# Stoneburner, Kurt

- ## DSC 650 - Assignment 10

Links to Deep Learning Sample Code: Word Embedding Examples:

- https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/first_edition/6.1-using-word-embeddings.ipynb (https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/first_edition/6.1-using-word-embeddings.ipynb)

RNN and LSTM Examples

- https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/first_edition/6.2-understanding-recurrent-neural-networks.ipynb (https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/first_edition/6.2-understanding-recurrent-neural-networks.ipynb)

ngram reference:

- https://www.analyticsvidhya.com/blog/2021/09/what-are-n-grams-and-how-to-implement-them-in-python/ (https://www.analyticsvidhya.com/blog/2021/09/what-are-n-grams-and-how-to-implement-them-in-python/)

**pad_sequence**: Used to transform lists within lists (2D array) to have a uniform inner dimension. Essentially, padding smaller arrays to the size of the largest, or trimming all arrays to max len

This function transforms a list (of length num_samples) of sequences (lists of integers) into a 2D Numpy array of shape (num_samples, num_timesteps). num_timesteps is either the maxlen argument if provided, or the length of the longest sequence in the list.

- https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/sequence/pad_sequences (https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/sequence/pad_sequences)

Working with tensors was a lost cause for this assignment. Keeping these for reference Convert Numpy Array to Tensor:

- https://www.projectpro.io/recipes/convert-numpy-array-tensor (https://www.projectpro.io/recipes/convert-numpy-array-tensor)

Convert Tensor to Numpy Array:

- https://www.delftstack.com/howto/numpy/python-convert-tensor-to-numpy-array/ (https://www.delftstack.com/howto/numpy/python-convert-tensor-to-numpy-array/)

Basic Text Preparation for modeling: Tokenize, ngram, Convert Text to numeric categorical (integer) value. Use Keras pad_sequences to convert integer lists to uniform lists. Feed into model.

In [1]:
```python
1   import os
2   from pathlib import Path
3   import sys
4   # //*** Imports and Load Data
5   import matplotlib.pyplot as plt
6   import numpy as np
7   import pandas as pd
8   import time
9
10  from tensorflow import keras
11  import tensorflow as tf
12  import datetime
13  from tensorflow.keras.optimizers import RMSprop
14  from tensorflow.keras import preprocessing
15  from tensorflow.keras.models import Sequential
16  from tensorflow.keras.layers import Flatten, Dense
17  from tensorflow.keras.layers import Embedding
18
19
20  #//*** Reusing Code from assignment 04
21  from chardet.universaldetector import UniversalDetector
22  from bs4 import BeautifulSoup
23
24
25  import re
26
27  #//*** Use the whole window in the IPYNB editor
28  from IPython.core.display import display, HTML
29  display(HTML("<style>.container { width:100% !important; }</style>"))
30
31  #//*** Maximize columns and rows displayed by pandas
32  pd.set_option('display.max_rows', 100)
33  pd.set_option('display.max_columns', None)
```

In [2]:
```python
1   #//***********************************************************
2   #//*** Plot a Fitted Models History of Loss and Accuracy
3   #//***********************************************************
4   def plot_model_history(input_history):
5       loss_key, acc_key = list(input_history.history.keys())[:2]
6       val_loss_key, val_acc_key = list(input_history.history.keys())[2:
7
8       acc = input_history.history[acc_key]
9       loss = input_history.history[loss_key]
10
11      val_loss = input_history.history[val_loss_key]
12      val_acc = input_history.history[val_acc_key]
13
14      epochs = range(1, len(loss) + 1)
15      plt.plot(epochs, acc, "b", label="Training Accuracy")
16      plt.title("Training Accuracy\nAccuracy should go up")
17      plt.xlabel("Epochs")
18      plt.ylabel("Loss")
19      plt.legend()
20      plt.show()
21
```

```
22        plt.plot(epochs, loss, "bo", label="Training Loss")
23
24        plt.title("Training Loss \nLoss should go down")
25        plt.xlabel("Epochs")
26        plt.ylabel("Loss")
27        plt.legend()
28        plt.show()
29
30
31        plt.plot(epochs, loss, "bo", label="Training loss")
32        plt.plot(epochs, val_loss, "b", label="Validation loss")
33        plt.title("Training and validation loss")
34        plt.xlabel("Epochs")
35        plt.ylabel("Loss")
36        plt.legend()
37        plt.show()
38
39        #//*** Plot the Validation Set Accuracy
40        plt.clf()
41
42        plt.plot(epochs, acc, "bo", label="Training accuracy")
43        plt.plot(epochs, val_acc, "b", label="Validation accuracy")
44        plt.title("Training and validation accuracy")
45        plt.xlabel("Epochs")
46        plt.ylabel("Accuracy")
47        plt.legend()
48        plt.show()
49
```

In [3]:
```
1  #//*** Get Working Directory
2  current_dir = Path(os.getcwd()).absolute()
3
4  #//*** Go up Two folders
5  project_dir = current_dir.parents[2]
6
7  #//*** IMDB Data Path
8  imdb_path = project_dir.joinpath("dsc650/data/external/imdb/aclImdb")
9
10 file_path = imdb_path.joinpath("train/pos")
11
12 #//*** Grab the first positive review text for testing
13 file_path = file_path.joinpath(os.listdir(file_path)[0])
14
15 with open(file_path,'r') as f:
16     sample_text = f.read()
17
18 print(sample_text)
```

Bromwell High is a cartoon comedy. It ran at the same time as some oth
er programs about school life, such as "Teachers". My 35 years in the
teaching profession lead me to believe that Bromwell High's satire is
much closer to reality than is "Teachers". The scramble to survive fin

In [4]:
```python
#//*** Randomly assign 20% of the training Data and copy to a validat
import os, pathlib, shutil, random

val_dir = imdb_path.joinpath("val")
train_dir = imdb_path.joinpath("train")
test_dir = imdb_path.joinpath("test")

for category in ("neg", "pos"):
    #//*** Skip if val folder exists (Delete Folder to resample)
    if os.path.exists(val_dir.joinpath(category)):
        break

    os.makedirs(val_dir.joinpath(category))
    files = os.listdir(train_dir.joinpath(category))
    random.Random(1337).shuffle(files)
    num_val_samples = int(0.2 * len(files))
    val_files = files[-num_val_samples:]
    for fname in val_files:
        shutil.move(train_dir / category / fname,
                    val_dir / category / fname)

