

تمرین اول پایتون

شبکه‌های عصبی بهار ۱۴۰۲

آریا ابراهیمی ۹۸۲۲۷۶۲۱۷۵

هدف از این تمرین، بررسی تاثیر تغییر تعداد لایه‌ها، تعداد نوروها و تغییر optimizer در یک شبکه عصبی می‌باشد. دیتاست استفاده شده در این تمرین، دیتاست MNIST می‌باشد و برای راحتی کار، کلاسی برای شبکه عصبی چندلایه تعریف شده است که مشخصات شبکه (نرخ یادگیری، تعداد لایه‌ها، تعداد نوروهای هر لایه و ...) را از طریق یک فایل json دریافت کرده و شبکه‌ای با مشخصات داده شده ایجاد می‌کند. سپس در فایل main.py یک شی از کلاس MLP برای کانفیگ انتخابی (می‌تواند نام یک کانفیگ یا تمامی کانفیگ‌ها باشد؛ برای آموزش تمامی کانفیگ‌ها باید از مقدار ورودی all استفاده کرد) ساخته شده و یادگیری می‌شود. مقادیر loss و accuracy در هر ۱۰۰ گام با استفاده از SummaryWriter ذخیره می‌شوند (هم برای داده‌های train و هم برای داده‌های test که در نمودارهای نمایش داده شده از مقادیر برای داده‌های test استفاده شده است) تا بتوان نتایج را در Tensorboard مشاهده کرد.

۱- بخش اول

۱.۱- بررسی تغییر تعداد لایه‌ها

میتوان در شکل ۱ مشاهده کرد که با افزایش تعداد لایه‌ها، روند یادگیری کندتر شده است. برای شبکه‌های سه‌لایه، روند یادگیری تقریباً بعد از ۵۰۰ گام converge شده است ولی شبکه عمیق‌تر با ده لایه و با استفاده از RMSProp optimizer تقریباً در گام‌های نهایی converge شده است.

با افزایش تعداد لایه‌ها، توانایی مدل برای شناسایی الگوهای پیچیده‌تر بیش‌تر می‌شود، ولی از طرف دیگر ریسک overfit شدن را افزایش می‌دهد و همچنین باعث می‌شود که شبکه سخت‌تر optimize شود.

همچنین ممکن است مشکل vanishing gradients نیز پیش بیاید. این مشکل بیان‌گر این است که در شبکه‌های عمیق، اگر گرادیان‌ها کوچک باشند، در back propagation هرچه قدر به عقب بروند کوچک‌تر می‌شوند که باعث می‌شود وزن‌های لایه‌های ابتدایی به خوبی تغییر نکنند. از آنجایی که در شبکه استفاده شده در این تمرین از تابع فعال‌ساز ReLU استفاده شده است، احتمال رخداد این مشکل کمتر است و احتمالاً دیر converge شدن به‌خاطر زیاد بودن پارامترها است.

میتوان مشاهده کرد که شبکه‌های کوچک‌تر هم سریع‌تر یادگیری شده‌اند و هم نتایج بهتر یا تقریباً برابری را می‌دهند که می‌تواند برای این باشد که داده‌های استفاده شده، پیچیدگی

زیادی ندارند و میتوان آنها را در یک فضای کوچکتر جداسازی کرد.(diminishing returns)

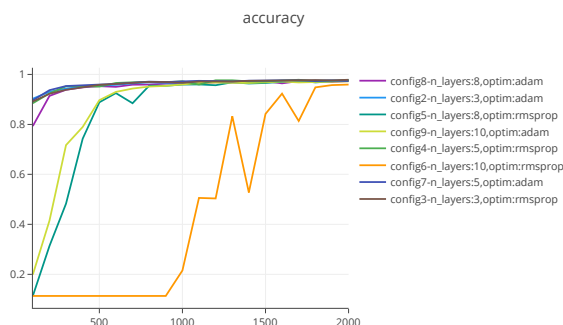


Figure 1: Accuracy of MLP for different number of layers.

