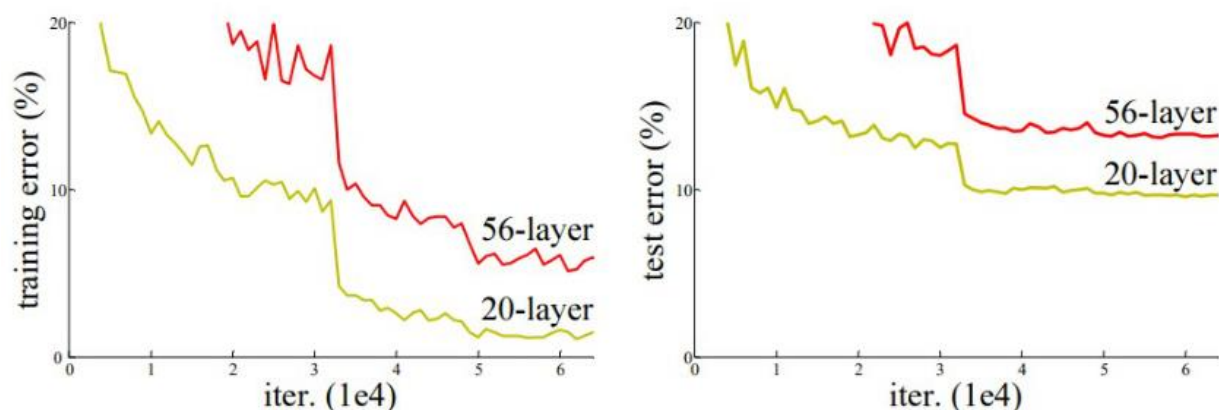


تشخیص حروف الفبای انگلیسی به وسیله شبکه های عصبی کانولوشنی و رزنت

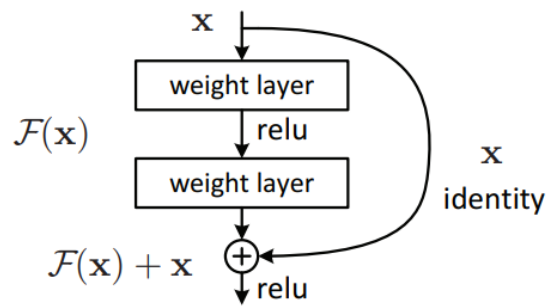
معماری شبکه عصبی ResNet

ResNet که در سال 2015 توسط محققان مایکروسافت ریسرچ پیشنهاد شد، معماری جدیدی به نام Residual Network را معرفی کرد. در معماری های شبکه عصبی هرچه لایه های بیشتری در یک شبکه عصبی عمیق استفاده می کند میزان خطا را کاهش دهد. این برای تعداد لایه های کمتری کار می کند، اما وقتی تعداد لایه ها را افزایش می دهیم، یک مشکل رایج در یادگیری عمیق مرتبط با آن وجود دارد که گرادینان Vanishing/Exploding نامیده می شود. این باعث می شود که گرادینان 0 یا خیلی بزرگ شود. بنابراین وقتی تعداد لایه ها را افزایش می دهیم، میزان خطای آموزش و آزمون نیز افزایش می یابد.



در نمودار بالا، می توانیم مشاهده کنیم که یک CNN 56 لایه نسبت به معماری 20 لایه CNN، میزان خطای بیشتری را در مجموعه داده های آموزشی و آزمایشی ارائه می دهد. پس از تجزیه و تحلیل بیشتر در مورد میزان خطا، نویسندگان توانستند به این نتیجه برسند که علت آن vanishing/exploding gradient است.

به منظور حل مشکل این معماری مفهومی به نام Residual Blocks را معرفی کرد. در این شبکه از تکنیکی به نام skip connections استفاده میکنیم. که لایه ها را با این تکنیک به لایه بعدی متصل میکنند و Resnet ها با چین این بلوک های باقیمانده در کنار هم ساخته می شوند.



