

به نام خدا



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

درس شبکه‌های عصبی
استاد صفابخش

تمرین دوم

علیرضا مازوچی

۴۰۰۱۳۱۰۷۵

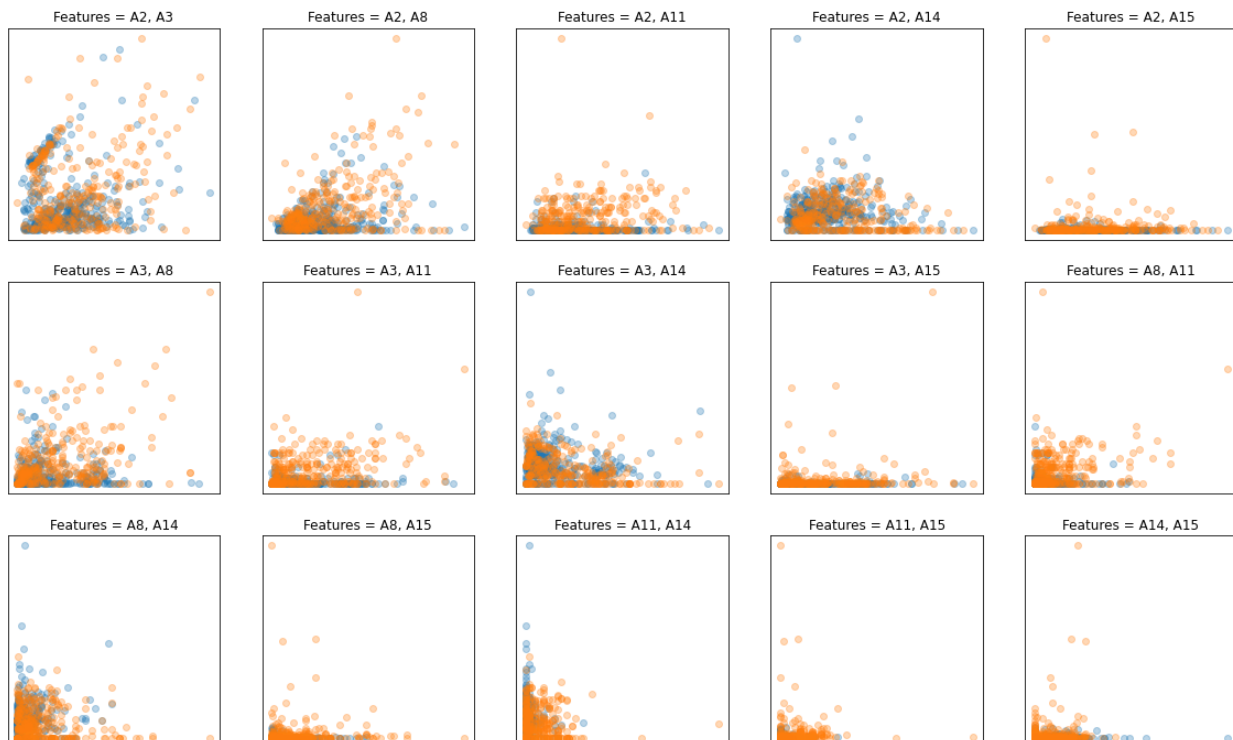
سوال ۱

مطابق با توضیحات موجود در مجموعه داده در حدود ۵ درصد داده‌ها یک یا چند ویژگی دارای مقدار تهی است. برای ستون‌هایی که دارای داده‌های پیوسته هستند از میانگین سایر داده‌ها و برای ستون‌هایی که دارای مقادیر گسسته هستند از ماکسیم مقدار برای پر کردن مقدار از دست رفته بهره می‌گیریم. مقادیر + و - ستون کلاس را به مقدار ۰ و ۱ تبدیل کردم. برای تمام ستون‌های ویژگی یک نرمال‌سازی مطابق با فرمول زیر هم به کار می‌گیریم تا رنج مقادیر برای تمام این ستون‌های مشابه با یکدیگر باشد:

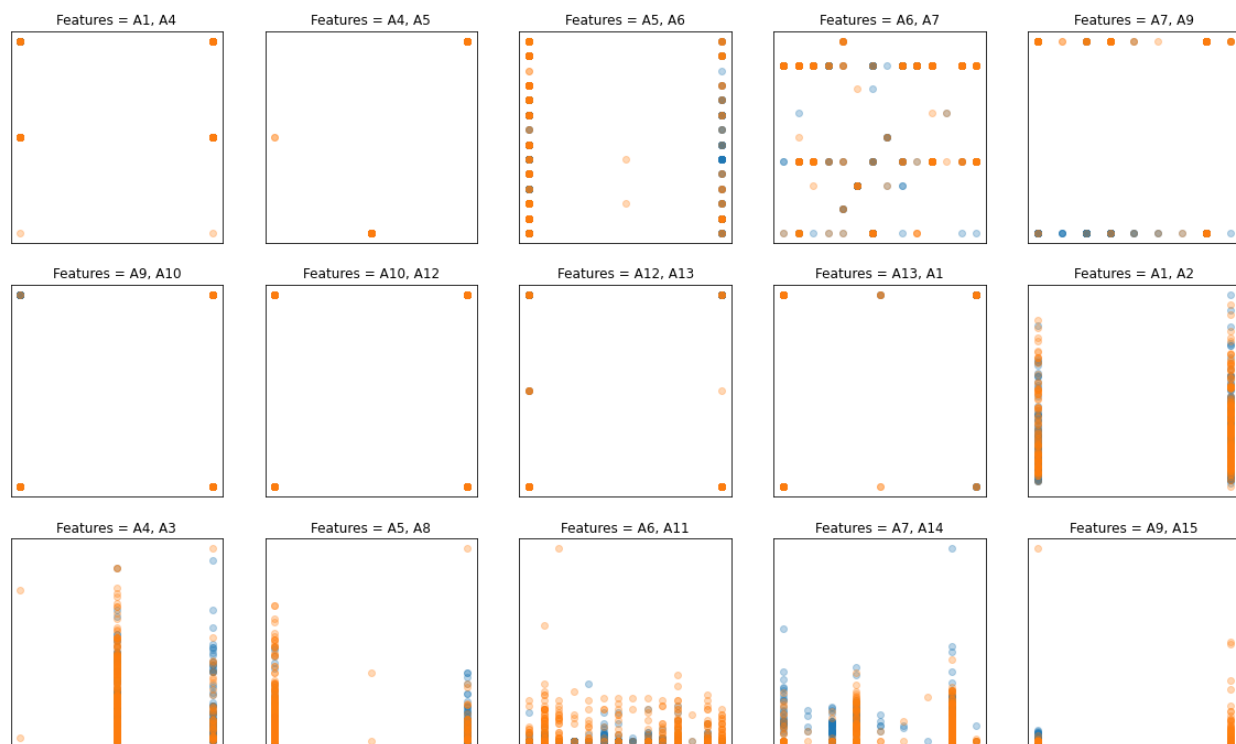
$$x_{normal} = \frac{x - \text{mean}(x)}{\text{std}(x)}$$

