In [38]:

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from keras.models import Model, Sequential
from keras.layers import LSTM, Activation, Dense, Dropout, Input, Embedding
from keras.optimizers import Adam
from keras.preprocessing.text import Tokenizer
from keras.preprocessing import sequence
from keras.utils import pad_sequences
from keras.utils import to_categorical
from keras.callbacks import EarlyStopping
%matplotlib inline
```

In [5]:

```python
df = pd.read_csv('/content/spam.csv',encoding='latin-1')
df.head()
```

Out[5]:

|   | v1 | v2 | Unnamed: 2 | Unnamed: 3 | Unnamed: 4 |
|---|-----|-----|-----|-----|-----|
| 0 | ham | Go until jurong point, crazy.. Available only ... | NaN | NaN | NaN |
| 1 | ham | Ok lar... Joking wif u oni... | NaN | NaN | NaN |
| 2 | spam | Free entry in 2 a wkly comp to win FA Cup fina... | NaN | NaN | NaN |
| 3 | ham | U dun say so early hor... U c already then say... | NaN | NaN | NaN |
| 4 | ham | Nah I don't think he goes to usf, he lives aro... | NaN | NaN | NaN |

In [6]:

```python
df.drop(['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'],axis=1,inplace=True)
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 2 columns):
 #   Column  Non-Null Count  Dtype
---  ------  --------------  -----
 0   v1      5572 non-null   object
 1   v2      5572 non-null   object
dtypes: object(2)
memory usage: 87.2+ KB
```

In [7]:

```python
sns.countplot(df.v1)
plt.xlabel('Label')
plt.title('Number of ham and spam messages')
```
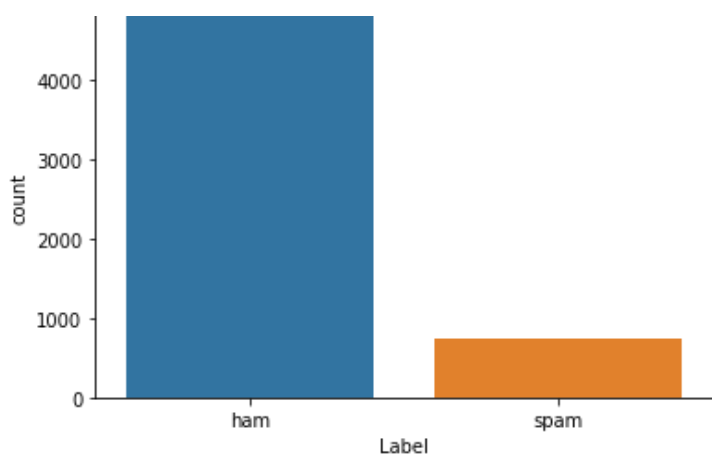
```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the
following variable as a keyword arg: x. From version 0.12, the only valid positional argu
ment will be `data`, and passing other arguments without an explicit keyword will result
in an error or misinterpretation.
  FutureWarning
```

Out[7]:

```
Text(0.5, 1.0, 'Number of ham and spam messages')
```

Number of ham and spam messages

5000

```
X = df.v2
Y = df.v1
le = LabelEncoder()
Y = le.fit_transform(Y)
Y = Y.reshape(-1,1)
```

In [17]:

```
X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=0.15)
```

In [30]:

```
max_words = 1000
max_len = 150
tok = Tokenizer(num_words=max_words)
tok.fit_on_texts(X_train)
sequences = tok.texts_to_sequences(X_train)
sequences_matrix = pad_sequences(sequences,maxlen=max_len)
```

In [35]:

```
tok.index_word
```

Out[35]:

```
{1: 'i',
 2: 'to',
 3: 'you',
 4: 'a',
 5: 'the',
 6: 'u',
 7: 'and',
 8: 'in',
 9: 'is',
 10: 'me',
 11: 'my',
 12: 'for',
 13: 'your',
 14: 'it',
 15: 'of',
 16: 'call',
 17: 'have',
 18: 'on',
 19: '2',
 20: 'that',
 21: 'now',
 22: 'are',
 23: 'so',
 24: 'but',
 25: 'not',
 26: 'can',
 27: 'do',
 28: 'at',
 29: 'or',
```

