

1. (1%)請比較有無**normalize(rating)**的差別。並說明如何**normalize**。

Normalize Rating 的有無在validation Loss的差距非常明顯，我如果用 Normalize Rating 配上latent dimension = 20 再加上 bias，loss最低可以到 0.6128。然而，我如果沒使用Normalize Rating(latent dimension, bias 條件一樣)，loss 最低是0.7771。

做法：首先，先計算rating 的mean跟std。再將所有的rating減掉mean並除以Std，這樣可以確保rating的平均為0，標準差為1。而在跑testing的時候，記得把training 的std & mean 用pickle等檔型預存下來，將預測出的 testing data \* std\_train + mean\_train。

2. (1%)比較不同的**latent dimension**的結果。

在Bias & Normalization都有設定的前提下，我測試了Latent Dimension = 15, 20, 25, 30, 35, 40 六種的狀況，結果如下：

Dimension	15	20	25	30	35	40
Val_loss	0.6133	0.6156	0.6153	0.6171	0.6139	0.6128

從我自己的測試中，這些不同的Dimension對Val\_Loss並沒有太大的影響。另外，若在重複做同樣的測試，我發現Loss最低的Dimension從40變成25，40反而變成次高的Loss，推論在這個區間下Dimension跟Loss並沒有太大的相關性。

3. (1%)比較有無**bias**的結果。

我比較Latent Dimension = 20,40、做Normalization之下Bias有無對於Val\_Loss的影響，結果如下：

Bias Or Not(20)	Yes	No
Loss	0.6156	0.6220

Bias Or Not(40)	Yes	No
Loss	0.6128	0.6204

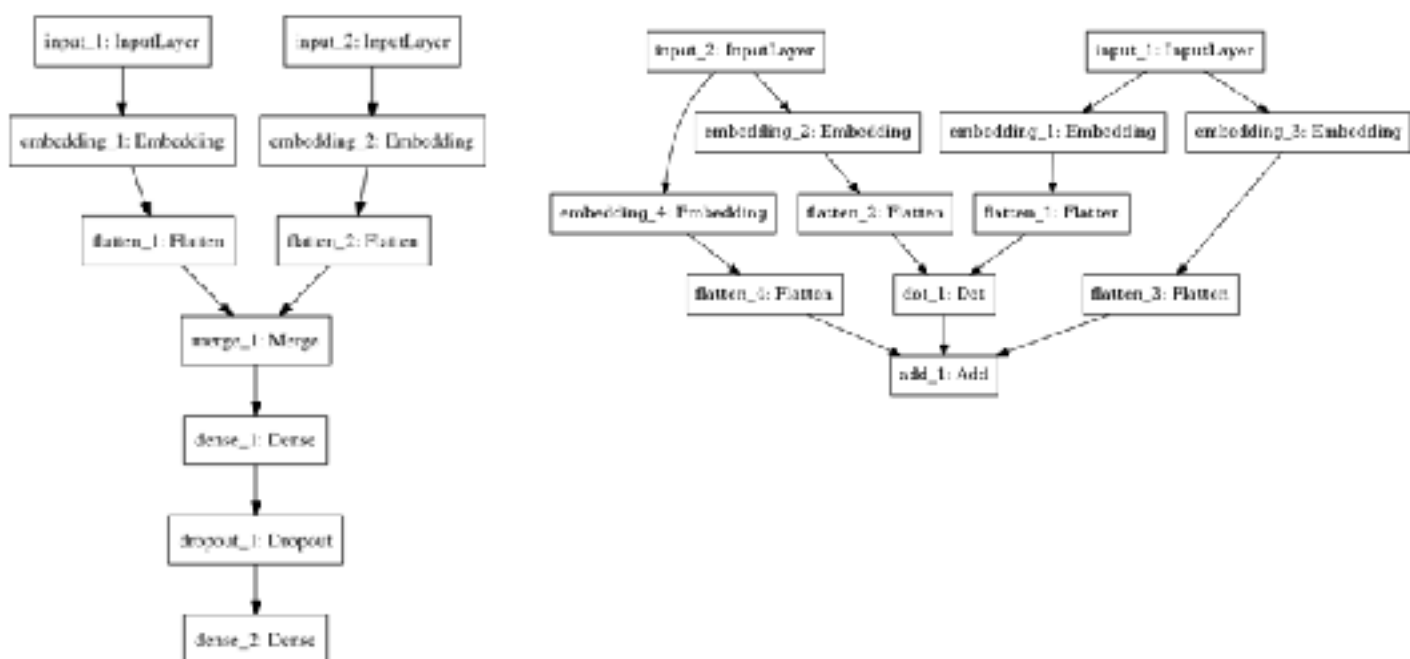
根據測試結果，在沒有Bias的情況下Val\_Loss是較高的，因為Bias有考慮使用者的評分傾向(ex:習慣評低分或高分)，能比較精確的呈現預測結果。

4. (1%)請試著用DNN來解決這個問題，並且說明實做的方法(方法不限)。並比較MF和NN的結果，討論結果的差異。

我Loss最低的DNN是用Linear Regression，因Classification的結果並沒有想像中來得好。另外，我只有放一層Dense(Unit = 256)，User & Movie從InputLayer輸入之後經過Embedding、Flatten、Merge(Concat merge)的處理之後就直接輸出，根據Mean Square Error 找最低的Loss，並用Batch = 128，patience = 10。在Kaggle上的Public的結果是0.859。(架構如左下)

跟DNN比起來，我MF的方法其實就是用助教的Code，在嘗試不同Latent Dimension之後選用Dimension= 40，Kaggle上的成績是0.875。(架構如右下)，架構上MF看起來比較胖是因為Bias也做了Embedding跟Flatten。