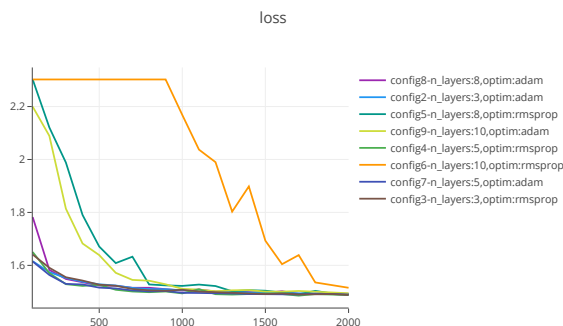


Figure 2: Loss of MLP for different number of layers.

در نتیجه استفاده از شبکه با سه لایه می‌تواند گزینه بهتری باشد زیرا هم دقت خوبی برای جداسازی داده‌های MNIST ارائه می‌کند و هم مشکلاتی که گفته شد را ایجاد نمی‌کند.

که باعث سرعت بخشیدن به یادگیری می‌شود ولی ممکن است باعث نویزی شدن آپدیت‌ها شود.

۲. **Adagrad**: هدف در این الگوریتم این است که نرخ یادگیری را adaptive کند، به این معنی که براساس شیب جهت هر پارامتر، یک scaling انجام دهد و می‌توان گفت درواقع آپدیت‌ها را نرمال می‌کند. می‌توان مشاهده کرد با استفاده از این الگوریتم، پارامترهایی که گرادیان بزرگ‌تری دارند، آپدیت کوچک‌تری خواهند داشت و گرادیان‌های کوچک‌تر آپدیت‌هایی بزرگ‌تر.

$$v_t^w = v_{t-1}^w + (\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t^w + \epsilon}} * \nabla w_t$$

$$v_t^b = v_{t-1}^b + (\nabla b_t)^2$$

$$b_{t+1} = b_t - \frac{\alpha}{\sqrt{v_t^b + \epsilon}} * \nabla b_t$$

۳. **RMSprop**: در Adagrad، پارامترها با شدت زیادی آپدیت می‌شوند. ایده RMSprop این است که به جای استفاده از مربع گرادیان‌ها از روش دیگری به نام exponentially decaying average مربع گرادیان‌ها استفاده شود.

$$S_t^w = \beta * v_{t-1}^w + (1 - \beta)(\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{S_t^w + \epsilon}} * \nabla w_t$$

$$S_t^b = \beta * v_{t-1}^b + (1 - \beta)(\nabla b_t)^2$$

$$b_{t+1} = b_t - \frac{\alpha}{\sqrt{S_t^b + \epsilon}} * \nabla b_t$$

۴. **Adam**: ترکیبی از دو optimizer مومنتموم و RMSprop می‌باشد. در الگوریتم مومنتموم از exponentially weighted averages استفاده می‌کنیم تا آپدیت‌ها smooth تر شوند. الگوریتم Adam به صورت زیر است:

۱.۲- بررسی تغییر تعداد نورن‌ها

همان‌طور که در شکل ۳ مشاهده می‌شود، در این مسئله با افزایش تعداد نورن‌ها به دقت بیشتر قابل توجهی نمی‌رسیم. با افزایش تعداد نورن‌ها در لایه‌های مخفی، می‌توان توابع پیچیده‌تری را تخمین زد ولی در این مسئله نیازی به این کار نیست زیرا شبکه‌های کوچک هم دقتی نزدیک به شبکه‌های بزرگ‌تر خواهند داشت و استفاده از شبکه بزرگ فقط زمان یادگیری را افزایش می‌دهد و باعث می‌شود که نیاز به توان محاسباتی بیشتری داشته باشیم.