```
30: 'get',
31: 'with',
32: 'if',
33: 'ur',
34: "i'm",
35: 'be',
36: 'will',
37: 'no',
38: 'just',
39: 'we',
40: 'this',
41: 'gt',
42: 'lt',
43: '4',
44: 'up',
45: 'ok',
46: 'free',
47: 'when',
48: 'out',
49: 'go',
50: 'how',
51: 'from',
52: 'what',
53: 'all',
54: 'know',
55: 'then',
56: 'got',
57: 'like',
58: 'good',
59: 'come',
60: 'was',
61: 'its',
62: 'am',
63: 'only',
64: 'time',
65: 'day',
66: 'he',
67: 'love',
68: 'send',
69: 'there',
70: 'as',
71: 'text',
72: 'want',
73: 'one',
74: 'going',
75: 'need',
76: 'by',
77: "i'll",
78: 'about',
79: 'r',
80: 'txt',
81: 'stop',
82: 'still',
83: 'home',
84: 'see',
85: 'sorry',
86: 'lor',
87: 'she',
88: 'mobile',
89: 'k',
90: 'reply',
91: 'dont',
92: 'today',
93: 'back',
94: 'n',
95: 'da',
96: 'any',
97: 'take',
98: 'our',
99: 'later',
100: 'new',
101: 'they',
```

```
102: 'tell',
103: "don't",
104: 'think',
105: 'pls',
106: 'hi',
107: 'ì',
108: 'please',
109: 'been',
110: 'some',
111: 'phone',
112: 'her',
113: 'hope',
114: 'much',
115: 'well',
116: 'did',
117: 'too',
118: 'him',
119: 'hey',
120: '1',
121: 'here',
122: 'd',
123: 'make',
124: 'has',
125: 'week',
126: 'night',
127: 'who',
128: 'oh',
129: 'msg',
130: 'where',
131: 'an',
132: 'claim',
133: 'great',
134: 'more',
135: 'dear',
136: 'happy',
137: 'yes',
138: 'way',
139: 'c',
140: 'wat',
141: 'had',
142: 'www',
143: 'work',
144: 'already',
145: 'number',
146: 'give',
147: 'say',
148: "it's",
149: 'should',
150: 'tomorrow',
151: 'right',
152: 'e',
153: 'cash',
154: 'amp',
155: 'doing',
156: 'meet',
157: 'message',
158: 'after',
159: 'im',
160: 'life',
161: 'prize',
162: 'said',
163: '3',
164: 'really',
165: 'ask',
166: 'them',
167: 'babe',
168: 'last',
169: 'very',
170: 'why',
171: 'morning',
172: 'find',
173: 'win',
```

```
174: 'yeah',
175: 'b',
176: 't',
177: 'also',
178: 'sure',
179: 'cos',
180: 'let',
181: 'lol',
182: 'uk',
183: 'pick',
184: 'miss',
185: 'urgent',
186: 'thanks',
187: 'anything',
188: 'care',
189: 'sent',
190: 'nokia',
191: 'again',
192: 'something',
193: '150p',
194: 'com',
195: 'cant',
196: 'min',
197: "i've",
198: 'won',
199: 'keep',
200: 'would',
201: 's',
202: 'went',
203: 'contact',
204: 'nice',
205: 'his',
206: 'buy',
207: 'wait',
208: 'place',
209: 'every',
210: 'money',
211: 'over',
212: 'us',
213: 'were',
214: 'even',
215: 'which',
216: 'co',
217: 'someone',
218: 'gonna',
219: 'soon',
220: 'down',
221: 'ya',
222: 'sms',
223: 'first',
224: '5',
225: 'chat',
226: 'thing',
227: "that's",
228: 'late',
229: 'could',
230: 'before',
231: 'leave',
232: 'sleep',
233: 'help',
234: 'dun',
235: 'feel',
236: 'special',
237: 'next',
238: 'wan',
239: 'customer',
240: "can't",
241: 'things',
242: 'service',
243: 'gud',
244: 'off',
245: 'always',
```

