



دانشگاه تهران

پردیس دانشکده‌های فنی

دانشکده مهندسی برق و کامپیوتر

گزارش مینی پروژه اول درس شبکه عصبی

یلدا فروتن

۸۱۰۱۹۶۲۶۵

استاد درس

جناب آقای دکتر کلهر

پاییز ۹۸

۱. مرور مسائل پیرامون CNN

۱,۱ انتخاب توابع و روش‌های بهینه‌سازی

۱,۱,۱ مقایسه توابع فعال‌ساز

از آنجایی که مسائل مطرح‌شده در شبکه‌های عصبی از نوع Logistic Regression نیست، عموماً در لایه‌های مخفی از توابع غیرخطی استفاده می‌شود و توابع خطی پاسخگو نخواهد بود. دلیل این امر آن است که استفاده از تابع خطی در لایه‌های پنهان بی‌فایده است؛ زیرا در هنگام Update کردن پارامترها، می‌توان بدون استفاده از توابع خطی و با اعمال ضربایی به وزن‌ها، توابع خطی را پیاده‌سازی کرد. تنها استفاده از توابع فعال‌ساز خطی در شبکه‌های عصبی، در لایه آخر مسائل مربوط به رگرسیون است. حتی در این نوع مسائل نیز، در لایه‌های پنهان از توابع غیرخطی همچون ReLU, Tanh و Sigmoid استفاده می‌شود.

برای اعمال روش‌های بهینه‌سازی همچون روش Gradient Descent لازم است تابع فعال‌ساز مشتق‌پذیر باشد. بنابراین تابع در یک بازه خاص نباید دارای مقدار 0 باشد. در ادامه سه تابع فعال‌ساز غیرخطی ReLU, Tanh و Sigmoid با یکدیگر مقایسه می‌شوند.

• تابع ReLU

تابع ReLU یا Rectified Linear Unit یک تابع ساده و کار راه‌انداز است. این تابع شبیه Ramp عمل می‌کند اما دیگر خطی نیست. همچنین همان‌طور که در ادامه آمده است تابع ReLU در مقادیر مثبت همچون تابع همانی است اما برای مقادیر منفی دارای مقدار صفر است.

$$ReLU: a = \max(0, z)$$

$$Leaky ReLU: a = \max(0.01z, z)$$

مزایا

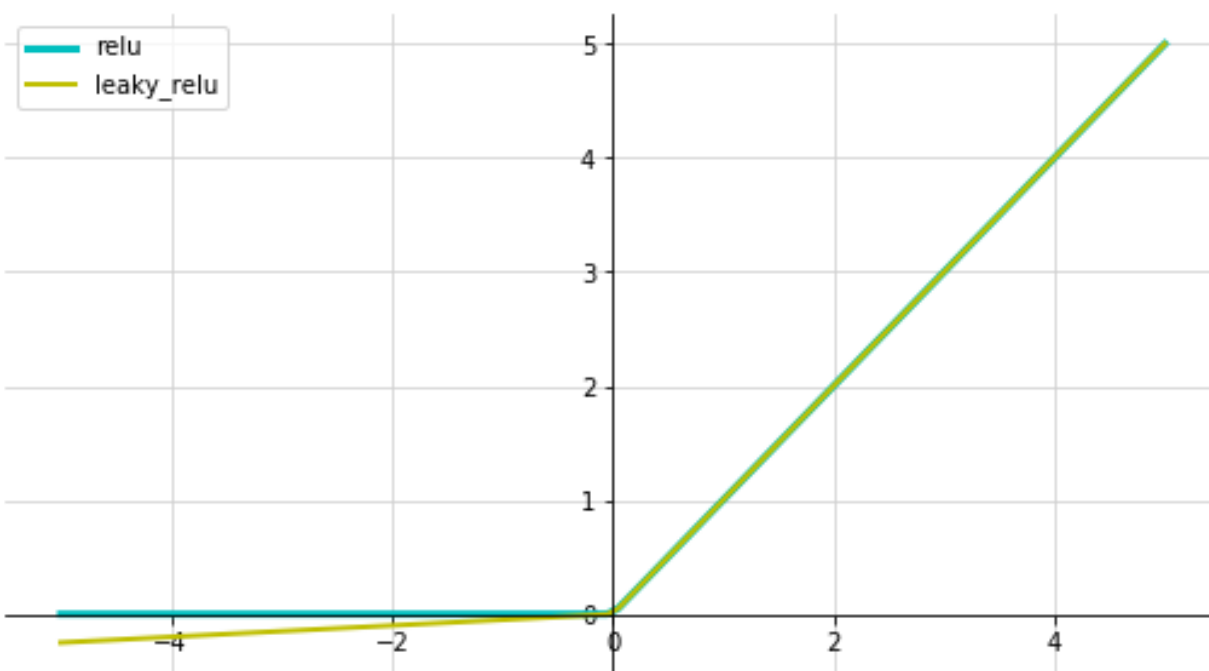
- مشتق تابع ReLU به ازای مقادیر مثبت، غیر صفر و یک‌به‌یک است؛ بنابراین مشکل Gradient Vanishing، که در قسمت بعد تعریف می‌شود، ندارد.
- تابع ReLU نسبت به سایر توابع فعال‌ساز ساده‌تر است.