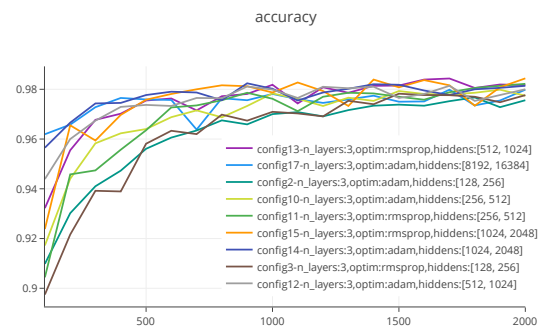


Figure 3: Accuracy of MLP for different number of neurons in hidden layers.

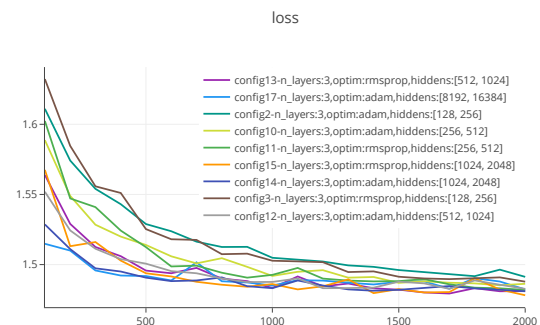


Figure 4: Loss of MLP for different number of neurons in hidden layers.

۱.۳- بررسی تغییر optimizer

در این قسمت به بررسی تاثیر optimizer های مختلف می‌پردازیم. برای این کار از optimzier های زیر استفاده شده است:

۱. **SGD**: در Stochastic Gradient Descent به جای استفاده از کل داده‌ها یا یک batch کوچک از داده‌ها برای آپدیت وزن‌ها، از تنها یک داده استفاده می‌کنیم

همان‌طور که مشاهده می‌شود، Adam و RMSprop تقریباً همانند یک‌دیگر عمل کرده‌اند و نتیجه بهتری دارند ولی SGD و Adagrad در iteration های بیشتری converge شده‌اند.

در نتیجه بهترین شبکه‌ای که trade off بین توان محاسباتی و دقت را رعایت کند می‌تواند به صورت زیر باشد:

برای optimizer از adam استفاده می‌کنیم چون در شبکه سه‌لایه به همراه RMSprop بهترین عملکرد را داشت و در شکل ۱ هم مشاهده شد که اگر شبکه عمیق‌تر شود بهتر از RMSprop عمل می‌کند.

برای عمق شبکه همان‌طور که اشاره شد، شبکه عمیق نیازی نداریم زیرا برای داده های MNIST نیازی به پیچیدگی زیادی نداریم و داده ها به یک شبکه ۳ لایه قابل جداسازی هستند. برای تعداد نورون ها در هر لایه هم می‌توان در شکل ۷ مشاهده کرد که با افزایش آنها accuracy خیلی افزایش پیدا نمی‌کند برای همین از کانفیگ شماره ۱۰ به عنوان شبکه انتخابی استفاده می‌کنیم که ۳ لایه دارد. لایه اول شامل ۲۵۶ نورون است و لایه دوم ۵۱۲ نورون دارد و optimizer استفاده شده Adam است.

$$v_{t,w} = \beta_1 * v_{t,w} + (1 - \beta_1) * (\nabla w_t)$$

$$S_{t,w} = \beta_2 * S_{t,w} + (1 - \beta_2) * (\nabla w_t)^2$$

$$v_{t,w}^{corrected} = \frac{v_{t,w}}{1 - \beta_1^t}$$

$$S_{t,w}^{corrected} = \frac{S_{t,w}}{1 - \beta_2^t}$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{S_{t,w}^{corrected} + \epsilon}} * v_{t,w}^{corrected}$$

$$v_{t,b} = \beta_1 * v_{t,b} + (1 - \beta_1) * (\nabla w_t)$$

$$S_{t,b} = \beta_2 * S_{t,b} + (1 - \beta_2) * (\nabla w_t)^2$$

$$v_{t,b}^{corrected} = \frac{v_{t,b}}{1 - \beta_1^t}$$

$$S_{t,b}^{corrected} = \frac{S_{t,b}}{1 - \beta_2^t}$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{S_{t,b}^{corrected} + \epsilon}} * v_{t,b}^{corrected}$$

حال به بررسی تفاوت optimizer ها در مسئله طبقه‌بندی MNIST با استفاده از MLP می‌پردازیم.

