, colormap, 0)[ :, :, 0]  
91 output = cv2.addWeighted(img\_cv2, alpha, heatmap, 1 - alpha, 0)  
92 final\_output = (heatmap, output)  
93  
94 plt.imsave(img\_file.replace('/img', '').replace('.jpg', '') + '\_gradcam.png', output, dpi = 500)  
95 plt.imshow(output)  
96 plt.close()

gradcam.py hosted with ❤ by GitHub

[view raw](#)

## Final Architecture — Load fine-tuned weights and create final hybrid architecture for training



Here is the python code to load all our fine-tuned neural network architecture and then create the final architecture:

```
1 import pandas as pd
2 import numpy as np
3 import pickle as pkl
4
5 from keras.preprocessing.image import ImageDataGenerator
6 from keras.preprocessing.text import Tokenizer
7 from keras.preprocessing.sequence import pad_sequences
8
9 from keras.applications.inception_v3 import InceptionV3
10 from keras.applications.inception_v3 import preprocess_input
11 from keras.layers import BatchNormalization, Conv2D, Conv1D, MaxPooling2D, MaxPooling1D, GlobalMaxPooling1D
12 from keras.models import Sequential, Model
13 from keras.regularizers import l2, l1
14 from tensorflow.keras.optimizers import Adam, SGD, RMSprop
15 from keras.callbacks import TensorBoard, EarlyStopping, ModelCheckpoint, History, LearningRateScheduler
16 from keras.losses import MeanAbsoluteError, MeanAbsolutePercentageError, MeanSquaredError
17 #from keras.utils import plot_model, Sequence
18 from keras.initializers import Constant
19 from tensorflow.keras.optimizers.schedules import ExponentialDecay
20
21 from keras.preprocessing.image import load_img
22 from keras.preprocessing.image import img_to_array
23
24 from attention import Attention
25 import keras
26
27 X_train_structured_std = pd.read_csv('/data/X_train_structured.csv', sep = ';')
28 X_test_structured_std = pd.read_csv('/data/X_test_structured.csv', sep = ';')
29
30 X_train_image = pd.read_csv('/data/X_train_image.csv', sep = ';')
31 X_test_image = pd.read_csv('/data/X_test_image.csv', sep = ';')
32
33 X_train_text = pd.read_csv('/data/X_train_text.csv', sep = ';')
34 X_test_text = pd.read_csv('/data/X_test_text.csv', sep = ';')
35
36 del X_train_text['Unnamed: 0']
37 del X_test_text['Unnamed: 0']
38
39
40
41 y_train_std = pd.read_csv('/data/y_train_std.csv', sep = ';')
42 y_test_std = pd.read_csv('/data/y_test_std.csv', sep = ';')
43
44
45 print('TRATN' len(X_train_structured_std) len(X_train_image) len(X_train_text) len(y_train_std)
```

12/24, 8:44 PM Hybrid (multimodal) neural network architecture : Combination of tabular, textual and image inputs to predict house prices. | by Dav...

```
46 print('TEST', len(X_test_structured_std), len(X_test_image), len(X_test_text), len(y_test_std))

47

48 #file = open("/data/model/X_numerical_scaler.pkl",'rb')
49 #X_numerical_scaler = pkl.load(file)
50 #file.close()

51

52 file = open("/data/model/y_numerical_scaler.pkl",'rb')
53 y_numerical_scaler = pkl.load(file)
54 file.close()

55

56 X_train_image['homeImage'] = X_train_image['homeImage'].apply(lambda x: '/data/img/'+str(x))
57 X_train_image['price'] = y_train_std.values
58 X_train_image

59

60 X_test_image['homeImage'] = X_test_image['homeImage'].apply(lambda x: '/data/img/'+str(x))
61 X_test_image['price'] = y_test_std.values
62 X_test_image

63

64

65

66

67 def load_house_images(df, image_size):
68     # initialize our images array (i.e., the house images themselves)
69     images = []
70     # loop over the indexes of the houses
71     for i, image_file in enumerate(df['homeImage'].values):
72
73         # find the four images for the house and sort the file paths,
74         # ensuring the four are always in the *same order*
75         # initialize our list of input images along with the output image
76         # after *combining* the four input images
77         inputImages = []
78         # loop over the input house paths
79         img = load_img(image_file, target_size=(image_size, image_size))
80         img = img_to_array(img)
81         img = img.reshape((1, img.shape[0], img.shape[1], img.shape[2]))
82
83         if i%1000 == 0:
84             print(i, image_file, img.shape)
85         elif i == 0:
86             print(i, image_file, img.shape)
87
88         images.append(img)
89         # tile the four input images in the output image such the first
```