سوال ۲

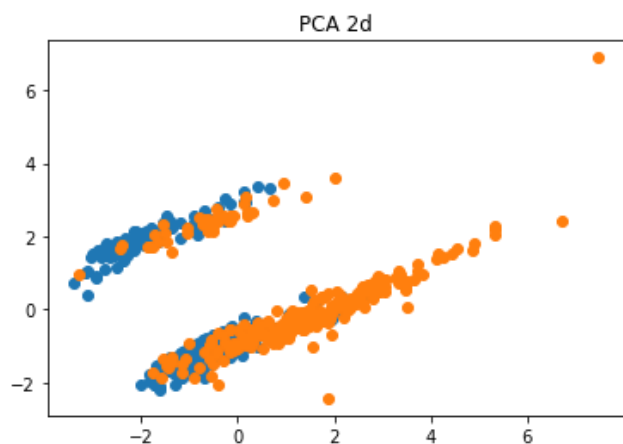
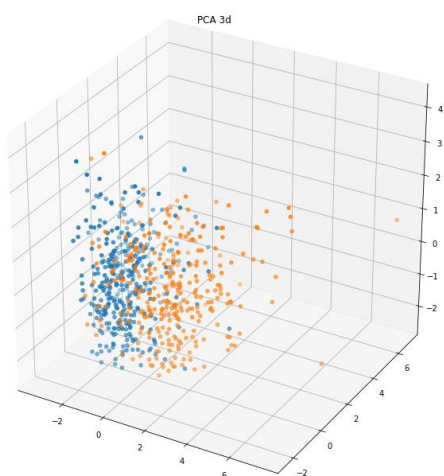
برای بررسی این مورد، چندین بررسی را انجام دادیم. اول آنکه ویژگی‌ها پیوسته را به صورت دو به دو در یک نمودار ترسیم کردیم ولی در هیچ کدام از نمودارها داده‌های دو کلاس از یکدیگر به صورت خطی جدا نشد:



ویژگی‌های پیوسته با توجه به رنج وسیع‌تری که نسبت به ویژگی‌های گسسته دارند احتمال بیشتری برای جداکردن کلاس‌ها دارند ولی با این حال ممکن است ویژگی‌های گسسته داده‌ها را جداکنند. لذا تعدادی از ترکیب‌های ویژگی‌ها گسسته با هم و با ویژگی‌های پیوسته را بررسی کردیم ولی در این حالت هم موفقیتی حاصل نشد:



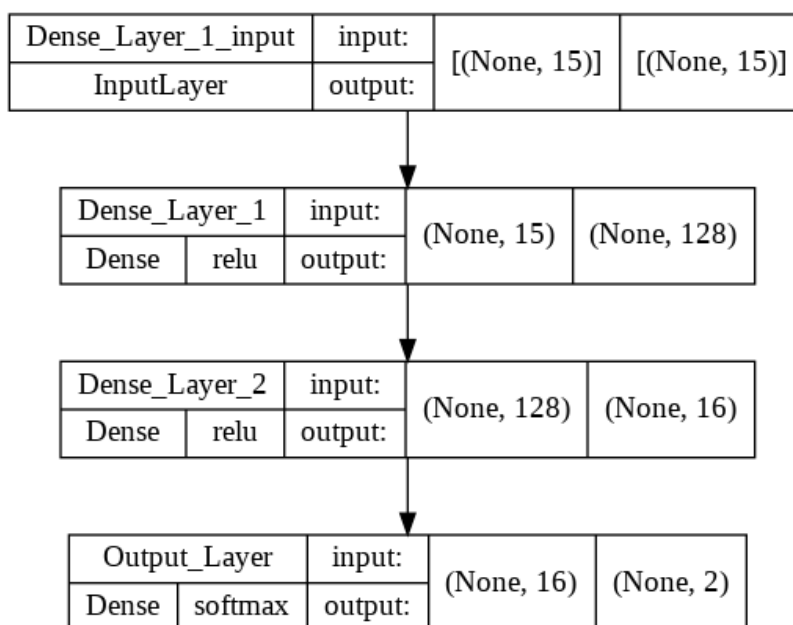
بار دیگر از روش PCA استفاده کردیم ولی این بار هم جداسازی داده‌ها انجام نشد:



حال که این روش‌های اولیه جواب ندادند، نوبت به بررسی یک روش پیچیده‌تر ولی با جواب قطعی می‌رسد. این بار از روش SVM با حاشیه سخت استفاده می‌کنم. اگر بتوان یک مدل SVM با حاشیه سخت پیدا کرد که به دقت ۱۰۰٪ برسد یعنی داده‌ها خطی جداپذیرند و اگر چنین مدلی وجود نداشته باشد یعنی خطی جداپذیر نیستند. در بررسی‌ای که من انجام دادم چنین مدلی پیدا نشد و لذا داده‌ها خطی جدا ناپذیرند.

