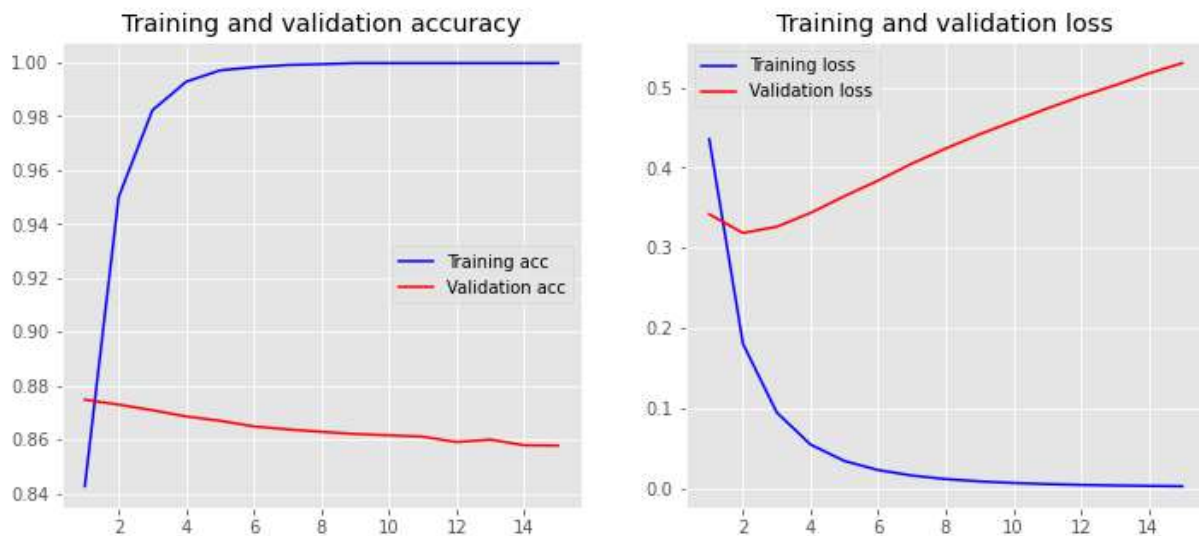


در این مجموعه دادگان دو تا دسته ی تست و آموزش هست و هر کدام شامل دو دسته فایل txt. هستند که یک مجموعه شان مثبت و دیگری منفی ست این دسته بندی در حکم متغیر هدف است. هر متن در یک فایل جداگانه ی txt. ذخیره شده. برای خواندن تمام این متن ها و تشکیل یک دیتافریم خام اولیه، ادرس و نام تمام فایل ها در هر فولدر pos و neg را در لیست mylist به کمک glob قرار میدهم. سپس هر کدام از این ها را با برچسب ۰ یا ۱ به ترتیب برای فایل های پوشه ی neg و pos در یک تابل در دو لیست poss و negs قرار می دهم و نهایتا تمام این ها را به یک دیتافریم با دو ستون sentences و target تبدیل می کنم. تمام این مراحل برای داده های تست و آموزش به طور جداگانه و یکسان انجام شده. دیتافریم ها شافل شده اند تا داده های با برچسب ۰ و ۱ درهم تلفیق شوند.

از sklearn با استفاده از CountVectorizer یک ساک کلمات از تمام جملات df\_train می سازم. این قابلیت در این مدل وجود دارد که برای متون انگلیسی کلمات stopword و علائم نگارشی را در ساک کلمه لحاظ نکند و نیاز به پیش پردازش های این چینی از نداریم. جملات دیتافریم ها را به صورت نمایشی از مدل ساک کلمات ساخته شده در X\_train و X\_test ذخیره می کنم. این ها باید به صورت آرایه های numpy تبدیل شوند تا بتوان در شبکه ی عصبی بخش بعد ازشان استفاده کرد. من به دلیل اینکه در کامپیوتر شخصی خودم حافظه ی کافی پردازش کل این داده ها را به صورت بردار نداشتم و امکان استفاده از پردازش های آنلاین هم برای فراهم نبود، هر دو مجموعه تیت و آموزش را از ۲۵۰۰۰ به ۱۰۰۰۰ کوچک تر کردم و بردارهای X\_train\_arr و X\_test\_arr برای ورود به شبکه عصبی مرحله بعد آماده اند.