مزیت افزودن این نوع اتصال پرش این است که اگر هر لایه ای به عملکرد معماری لطمه بزند، با تنظیم کردن از آن عبور می کند. در صورتیکه نگاشت همانی، نگاشت بهینه مورد نظر باشد، نزدیک کردن تابع residual به صفر خیلی ساده تر از بدست آوردن نگاشت همانی توسط چندین لایه کانولوشنی پشت سر هم با توابع غیرخطی هست. در واقع اینطور میتوان بیان کرد که شبکه هر جا لازم بدون لایه های وزن دار رو در نظر میگیره و اثر اون رو تو خروجی میاره و هر جا نیاز به نگاشت همانی بود، ضرایب کانولوشنی رو به سمتی میره که مثل نگاشت همانی عمل کنه. گرادینان از هر دو مسیر عبور میکنه از skip ها عبور میکنه تا گرادینان به لایه های ابتدایی هم برسه تا وزن بهتر و دقیق تر آپدیت بشن و کارایی شبکه بالا بره از لایه ها هم باید عبور کنه تا هر جا لازم بود ضرایب رو جوری تنظیم و آپدیت کنه که تابع residual رو صفر کنن هر جا هم لازم نبود کار خودشو میکنه و اثرش میاد تو خروجی در هر حال اون وزن ها آپدیت میشوند.

فرمول نگاشت به صورت زیر است

$$F(x) := H(x) - x \text{ which gives } H(x) := F(x) + x.$$

Architecture

| layer name | output size | 18-layer | 34-layer | 50-layer | 101-layer | 152-layer |
|------------|-------------|---|---|---|--|--|
| conv1 | 112×112 | 7×7, 64, stride 2 | | | | |
| conv2_x | 56×56 | 3×3 max pool, stride 2 | | | | |
| | | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$ |
| conv3_x | 28×28 | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$ | $\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$ |
| conv4_x | 14×14 | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$ | $\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$ |
| conv5_x | 7×7 | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$ | $\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$ |
| | 1×1 | average pool, 1000-d fc, softmax | | | | |
| FLOPs | | 1.8×10^9 | 3.6×10^9 | 3.8×10^9 | 7.6×10^9 | 11.3×10^9 |

در بحث BP این قاعده که خطا و ارور از انتها برمیگردد به ابتدا و همه لایه ها را آپدیت می‌کند وجود دارد حال این خطا که بر می‌گردد از تمام مسیر ها شبکه به لایه قبلی برمیگردد . وقتی skip connection رو قرار می‌دهیم زمانی که گرادینان یکی از وزن ها صفر شد از مسیر آن گرادینان حرکت نمی‌کند و با skip connection عبور می‌کند از لایه ای که وزن آن صفر شده است . وقتی skip connection نبود آن لایه نمیتوانست نقش همانی پیدا کند و تا اگر هم وزنی ندارد ورودی و خروجی یکی شوند .

مورد بعد راجع به ساختار لایه های کانولوشنی اینکه همانطور که در شکل بالا میبینیم ما کانولوشن های $1*1$ را در لایه قرار می‌دهیم و این باعث میشود که تعداد فیچرهای در اون لایه کم یا زیاد بکنیم مثلا 64 تا فیچر داریم و میرسونیم به 128 فیچر البته آن عدد 128 نشان دهنده 128 نوع مختلف فیلترها هستند . در واقع تک تک صفحه هات ورودی را در یک وزن ضرب کرده و این باعث افزایش فیچر میشود و سائز ورودی را کاهش نمیدهد .

ساختار دیتاست

دیتا ست ما شامل سه فایل train , validation , test میباشد که این تصاویر شامل حروف الفبای انگلیسی که به صورت حرکات دست برای ناشنویان نشان میدهد میباشد . که ما این تصاویر و در گوگل درایو آپلود کرده و سپس با اتصال کولب به آن به دیتا دسترسی پیدا میکنیم .