سوال ۳

به عنوان اولین تلاش شبکه عصبی زیر را ایجاد کردم:

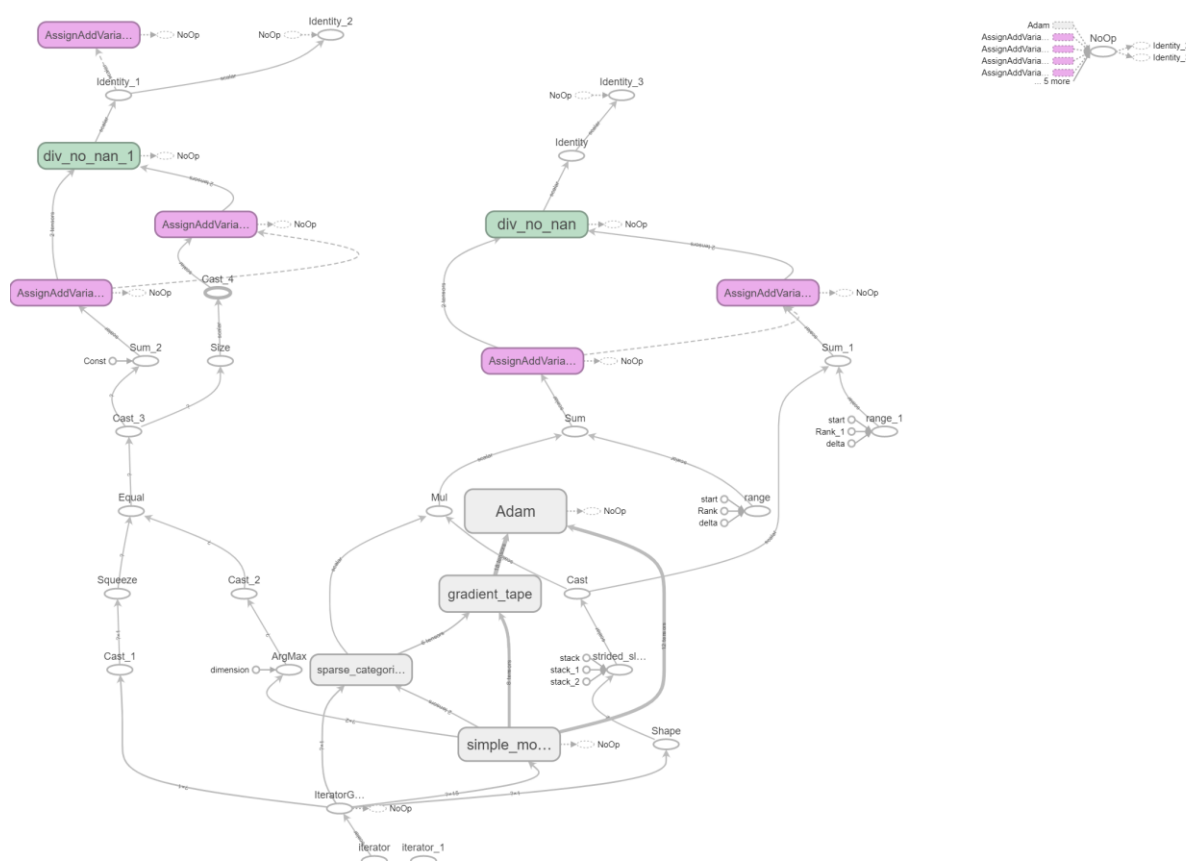


در این شبکه لایه ورودی شامل ۱۵ ویژگی است. سپس از این ۱۵ ویژگی توسط یک لایه Dense ۱۲۸ ویژگی استخراج می‌شود. در لایه Dense بعد ۱۶ ویژگی ترکیبی و پیچیده حاصل می‌شود و نهایتاً در لایه خروجی دو ویژگی ایجاد می‌شود که هر کدام متناسب با احتمال تعلق داده به یکی از دو کلاس موجود است. از Optimizer آدام و از تابع خطای Sparse Categorical Cross Entropy کمک گرفته‌ام.

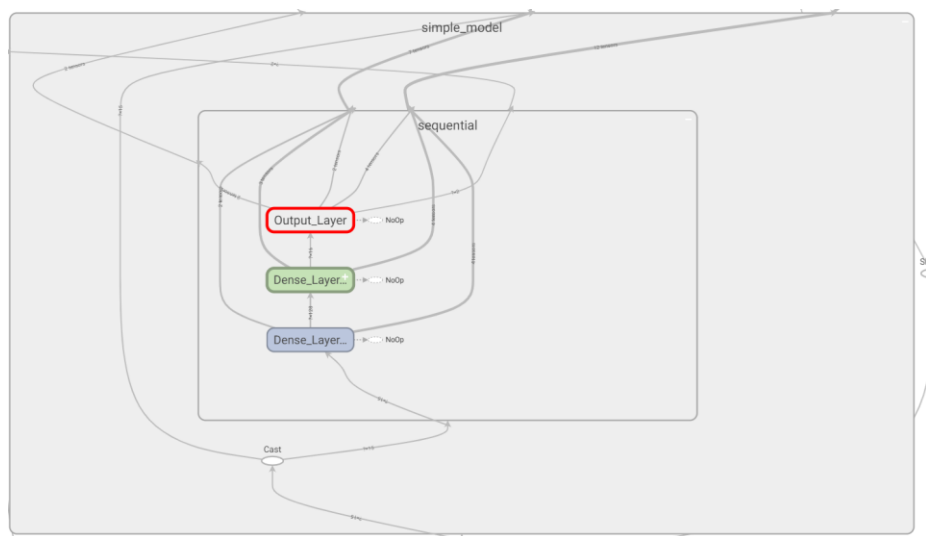
در لایه‌های Dense میانی از ReLU به عنوان تابع فعال‌سازی استفاده شده است تا عملکرد غیرخطی به مدل داده شود و در آخرین لایه که لایه خروجی باشد از یک لایه

Softmax استفاده شده است تا خروجی از جنس احتمال باشد و با تابع خطا استفاده می‌شود سازگار باشد. تعداد گام هم برابر با ۲۰ در نظر گرفته شده است.