```

## Load IMDB Dataset

In [5]:
```python
#//*** Use Universal Detector to determine file encoding.
#//*** Borrowed from Assignment04
def read_file_with_encoding(filepath):

    detector = UniversalDetector()

    try:
        with open(filepath) as f:
            return f.read()
    except UnicodeDecodeError:
        detector.reset()
        with open(filepath, 'rb') as f:
            for line in f.readlines():
                detector.feed(line)
                if detector.done:
                    break
        detector.close()
        encoding = detector.result['encoding']
        with open(filepath, encoding=encoding) as f:
            return f.read()

#//*** Borrowed from Assignment04
def parse_html_payload(payload):
    """
    This function uses Beautiful Soup to read HTML data
```

```
26          and return the text.  If the payload is plain text, then
27          Beautiful Soup will return the original content
28          """
29          soup = BeautifulSoup(payload, 'html.parser')
30          return str(soup.get_text()).encode('utf-8').decode('utf-8')
31
32  def load_dataset(dir_path):
33
34          text = []
35          targets = []
36
37          #//*** Crawl the neg and pos folders
38          for category in ("neg", "pos"):
39              files = os.listdir(dir_path.joinpath(category))
40
41              #//*** Loop through each file in the folder
42              for file in files:
43                  try:
44                      #//*** Add processed file to text
45                      text.append(
46                          #//*** Strip HTML Tags
47                          parse_html_payload(
48                              #//*** Read File from disk. Function uses Uni
49                              read_file_with_encoding(
50                                  dir_path.joinpath(category).joinpath(file
51
52                      #//*** Append Target Value
53                      if category == 'neg':
54                          targets.append(0)
55                      else:
56                          targets.append(1)
57                  except:
58                      print(f"Dropping File: {file} due to decoding issues"
59
60          #//*** Targets returned as Numpy float32 array
61          return text, np.asarray( targets).astype('float32')
62
```