12/24, 8:44 PM Hybrid (multimodal) neural network architecture : Combination of tabular, textual and image inputs to predict house prices. | by Dav...

```
90     # image goes in the top-right corner, the second image in the
91     # top-left corner, the third image in the bottom-right corner,
92     # and the final image in the bottom-left corner
93     # return our set of images
94 images = np.array(images)
95 images = images.reshape((len(images), image_size, image_size, 3))
96 return images
97
98
99
100 X_train_image_array = load_house_images(X_train_image, image_size = 128)
101 X_test_image_array = load_house_images(X_test_image, image_size = 128)
102
103 pkl.dump(X_train_image_array, open('/data/model/X_train_image_array.pkl','wb'))
104 pkl.dump(X_test_image_array, open('/data/model/X_test_image_array.pkl','wb'))
105
106
107 file = open("/data/model/X_train_image_array.pkl",'rb')
108 X_train_image_array = pkl.load(file)
109 file.close()
110
111 file = open("/data/model/X_test_image_array.pkl",'rb')
112 X_test_image_array = pkl.load(file)
113 file.close()
114
115
116 def prepare_text_data_train(dataf):
117     for col in dataf.columns:
118         dataf[col] = dataf[col].str.lower()
119
120         tokenizer = Tokenizer(num_words=2096) #5000
121         tokenizer.fit_on_texts(dataf[col])
122         pkl.dump(tokenizer, open('/data/model/text_tokenizer.pkl','wb'))
123
124         vocab_size = len(tokenizer.word_index) + 1
125         print(vocab_size)
126
127         pkl.dump(vocab_size, open('/data/model/text_vocab_size.pkl','wb'))
128
129         maxlen = int(np.mean(dataf[col].str.len()))
130         pkl.dump(maxlen, open('/data/model/text_maxlen.pkl','wb'))
131
132 def prepare_text_data_inference(dataf):
133     for col in dataf.columns:
134         dataf[col] = dataf[col].str.lower()
```

```
135
136     file = open("/data/model/text_tokenizer.pkl",'rb')
137     tokenizer = pkl.load(file)
138     file.close()
139
140     data_tokens = tokenizer.texts_to_sequences(dataf[col])
141
142     file = open("/data/model/text_vocab_size.pkl",'rb')
143     vocab_size = pkl.load(file)
144     file.close()
145
146     file = open("/data/model/text_maxlen.pkl",'rb')
147     maxlen = pkl.load(file)
148     file.close()
149
150     data_tokens_pad = pad_sequences(data_tokens, padding='post', maxlen=max(maxlen, 300), tr
151
152     return data_tokens_pad, data_tokens, vocab_size, max(maxlen, 300)
153
154     prepare_text_data_train(X_train_text)
155     train_tokens_pad, train_tokens, vocab_size, maxlen = prepare_text_data_inference(X_train_text)
156     test_tokens_pad, test_tokens, vocab_size, maxlen = prepare_text_data_inference(X_test_text)
157
158     file = open("/data/model/text_tokenizer.pkl",'rb')
159     tokenizer = pkl.load(file)
160     file.close()
161
162
163
164     fine_tuned_inceptionv3 = keras.models.load_model('/data/model/best_iv3_model_simple.hdf5', compi
165     fine_tuned_glove_bilstm_att = keras.models.load_model('/data/model/best_nlp_model.hdf5', compil
166     fine_tuned_structured = keras.models.load_model('/data/model/best_structured_model.hdf5', compil
167
168     fine_tuned_inceptionv3_layer = Model(inputs=fine_tuned_inceptionv3.input, outputs=fine_tuned_in
169     fine_tuned_glove_bilstm_att_layer = Model(inputs=fine_tuned_glove_bilstm_att.input, outputs=fine
170     fine_tuned_structured_layer = Model(inputs=fine_tuned_structured.input, outputs=fine_tuned_struct
171
172     for layer in fine_tuned_structured_layer.layers:
173         layer._name = layer.name + '_structured'
174
175     for layer in fine_tuned_glove_bilstm_att_layer.layers:
176         layer._name = layer.name + '_nlp'
177
178     for layer in fine_tuned_inceptionv3_layer.layers:
```