DNN除了Loss較小之外，它參數調整的彈性較高，可以再去調整Dense的層數、Dropout跟是否要做Batch Normalization，這些都可以讓DNN的結果遠高於MF。但MF的優點是它大概花2~3個Epoch就可以收斂到最小值了，而DNN會需要15~20個Epoch才可以達到最小的Loss。



5. (1%)請試著將movie的embedding用t降維後，將movie category當作label來作圖。

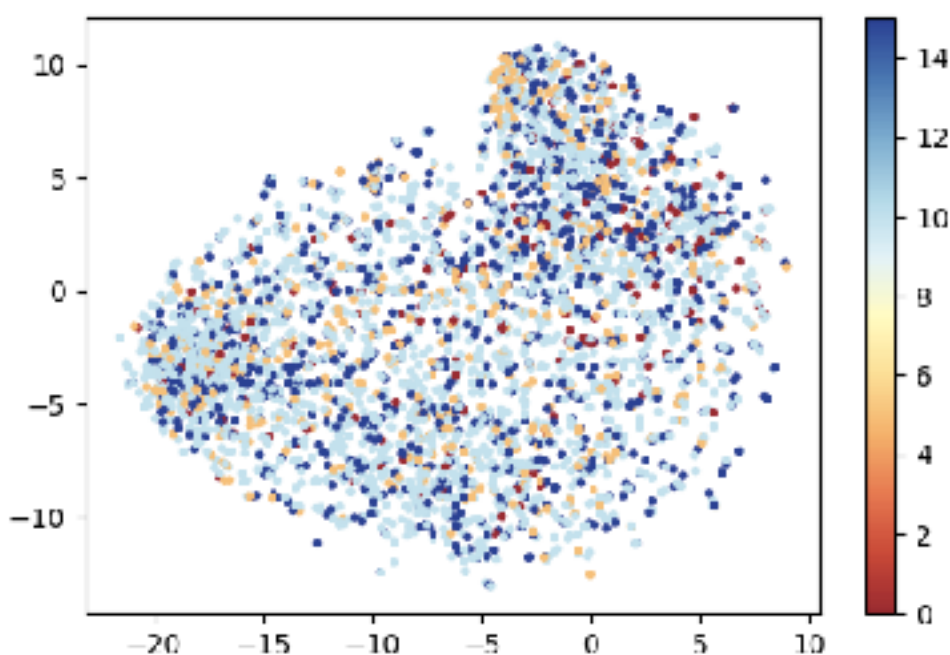
我主要將movie category分成4大類，如下表：

Category	Type
0	Animation, Children's
1	Crime, Horror, Thriller
2	Drama, Comedy, Musical
3	Adventure, Action

我得到的降為  
圖片如右，可以發  
現Category幾乎都  
混在在一起，沒有  
很明顯的呈現區別。

從這樣訓練過  
後的embedding降維  
對應到不同的  
Movie category，它  
的群聚關係不甚明  
顯。此外，在得到

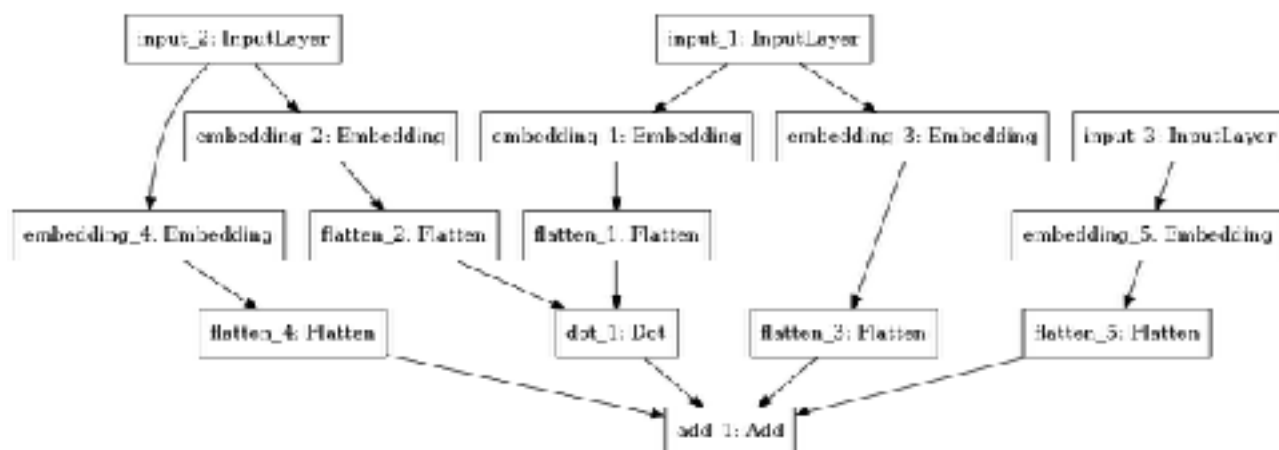
這樣的結果後我也重新檢視了一下movie.csv不同影片的分類，助教建議  
的影片分類也不全然正確，有些電影是同時包含在兩個類別裡面(ex: ID =  
2027 的 Mafia同時有Crime & Comedy)，這也影響了視覺上影片類別的辨  
認。



6. (BONUS)(1%)試著使用除了rating以外的feature, 並說明你的作法和結果，結果好壞不會影響評分。

參照第五題的做法，我把Movie Category當成Bias也加進去了。分類參考上題4大類 + 「其他類別」 + Unknown Movie共六大類，Latency Dimension = 20、有做Normalization。

架構如下圖：



最後，我的Val\_Loss是0.6095，比第二題同樣參數下的結果還低，Kaggle上的Error是0.87059 (過Strong Baseline)。