معایب

- در هنگام استفاده از روش‌های Gradient Search، به ازای مقادیر منفی، مشتق تابع ReLU صفر خواهد شد و مقادیر وزن‌ها update نخواهد شد. راه‌حل آن قرار دادن مشتق تقریبی برای مقادیر منفی است که به این تابع جدید Leaky ReLU می‌گویند.
- فقط به‌عنوان تابع فعال‌ساز در لایه‌های میانی شبکه می‌توان از آن استفاده کرد.
- تابع ReLU غیرخطی است اما درواقع از دو بخش خطی ساخته شده است.

کاربرد

- تابع ReLU را می‌توان به‌عنوان Default Function نام‌گذاری کرد؛ در مواردی که رفتار تابع مشخص نیست، می‌توان از ReLU یا Leaky ReLU استفاده کرد. به‌طور کلی تابع ReLU یک General Approximator است. تابع فعال‌ساز Leaky ReLU بهتر از تابع ReLU است اما کاربرد آن کمتر است.



شکل ۱ - توابع فعال‌ساز ReLU و Leaky ReLU

• تابع Sigmoid

تابع Sigmoid شبیه به تابع Sign نرم شده است که مقادیر مختلف ورودی را به بازه ۰ تا ۱ می‌برد.

$$\text{Sigmoid: } a = \frac{1}{1 + e^{-z}}$$

مزایا

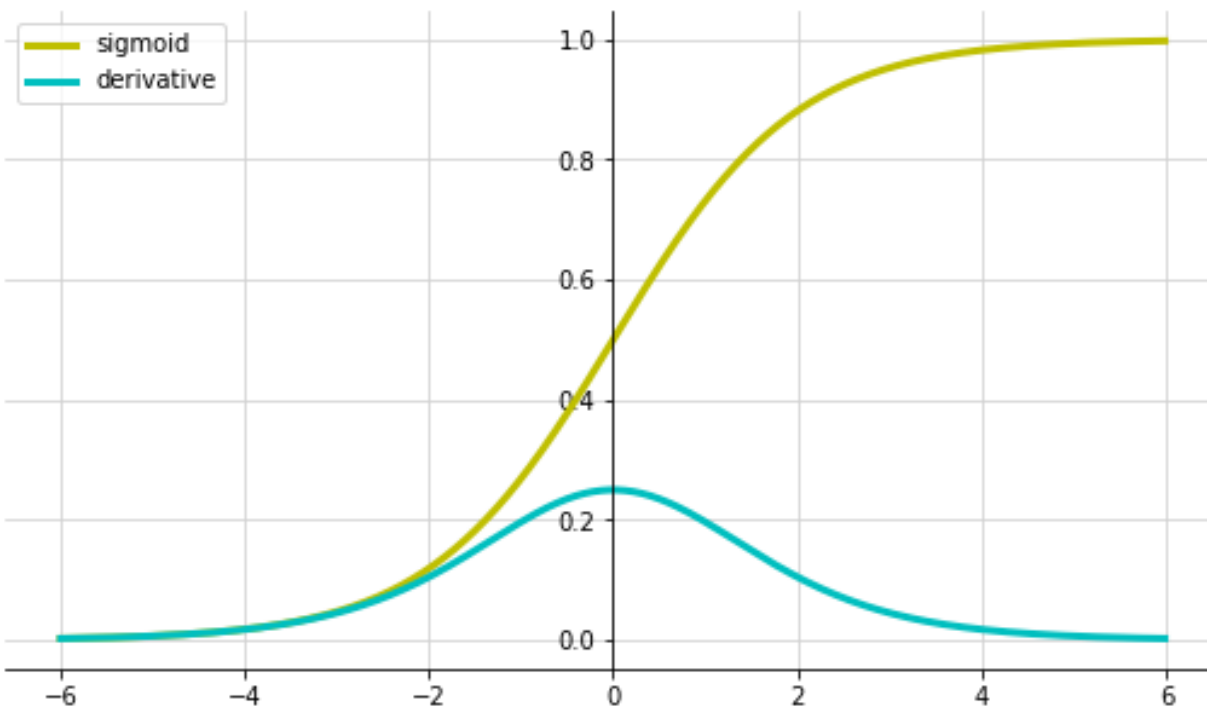
- شاید مهم‌ترین فایده تابع Sigmoid آن است که تابع پله را به صورت غیرخطی پیاده‌سازی می‌کند. همچنین در بازه ۳- تا ۳، تغییرات کوچک در ورودی تابع، منجر به تغییرات بزرگ در خروجی می‌گردد.
- تابع Sigmoid و حتی بخش‌هایی از آن همچنان غیرخطی است.
- تابع Sigmoid و مشتق آن Smooth هستند.