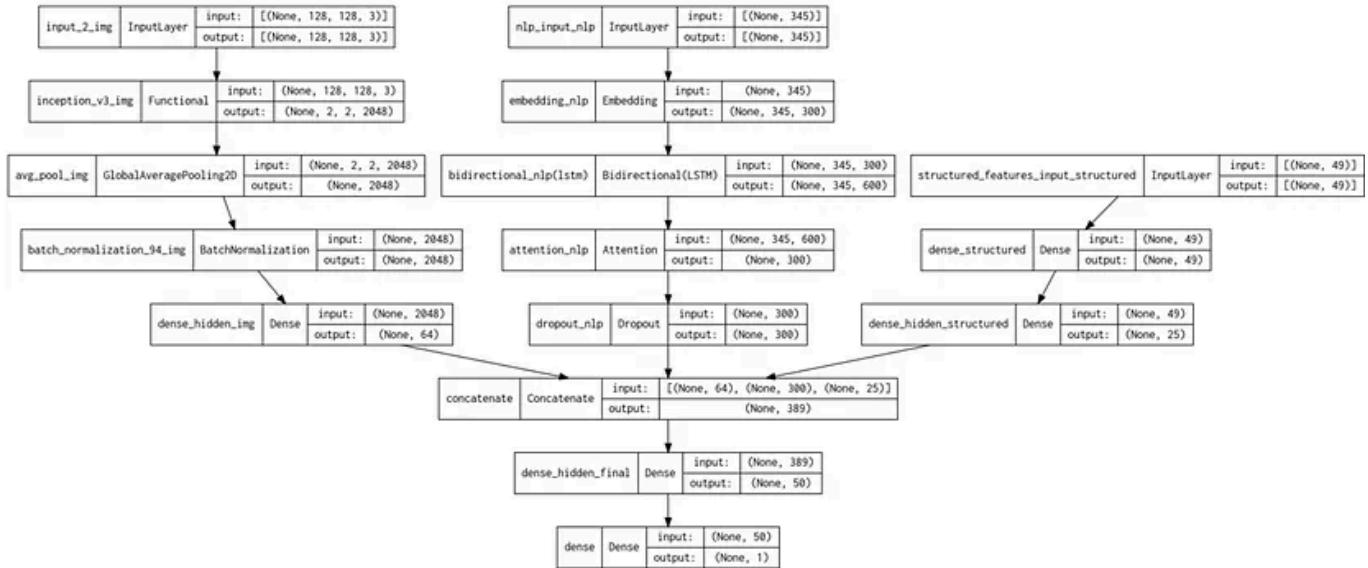
```
179     layer._name = layer.name + '_img'
180
181     x = keras.layers.concatenate([fine_tuned_inceptionv3_layer.output,
182                                     fine_tuned_glove_bilstm_att_layer.output,
183                                     fine_tuned_structured_layer.output])
184     x = Dense(50, activation='relu', name='dense_hidden_final')(x)
185     last_layer = Dense(1, activation='linear')(x)
186
187
188     model = keras.Model(inputs=[fine_tuned_inceptionv3_layer.input,
189                                 fine_tuned_glove_bilstm_att_layer.input,
190                                 fine_tuned_structured_layer.input
191                               ],
192                           outputs=[last_layer])
193
194     initial_learning_rate = 0.1
195
196     ''
197     lr_schedule = ExponentialDecay(
198         initial_learning_rate,
199         decay_steps=100000,
200         decay_rate=0.96,
201         staircase=True)
202     ''
203
204     model.compile(optimizer=Adam(learning_rate=initial_learning_rate, epsilon=1), loss="mean_squared_error")
205     print(model.summary())
206
207     from tensorflow.keras.utils import plot_model
208     plot_model(model, to_file="final_model.png",
209                 show_shapes=True,
210                 show_dtype=False,
211                 show_layer_names=True,
212                 dpi=350)
213
214     print('saving...')
215     model.save('/data/model/best_final_model_simple.h5')
216     print('saved DONE!')
217
218     print('loading...')
219     reconstructed_model = keras.models.load_model("/data/model/best_final_model_simple.h5", compile=False)
220     print('loaded DONE!')
221
222     plot_model(reconstructed_model, show_shapes=True, to_file='final_model_image.png')
223
```

```
224     reconstructed_model.compile(optimizer=Adam(learning_rate=initial_learning_rate, epsilon=1), loss='mse')
225     print(reconstructed_model.summary())
226
227     stop = EarlyStopping(monitor="val_loss", patience=15, restore_best_weights=True, mode='min', verbose=1)
228     best = ModelCheckpoint(filepath='/data/model/best_final_model_weights_only.h5',
229                           save_best_only=True,
230                           save_weights_only=False,
231                           monitor='val_loss',
232                           mode='min', verbose=1)
233
234     results = reconstructed_model.fit([X_train_image_array,
235                                         train_tokens_pad,
236                                         X_train_structured_std],
237                                         y_train_std,
238                                         epochs=500,
239                                         batch_size = 250,
240                                         validation_data=([X_test_image_array,
241                                             test_tokens_pad,
242                                             X_test_structured_std],
243                                             y_test_std),
244                                         callbacks=[stop, best],
245                                         )
246
247     results.save('/data/model/final_hybrid_model.h5')
```