## Assignment 10.1

```
In [27]:   1  #//*** Vectorize a corpus
           2  class Vectorizer:
           3      def __init__(self,**kwargs):
           4          self.corpus_tokens = []
           5          self.corpus_ngrams = []
           6
           7          self.max_tokens = None
           8          self.ngram_size = 1
           9          self.tidyup = True
          10
          11          self.max_element_count = -1
          12
          13          for key,value in kwargs.items():
          14              if key =="max_tokens":
          15                  self.max_tokens = value
```

```
16
17                    if key == "ngrams":
18                        self.ngram_size = value
19
20                    if key == "tidyup":
21                        self.tidyup = value
22
23
24            #//*** One Hot Encoding Dictionaries
25            #//*** Key = Token Index, Value = Word
26            self.ngram_index = {}
27
28            #//*** Key = Word, Value = Token Index
29            self.vocabulary_index = {}
30
31        def tokenize(self,raw_text):
32            #//*** Initialize Output Tokens
33            tokens = []
34
35            #//*** Split Text into words
36            for x in re.split("\s",raw_text):
37
38                #//*** Findall Non text characters in each word
39                non_text = re.findall("\W",x)
40
41                #//*** Remove non_text Characters
42                for i in non_text:
43                    x = x.replace(i,"")
44
45                #//*** If X has length, append out
46                if len(x) > 0:
47                    tokens.append(x.lower())
48            return tokens
49
50        def build_ngrams(self):
51            if self.ngram_size <= 0:
52                print("Ngram size must be an integer > 0")
53                print("Quitting!")
54                return None
55
56            #//*** Using unigrams, use tokens
57            if self.ngram_size == 1:
58                self.corpus_ngrams = self.corpus_tokens
59                return
60
61            self.corpus_ngrams = []
62
63            #//*** Get each token group from corpus_tokens
64            for token in self.corpus_tokens:
65
66                loop_ngram = []
67
68                #//*** Use an index based range to loop through tokens
69                for x in range(0,len(token) ):
70
71                    #//*** Check if index + ngram_size exceeds the length
```

```python
 72                         if x+self.ngram_size <= len(token):
 73
 74                             result = ""
 75
 76                             #//*** Build the ngram
 77                             for y in range(self.ngram_size):
 78                                 #print(self.tokens[x+y])
 79                                 result += token[x+y] + " "
 80
 81                             loop_ngram.append(result[:-1])
 82
 83                         else:
 84                             break
 85
 86                 #//*** Grab Token Element Count, Keep the greatest count
 87                 if len(loop_ngram) > self.max_element_count:
 88                     self.max_element_count = len(loop_ngram)
 89
 90
 91                 #//*** Token group ngram is built. Add loop_ngram to corp
 92                 self.corpus_ngrams.append(loop_ngram)
 93
 94
 95
 96     def build_vocabulary(self,corpus):
 97         if not isinstance(corpus,list) :
 98             print("Vectorizer Requires a corpus (list of text):")
 99             return None
100
101         self.tokens = []
102
103         print("Tokenizing...")
104         #//*** Tokenize each text entry in the corpus
105         for raw_text in corpus:
106             self.corpus_tokens.append(self.tokenize(raw_text))
107
108         print("Building ngrams...")
109         #//*** Build ngrams (Defaults to unigrams)
110         self.build_ngrams()
111
112         word_freq = {}
113
114         print("Building Vocabulary...")
115         #//*** Build dictionary of unique words
116         #//*** Loop through each element of the corpus
117         for element in self.corpus_ngrams:
118
119
120             #//*** Grab Token Element Count, Keep the greatest count
121             if len(element) > self.max_element_count:
122                 self.max_element_count = len(element)
123
124             #//*** Process each individual ngram
125             for ngram in element:
126
127
```

```python
128
129                          #//*** Add unique words to dictionaries
130                          if ngram not in self.vocabulary_index.keys():
131                              index = len(self.ngram_index.values())
132                              self.ngram_index[ index ] = ngram
133                              self.vocabulary_index [ ngram ] = index
134
135                              #//*** Initialize Word Frequency
136                              word_freq[ ngram ] = 1
137                          else:
138                              #//*** Increment Word Frequency
139                              word_freq[ ngram ] += 1
140
141          #//*** END for element in self.corpus_ngrams:
142          if self.max_tokens != None:
143
144              #//*** Check if token count exceeds max tokens
145              if self.max_tokens < len(self.ngram_index.items()):
146
147                  print("Sorting Word Frequency...")
148                  #//*** Sort the Word Frequency Dictionary. Keep the l
149                  word_freq = dict(sorted(word_freq.items(), key=lambda
150
151                  print("Building Token Dictionary")
152                  #//*** Get list of keys that are lowest frequency
153                  for key in list(word_freq.keys())[self.max_tokens:]:
154                      #//*** Delete Low Frequency ngrams
155                      del word_freq[ key ]
156
157                  self.ngram_index = {}
158                  self.vocabulary_index = {}
159
160                  print("Rebuilding Vocabulary")
161                  #//*** Rebuild ngram_index & vocabulary_index
162                  for ngram in word_freq.keys():
163                      index = len(self.ngram_index.values())
164                      self.ngram_index[ index ] = ngram
165                      self.vocabulary_index [ ngram ] = index
166
167              #//*** END Trim Low Frequency ngrams
168          self.word_freq = word_freq
169
170      #//*** One Hot encode the corpus.
171      #//*** Handling the corpus as a whole increases processing speed
172      #//*** Hot encode to a sparse tensor to for increased encoding sp
173      def one_hot_encode(self,corpus):
174
175          #//*** Encoded Results
176          results = []
177
178          #//*** Set the Max array size to the total number of items in
179          array_size = len(self.ngram_index.keys())
180
181
182          start_time = datetime.datetime.now()
183          count = 0
```

```python
184
185
186
187          for element in corpus:
188              #//*** hot encode each ngram
189              result = []
190              for ngram in element:
191
192                  #//*** Skip words not in self.vocabulary_index
193                  #//*** These are skipped due to max_tokens limitation
194                  if ngram not in self.vocabulary_index.keys():
195                      continue
196
197                  sparse_tensor = tf.SparseTensor(indices=[[0,self.voca
198                  #index = self.vocabulary_index[ngram]
199
200                  #base_array = np.zeros(array_size, dtype=int)
201
202                  #base_array [index] = 1
203
204
205                  #//*** Add the one-hot-encoded word to encoded text
206                  result.append(sparse_tensor)
207
208
209              #//*** END for ngram in tokens:
210
211              result = tf.sparse.concat(axis=1, sp_inputs=result)
212              #//*** concat Sparse Matrix
213              results.append( result )
214
215              count += 1
216
217
218
219              #//*** Print a status update every 1000 items
220              if count % 100 == 0:
221                  print(f"{count} / {len(corpus)} Encoded: {datetime.da
222
223          #//*** Concat List of Sparse Matrixes into a sparse matrix
224          #results =  tf.sparse.concat(axis=1, sp_inputs=results)
225
226          print(f"Encoding Complete: {datetime.datetime.now() - start_t
227
228          return results
229
230      #//***
231      def integer_encode(self,corpus):
232          #//*** Encoded Results
233          results = []
234
235          #//*** Set the Max array size to the total number of items in
236          array_size = len(self.ngram_index.keys())
237
238
239          start_time = datetime.datetime.now()
```

```python
240            count = 0
241
242
243
244            for element in corpus:
245                #//*** hot encode each ngram
246                result = []
247                for ngram in element:
248
249                    #//*** Skip words not in self.vocabulary_index
250                    #//*** These are skipped due to max_tokens limitation
251                    if ngram not in self.vocabulary_index.keys():
252                        continue
253
254                    #//*** Get integer value of ngram from dictionary.
255                    #//*** Add to result
256                    result.append(self.vocabulary_index[ngram])
257
258
259                #//*** END for ngram in tokens:
260
261                #//*** result is a complete encoded element
262                results.append( np.array(result).astype(np.float32) )
263
264                count += 1
265
266
267
268                #//*** Print a status update every 1000 items
269                if count % 5000 == 0:
270                    print(f"{count} / {len(corpus)} Encoded: {datetime.da
271
272            print(f"Encoding Complete: {datetime.datetime.now() - start_t
273
274            #//*** results is a collection of encoded elements
275            return np.array(results,dtype=object)
276
277
278        def encode(self,corpus,encoding='int'):
279
280            if not isinstance(corpus,list) :
281                print("Vectorizer Requires a corpus (list of text):")
282                return None
283
284            self.corpus_tokens = []
285            self.corpus_ngrams = []
286            print("Tokenizing...")
287            #//*** Tokenize each text entry in the corpus
288            for raw_text in corpus:
289                self.corpus_tokens.append(self.tokenize(raw_text))
290
291            print("Building ngrams...")
292            #//*** Build ngrams (Defaults to unigrams)
293            self.build_ngrams()
294
295            if encoding == 'onehot':
```

```
296                     print("One Hot Coding....")
297
298                     #//*** One Hot Encode Values. These are actually sparse t
299                     encoded = self.one_hot_encode(self.corpus_ngrams)
300
301             if encoding == 'int':
302                     print("Interger encoding....")
303
304                     #//*** Convert ngrams to integers. These are actually spa
305                     encoded = self.integer_encode(self.corpus_ngrams)
306
307                     #//*** Convert lists to Numpy array of float 32 type. Thi
308                     #encoded = np.asarray(encoded).astype('float32')
309
310             #//*** TidyUp (Delete) ngrams and Tokens
311             if self.tidyup:
312                     self.corpus_tokens = []
313                     self.corpus_ngrams = []
314
315
316             return encoded
317
318       #//*** Convert One-Hot-Encoding to text
319       def decode(self,elements):
320
321             results = []
322
323             #//*** For Each element in Corpus
324
325             decoded = ""
326
327             #//*** For Each ngram (word(s)) in Elements
328             for ngram in elements:
329
330                     #//*** Grab Index of 1 from sparse tensor
331                     index = ngram.indices[0].numpy()[1]
332
333                     #ngram = list(ngram.numpy())
334
335                     decoded += self.ngram_index[ index ] + " "
336
337             #//*** END for ngram in elements:
338             results.append( decoded[:-1])
339
340             #//*** END for elements in corpus:
341             return results
342
343   print("Loading Raw Validation Set")
344   val_x_train, val_y_train = load_dataset(val_dir)
345
346   #//*** Test the Vectorizer with some sample data
347   vectorizer = Vectorizer(max_tokens=100,ngrams=2, tidyup=False)
348   vectorizer.build_vocabulary(val_x_train[:5])
349   start_time = datetime.datetime.now()
350
351   temp_vals = vectorizer.encode(val_x_train[:5],encoding='onehot')
```

```
352
353    print(f"Run Time: {datetime.datetime.now() - start_time}")
354
355
356    int_vals = vectorizer.encode(val_x_train[:5],encoding='int')
357    print("Integer Encoding:")
358    print(int_vals)
359    print()
360    print()
361
362    print("Sample Text: (First 500 Chars)")
363    for element in val_x_train[:5]:
364        print(element[:500])
365        print("====")
366    print()
367    print()
368
369    print("Tokens: (First 100 tokens)")
370    for token in vectorizer.corpus_tokens:
371        print(token[:100])
372        print("====")
373    print()
374    print()
375
376    print("ngrams: (First 50 tokens)")
377    for token in vectorizer.corpus_ngrams:
378        print(token[:100])
379        print("====")
380    print()
381    print()
382    print("Small one hot encoded Sample:")
383    print(temp_vals)
384    print()
385    print()
386    print("Encoded Vocabulary")
387    print(vectorizer.vocabulary_index)
388    print()
389    print()
390    print("Decoded Text from vocabulary (limited by max tokens)")
391    for result in vectorizer.decode(temp_vals):
392        print(result)
393        print()
394
395    del temp_vals
396    del int_vals
397    del vectorizer
```

```
Loading Raw Validation Set
Tokenizing...
Building ngrams...
Building Vocabulary...
Sorting Word Frequency...
Building Token Dictionary
Rebuilding Vocabulary
Tokenizing...
Building ngrams...
One Hot Coding....
Encoding Complete: 0:00:00.046011
Run Time: 0:00:00.051012
Tokenizing...
```

```
Building ngrams...
Interger encoding....
Encoding Complete: 0:00:00
Integer Encoding:
[array([ 2., 76., 77.,  5., 78., 79., 80., 81., 82., 21., 83., 84., 8
5.,
```