با استفاده از keras از tensorflow یک شبکه عصبی با یک لایه ورودی ۱۰۰۰۰ تایی، با یک لایه پنهان ۱۰ تایی با تابع فعالسازی relu و یک تک واحد خروجی با تابع فعالسازی سیگموئید تعریف کردم. این شبکه با معیار accuracy آموزش می بیند و برای تابع loss از binary\_crossentropy ساخته می شود و با ۱۵ ایپاک و در دسته های ۶۴ تایی آپدیت وزن و آموزش صورت می گیرد.

نمودار تست و آموزش روی loss و accuracy آمده و نتیجه ی بهترین مدل به دست آمده و نهایی روی داده های تست و آموزش هم چاپ شده. این نتیجه روی ایپاک ۲ در ابتدای کار به دست آمده (با معیار loss) و بعد هرچه آموزش ادامه پیدا می کند بیش برازش بیشتر می شود و مدل قابلیت generalization از دست می دهد. دقت ۸۶ دسته بندی درصد



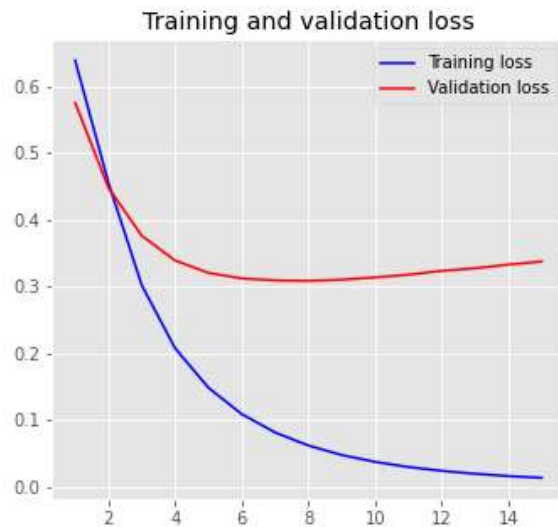
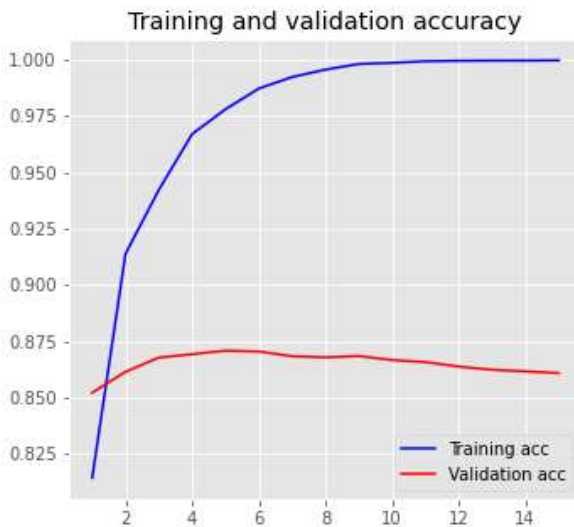
Training Accuracy: 0.9997  
Training Loss: 0.0021

Testing Accuracy: 0.8578  
Testing Loss: 0.5302

یک مدل دوم می‌سازم که به جای داده‌های ساک کلمات از روش tfidf استفاده کند. با استفاده از TfidfTransformer از sklearn فرکانس تکرار کلمات را در هر متن پیدا می‌کند. دوباره داده‌ها را به صورت آرایه‌ی numpy تبدیل می‌کنم و به همان علت محدودیت سیستم پردازش که دارم با داده‌های ۱۰۰۰۰ تایی کار می‌کنم. همان معماری شبکه‌ی قبلی را روی این داده‌ها آموزش می‌دهم. نحوه‌ی آموزش هم همان است.

نمودار تست و آموزش آورده شد. بهترین نتیجه در حدود ایپاک ۵ به دست آمده و بعد افت نتیجه روی تست را داریم. (بیش برارش) تغییری با مدل قبلی در نتیجه و عملکرد مدل به دست نیامد.

Training Accuracy: 0.9997  
Training Loss: 0.0116  
Testing Accuracy: 0.8608  
Testing Loss: 0.3376



مدل نهایی به دست آمده از این شبکه با دقت ۸۶ درصد قابل استفاده ست.