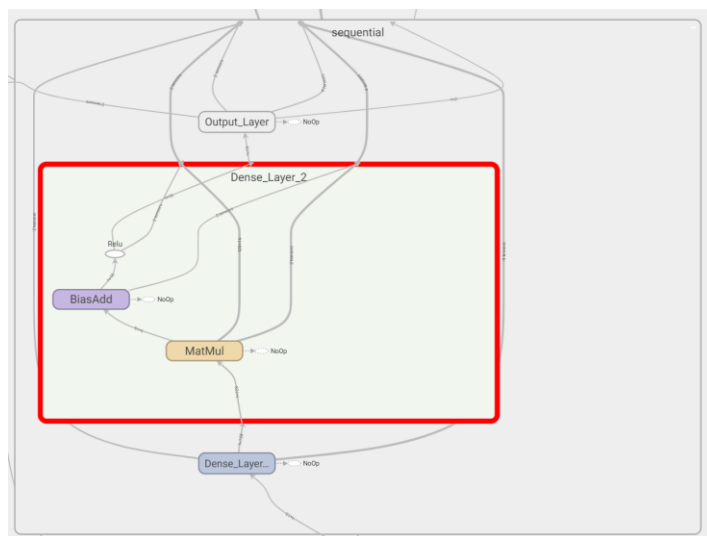
بعد از آموزش به صحت ۹۲٪ روی داده‌های آموزش و ۸۸٪ روی داده‌های اعتبارسنجی رسیدیم. همچنین اگر تمایل داشته باشیم می‌توانیم اطلاعاتی را از روی tensor board هم ببینیم. مثلاً گراف شبکه عبارت است از:



در قسمت پایین گراف کل مدل قرار گرفته است «simple_mo...» که با کلیک بر روی آن اطلاعات بیشتری راجع به خود مدل دریافت می‌کنیم:

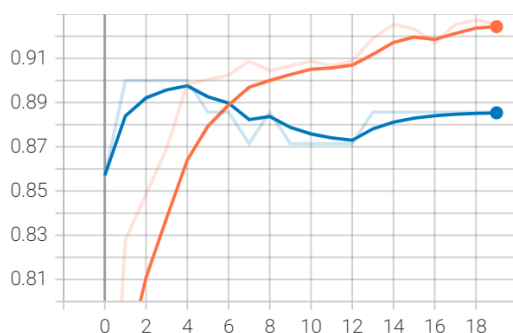


چنانچه روی یکی از لایه‌ها نظیر لایه وسط سبز رنگ کلیک کنیم اطلاع بیشتری راجع به خود لایه بدست می‌آوریم:

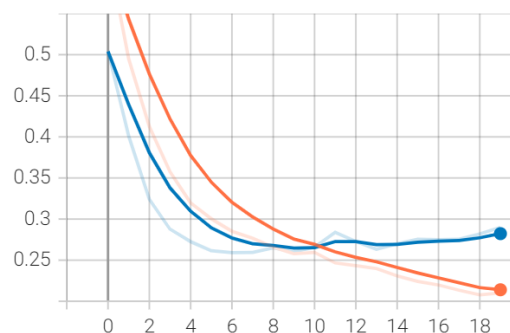


می‌توانیم بخش Scalars را هم بررسی کنیم:

epoch_accuracy

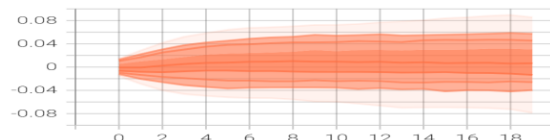
epoch_accuracy
tag: epoch_accuracy

epoch_loss

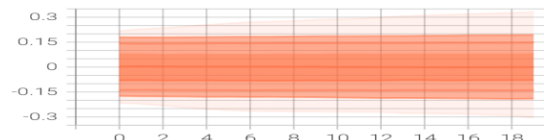
epoch_loss
tag: epoch_loss

متناسب با نتایج این قسمت به نظر می‌رسد که از گام ۱۲ به بعد مدل دچار بیش‌برازش شده است و بهتر بود که در همان قسمت یادگیری مدل متوقف شود. در قسمت بعدی tensor board یعنی Distribution نتایج زیر را می‌توان دید. نسبتاً فراوانی مقدار پارامترها قابل قبول است و از این بابت مشکل جدی‌ای دیده نمی‌شود.

Dense_Layer_1

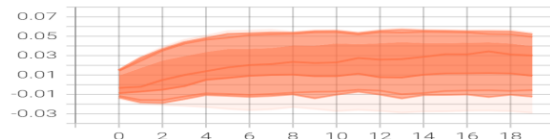
Dense_Layer_1/bias_0
tag: Dense_Layer_1/bias_0

train

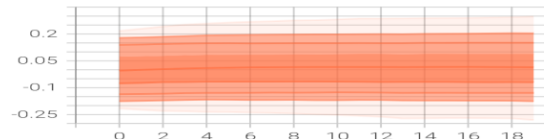
Dense_Layer_1/kernel_0
tag: Dense_Layer_1/kernel_0

train

Dense_Layer_2

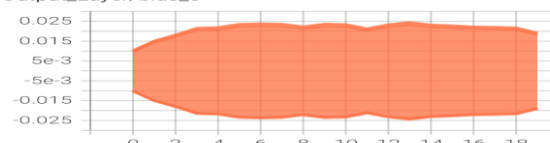
Dense_Layer_2/bias_0
tag: Dense_Layer_2/bias_0

train

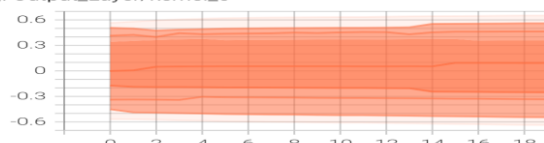
Dense_Layer_2/kernel_0
tag: Dense_Layer_2/kernel_0