معایب

- یکی از معایب تابع Sigmoid آن است که تابع بین ۰ تا ۱ بوده و نسبت به ۰ متقارن نیست؛ بنابراین تمام مقادیر آن مثبت هستند.
- مشتق تابع فعال‌ساز Sigmoid همواره بین ۰ تا ۱ است. از آنجایی که در روش‌های کاهش خطا به صورت Gradient Search، مشتقات در زنجیره‌های طولانی در یکدیگر ضرب می‌شوند، احتمال کوچک کردن (نزدیک به صفر شدن) وزن‌های مربوط به لایه‌های ابتدایی شبکه زیاد است که به آن Gradient Vanishing می‌گویند زیرا اثر گرادیان در لایه‌های عقب‌تر را کم می‌کند و مؤلفه خطای بازگشتی منجر به یادگیری نمی‌شود.

کاربرد

- تابع Sigmoid برای کاربردهای Classification مناسب است و سریع‌تر از ReLU فرایند یادگیری را همگرا می‌سازد.
- تابع Sigmoid برای لایه خروجی مسائل Binary Classification که دارای دو کلاس هستند، مناسب است.



شکل ۲ - تابع فعال‌ساز Sigmoid و مشتق آن

• تابع Tanh

تابع Tanh درواقع شیفت‌یافته تابع Sigmoid است که به‌صورت زیر تعریف می‌شود.