```
246: '18',
247: 'around',
248: 'many',
249: 'friends',
250: 'v',
251: 'thk',
252: 'per',
253: 'getting',
254: 'tonight',
255: 'may',
256: 'hello',
257: 'told',
258: 'wish',
259: 'friend',
260: 'ìï',
261: 'year',
262: 'try',
263: 'other',
264: 'fine',
265: 'tone',
266: 'use',
267: '16',
268: 'people',
269: '50',
270: 'yet',
271: 'x',
272: 'lunch',
273: "you're",
274: 'same',
275: 'guaranteed',
276: 'stuff',
277: '6',
278: 'talk',
279: 'waiting',
280: 'name',
281: 'haha',
282: 'mins',
283: 'finish',
284: 'few',
285: 'days',
286: 'best',
287: 'holiday',
288: 'live',
289: 'class',
290: 'man',
291: "didn't",
292: 'job',
293: 'trying',
294: 'enjoy',
295: 'heart',
296: 'meeting',
297: 'being',
298: 'yup',
299: 'coming',
300: 'thought',
301: 'having',
302: 'ill',
303: 'long',
304: 'draw',
305: 'yo',
306: 'cool',
307: 'account',
308: 'line',
309: 'person',
310: 'better',
311: 'bit',
312: 'y',
313: 'never',
314: 'done',
315: 'car',
316: 'problem',
317: 'play',
```

```
318: 'nothing',
319: 'another',
320: 'smile',
321: 'cs',
322: '7',
323: 'wk',
324: 'big',
325: 'house',
326: 'mind',
327: 'real',
328: 'end',
329: 'receive',
330: 'camera',
331: 'thats',
332: 'eat',
333: 'word',
334: 'dat',
335: 'check',
336: 'chance',
337: 'ready',
338: 'latest',
339: 'lar',
340: 'dinner',
341: '150ppm',
342: 'awarded',
343: 'box',
344: 'shows',
345: '1st',
346: 'boy',
347: 'sir',
348: 'world',
349: 'bt',
350: 'into',
351: 'than',
352: 'start',
353: 'guess',
354: 'sweet',
355: 'room',
356: 'half',
357: 'ah',
358: 'girl',
359: 'bad',
360: 'jus',
361: 'maybe',
362: 'den',
363: 'å£1',
364: 'guys',
365: 'wanna',
366: 'look',
367: 'watching',
368: 'landline',
369: 'month',
370: 'code',
371: 'might',
372: 'once',
373: 'plan',
374: 'god',
375: 'sat',
376: 'called',
377: 'rate',
378: 'pa',
379: 'xxx',
380: 'shall',
381: 'because',
382: 'video',
383: 'ever',
384: 'aight',
385: 'remember',
386: 'kiss',
387: '9',
388: 'liao',
389: 'lot',
```

```
390: 'dunno',
391: 'collect',
392: 'actually',
393: 'early',
394: 'tv',
395: 'po',
396: 'little',
397: 'fun',
398: "he's",
399: 'å£1000',
400: 'speak',
401: 'minutes',
402: 'put',
403: 'enough',
404: 'left',
405: 'forgot',
406: 'part',
407: 'bed',
408: 'll',
409: 'does',
410: 'watch',
411: 'cost',
412: 'missing',
413: 'evening',
414: "how's",
415: 'everything',
416: 'probably',
417: 'tmr',
418: 'offer',
419: 'office',
420: 'working',
421: 'without',
422: 'made',
423: 'ringtone',
424: 'wife',
425: 'dis',
426: '10',
427: 'apply',
428: 'easy',
429: 'since',
430: 'award',
431: 'reach',
432: 'princess',
433: '8',
434: 'baby',
435: 'thank',
436: 'shit',
437: 'birthday',
438: 'details',
439: 'tones',
440: 'wif',
441: 'wot',
442: 'wont',
443: 'okay',
444: 'fuck',
445: 'didnt',
446: 'texts',
447: 'looking',
448: 'm',
449: 'dad',
450: 'hear',
451: 'thanx',
452: 'mail',
453: 'important',
454: 'orange',
455: 'between',
456: 'quite',
457: 'while',
458: 'wanted',
459: 'entry',
460: 'bus',
461: 'must',
```