In [7]:
```python
1  print("Loading Raw Validation Set")
2  raw_val_x_train, val_y_train = load_dataset(val_dir)
3
4  print("Loading Raw Train Data")
5  raw_x_train, y_train = load_dataset(train_dir)
6
7  print("Loading Raw Test Data")
8  raw_x_test, y_test = load_dataset(test_dir)
9  print("Done")
10
11 val_y_train = np.array(val_y_train,dtype=object)
12 y_train = np.array(y_train,dtype=object)
13 y_test = np.array(y_test,dtype=object)
14
15 y_train = np.asarray(y_train).astype(np.int)
```

```
Loading Raw Validation Set
Loading Raw Train Data
Dropping File: 7714_1.txt due to decoding issues
Dropping File: 11351_9.txt due to decoding issues
Dropping File: 8263_9.txt due to decoding issues
Loading Raw Test Data
Dropping File: 4414_1.txt due to decoding issues
Dropping File: 6973_1.txt due to decoding issues
Dropping File: 2464_10.txt due to decoding issues
Dropping File: 5281_10.txt due to decoding issues
Done
```

In [8]:
```python
1
2
3
4  #//*** Test the Vectorizer with some sample data
5  max_tokens = 20000
6  #maxlen = 1000 #//*** Limit reviews to this length, Leave blank to us
7
8  ngrams = 1
9
10 #//*** Initialize vectorizer
11 vectorizer = Vectorizer(max_tokens=max_tokens,ngrams=ngrams)
12
13 #//*** Build Vocabulary based on the training text
14 vectorizer.build_vocabulary(raw_x_train)
15
16 #//*** maxlen required: This is maximum number of tokens/ngrams to us
17 #//*** pads all articles to the same word count. This required to hav
18
19 maxlen = vectorizer.max_element_count #//*** Sets maxlen to the large
20
21
```

```
22  #//*** Encode Validation, training and test data
23
24  print("Encoding Validation Data...")
25  val_x_train = vectorizer.encode(raw_val_x_train)
26  print(vectorizer.max_element_count)
27  print("Encoding Training Data...")
28  x_train = vectorizer.encode(raw_x_train)
29
30  print("Encoding Test Data...")
31  x_test = vectorizer.encode(raw_x_test)
32
33
34  print("Padding Validation...")
35  val_x_train = preprocessing.sequence.pad_sequences(val_x_train, maxle
36
37  print("Padding Training Data...")
38  x_train = preprocessing.sequence.pad_sequences(x_train, maxlen=maxlen
39
40  print("Padding Test Data...")
41  x_test = preprocessing.sequence.pad_sequences(x_test, maxlen=maxlen)
42
43  val_y_train = np.asarray(val_y_train).astype(np.int)
44  x_test = np.asarray(x_test).astype(np.int)
45  y_test = np.asarray(y_test).astype(np.int)
46
47  print("Done")
48
49  print(type(val_x_train),type(val_x_train[0]),x_train.shape)
50  print(type(val_y_train),type(val_y_train[0]),y_train.shape)
51  print(type(x_test),type(x_test[0]))
52  print(type(y_test),type(y_test[0]))
53  print(x_test.shape,y_test.shape)
```

```
Tokenizing...
Building ngrams...
Building Vocabulary...
Sorting Word Frequency...
Building Token Dictionary
Rebuilding Vocabulary
Encoding Validation Data...
Tokenizing...
Building ngrams...
Interger encoding....
5000 / 5000 Encoded: 0:00:00.496111
Encoding Complete: 0:00:00.496111
2450
Encoding Training Data...
Tokenizing...
Building ngrams...
Interger encoding....
5000 / 19997 Encoded: 0:00:00.484107
10000 / 19997 Encoded: 0:00:00.977218
15000 / 19997 Encoded: 0:00:01.466329
Encoding Complete: 0:00:01.963440
Encoding Test Data...
Tokenizing...
```

```
Building ngrams...
Interger encoding....
5000 / 24996 Encoded: 0:00:00.477108
10000 / 24996 Encoded: 0:00:00.949213
15000 / 24996 Encoded: 0:00:01.431323
20000 / 24996 Encoded: 0:00:01.913430
Encoding Complete: 0:00:02.384537
Padding Validation...
Padding Training Data...
Padding Test Data...
Done
<class 'numpy.ndarray'> <class 'numpy.ndarray'> (19997, 2450)
<class 'numpy.ndarray'> <class 'numpy.int32'> (19997,)
<class 'numpy.ndarray'> <class 'numpy.ndarray'>
<class 'numpy.ndarray'> <class 'numpy.int32'>
(24996, 2450) (24996,)
```

In [9]:
```python
1  print(type(val_x_train),type(val_x_train[0]),x_train.shape)
2  print(type(val_y_train),type(val_y_train[0]),y_train.shape)
3  print(type(x_test),type(x_test[0]))
4  print(type(y_test),type(y_test[0]))
5  print(x_test.shape,y_test.shape)
```

```
<class 'numpy.ndarray'> <class 'numpy.ndarray'> (19997, 2450)
<class 'numpy.ndarray'> <class 'numpy.int32'> (19997,)
<class 'numpy.ndarray'> <class 'numpy.ndarray'>
<class 'numpy.ndarray'> <class 'numpy.int32'>
(24996, 2450) (24996,)
```

https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/first_edition/6.1-using-word-embeddings.ipynb (https://github.com/fchollet/deep-learning-with-python-notebooks/blob/master/first_edition/6.1-using-word-embeddings.ipynb)

```
In [10]:    1  """
            2
            3  from keras.datasets import imdb
            4  from keras import preprocessing
            5
            6  # Number of words to consider as features
            7  max_features = 10000
            8  # Cut texts after this number of words
            9  # (among top max_features most common words)
           10  maxlen = 20
           11
           12  # Load the data as lists of integers.
           13  (x_train, y_train), (x_test, y_test) = imdb.load_data(num_words=max_f
           14
           15  # This turns our lists of integers
           16  # into a 2D integer tensor of shape `(samples, maxlen)`
           17  #x_train = preprocessing.sequence.pad_sequences(x_train, maxlen=maxle
           18  #x_test = preprocessing.sequence.pad_sequences(x_test, maxlen=maxlen)
           19  print(type(x_train),type(x_train[0]))
           20  print(type(y_train),type(y_train[0]))
           21  """
           22  print()
```

**pad_sequence**: Used to transform lists within lists (2D array) to have a uniform inner dimension. Essentially, padding smaller arrays to the size of the largest, or trimming all arrays to max len

This function transforms a list (of length num_samples) of sequences (lists of integers) into a 2D Numpy array of shape (num_samples, num_timesteps). num_timesteps is either the maxlen argument if provided, or the length of the longest sequence in the list.