$$\text{Tanh}: a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

مزایا

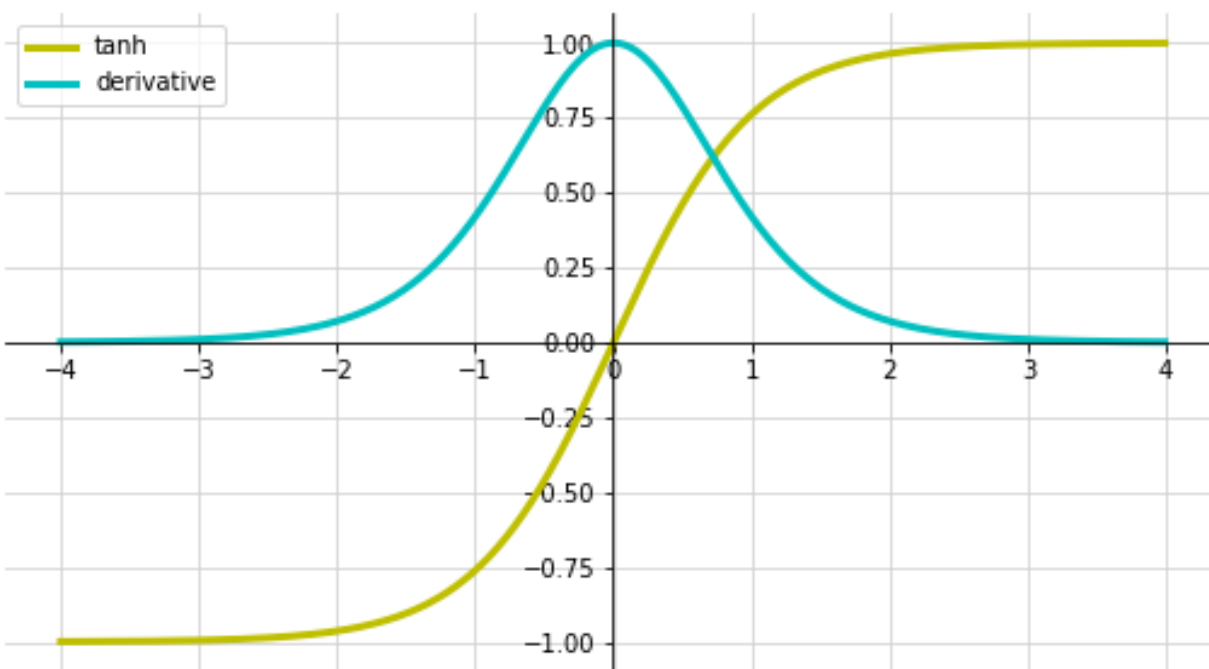
- تابع Tanh عملکردی شبیه به Sigmoid دارد با این تفاوت که نسبت به صفر متقارن است. بنابراین فرایند یادگیری برای آن آسان‌تر است.
- گرادیان تابع Tanh نسبت به گرادیان تابع Sigmoid، بزرگ‌تر بوده که فرایند بهینه‌سازی را سریع‌تر می‌کند.

معایب

- تابع Tanh همانند Sigmoid منجر به Gradient Vanishing می‌شود.

کاربرد

- در شبکه‌های عصبی به عنوان تابع فعال‌ساز لایه‌های پنهان مناسب است اما برای لایه‌های خروجی مناسب نیست.



شکل ۳ - تابع فعال‌ساز Tanh و مشتق آن

بنابراین تابع ReLU به دلیل داشتن مقادیر متفاوت از صفر برای مقادیر مثبت، نسبت به دو تابع Sigmoid و Tanh ارجحیت دارد.

۲,۱,۱ توابع Loss

توابع Loss Function در راستای آموزش شبکه عصبی مفید هستند؛ به گونه‌ای که با داشتن ورودی‌ها و خروجی‌های ایده‌آل (Target)، می‌توان میزان خطا بین خروجی شبکه (Predicted) و Target را محاسبه کرد. توابع خطا را می‌توان براساس نوع مسئله شبکه عصبی دسته‌بندی کرد.

- **مسائل Regression:** در این مسائل خروجی Target به صورت پیوسته است و شبکه عصبی یک تابع را تقریب می‌زند. تابع خطا رایج برای این گونه مسائل تابع Mean Square Error یا MSE است.

- مسائل **Classification**: مقادیر خروجی در مسائل Classification به صورت گسسته بوده و نمایانگر Label ورودی‌ها است. تابع خطا مرسوم برای این نوع مسائل، Cross-entropy است.
 - مسائل **Embedding**: در این نمونه مسائل، مقایسه‌ای بین دو مقدار ورودی صورت می‌گیرد که آیا شبیه به هم هستند یا نه. تابع خطا رایج برای این مسائل، Hinge Loss است.
- در ادامه چهار نمونه تابع خطا رایج به همراه مزایا و کاربرد آن‌ها بررسی شده است.

• Hinge Loss

یکی از جایگزین‌های تابع Cross-entropy برای مسائل **Binary Classification** تابع خطا Hinge است که برای مدل‌های طبقه‌بندی در ماشین‌های بردار پشتیبانی (SVMs) کاربرد دارد. درواقع تابع Hinge برای مسائل **Binary Classification** هنگامی که Target دارای مقادیر ۱ است، به کار گرفته می‌شود تا نمونه‌ها را به سمت داشتن علامت مناسب مثبت یا منفی بکشاند. در برخی موارد استفاده از تابع خطا Hinge باعث بهبود کارایی شبکه نسبت به استفاده از خطا Cross-entropy در مسائل **Binary Classification** شده است.