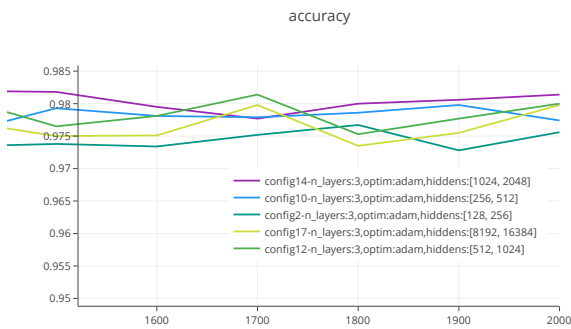


Figure 7: Accuracy of Adam optimizer for 3-layer MLPs with different number of neurons in hidden layers.

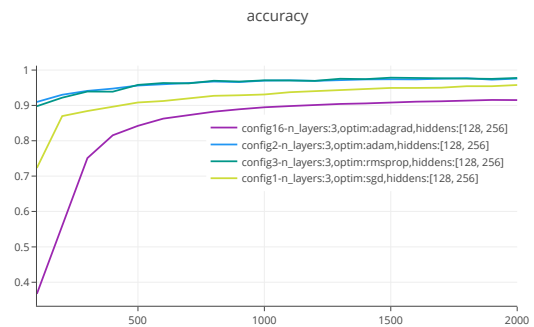


Figure 5: Accuracy of MLP for different optimizers.

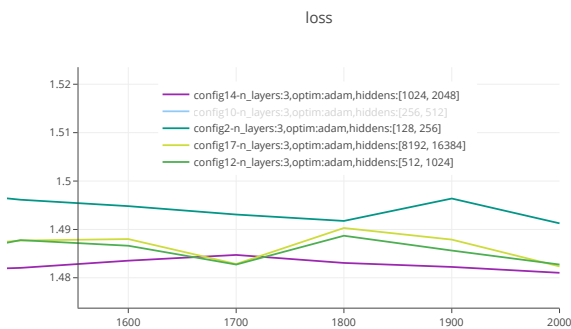


Figure 8: Loss of Adam optimizer for 3-layer MLPs with different number of neurons in hidden layers.

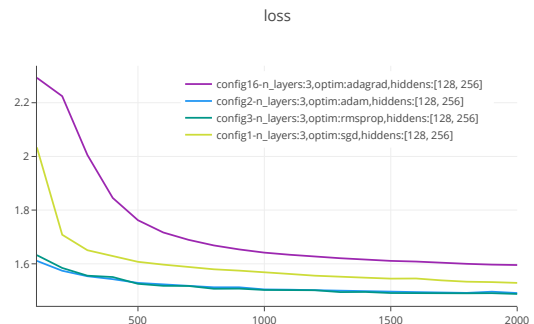


Figure 6: Loss of MLP for different optimizers.

برای رفع این مشکل اگر برای آموزش داده‌های ۵ تا ۹ را نداشته باشیم می‌توان از روش‌های unsupervised برای تشخیص داده‌های دورافتاده استفاده کرد. برای مثال می‌توان نماینده هر کلاس را در حین یادگیری پیدا کرد و در نهایت در زمان تست اگر داده‌ای به شبکه داده شد، نزدیک بودن به مرکز کلاس تخمین زده شده نیز بررسی شود و اگر فاصله از یک threshold بیشتر بود می‌توان داده را به عنوان داده دورافتاده اعلام کرد.