- https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/sequence/pad_sequences (https://www.tensorflow.org/api_docs/python/tf/keras/preprocessing/sequence/pad_sequences)

# 10.2

Using listings 6.16, 6.17, and 6.18 in Deep Learning with Python as a guide, train a sequential model with embeddings on the IMDB data found in data/external/imdb/. Produce the model performance metrics and training and validation accuracy curves within the Jupyter notebook.

```
In [11]:    1
            2
            3  def get_sequential_embedded_model():
            4      model = Sequential()
            5      # We specify the maximum input length to our Embedding layer
            6      # so we can later flatten the embedded inputs
            7      model.add(Embedding(max_tokens, 8, input_length=maxlen))
            8      # After the Embedding layer,
            9      # our activations have shape `(samples, maxlen, 8)`.
           10
```

```
11        # We flatten the 3D tensor of embeddings
12        # into a 2D tensor of shape `(samples, maxlen * 8)`
13        model.add(Flatten())
14
15        # We add the classifier on top
16        model.add(Dense(1, activation='sigmoid'))
17        model.compile(optimizer='rmsprop', loss='binary_crossentropy', me
18
19        return model
20
21  model = get_sequential_embedded_model()
22  model.summary()
23
24  #//*** Code to check the data type expected for each model layer
25  [print(i.shape, i.dtype) for i in model.inputs]
26  [print(o.shape, o.dtype) for o in model.outputs]
27  [print(l.name, l.input_shape, l.dtype) for l in model.layers]
28
29  #//*** Reference for recasting lists as an np.array of float32 type
30  #x = np.asarray(x).astype('float32')
31
```

```
Model: "sequential"

_____
Layer (type)                 Output Shape              Param #
=================================================================
embedding (Embedding)        (None, 2450, 8)           160000
_____
flatten (Flatten)            (None, 19600)             0
_____
dense (Dense)                (None, 1)                 19601
=================================================================
Total params: 179,601
Trainable params: 179,601
Non-trainable params: 0
_____
(None, 2450) <dtype: 'float32'>
(None, 1) <dtype: 'float32'>
embedding (None, 2450) float32
flatten (None, 2450, 8) float32
dense (None, 19600) float32
```

Out[11]: [None, None, None]

In [12]:
```
1  history = model.fit(val_x_train,val_y_train,
2                      epochs=10,
3                      batch_size=32,
4                      validation_split=0.2
5                      )
6
7  score = model.evaluate(np.asarray(x_test).astype(np.int), np.asarray(
8
9  print('Test loss:', score[0])
10  print('Test accuracy:', score[1])
11
12  plot_model_history(history)
```
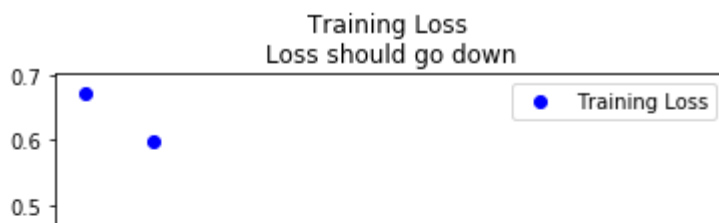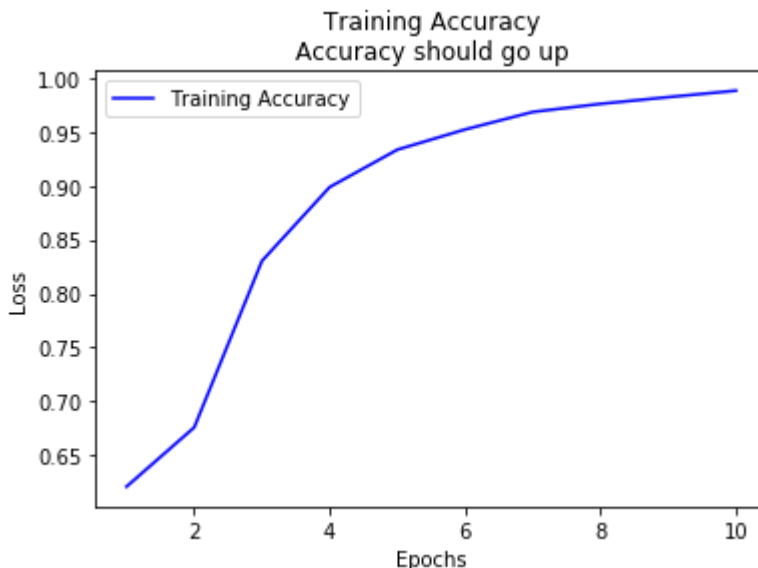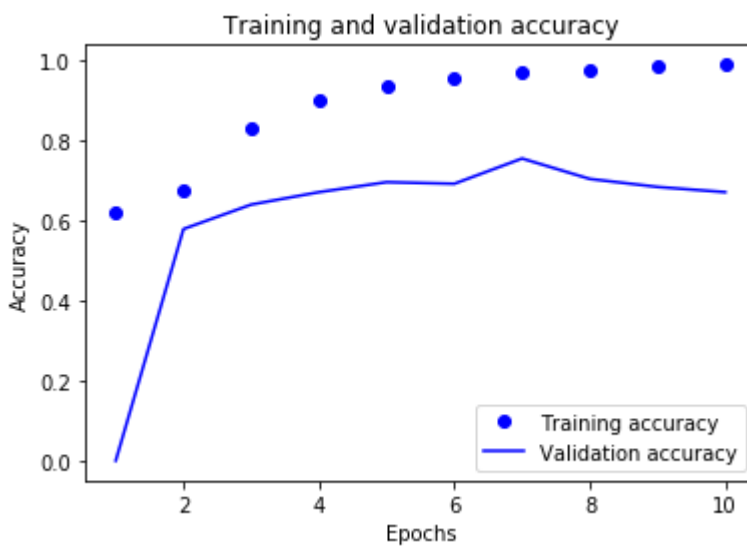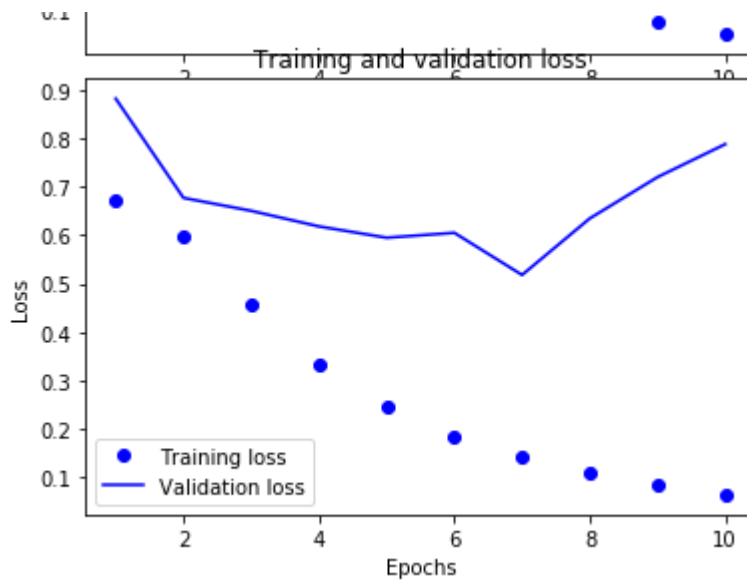
```
Epoch 1/10
125/125 [==============================] - 2s 9ms/step - loss: 0.6708
- acc: 0.6208 - val_loss: 0.8819 - val_acc: 0.0000e+00
Epoch 2/10
125/125 [==============================] - 1s 7ms/step - loss: 0.5973
- acc: 0.6758 - val_loss: 0.6769 - val_acc: 0.5780
Epoch 3/10
125/125 [==============================] - 1s 8ms/step - loss: 0.4582
- acc: 0.8303 - val_loss: 0.6502 - val_acc: 0.6390
Epoch 4/10
125/125 [==============================] - 1s 7ms/step - loss: 0.3333
- acc: 0.8988 - val_loss: 0.6180 - val_acc: 0.6700
Epoch 5/10
125/125 [==============================] - 1s 7ms/step - loss: 0.2458
- acc: 0.9335 - val_loss: 0.5947 - val_acc: 0.6950
Epoch 6/10
125/125 [==============================] - 1s 7ms/step - loss: 0.1850
- acc: 0.9523 - val_loss: 0.6048 - val_acc: 0.6910
Epoch 7/10
125/125 [==============================] - 1s 7ms/step - loss: 0.1410
- acc: 0.9685 - val_loss: 0.5178 - val_acc: 0.7540
Epoch 8/10
125/125 [==============================] - 1s 7ms/step - loss: 0.1088
- acc: 0.9760 - val_loss: 0.6345 - val_acc: 0.7030
Epoch 9/10
125/125 [==============================] - 1s 7ms/step - loss: 0.0836
- acc: 0.9822 - val_loss: 0.7198 - val_acc: 0.6830
Epoch 10/10
125/125 [==============================] - 1s 7ms/step - loss: 0.0642
- acc: 0.9883 - val_loss: 0.7878 - val_acc: 0.6700
Test loss: 0.4089915454387665
Test accuracy: 0.8279724717140198
```