final\_hybrid\_architecture.py hosted with ❤ by GitHub

[view raw](#)

The final architecture:



Final Hybrid (Multimodal) Neural Netowrk Architecture

We can see the concatenate layer in the center that concatenate our 3 previous neural network architectures; image, text and structured data respectively.

## Results and Conclusions

If we take a quick look at our first results:

```
1 import pandas as pd
2 import pickle as pkl
3 from attention import Attention
4 import keras
5 from sklearn import metrics
6
7 print('loading...')
8 final_model = keras.models.load_model("/data/model/best_final_model_weights_only.h5", compile=False)
9 fine_tuned_inceptionv3 = keras.models.load_model('/data/model/best_iv3_model_simple.hdf5', compile=False)
10 fine_tuned_glove_bilstm_att = keras.models.load_model('/data/model/best_nlp_model.hdf5', compile=False)
11 fine_tuned_structured = keras.models.load_model('/data/model/best_structured_model.hdf5', compile=False)
12 fine_tuned_structured_nlp = keras.models.load_model('/data/model/best_model_struct_nlp.h5', compile=False)
13 print('loaded DONE!')
14
15 X_train_structured_std = pd.read_csv('/data/X_train_structured.csv', sep = ';')
16 X_test_structured_std = pd.read_csv('/data/X_test_structured.csv', sep = ';')
17
18 file = open("/data/model/X_train_image_array.pkl",'rb')
19 X_train_image_array = pkl.load(file)
20 file.close()
21
22 file = open("/data/model/X_test_image_array.pkl",'rb')
23 X_test_image_array = pkl.load(file)
24 file.close()
25
26 X_train_text = pd.read_csv('/data/X_train_text.csv', sep = ';')
27 X_test_text = pd.read_csv('/data/X_test_text.csv', sep = ';')
28
29 del X_train_text['Unnamed: 0']
30 del X_test_text['Unnamed: 0']
31
32 y_train_std = pd.read_csv('/data/y_train_std.csv', sep = ';')
33 y_test_std = pd.read_csv('/data/y_test_std.csv', sep = ';')
34
35 file = open("/data/model/y_numerical_scaler.pkl",'rb')
36 y_numerical_scaler = pkl.load(file)
37 file.close()
38
39 def prepare_text_data_train(dataf):
40     for col in dataf.columns:
41         dataf[col] = dataf[col].str.lower()
42
43         tokenizer = Tokenizer(num_words=2096) #5000
44         tokenizer.fit_on_texts(dataf[col])
45         np1 = np1.append(tokenizer, open('/data/model/text_tokenizer.pkl', 'wb'))
```

```
46
47     vocab_size = len(tokenizer.word_index) + 1
48     print(vocab_size)
49
50     pk1.dump(vocab_size, open('/data/model/text_vocab_size.pkl','wb'))
51
52     maxlen = int(np.mean(dataf[col].str.len()))
53     pk1.dump(maxlen, open('/data/model/text_maxlen.pkl','wb'))
54
55 def prepare_text_data_inference(dataf):
56     for col in dataf.columns:
57         dataf[col] = dataf[col].str.lower()
58
59         file = open("/data/model/text_tokenizer.pkl",'rb')
60         tokenizer = pk1.load(file)
61         file.close()
62
63         data_tokens = tokenizer.texts_to_sequences(dataf[col])
64
65         file = open("/data/model/text_vocab_size.pkl",'rb')
66         vocab_size = pk1.load(file)
67         file.close()
68
69         file = open("/data/model/text_maxlen.pkl",'rb')
70         maxlen = pk1.load(file)
71         file.close()
72
73         data_tokens_pad = pad_sequences(data_tokens, padding='post', maxlen=max(maxlen, 300), tr
74
75     return data_tokens_pad, data_tokens, vocab_size, max(maxlen, 300)
76
77 prepare_text_data_train(X_train_text)
78 train_tokens_pad, train_tokens, vocab_size, maxlen = prepare_text_data_inference(X_train_text)
79 test_tokens_pad, test_tokens, vocab_size, maxlen = prepare_text_data_inference(X_test_text)
80