پیاده سازی معماری Resnet-50

الف) با استفاده از ImageDataGenerator برای تولید دسته ای از داده های تصویر بصورت عدد که برای آموزش یک مدل یادگیری عمیق شروع به دسته بندی داده ها میکنیم .

کد چندین پارامتر را برای تقویت داده ها مشخص می کند:

- 1- preprocessing_function=preprocess_input نشان می دهد که یک تابع پیش پردازش به نام preprocess_input برای هر تصویر اعمال می شود.
- 2- shear_range=0.2 محدوده ای را مشخص می کند که در آن تبدیل های برشی تصادفی می توانند روی تصاویر اعمال شوند.
- 3- horizontal_flip=True نشان می دهد که چرخش های افقی تصادفی را می توان روی تصاویر اعمال کرد.
- 4- validation_split=0.4 مشخص می کند که 40 درصد از داده ها برای اعتبارسنجی استفاده می شود، در حالی که 60 درصد باقی مانده برای آموزش استفاده می شود.

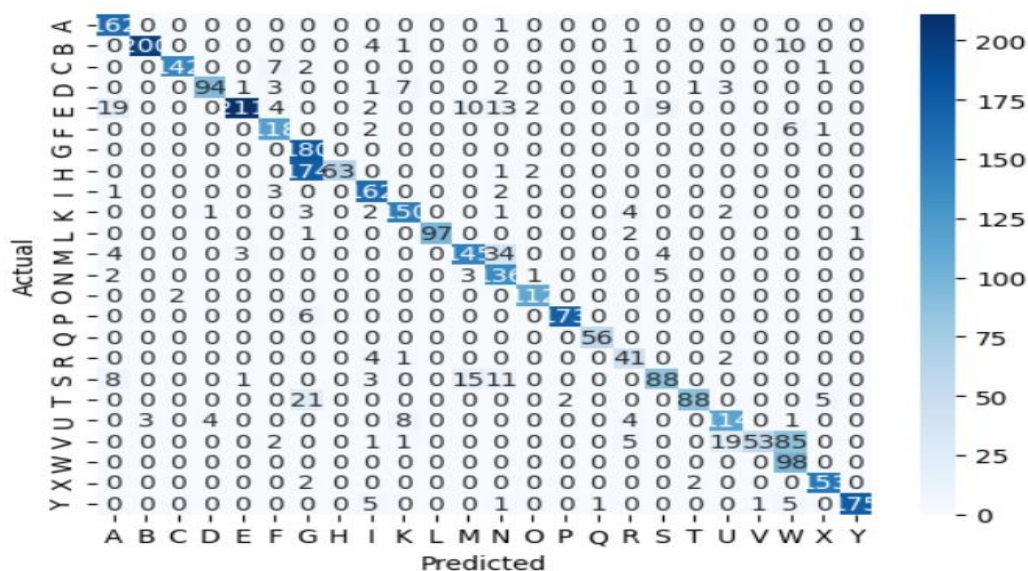
ب) با استفاده از تابع `flow_from_directory` از شی `train_datagen` یک `generator` برای داده های آموزشی ایجاد می کند که :

- 1- `train_data_dir`: این مسیر دایرکتوری است که داده های آموزشی در آن قرار دارند.
- 2- `target_size`: اندازه مورد نظر تصاویر را در داده های آموزشی مشخص می کند. چند برابر `(img_width, img_height)` طول می کشد.
- 3- `batch_size`: تعداد نمونه های هر دسته را در طول آموزش مشخص می کند.
- 4- `class_mode`: این نوع برجسب هایی که باید تولید شوند را مشخص می کند. در این مورد، روی "categorical" تنظیم می شود

به همین ترتیب `test_datagen` و `validation_datagen` را محاسبه میکنیم . پس از آن یک لایه `resnet 50` را از API `Keras` دریافت کرده و بعد یک `MaxPooling` برای کاهش ابعاد تصاویر اعمال میکنیم و بعد یک لایه `mlp` با اکتیویشن `relu` و بعد از آن `softmax` اعمال خواهیم کرد .