Training Accuracy
Accuracy should go up



Training Loss
Loss should go down

### Training and validation loss



### Training and validation accuracy



```
In [13]:    1  history = model.fit(x_train,y_train,
            2                       epochs=8,
            3                       batch_size=32,
            4                       validation_split=0.2
            5                      )
            6
            7  score = model.evaluate(np.asarray(x_test).astype(np.int), np.asarray(
            8
            9  print('Test loss:', score[0])
           10  print('Test accuracy:', score[1])
           11
```

```
Epoch 1/8
500/500 [==============================] - 4s 8ms/step - loss: 0.2878
- acc: 0.8827 - val_loss: 0.4963 - val_acc: 0.7975
Epoch 2/8
500/500 [==============================] - 4s 8ms/step - loss: 0.2160
- acc: 0.9155 - val_loss: 0.5547 - val_acc: 0.7803
Epoch 3/8
500/500 [==============================] - 4s 8ms/step - loss: 0.1730
- acc: 0.9362 - val_loss: 0.4601 - val_acc: 0.8198
Epoch 4/8
500/500 [==============================] - 4s 8ms/step - loss: 0.1423
- acc: 0.9492 - val_loss: 0.4608 - val_acc: 0.8290
Epoch 5/8
500/500 [==============================] - 4s 8ms/step - loss: 0.1186
- acc: 0.9595 - val_loss: 0.6408 - val_acc: 0.7710
Epoch 6/8
500/500 [==============================] - 4s 8ms/step - loss: 0.0997
```

In [ ]:  1

## 10.3

Using listing 6.27 in Deep Learning with Python as a guide, fit the same data with an LSTM layer.
Produce the model performance metrics and training and validation accuracy curves within the
Jupyter notebook.