81 y_pred_img = fine_tuned_inceptionv3.predict(X_test_image_array)
82 y_pred_nlp = fine_tuned_glove_bilstm_att.predict(test_tokens_pad)
83 y_pred_structured = fine_tuned_structured.predict(X_test_structured_std)
84 y_pred_structured_nlp = fine_tuned_structured_nlp.predict([test_tokens_pad, X_test_structured_st
85 y_pred_all = final_model.predict([X_test_image_array, test_tokens_pad, X_test_structured_std])
86
87 y_true = y_numerical_scaler.inverse_transform(y_test_std)
88 y_true_train = y_numerical_scaler.inverse_transform(y_train_std)
89 y_pred_img = y_numerical_scaler.inverse_transform(y_pred_img)
```

```
90     y_pred_nlp = y_numerical_scaler.inverse_transform(y_pred_nlp)
91     y_pred_structured = y_numerical_scaler.inverse_transform(y_pred_structured)
92     y_pred_structured_nlp = y_numerical_scaler.inverse_transform(y_pred_structured_nlp)
93     y_pred_all = y_numerical_scaler.inverse_transform(y_pred_all)
94
95     y_pred_compare = pd.DataFrame(y_true)
96     y_pred_compare.columns = ['Y_TRUE']
97     y_pred_compare['Y_PRED_STRUCTURED'] = y_pred_structured
98     y_pred_compare['Y_PRED_NLP'] = y_pred_nlp
99     y_pred_compare['Y_PRED_IMG'] = y_pred_img
100    y_pred_compare['Y_PRED_HYBRID_NLP'] = y_pred_structured_nlp
101    y_pred_compare['Y_PRED_HYBRID_ALL'] = y_pred_all
102    y_pred_compare['MEAN'] = np.mean(y_true_train.ravel())
103    print('MEAN price:', np.mean(y_true_train.ravel()))
104
105
106    print('R2')
107    print('STRUCTURED', metrics.r2_score(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_STRUCTURED']))
108    print('NLP', metrics.r2_score(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_NLP']))
109    print('IMG', metrics.r2_score(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_IMG']))
110    print('HYBRID_ALL', metrics.r2_score(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_HYBRID_ALL']))
111    print('HYBRID_NLP', metrics.r2_score(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_HYBRID_NLP']))
112    print('DUMMY_MEAN_MODEL', metrics.r2_score(y_pred_compare['Y_TRUE'], y_pred_compare['MEAN']))
113    print('-----')
114    print('RMSE')
115    print('STRUCTURED', np.sqrt(metrics.mean_squared_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_STRUCTURED'])))
116    print('NLP', np.sqrt(metrics.mean_squared_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_NLP'])))
117    print('IMG', np.sqrt(metrics.mean_squared_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_IMG'])))
118    print('HYBRID_ALL', np.sqrt(metrics.mean_squared_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_HYBRID_ALL'])))
119    print('HYBRID_NLP', np.sqrt(metrics.mean_squared_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_HYBRID_NLP'])))
120    print('DUMMY_MEAN_MODEL', np.sqrt(metrics.mean_squared_error(y_pred_compare['Y_TRUE'], y_pred_compare['MEAN'])))
121    print('-----')
122    print('MAE')
123    print('STRUCTURED', metrics.mean_absolute_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_STRUCTURED']))
124    print('NLP', metrics.mean_absolute_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_NLP']))
125    print('IMG', metrics.mean_absolute_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_IMG']))
126    print('HYBRID_ALL', metrics.mean_absolute_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_HYBRID_ALL']))
127    print('HYBRID_NLP', metrics.mean_absolute_error(y_pred_compare['Y_TRUE'], y_pred_compare['Y_PRED_HYBRID_NLP']))
128    print('DUMMY_MEAN_MODEL', metrics.mean_absolute_error(y_pred_compare['Y_TRUE'], y_pred_compare['MEAN']))
```