train

Output_Layer

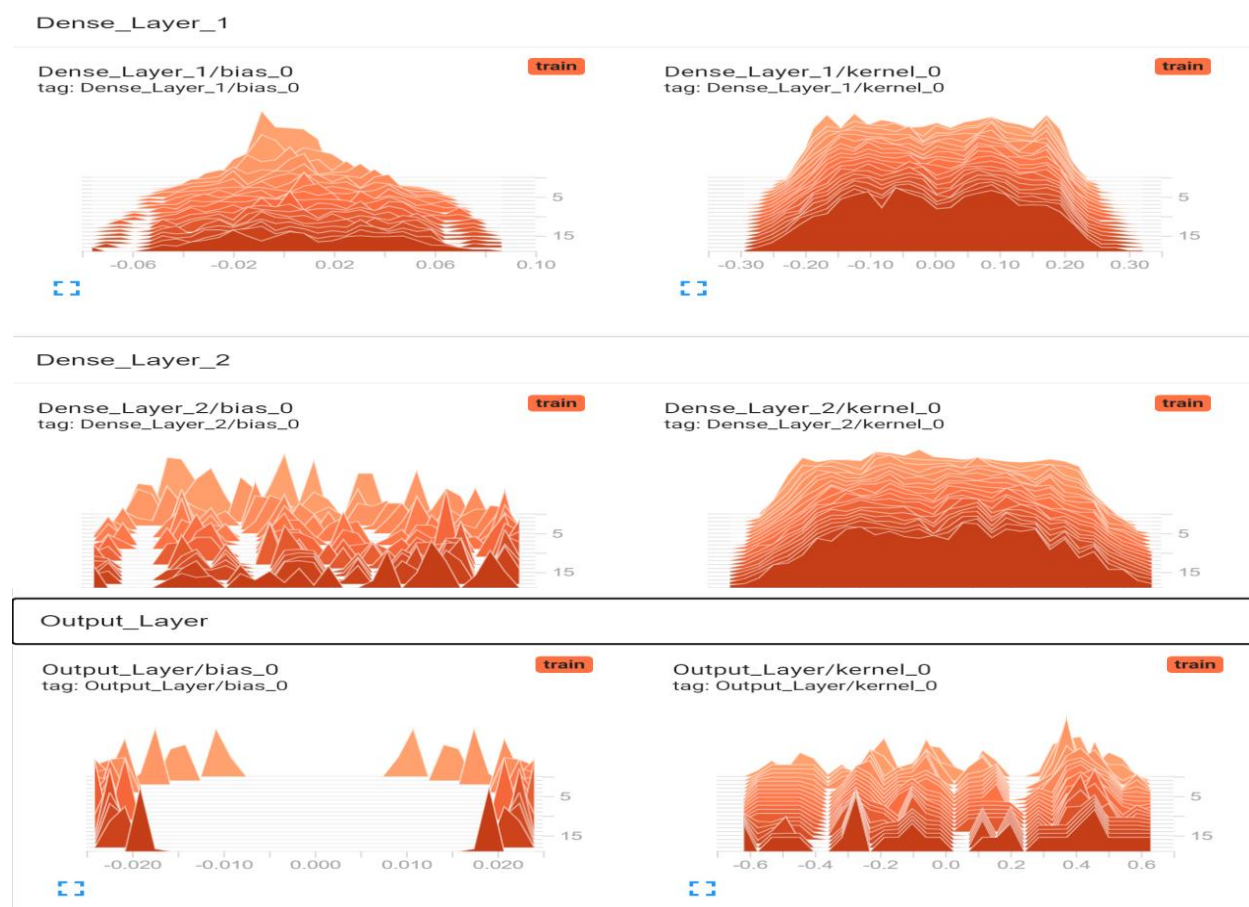
Output_Layer/bias_0
tag: Output_Layer/bias_0

train

Output_Layer/kernel_0
tag: Output_Layer/kernel_0

train

در قسمت Histogram نمودارهای زیر دیده می‌شود. به نظر bias هر سه لایه آموزش مناسبی دیده است و وزن‌های متنوعی را اختیار کرده‌اند. وزن‌های مربوط به کرنل خروجی هم قابل قبول است اما وزن‌های مربوط به کرنل دو لایه مخفی تغییر جدی‌ای پیدا نکرده‌اند که نکته مثبتی نیست.



سوال ۴

در شبکه ساده قسمت قبل دو لایه میانی با تعداد نرون ۱۲۸ و ۱۶ داشتیم. از شبکه قبل تنظیمات لایه خروجی، Optimizer، تابع خطا، نوع توابع فعال‌ساز لایه‌های میانی و تعداد گام آموزش را بدون تغییر نگه می‌دارم. تعداد لایه‌های میانی و تعداد نرون هر لایه را به عنوان پارامترهای ورودی در نظر می‌گیرم. برای این دو پارامتر طبیعتاً بی‌نهایت عدد را می‌توان تست کرد و نمی‌توان به پارامترهای بهینه رسید؛ اما می‌توان

به یکی از جواب‌های نزدیک به حالت بهینه دست پیدا کرد. در جدول زیر ۵۷ اجرای مختلف آورده شده است. تعداد لایه‌های مخفی از ۰ تا ۷ متغیر است. توجه کنید که برای هر حالت تنها یک بار اجرا انجام شده است و کاملاً محتمل است که در اجراهای بعد نتایج متفاوت باشد.