In [14]:
```python
1
2
3  from keras.models import Sequential
4  from keras.layers import Embedding, SimpleRNN,LSTM
5
6  def get_lstm_model():
7
8      model = Sequential()
9      model.add(Embedding(max_tokens, 32))
10     model.add(LSTM(32))
11     model.add(Dense(1, activation='sigmoid'))
12
13     model.compile(optimizer='rmsprop',
14                   loss='binary_crossentropy',
15                   metrics=['acc'])
16
17
18 history = model.fit(val_x_train, val_y_train,
19                 epochs=10,
20                 batch_size=128,
21                 validation_split=0.2)
22
23 score = model.evaluate(np.asarray(x_test).astype(np.int), np.asarray(
24
25 print('Test loss:', score[0])
26 print('Test accuracy:', score[1])
27
28 plot model history(history)
```

```
Epoch 1/10
32/32 [==============================] - 1s 27ms/step - loss: 0.0881 -
acc: 0.9685 - val_loss: 0.7499 - val_acc: 0.7710
Epoch 2/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0703 -
acc: 0.9770 - val_loss: 0.7489 - val_acc: 0.7780
Epoch 3/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0569 -
acc: 0.9830 - val_loss: 0.8917 - val_acc: 0.7380
Epoch 4/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0464 -
acc: 0.9875 - val_loss: 0.8226 - val_acc: 0.7510
Epoch 5/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0377 -
acc: 0.9902 - val_loss: 0.8585 - val_acc: 0.7510
Epoch 6/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0313 -
acc: 0.9933 - val_loss: 1.0313 - val_acc: 0.7050
Epoch 7/10
32/32 [==============================] - 1s 27ms/step - loss: 0.0255 -
acc: 0.9950 - val_loss: 0.9157 - val_acc: 0.7430
Epoch 8/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0208 -
acc: 0.9965 - val_loss: 0.9251 - val_acc: 0.7420
Epoch 9/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0171 -
acc: 0.9975 - val_loss: 1.0651 - val_acc: 0.7080
```
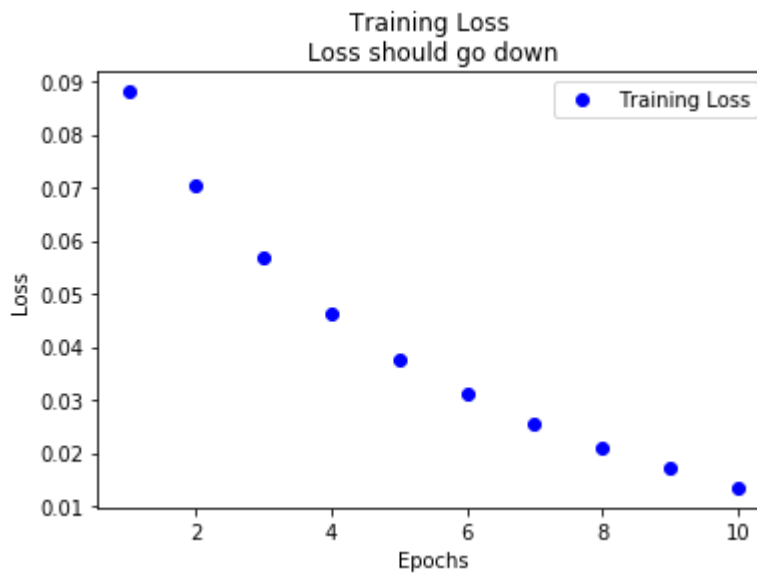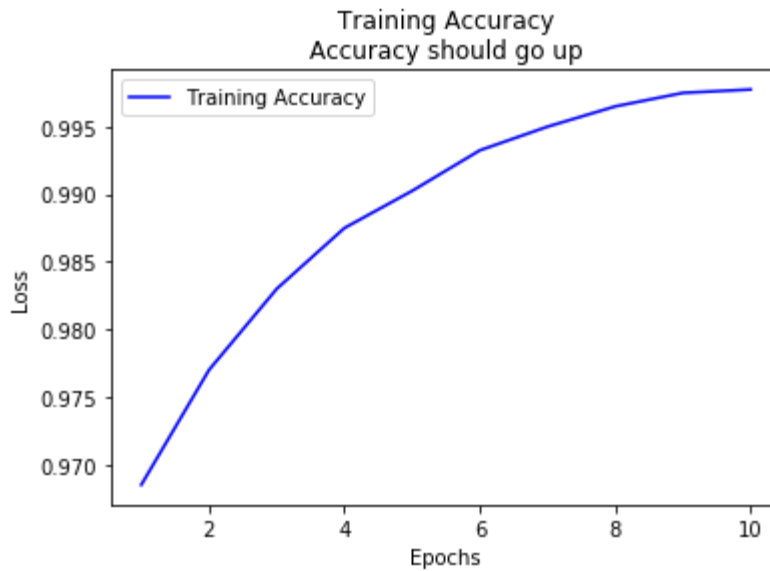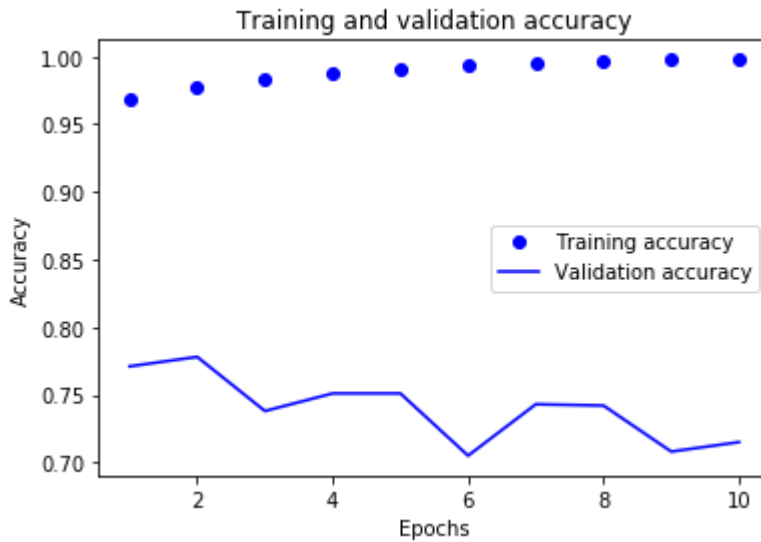
```
Epoch 10/10
32/32 [==============================] - 1s 26ms/step - loss: 0.0135 -
acc: 0.9977 - val_loss: 1.0511 - val_acc: 0.7150
Test loss: 0.4457803964614868
Test accuracy: 0.8531364798545837
```

In [15]:
```python
1  history = model.fit(x_train, y_train,
2                      epochs=8,
3                      batch_size=128,
4                      validation_split=0.2)
5
6  score = model.evaluate(np.asarray(x_test).astype(np.int), np.asarray(
7
8  print('Test loss:', score[0])
```

```
Epoch 1/8
125/125 [==============================] - 3s 27ms/step - loss: 0.0648
- acc: 0.9773 - val_loss: 0.8485 - val_acc: 0.7442
Epoch 2/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0523
- acc: 0.9834 - val_loss: 0.7618 - val_acc: 0.7663
Epoch 3/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0432
- acc: 0.9874 - val_loss: 0.8154 - val_acc: 0.7508
Epoch 4/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0359
- acc: 0.9905 - val_loss: 0.7709 - val_acc: 0.7663
Epoch 5/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0297
- acc: 0.9924 - val_loss: 0.8048 - val_acc: 0.7607
Epoch 6/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0247
- acc: 0.9942 - val_loss: 0.9268 - val_acc: 0.7372
Epoch 7/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0203
- acc: 0.9955 - val_loss: 0.9174 - val_acc: 0.7420
Epoch 8/8
125/125 [==============================] - 3s 26ms/step - loss: 0.0167
- acc: 0.9966 - val_loss: 0.8256 - val_acc: 0.7688
Test loss: 0.46850863099098206
Test accuracy: 0.8548167943954468
```

# 10.4

Using listing 6.46 in Deep Learning with Python as a guide, fit the same data with a simple 1D convnet. Produce the model performance metrics and training and validation accuracy curves within the Jupyter notebook.