Model	R2 Score	RMSE	MAE
Dummy Mean House Price model (calculated on train set)	0.00	280 790.14\$	187 199.09\$
Structured Data NN Architecture alone	<b>0.808</b>	<b>123 102.14\$</b>	<b>79 105.03\$</b>
Text Data NN Architecture alone	0.390	219 223.63\$	143 342.99\$
Image Data NN Architecture alone	<b>0.076</b>	<b>269 962.43\$</b>	<b>176 746.92\$</b>
Hybrid [Structured+Text] NN Architecture	0.809	122 700.20\$	77 911.38\$
Hybrid [Structured+Text+Image] NN Architecture	<b>0.783</b>	<b>130 720.34\$</b>	<b>85 008.93\$</b>

Model benchmark

The worst model of all is the dummy mean house price model where the price of each house = average price of the whole training set.

The second worst model of all is the **Image Data NN Architecture** alone where there seems to be very slight house price predictive power, meaning that house images doesn't correlate very much with the house prices.

Does it make sense ?

The price of a house is not limited to the appearance of the exterior. The local real estate market seems a lot more important. For example, if we take the same house and simply move it to another **regionCluster**, the price could be very different. Therefore, it should not be surprising that our Image Data NN Architecture alone trained on image didn't perform as well as the Structured Data NN Architecture trained on the attributes containing much more context about the local real estate market. Since our result is bad, we could use some strategies to improve our image data neural network and dataset:

- Having a dataset where photos are more consistent and taken from the same angles could be helpful.

- It is difficult to represent a full house using a single image. Having a dataset with multiple photos for each house (interior, exterior, kitchen, bathroom, etc.) could have helped our model a lot. Multiple images can be done by creating a montage that combines all images into a single image and then pass the montage through the CNN input.
- Direct regression task on images maybe not the best choice. Using bounding boxes to detect objects in images like different types of garages, windows, doors, patio, swimming pool, etc. could be an interesting alternative and surely reduced some unwanted pixel noise in images but requires much more work and time to prepare datasets.
- Maybe more epochs could have helped... choosing a different model than InceptionV3... hyper-parameter tuning... etc.

Our **Text Data NN Architecture** alone performed better than the dummy mean house price model. There seems to be predictive information / correlations on house prices when using the text description.

However, adding the textual description with our structured data (**Hybrid [Structured+Text] NN architecture**) seems to have only added noises because the performance is almost the same as structured data alone. One explanation could be that the house textual description contains redundant information from our structured features.

Third place is our **Final Hybrid neural network architecture** where images and text seem to have added noise to the structured features.

The 2 winners (almost tied) are the **Structured Data Neural Network** with a MAE of **79 105\$** and the **Hybrid (Structured + Text) Neural Network** with a

MAE of **77 911\$** (Sounds high but our Austin houses present in the dataset are expensives with an average price of **499 635\$**)

## Final words

In a context where we have better informative images and better informative texts correlated to our target, it would be very interesting to experiment again with this kind of hybrid architecture and see how it goes !

Also, keep in mind that in practice, it is difficult to combine different noise levels and conflicts between modalities. Moreover, modalities have different quantitative influence on the prediction output.

The most common method is to combine high-level embeddings from the different inputs by concatenating them, just like we did. The problem with this approach is that it would give an equal importance to all the sub-networks / modalities which is highly unlikely in real-life situations. For example, using a weighted combination of networks could be an alternative.

Multimodal

Deep Learning

NLP

Computer Vision

Neural Networks



Written by **Dave Cote, M.Sc.**

259 Followers · 34 Following

Follow

Data Scientist in an insurance company. More than 10 years in Data Science and for delivering actionable « Data-Driven » solutions.