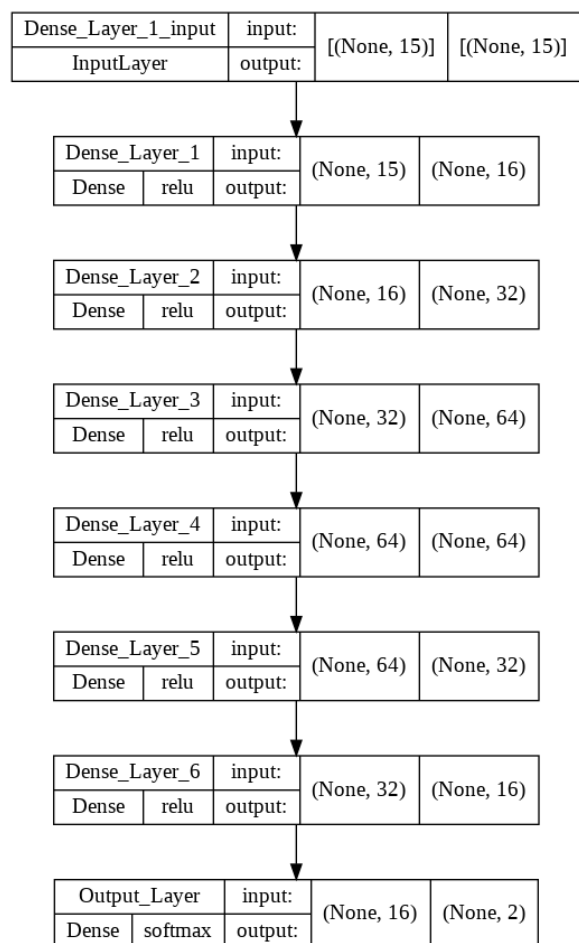
| ردیف | تعداد لایه مخفی | تعداد نورون لایه‌های مخفی | صحت اعتبارسنجی |
|------|-----------------|---------------------------|----------------|
| ۱ | ۰ | - | ۸۸/۵۷٪ |
| ۲ | ۱ | ۱۶ | ۸۸/۵۷٪ |
| ۳ | ۱ | ۳۲ | ۸۴/۲۹٪ |
| ۴ | ۱ | ۶۴ | ۸۴/۲۹٪ |
| ۵ | ۱ | ۱۲۸ | ۸۸/۵۷٪ |
| ۶ | ۱ | ۲۵۶ | ۸۵/۷۱٪ |
| ۷ | ۲ | ۱۶-۱۶ | ۸۸/۵۷٪ |
| ۸ | ۲ | ۱۶-۶۴ | ۸۴/۲۹٪ |
| ۹ | ۲ | ۱۶-۱۲۸ | ۸۴/۲۹٪ |
| ۱۰ | ۲ | ۶۴-۱۶ | ۸۵/۷۱٪ |
| ۱۱ | ۲ | ۶۴-۶۴ | ۸۷/۱۴٪ |
| ۱۲ | ۲ | ۶۴-۱۲۸ | ۸۴/۲۹٪ |
| ۱۳ | ۲ | ۱۲۸-۱۶ | ۸۵/۷۱٪ |
| ۱۴ | ۲ | ۱۲۸-۶۴ | ۸۲/۸۶٪ |
| ۱۵ | ۲ | ۱۲۸-۱۲۸ | ۸۸/۵۷٪ |
| ۱۶ | ۳ | ۲۵۶-۱۲۸-۶۴ | ۸۵/۷۱٪ |
| ۱۷ | ۳ | ۱۲۸-۶۴-۳۲ | ۸۵/۷۱٪ |
| ۱۸ | ۳ | ۶۴-۳۲-۱۶ | ۸۸/۵۷٪ |
| ۱۹ | ۳ | ۶۴-۱۲۸-۲۵۶ | ۸۵/۷۱٪ |
| ۲۰ | ۳ | ۳۲-۶۴-۱۲۸ | ۸۷/۱۴٪ |
| ۲۱ | ۳ | ۱۶-۳۲-۶۴ | ۸۵/۷۱٪ |
| ۲۲ | ۳ | ۱۲۸-۱۲۸-۱۲۸ | ۸۸/۵۷٪ |
| ۲۳ | ۳ | ۱۶-۱۶-۱۶ | ۸۴/۲۹٪ |

| | | | |
|--------|---------------------|---|----|
| ጸ፩/፱% | ፲፯-፲፭-፲፯ | ፯ | ፯፭ |
| ጸ፱/፲፭% | ፲፭-፲፯-፲፭ | ፯ | ፯፩ |
| ጸ፭/፯፱% | ፲፯-፭፭-፯፯-፲፭ | ፯ | ፯፭ |
| ጸ፯/ጸ፭% | ፯፩፭-፲፯-፭፭-፯፯ | ፭ | ፯፱ |
| ጸ፱/፲፭% | ፯፩፭-፭፭-፯፯-ጸ | ፭ | ፯ጸ |
| ጸ፩/፱% | ፲፭-፯፯-፭፭-፲፯ | ፭ | ፯፱ |
| ጸ፯/ጸ፭% | ፯፯-፭፭-፲፯-፯፩፭ | ፭ | ፯፬ |
| ጸ፩/፱% | ጸ-፯፯-፭፭-፯፩፭ | ፭ | ፯፱ |
| ጸ፩/፱% | ፯፩፭-፯፩፭-፯፩፭-፯፩፭ | ፭ | ፯፯ |
| ጸጸ/፩፱% | ፲፭-፲፭-፲፭-፲፭ | ፭ | ፯፯ |
| ጸ፲/፭፯% | ፲፯-፭፭-፯፯-፲፭-ጸ | ፩ | ፯፭ |
| ጸ፭/፯፱% | ፯፩፭-፲፯-፭፭-፯፯-፲፭ | ፩ | ፯፩ |
| ጸ፩/፱% | ፯፩፭-፲፯-፭፭-፭፭-፲፭-ጸ | ፩ | ፯፭ |
| ጸ፯/ጸ፭% | ጸ-፲፭-፯፯-፭፭-፲፯ | ፩ | ፯፱ |
| ጸ፩/፱% | ፲፭-፯፯-፭፭-፲፯-፯፩፭ | ፩ | ፯ጸ |
| ጸ፩/፱% | ጸ-፲፭-፭፭-፲፯-፯፩፭ | ፩ | ፯፱ |
| ጸ፱/፲፭% | ፲፯-፲፯-፲፯-፲፯-፲፯ | ፩ | ፭፬ |
| ጸ፭/፯፱% | ፲፭-፲፭-፲፭-፲፭-፲፭ | ፩ | ፭፱ |
| ጸ፱/፲፭% | ፲፭-፯፯-፭፭-፯፯-፲፭ | ፩ | ፭፯ |
| ፱፬/፬% | ፭፭-፯፯-፲፭-፯፯-፭፭ | ፩ | ፭፯ |
| ጸ፱/፲፭% | ፲፯-፭፭-፯፯-፲፭-ጸ-፭ | ፭ | ፭፭ |
| ጸ፯/ጸ፭% | ፯፩፭-፲፯-፭፭-፯፯-፲፭-ጸ | ፭ | ፭፩ |
| ጸ፩/፱% | ፭-ጸ-፲፭-፯፯-፭፭-፲፯ | ፭ | ፭፭ |
| ጸ፱/፲፭% | ጸ-፲፭-፯፯-፭፭-፲፯-፯፩፭ | ፭ | ፭፱ |
| ፱፲/፭፯% | ፲፭-፯፯-፭፭-፭፭-፯፯-፲፭ | ፭ | ፭ጸ |
| ጸ፱/፲፭% | ፭፭-፯፯-፲፭-፲፭-፯፯-፭፭ | ፭ | ፭፱ |
| ጸ፱/፲፭% | ፲፯-፲፯-፲፯-፲፯-፲፯-፲፯ | ፭ | ፩፬ |
| ጸ፩/፱% | ፲፭-፲፭-፲፭-፲፭-፲፭-፲፭ | ፭ | ፩፱ |
| ጸጸ/፩፱% | ፯፩፭-፲፯-፭፭-፯፯-፲፭-ጸ-፭ | ፱ | ፩፯ |
| ጸ፬/፬% | ፭-ጸ-፲፭-፯፯-፭፭-፲፯-፯፩፭ | ፱ | ፩፯ |