In [24]:
```python
from keras import layers

def get_conv1d_model():
    model = Sequential()
    model.add(layers.Embedding(max_tokens, 128, input_length=maxlen))
    model.add(layers.Conv1D(32, 7, activation='relu'))
    model.add(layers.MaxPooling1D(5))
    model.add(layers.Conv1D(32, 7, activation='relu'))
    model.add(layers.GlobalMaxPooling1D())
    model.add(layers.Dense(1))

    model.compile(optimizer=RMSprop(learning_rate=1e-4),
                  loss='binary_crossentropy',
                  metrics=['acc'])

    model.summary()
    return model


model = get_conv1d_model()
history = model.fit(val_x_train, val_y_train,
                    epochs=30,
                    batch_size=128,
                    validation_split=0.2)

score = model.evaluate(np.asarray(x_test).astype(np.int), np.asarray(

print('Test loss:', score[0])
print('Test accuracy:', score[1])

```

```
Model: "sequential_6"
_____
Laver (tvpe)                  Output Shape              Param #
```

In [25]:
```python
1  history = model.fit(x_train, y_train,
2                      epochs=10,
3                      batch_size=128,
4                      validation_split=0.2)
5
6  score = model.evaluate(np.asarray(x_test).astype(np.int), np.asarray(
7
8  print('Test loss:', score[0])
9  print('Test accuracy:', score[1])
```

```
Epoch 1/10
125/125 [==============================] - 99s 789ms/step - loss: 3.13
18 - acc: 0.6251 - val_loss: 1.2540 - val_acc: 0.0000e+00
Epoch 2/10
125/125 [==============================] - 97s 775ms/step - loss: 0.65
54 - acc: 0.6251 - val_loss: 0.9645 - val_acc: 0.0000e+00
Epoch 3/10
125/125 [==============================] - 99s 789ms/step - loss: 0.63
06 - acc: 0.6251 - val_loss: 0.9557 - val_acc: 2.5000e-04
Epoch 4/10
125/125 [==============================] - 101s 809ms/step - loss: 0.5
968 - acc: 0.6288 - val_loss: 0.9135 - val_acc: 0.0207
Epoch 5/10
125/125 [==============================] - 98s 786ms/step - loss: 0.54
06 - acc: 0.6777 - val_loss: 0.9148 - val_acc: 0.1615
Epoch 6/10
125/125 [==============================] - 99s 794ms/step - loss: 0.44
34 - acc: 0.8117 - val_loss: 0.8259 - val_acc: 0.4803
Epoch 7/10
125/125 [==============================] - 98s 781ms/step - loss: 0.34
56 - acc: 0.8776 - val_loss: 0.4988 - val_acc: 0.7897
Epoch 8/10
125/125 [==============================] - 97s 779ms/step - loss: 0.28
25 - acc: 0.9052 - val_loss: 0.6948 - val_acc: 0.7120
Epoch 9/10
125/125 [==============================] - 97s 779ms/step - loss: 0.23
84 - acc: 0.9191 - val_loss: 0.7116 - val_acc: 0.7322
Epoch 10/10
125/125 [==============================] - 100s 798ms/step - loss: 0.2
044 - acc: 0.9331 - val_loss: 0.7217 - val_acc: 0.7487
Test loss: 0.4732949435710907
Test accuracy: 0.8379340767860413
```
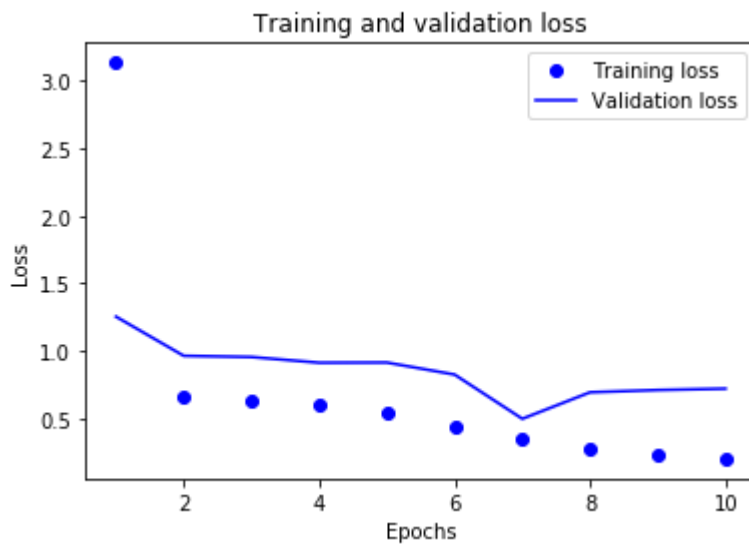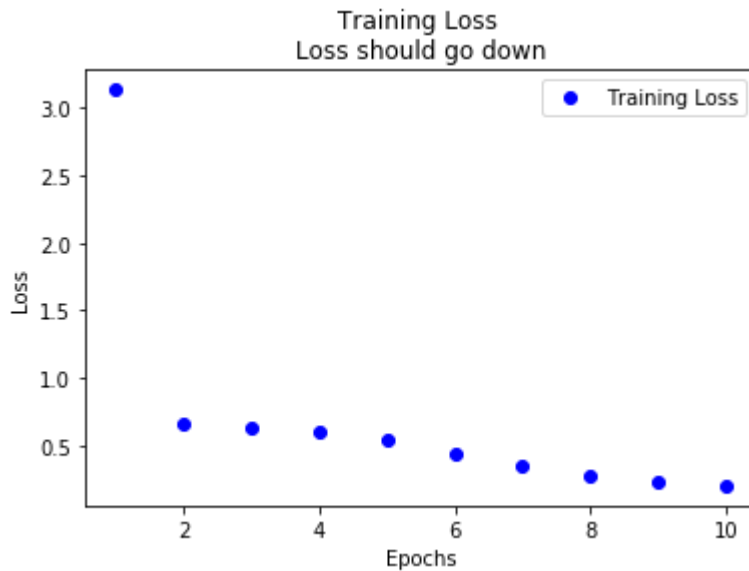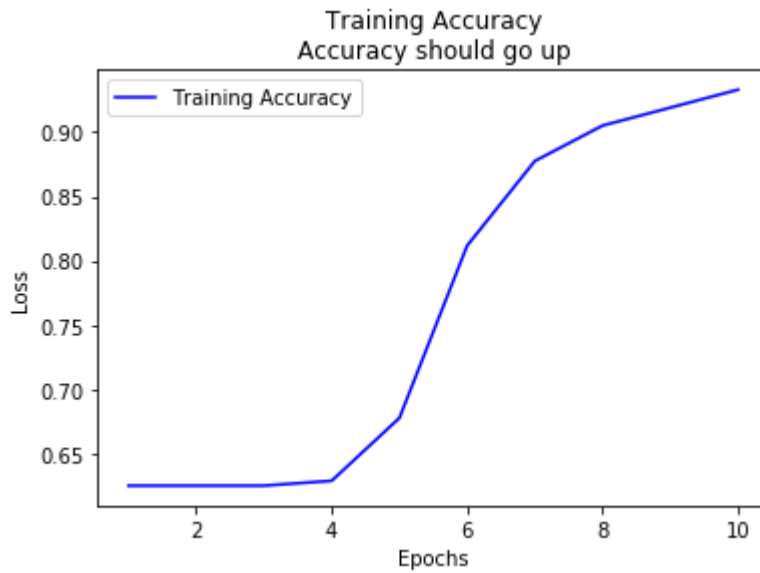
In [26]: 1 plot_model_history(history)

### Training Accuracy
### Accuracy should go up



### Training Loss
### Loss should go down



### Training and validation loss



Training and validation accuracy

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:    1    # //*** CODE HERE