```
462: "there's",
463: 'those',
464: '000',
465: 'weekend',
466: 'until',
467: 'plus',
468: 'says',
469: 'true',
470: 'bring',
471: 'food',
472: 'wake',
473: 'network',
474: 'stay',
475: 'anyway',
476: 'shopping',
477: 'most',
478: 'selected',
479: 'update',
480: 'sexy',
481: 'alright',
482: 'pain',
483: 'else',
484: 'leh',
485: 'worry',
486: 'though',
487: 'school',
488: 'yours',
489: '2nd',
490: 'price',
491: 'collection',
492: 'pay',
493: 'means',
494: 'asked',
495: '500',
496: 'xmas',
497: 'bored',
498: 'weekly',
499: 'show',
500: 'decimal',
501: 'lei',
502: 'yesterday',
503: 'came',
504: 'hour',
505: 'gift',
506: 'havent',
507: 'away',
508: 'music',
509: 'able',
510: 'hav',
511: 'luv',
512: 'valid',
513: 'join',
514: '10p',
515: 'abt',
516: 'hurt',
517: 'movie',
518: 'coz',
519: 'hot',
520: "we're",
521: 'nite',
522: 'dreams',
523: 'messages',
524: 'date',
525: 'missed',
526: 'attempt',
527: 'g',
528: 'mob',
529: '8007',
530: 'afternoon',
531: 'town',
532: 'family',
533: 'huh',
```

```
534: 'making',
535: 'game',
536: 'address',
537: 'till',
538: 'two',
539: 'sch',
540: 'guy',
541: 'online',
542: 'sleeping',
543: 'both',
544: 'plz',
545: 'news',
546: 'tot',
547: 'points',
548: 'wants',
549: 'oso',
550: 'run',
551: 'pic',
552: 'haf',
553: 'de',
554: "what's",
555: 'sad',
556: 'id',
557: 'years',
558: "haven't",
559: 'noe',
560: 'å£500',
561: 'set',
562: 'wen',
563: 'together',
564: 'aft',
565: 'words',
566: 'mean',
567: 'vouchers',
568: 'these',
569: 'shop',
570: 'club',
571: 'gr8',
572: 'dude',
573: 'old',
574: 'private',
575: 'calls',
576: 'double',
577: "won't",
578: 'pics',
579: 'times',
580: 'yourself',
581: 'either',
582: 'lose',
583: 'final',
584: 'til',
585: 'wid',
586: 'busy',
587: 'national',
588: 're',
589: 'post',
590: 'goes',
591: 'change',
592: 'calling',
593: 'story',
594: 'pounds',
595: 'driving',
596: 'friendship',
597: 'colour',
598: 'unsubscribe',
599: 'å£100',
600: 'goin',
601: 'trip',
602: 'ard',
603: "'",
604: 'http',
605: 'å£2000',
```

```
606: 'juz',
607: 'walk',
608: 'started',
609: 'available',
610: '86688',
611: 'found',
612: 'face',
613: 'email',
614: 'gd',
615: 'congrats',
616: 'games',
617: 'hair',
618: 'bonus',
619: 'feeling',
620: 'top',
621: 'drop',
622: 'rite',
623: 'saying',
624: 'tried',
625: 'delivery',
626: 'chikku',
627: 'drink',
628: 'brother',
629: 'question',
630: 'answer',
631: 'charge',
632: 'beautiful',
633: 'msgs',
634: 'taking',
635: 'sae',
636: 'took',
637: 'saw',
638: 'close',
639: 'statement',
640: 'expires',
641: 'head',
642: 'happen',
643: 'simple',
644: 'å£5000',
645: 'wil',
646: 'cause',
647: 'leaving',
648: 'fast',
649: 'far',
650: 'makes',
651: 'full',
652: '11',
653: 'believe',
654: "she's",
655: 'mum',
656: 'book',
657: 'takes',
658: 'kind',
659: 'choose',
660: 'okie',
661: 'order',
662: 'finished',
663: 'services',
664: 'whats',
665: 'thinking',
666: 'tomo',
667: 'await',
668: "c's",
669: 'company',
670: 'identifier',
671: 'visit',
672: 'todays',
673: 'sun',
674: "we'll",
675: 'happened',
676: 'w',
677: 'don',
```