دقت مدل و confusion matrix

مقدار دقت مدل برابر است با 0.83 درصد و مقدار دقت تست مدل برابر با 0.833333 درصد میباشد که البته با تنظیم مقدار `epoch` مناسب میتوان به دقت مناسبی دست یافت . همچنین میتوان `confusion matrix` را برای مدل بالا محاسبه نمود که برابر مقدار زیر است .



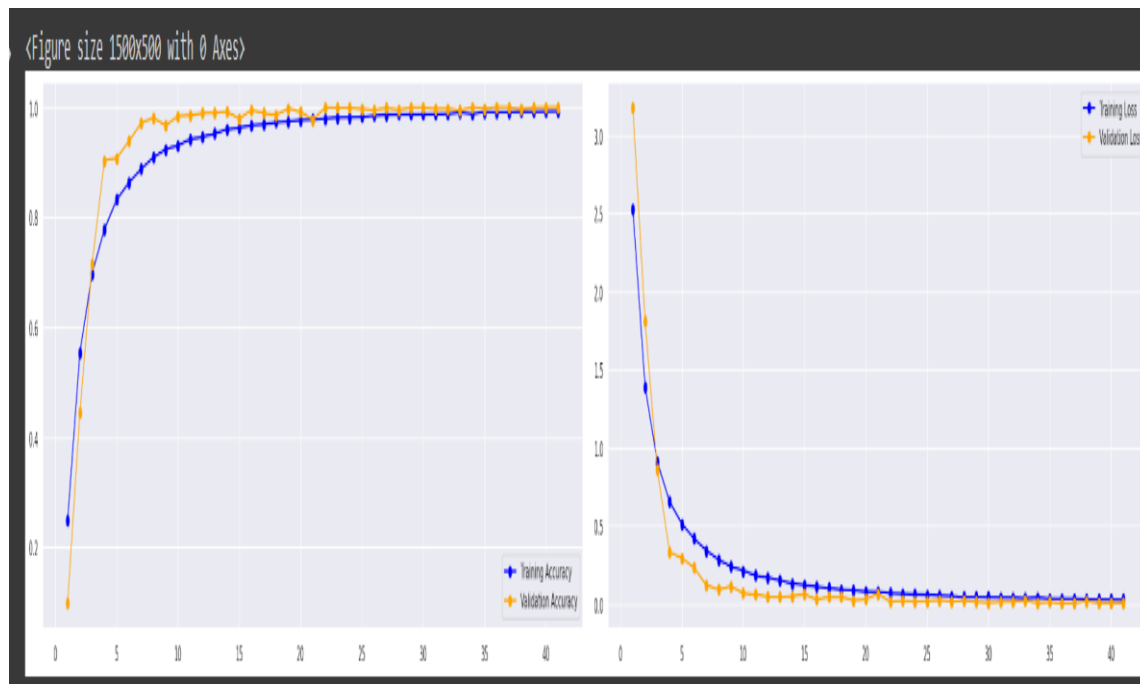
پیاده سازی با استفاده از CNN

طبق روال قبلی چون داده هامون عکس هستند ابتدا با استفاده از تابع `ImageDataGenerator` از `API Kera` عکسامون و به عدد تبدیل میکنیم. سپس یک لایه کانولوشن با اکتیویشن فانکشن `relu` و یک لایه `BatchNormalization` برای نرمال سازی خروجی های بدست آمده استفاده میکنیم و در آخر با `MaxPool2D` ابعاد لایه رو کاهش میدهیم و البته به تعداد لایه ها اضافه میشود. پس از فلت کردن تصاویر با اعمال یک لایه `dropout` و بعد یک لایه `mlp` دیگر به ساختار مناسبی میرسم.