همان‌طور که گفته شد، لازم است خروجی مربوط به Target دارای مقادیر -۱ یا ۱+ باشد. همچنین خروجی پیش‌بینی شده باید مقادیری در بازه [-1,1] دارا باشد. بنابراین استفاده از تابع فعال‌ساز Tanh مناسب است زیرا خروجی آن در بازه مدنظر قرار دارد.

$$\text{Hinge Loss: } J(y) = \sum_{i=1}^n \max(0, 1 - t_i \cdot h_i)$$

در رابطه بالا منظور از t، خروجی Target و منظور از h خروجی پیش‌بینی شده است.

• Softmax Cross Entropy

منظور از Softmax Cross Entropy، تابع فعال‌ساز Softmax و تابع هزینه Cross entropy است. درواقع تابع Softmax، یک حالت Soft از تابع ماکزیمم است که احتمال ماکزیمم و نزدیک به آن را می‌دهد. ورودی تابع Softmax یک بردار N بعدی و خروجی آن یک بردار بین ۰ تا ۱ است که در ادامه آمده است.

$$P_i = \frac{e^{a_i}}{\sum_{k=1}^N e_k^a}$$

از آنجایی که تابع Softmax یک توزیع احتمالی به عنوان خروجی می دهد، در لایه خروجی شبکه های عصبی از آن استفاده می گردد. در هنگام Backpropagation لازم است مشتق و یا گرادیان Softmax محاسبه گردد که در ادامه آمده است.

$$\frac{\partial p_i}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} = \dots = p_i(\delta_{ij} - p_j)$$

تابع هزینه Cross Entropy فاصله بین توزیع واقعی مدل و توزیعی که شبکه به دست آورده را مشخص می کند که در ادامه رابطه آن آمده است.

$$H(y, p) = - \sum_i y_i \log(p_i)$$

در نهایت از مشتق تابع فعال ساز Softmax برای محاسبه مشتق تابع خطا Cross entropy استفاده می شود.

• Mean Squared Error

مجموع مربعات خطا، تابع پیش فرض مسائل Regression است. این خطا به صورت میانگین مجذور اختلاف بین خروجی های پیش بینی شده و خروجی Target محاسبه می گردد. MSE جدای از علامت Predict و Target، همواره دارای مقدار مثبت است و بهترین مقدار خطا برای آن صفر خواهد بود. استفاده از مجذور خطا نمایانگر آن است که اشتباهات بزرگ تر منجر به خطاهای بیش تر می شود.

$$MSE \text{ Loss: } J(y) = \frac{1}{n} \sum_{i=1}^n (t_i - h_i)^2$$

• Log Loss

تابع Log Loss یا Binary Cross Entropy برای مسائل Classification کاربرد دارد. برای نمونه اگر خروجی های مسئله به صورت دو کلاس باشند، از تابع Binary Cross Entropy استفاده می شود و رابطه آن به صورت زیر است:

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

به عنوان مثال اگر دو کلاس وجود داشته باشد و برچسب های صفر و یک را به آن ها اختصاص داده شده باشد، تابع loss به صورت زیر محاسبه می گردد:

$$y = 0 \rightarrow H_p(q) = -\frac{1}{N} (\log(1 - p(y))) \rightarrow H_p(q) = \begin{cases} 0 \\ \infty \end{cases}$$

$$y = 1 \rightarrow H_p(q) = -\frac{1}{N} \log(p(y)) \rightarrow H_p(q) = \begin{cases} 0 \\ \infty \end{cases}$$

بنابراین اگر target مربوط به یک کلاس ۰ باشد و سیستم به درستی ۰ پیش بینی کرده باشد، در این صورت مقدار خطا صفر خواهد بود. اما اگر شبکه به اشتباه ۱ تشخیص دهد، مقدار خطا بینهایت می شود. همین شیوه برای سایر label ها نیز در نظر گرفته شده است.