| | | | |
|--------|-----------------------------|---|----|
| ۸۷/۱۴٪ | ۱۶-۳۲-۶۴-۱۲۸-۶۴-۳۲-۱۶ | ۷ | ۵۴ |
| ۹۰/۰٪ | ۱۲۸-۶۴-۳۲-۱۶-۳۲-۶۴-۱۲۸ | ۷ | ۵۵ |
| ۸۷/۱۴٪ | ۱۲۸-۱۲۸-۱۲۸-۱۲۸-۱۲۸-۱۲۸-۱۲۸ | ۷ | ۵۶ |
| ۹۰/۰٪ | ۱۶-۱۶-۱۶-۱۶-۱۶-۱۶-۱۶ | ۷ | ۵۷ |

در این جدول ۴ شبکه برتر عبارت است از:

| ردیف | تعداد لایه مخفی | تعداد نوروں لایه‌های مخفی | صحت اعتبارسنجی |
|------|-----------------|---------------------------|----------------|
| ۴۸ | ۶ | ۱۶-۳۲-۶۴-۶۴-۳۲-۱۶ | ۹۱/۴۳٪ |
| ۴۳ | ۵ | ۶۴-۳۲-۱۶-۳۲-۶۴ | ۹۰/۰٪ |
| ۵۵ | ۷ | ۱۲۸-۶۴-۳۲-۱۶-۳۲-۶۴-۱۲۸ | ۹۰/۰٪ |
| ۵۷ | ۷ | ۱۶-۱۶-۱۶-۱۶-۱۶-۱۶-۱۶ | ۹۰/۰٪ |



معماری برترین شبکه یعنی مدل ۴۸-ام به صورت مقابل است. با این مدل می‌توان به صحت این مدل روی داده‌های آموزش برابر با ۸۵/۵۱٪ درصد است. که ۵ درصد کمتر از صحت اعتبارسنجی است ولی باز فاصله زیادی با آن ندارد.

با بررسی جدول نتایج موجود در گزارش می‌توان دید که مدل‌هایی که یک لایه مخفی دارند یا اصلاً لایه مخفی‌ای ندارند به طور کلی صحت پایینی دارند و مدل‌هایی که ۷ لایه دارند (بر خلاف انتظار من!) از نظر میانگین بهترین نتایج را داشته‌اند. اولین لایه برخی از مدل‌ها دارای تعداد نوروں پایینی مانند ۴ نوروں است (مدل ۴۶ یا ۵۳). در

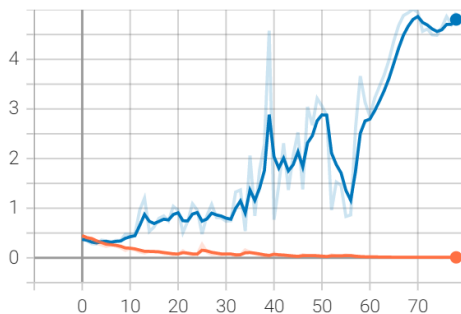
این موارد بخش مهمی از دانش شبکه در اولین لایه از بین می‌رود و خروجی هم مناسب نخواهد بود. همچنین برخی از مدل‌های نسبتاً پیچیده مانند مدل ۳۲ تنها دارای صحت ۸۵٪ روی مجموعه اعتبارسنجی هستند ولی صحت ۹۹٪ی روی مجموعه آموزشی دارند! این نشان می‌دهد که این مدل‌های پیچیده بر روی مجموعه آموزشی بیش‌برازش شده است که در سوال بعدی دقیق‌تر بررسی می‌شوند. نهایتاً باید توجه کرد که حتی مدل ۷ یا ۱۵ با داشتن دو لایه مخفی صحتی نزدیک به بهترین صحت داشته است (۸۸/۵۷٪). این نشان می‌دهد که برای این مسئله حتی با مدل‌های ساده هم می‌توان به صحت قابل قبولی رسید و اگر قرار باشد در مصالحه پیچیدگی مدل، زمان اجرا و صحت یک مدل را انتخاب کنیم، قطعاً این مدل دو لایه بهتر از مدل‌های شش هفت لایه خواهد بود.