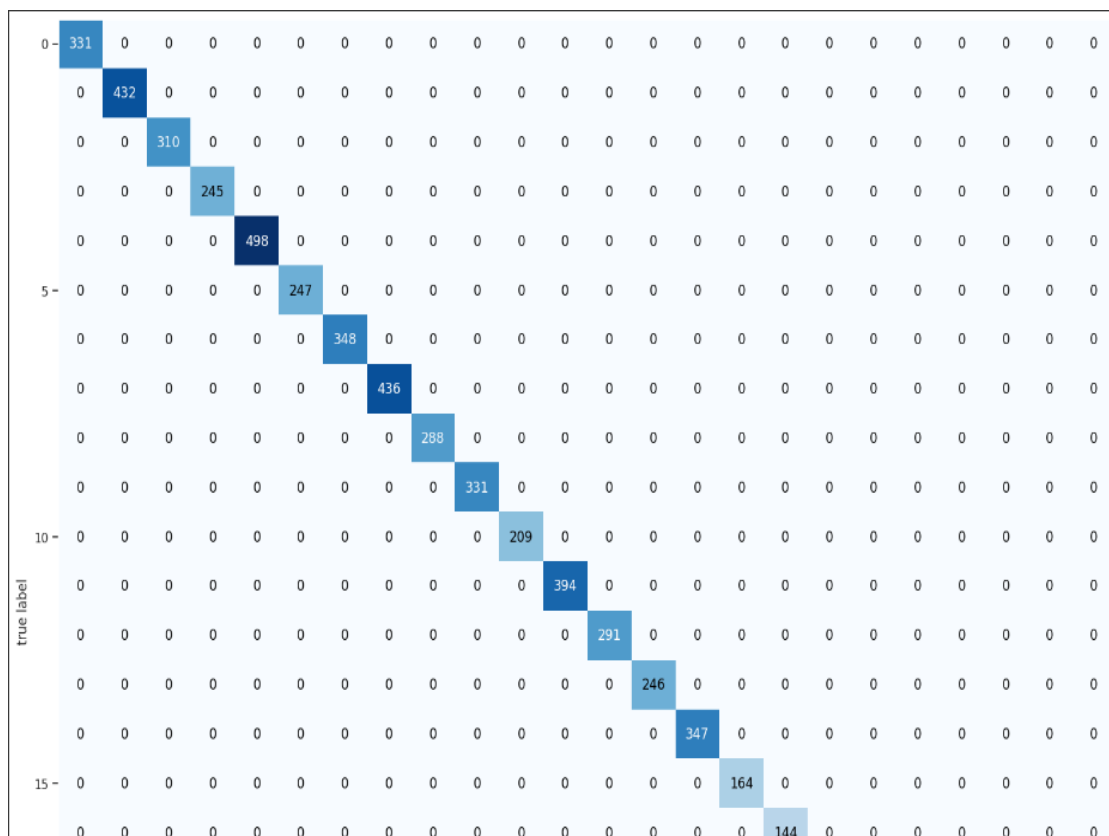
برای تشریح بیشتر اشاره ای به لایه `dropout(0.2)` میکنیم. این لایه حین آموزش این نورون ها، از تعدادی از آن ها به صورت تصادفی چشمپوشی شود. چشمپوشی یعنی اینکه آن نورون های خاص، در مسیر رفت یا برگشت در نظر گرفته نمی شوند. که با قرار دادن مقدار `0.2` نشان میده در طول فرآیند آموزش مدل مقدار `20` درصد از آن در نظر گرفته نمیشود.

دقت مدل

دقت مدل `0.99` درصد و دقت تست آن `0.9998` میباشد و نمودار در زیر نشان دهنده تغییرات دقت و مقدار `loss` را مشاهده میکنید.



Confusion matrix



به طور کلی برای مقایسه الگوریتم ها

- 1- مقدار ران تایم الگوریتم CNN نسبت به ResNet کمتر بوده .
- 2- دقت الگوریتم ResNet و CNN باهم برابر هستند و قابل تنظیم اند به وسیله callback_function و تغییرات مقدار epoch ها .
- 3- پیچیدگی زمانی الگوریتم CNN کمتر است .

پایان