۳.۱.۱ روش های بهینه سازی

• Gradient Descent

گرایان کاهشی مهم ترین تکنیک برای آموزش شبکه های عصبی و روشی برای بهینه سازی سیستم های هوشمند است. گرایان کاهشی عمدتاً در راستای آپدیت کردن وزن های شبکه برای مینیمم کردن تابع خطا استفاده می شود. می دانیم که شبکه عصبی با تکنیک های Backpropagation آموزش داده می شود؛ بدین صورت که شبکه ابتدا در مسیر forward، ورودی ها و وزن های متناظر با آن ها ضرب داخلی کرده و با یک بایاس جمع و سپس تابع فعال ساز را بر روی آن پیاده می کند. سپس در مسیر Backward، شبکه با استفاده از اپتیمایزر مقادیر وزن ها و بایاس ها را آپدیت می کند. در روش گرایان کاهشی، آپدیت کردن پارامترها با کمک گرایان تابع خطا در مسیر برگشت است. روابط مربوط به گرایان کاهشی در ادامه آمده است. لازم به ذکر است تابع خطا مستقل از اپتیمایزر بوده و همانند قسمت قبل باید انتخاب گردد.

$$w = w - \alpha \frac{\partial}{\partial w} J(w)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w)$$

علاوه بر انتخاب تابع هزینه، باید مقدار (که نشان گر نرخ یادگیری است تعیین گردد. اگر (کوچک انتخاب شود، همگرایی تضمین شده اما زمان آموزش شبکه طولانی می گردد. درحالی که با انتخاب (بزرگ، شبکه سریع تر یاد می گیرد اما ممکن است واگرا گردد. درواقع انتخاب نرخ یادگیری، الگوریتم مشخصی ندارد.

مدل استاندارد گرادیان کاهشی گرادیان همه داده ها را به صورت یکجا محاسبه می کند و همه داده ها را فقط یکبار آپدیت می کند. این در حالی است روش گرادیان برای دیتاست های بزرگ کند و غیرقابل کنترل می شود. بنابراین از روش Stochastic Gradient Descent برای حل مشکل گفته شده استفاده می شود. روش SGD برای هر نمونه از داده های Training، عمل آپدیت کردن پارامترها را انجام می دهد که منجر به سرعت بیشتر نسبت به GD می شود. بنابراین روابط گرادیان برای حالت Stochastic به صورت زیر می گردد. $x(i)$ و $y(i)$ بیان گر نمونه های Training set هستند.

$$\theta = \theta - \alpha \frac{\partial}{\partial \theta} J(\theta; x(i); y(i))$$

به دلیل به روزرسانی مکرر پارامترها در روش SGD، تابع خطا نوسان خواهد کرد. بنابراین احتمال یافتن مینیمم های محلی و بهتر نسبت به روش GD افزایش می یابد. از طرفی به روزرسانی های مکرر، همگرایی به یک نقطه مینیمم محلی را سخت تر کرده و منجر به Overshooting تابع هزینه می گردد. البته با کاهش نرخ یادگیری می توان همگرایی را همانند روش GD به دست آورد.

مشکل عدم ثبات همگرایی در روش SGD منجر به استفاده از روش Mini-Batch Gradient Descent می شود.

درواقع روش Mini-Batch GD مزیت هر دو روش GD و SGD را دارد که به ازای هر n نمونه از مجموعه Training، پارامترهای شبکه را آپدیت می کند. این روش یادگیری را به سمت همگرایی با ثبات بیشتر می برد.

روش گرادیان و حتی تغییرات آن (SGD و Mini-Batch GD) همچنان دارای مشکلاتی هستند. به عنوان مثال انتخاب Learning Rate مناسب همچنان مسئله است. همچنین همه پارامترها با نرخ یادگیری یکسان آپدیت می شوند. درحالی که برای پترن های کم تکرار نرخ یادگیری بزرگ تر مطلوب است. علاوه بر آن توابع هزینه به دلیل non-convex بودن، دارای چندین نقطه مینیمم هستند. بدین صورت یافتن بهترین نقطه (Global Optimum و داشتن Domain of Attraction) سخت خواهد بود چراکه معمولاً تابع در یک نقطه مینیمم

محلی گیر می‌کند. همچنین وجود نقاط Saddle Point مشکل‌ساز خواهد بود؛ درواقع این نقاط در یک بعد مینیمم بوده و در بعد دیگر ماکزیمم هستند.