```
678: 'test',
679: 'sounds',
680: 'awesome',
681: 'lots',
682: 'auction',
683: 'card',
684: 'second',
685: 'secret',
686: 'hard',
687: 'hit',
688: 'hours',
689: 'break',
690: '12hrs',
691: 'pub',
692: 'mine',
693: 'poly',
694: 'wit',
695: "doesn't",
696: 'smiling',
697: 'worth',
698: 'everyone',
699: 'voucher',
700: 'net',
701: 'neva',
702: 'open',
703: 'content',
704: 'needs',
705: 'winner',
706: 'minute',
707: 'ring',
708: 'used',
709: 'read',
710: 'lucky',
711: 'å£3',
712: 'sister',
713: 'blue',
714: 'comes',
715: 'loving',
716: 'welcome',
717: 'lets',
718: 'player',
719: 'goodmorning',
720: 'prob',
721: 'outside',
722: '750',
723: 'wrong',
724: 'land',
725: 'o',
726: 'tel',
727: 'bout',
728: 'sis',
729: 'pretty',
730: 'info',
731: 'pm',
732: 'tho',
733: 'gone',
734: '100',
735: 'hows',
736: 'hmm',
737: 'knw',
738: 'nt',
739: 'lesson',
740: 'operator',
741: 'angry',
742: 'course',
743: 'project',
744: 'finally',
745: 'whatever',
746: 'row',
747: 'search',
748: 'mobileupd8',
749: 'seeing',
```

```
750: "you'll",
751: 'luck',
752: 'cut',
753: 'ha',
754: 'fucking',
755: 'ltd',
756: 'frm',
757: 'college',
758: 'meant',
759: 'opt',
760: 'darlin',
761: 'fone',
762: '08000930705',
763: 'boytoy',
764: 'drive',
765: 'alone',
766: 'each',
767: 'anytime',
768: 'decided',
769: 'case',
770: 'balance',
771: 'lovely',
772: "you've",
773: 'carlos',
774: 'parents',
775: 'savamob',
776: 'sea',
777: 'saturday',
778: 'girls',
779: 'anyone',
780: 'oredi',
781: '30',
782: 'download',
783: 'etc',
784: 'hand',
785: 'ac',
786: 'telling',
787: 'mah',
788: 'frnd',
789: '2003',
790: '800',
791: 'ni8',
792: 'invited',
793: 'suite342',
794: '2lands',
795: 'th',
796: 'pobox',
797: 'light',
798: 'bslvyl',
799: 'friday',
800: 'tc',
801: 'their',
802: 'b4',
803: 'mates',
804: 'side',
805: 'mate',
806: 'treat',
807: 'bank',
808: 'reveal',
809: 'thinks',
810: '\x89û',
811: 'mrng',
812: 'camcorder',
813: 'sunday',
814: 'smth',
815: 'ass',
816: 'reading',
817: 'st',
818: 'earlier',
819: 'father',
820: 'john',
821: 'knew',
```

```
822: 'felt',
823: 'least',
824: 'weeks',
825: 'press',
826: 'offers',
827: 'smoke',
828: 'lost',
829: 'dating',
830: 'snow',
831: 'hold',
832: 'forget',
833: 'wonderful',
834: 'frnds',
835: 'wkly',
836: '87066',
837: 'enter',
838: 'touch',
839: 'type',
840: 'seen',
841: 'fr',
842: "i'd",
843: 'pass',
844: 'chennai',
845: 'å£250',
846: 'crazy',
847: 'comp',
848: 'phones',
849: 'extra',
850: 'reason',
851: 'laptop',
852: 'f',
853: 'un',
854: 'redeemed',
855: '04',
856: 'mom',
857: 'numbers',
858: 'å£350',
859: 'sex',
860: 'understand',
861: 'unlimited',
862: 'support',
863: 'motorola',
864: '08000839402',
865: 'fri',
866: 'yar',
867: '20',
868: 'caller',
869: '03',
870: 'muz',
871: 'wow',
872: 'xx',
873: 'ipod',
874: 'whole',
875: 'hungry',
876: 'rental',
877: 'comin',
878: 'currently',
879: 'met',
880: 'gas',
881: 'semester',
882: 'hee',
883: 'quiz',
884: "wasn't",
885: 'std',
886: 'charged',
887: 'crave',
888: 'log',
889: 'rs',
890: 'yr',
891: 'point',
892: 'eve',
893: 'jay',
```