اگر به داده‌های دیگر (۵ تا ۹) در زمان آموزش دسترسی داشته باشیم می‌توان یک خروجی دیگر به شبکه اضافه کرد. (یک لایه Dense با یک نورون و تابع فعال‌سازی sigmoid که درواقع binary classification انجام می‌دهد) تسک مربوط به این خروجی جدید، جداسازی کلاس‌های ۰ تا ۴ از بقیه کلاس‌ها است به این صورت که اگر اعداد در بازه ۰ تا ۴ به شبکه داده شدند، این خروجی باید عددی نزدیک به ۱ و اگر داده دیگری داده شد باید عددی نزدیک به ۰ خروجی داده شود. می‌توان این خروجی جدید را با داده‌های MNIST یادگیری کرد به این صورت که به کلاس ۰ تا ۴ لیبل‌های ۱ و به کلاس‌های دیگر لیبل ۰ اختصاص داده می‌شود و با استفاده از این لیبل‌ها می‌توان تابع loss را محاسبه کرد و شبکه را آموزش داد.

در مرحله تست ابتدا قبل از اعلام لیبل پیش‌بینی شده برای ورودی، مقدار خروجی جدید را بررسی می‌کنیم و اگر مقداری نزدیک ۱ داشت می‌توان لیبل پیش‌بینی شده را گزارش کرد و اگر مقدار نزدیک ۰ بود می‌توان داده ورودی را به عنوان others پیش‌بینی کرد و به لیبل گزارش شده اهمیت نداد.

این مدل دوباره با وزن‌های ابتدایی random uniform ران شده و به نتایج زیر رسیده‌ایم: (در فایل test_properties.ipynb قابل مشاهده است)

Accuracy	Precision	Recall	F1 score
۰.۹۸۰	۰.۹۸۰	۰.۹۸۰	۰.۹۷۹

نمودار ROC در شکل ۹ و نمودار یادگیری در شکل ۱۰ نمایش داده شده‌اند.

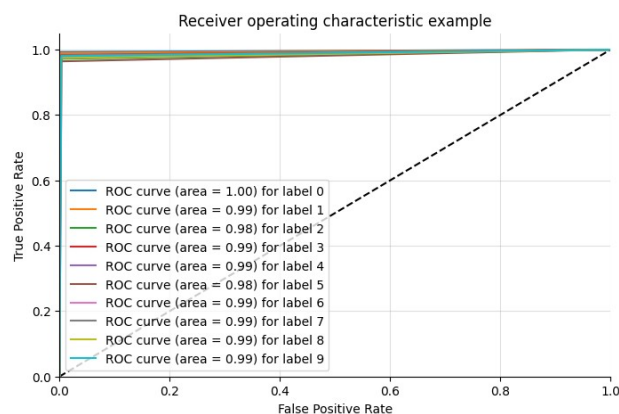


Figure 9: ROC curve of the proposed configuration.

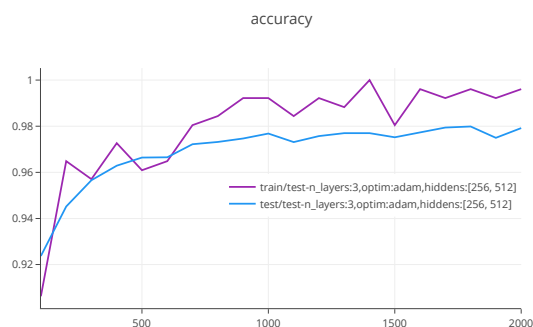


Figure 10: Learning curve of the proposed configuration.

۲- بخش دوم

با آموزش شبکه روی داده‌های ۰ تا ۴، عمل کرد شبکه برای همین کلاس‌ها در تست خوب است و به accuracy تقریباً برابر با ۹۹ درصد می‌رسیم ولی برای داده‌هایی که در مرحله آموزش مشاهده نشده‌اند عمل کرد خوبی ندارد و به اشتباه جز کلاس‌های ۰ تا ۴ تشخیص می‌دهد. (کد مربوط به این قسمت در test_zero_to_four.ipynb قابل مشاهده است)