سوال ۵

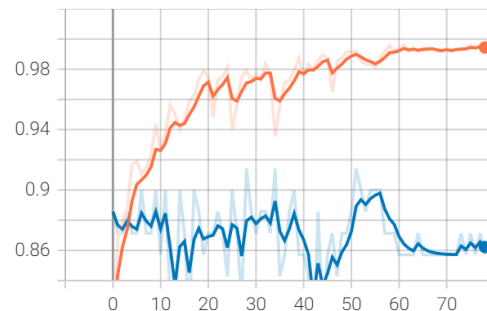
برای دستیابی به یک مدل بیش‌برازش شده باید یک شبکه پیچیده با تعداد پارامتر بسیار زیاد در نظر بگیریم تا به یک صحت زیاد روی داده‌های آموزش ولی صحت کم روی داده‌های اعتبارسنجی برسیم. مثلاً یک شبکه با هشت لایه مخفی که در هر لایه ۵۱۲ نورون داشته باشد به نظر مناسب می‌آید. سایر تنظیمات مانند دو سوال قبل است با این تفاوت که ۲۰۰ گام برای آموزش در نظر گرفته‌ام تا مدل زمان کافی برای بیش‌برازش شدن روی داده‌های آموزش را داشته باشد. نهایتاً آنکه از یک call back از نوع Early stopping روی صحت آموزشی استفاده کرده‌ام.

پس از آموزش، صحت مدل روی داده‌های آموزشی برابر با ۹۹/۱۷٪ است که عدد بسیار مناسبی است اما صحت همین مدل روی داده‌های اعتبارسنجی و تست که در زمان آموزش دیده نشده است به ترتیب عبارت است از ۸۷/۱۴٪ و ۸۱/۱۶٪. اگر مقدار خطا را در نظر بگیریم اختلاف شدیدتر هم می‌شود؛ مقدار خطا بر روی داده‌های آموزشی، اعتبارسنجی و تست به ترتیب برابر است با ۰/۰۱۳۰، ۳/۶۶۲۳ و ۵/۲۶۴۳. این‌ها همه حاکی از آن است که مدل بیش‌برازش شده است. چرا که یک مدل بیش‌برازش شده روی مجموعه داده آموزشی دقت بسیار خوبی دارد ولی روی مجموعه داده تست خیر. این مسئله در قسمت tensor board هم به خوبی عیان است:

epoch_loss
tag: epoch_loss



epoch_accuracy
tag: epoch_accuracy

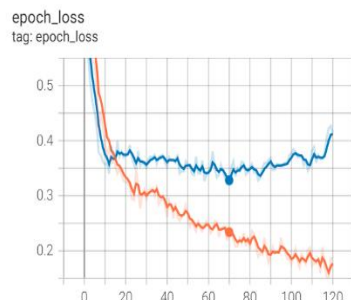


همچنین می‌دانیم مدل بیش‌برازش‌شده دارای تعداد پارامتر زیاد است. باتوجه به اینکه مدل حتی می‌تواند با دو لایه ساده با تعداد نورون نسبتاً کم مخفی می‌تواند به صحت مناسبی برسد طبیعی است که این تعداد از لایه‌های با تعداد نورون بالا از نظر تئوری مدل را بیش از حد پیچیده می‌کند.

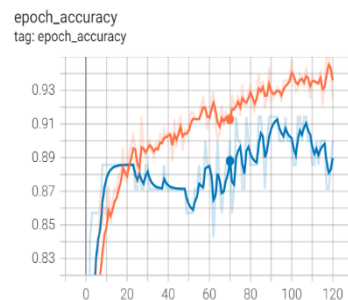
سوال ۶

برای آنکه تعمیم‌پذیری مناسب داشته باشیم شرطش آن است که پیچیدگی مدل مناسب باشد و صحت مدل روی مجموعه اعتبارسنجی بالا باشد. لذا یکی از مدل‌های مناسب سوال ۴ مانند مدل هفت لایه شماره ۵۵ را انتخاب کردم چراکه مطمئن هستیم دقت روی داده‌های اعتبارسنجی بالایی دارد. (علت عدم انتخاب مدل ۴۸ صحت پایین آن در اجرای مجدد بوده است!)؛ در عین حال در میان لایه‌های مخفی، لایه‌های Dropout برای تعمیم‌پذیری بیشتر اضافه کردم. از یک callback از نوع Early stopping و بر روی خطای اعتبارسنجی کمک گرفتم تا مطمئن شوم مدل زیاد از حد روی داده‌های آموزش پیش نمی‌رود.

پس از آموزش، مقدار accuracy بر روی سه مجموعه آموزشی، اعتبارسنجی و تست به ترتیب عبارت است از: ۹۴/۴۰٪، ۹۰٪ و ۸۴/۷۸٪. برای اطمینان به سراغ tensor board هم می‌رویم:



| Name | Smoothed Value | Value | Step | Time | Relative |
|------------|----------------|--------|------|---------------------|----------|
| train | 0.2339 | 0.2363 | 70 | Sun Mar 6, 12:10:28 | 12s |
| validation | 0.3282 | 0.3214 | 70 | Sun Mar 6, 12:10:28 | 12s |



| Name | Smoothed Value | Value | Step | Time | Relative |
|------------|----------------|--------|------|---------------------|----------|
| train | 0.9129 | 0.9108 | 70 | Sun Mar 6, 12:10:28 | 12s |
| validation | 0.888 | 0.9 | 70 | Sun Mar 6, 12:10:28 | 12s |

باتوجه به آنکه مقدار ۵۰ برای patient تنظیم شده است، خروجی فاز آموزش، خروجی مدل در گام ۷۰ خواهد بود. در این گام مقدار خطا در یک مینیمم مناسب قرار دارد و مقدار Accuracy هم قابل قبول است. اما اگر از callback استفاده نمی‌کردیم خروجی گام‌های آخر را می‌داشتیم که دچار بیش‌برازش شده است چراکه در گام ۱۲۰ خطای آموزش شدیداً کاهش یافته است ولی میزان خطای اعتبارسنجی در حال افزایش بوده است. در گام‌های آخر میزان accuracy داده‌های اعتبارسنجی هم دیگر از روال افزایشی خود خارج شده است.