```
894: 'mayb',
895: 'complimentary',
896: 'listen',
897: 'mr',
898: 'surprise',
899: "'",
900: 'congratulations',
901: 'bath',
902: 'nope',
903: 'gal',
904: 'red',
905: 'correct',
906: 'area',
907: 'otherwise',
908: 'cum',
909: 'ago',
910: 'checking',
911: 'couple',
912: 'slowly',
913: 'ten',
914: 'freemsg',
915: 'laugh',
916: 'mobiles',
917: 'information',
918: 'small',
919: 'party',
920: 'loads',
921: 'somewhere',
922: 'å£10',
923: 'ge',
924: 'slow',
925: 'march',
926: 'india',
927: 'admirer',
928: 'film',
929: 'loved',
930: 'within',
931: 'immediately',
932: 'die',
933: 'eg',
934: 'almost',
935: 'discount',
936: 'whenever',
937: 'spend',
938: 'hmv',
939: 'kate',
940: 'mu',
941: 'leaves',
942: 'across',
943: 'darren',
944: 'credit',
945: 'different',
946: 'rest',
947: 'dream',
948: 'reward',
949: 'christmas',
950: 'monday',
951: "isn't",
952: '0800',
953: 'gym',
954: 'asking',
955: 'confirm',
956: 'plans',
957: 'difficult',
958: 'rply',
959: 'drugs',
960: 'ends',
961: 'usf',
962: 'valued',
963: 'link',
964: 'exam',
965: 'park',
```

```
966: 'store',
967: 'empty',
968: 'eh',
969: 'supposed',
970: 'through',
971: 'askd',
972: 'police',
973: 'via',
974: 'uncle',
975: 'feels',
976: 'dnt',
977: 'computer',
978: 'gotta',
979: 'wana',
980: 'hospital',
981: 'picking',
982: 'heard',
983: 'completely',
984: 'sort',
985: 'reached',
986: 'nobody',
987: 'lovable',
988: 'short',
989: 'tonite',
990: 'rakhesh',
991: 'em',
992: 'weed',
993: 'clean',
994: 'poor',
995: 'hmmm',
996: '4u',
997: 'mon',
998: 'sony',
999: 'call2optout',
1000: 'studying',
...}
```

In [37]:

```python
TOT_SIZE = len(tok.word_index)+1
```

In [40]:

```python
lstm_model = Sequential()
lstm_model.add(Embedding(TOT_SIZE, 32, input_length=max_len))
lstm_model.add(LSTM(100))
lstm_model.add(Dropout(0.4))
lstm_model.add(Dense(20, activation="relu"))
lstm_model.add(Dropout(0.3))
lstm_model.add(Dense(1, activation = "sigmoid"))
lstm_model.compile(loss = "binary_crossentropy", optimizer = "adam", metrics = ["accurac
y"])
```

In [41]:

```python
lstm_model.summary()
```

```
Model: "sequential_1"
_____
 Layer (type)                Output Shape              Param #
=================================================================
 embedding_1 (Embedding)     (None, 150, 32)           262496

 lstm_1 (LSTM)               (None, 100)               53200

 dropout_2 (Dropout)         (None, 100)               0

 dense_2 (Dense)             (None, 20)                2020

 dropout_3 (Dropout)         (None, 20)                0

 dense_3 (Dense)             (None, 1)                 21
```

```
==============================================================
Total params: 317,737
Trainable params: 317,737
Non-trainable params: 0
_____
```

In [44]:

```python
lstm_model.fit(sequences_matrix,Y_train,batch_size=128,epochs=10,
          validation_split=0.2,
          workers = 10,
          callbacks=[EarlyStopping(monitor='val_loss',min_delta=0.0001)])
```

```
Epoch 1/10
30/30 [==============================] - 10s 329ms/step - loss: 0.0851 - accuracy: 0.9805
- val_loss: 0.0558 - val_accuracy: 0.9852
Epoch 2/10
30/30 [==============================] - 10s 326ms/step - loss: 0.0497 - accuracy: 0.9871
- val_loss: 0.0464 - val_accuracy: 0.9873
```

Out[44]:

```
<keras.callbacks.History at 0x7f6696448550>
```

In [45]:

```python
lstm_model.save('LSTM.h5')
```

In [47]:

```python
test_sequences = tok.texts_to_sequences(X_test)
test_sequences_matrix = pad_sequences(test_sequences,maxlen=max_len)
```

In [49]:

```python
acc = lstm_model.evaluate(test_sequences_matrix,Y_test)
```

```
27/27 [==============================] - 1s 32ms/step - loss: 0.0442 - accuracy: 0.9856
```

In [51]:

```python
print('Test set\n  Loss: {:0.3f}\n  Accuracy: {:0.3f}'.format(acc[0],acc[1]))
```

```
Test set
  Loss: 0.044
  Accuracy: 0.986
```