• Momentum

در روش SGD به دلیل افزایش واریانس، همگرایی مشکل می‌شود. استفاده از روش Momentum منجر به افزایش سرعت SGD می‌گردد؛ بدین‌صورت که Momentum منجر به کاهش تغییرات عمودی (مسیر نامربوط) و افزایش تغییرات افقی (مسیر مربوط) می‌شود. درواقع Momentum علاوه بر گرادیان قدم کنونی، کسری از گرادیان قدم‌های پیشین را برای آپدیت کردن بردار کنونی استفاده می‌کند.

$$V_{dw} = \beta V_{dw} + (1 - \beta)dw$$

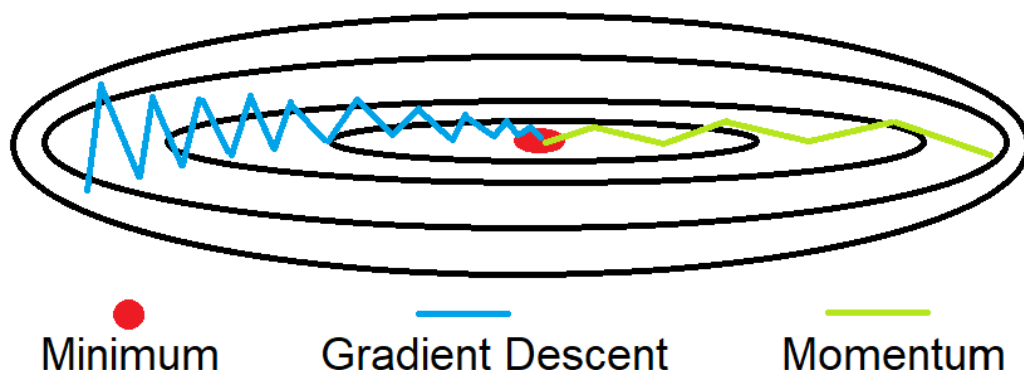
$$V_{db} = \beta V_{db} + (1 - \beta)db$$

رایج‌ترین مقدار برای پارامتر β ، ۰٫۹ است. بدین معنا که برای آپدیت کردن پارامترهای مدنظر از ۹ گرادیان قبلی میانگین می‌گیرد. درنهایت پارامترها به‌صورت زیر آپدیت می‌شوند.

$$w = w - \alpha V_{dw}$$

$$b = b - \alpha V_{db}$$

درنتیجه هنگام آپدیت شدن پارامترها با کمک momentum، همگرایی با سرعت بیشتر اتفاق می‌افتد و نوسانات را کاهش می‌دهد.



شکل ۴ - مقایسه همگرایی تابع خطا در روش Gradient Descent و Momentum

• AdaDelta

درواقع بهینه‌سازی تابع خطا به روش‌های Adaptive راه‌حلی برای انتخاب Learning Rate براساس پارامترهای شبکه هستند. به گونه‌ای که برای پارامترهای پرتغییر، آپدیت‌های کوچک و برای پارامترهای که تغییرات اندک دارند از آپدیت‌های بزرگ استفاده می‌کنند. بنابراین برای هر پارامتری براساس گرادیان قبلی، از نرخ یادگیری متفاوتی استفاده می‌شود. بنابراین در روش Adagrad از نرخ یادگیری متفاوت برای هر پارامتر (در هر قدم t استفاده می‌شود. در رابطه زیر، منظور از $g_{t,i}$ گرادیان تابع خطا است. بدیهی است که مزیت روش‌های Adaptive آن است که لازم نیست نرخ یادگیری به روش دستی تنظیم شود.

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\alpha}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}$$

در روش Adagrad، از آنجایی که مجذور گرادیان‌های قبلی جمع می‌شوند، Learning Rate کوچک و کوچک‌تر می‌گردد. کاهش Learning Rate، منجر به کاهش سرعت یادگیری شبکه می‌شود. برای حل این مشکل، الگوریتم AdaDelta پیشنهاد شده است. در AdaDelta به جای ذخیره کردن مجذور همه گرادیان‌های قبلی، بخشی از آن‌ها را ذخیره می‌کنند. به طوری که به جای ذخیره مجذور گرادیان همه پارامترهای w قبلی، مجموع گرادیان‌های قبلی به شکل بازگشتی به عنوان میانگین مجذور گرادیان پارامترهای قبلی تعریف می‌گردد. حال میانگین کنونی در $\text{step}=t$ به میانگین قبلی و گرادیان کنونی وابسته است که به صورت زیر تعریف می‌گردد. لازم به ذکر است (همان ترم momentum است که در حدود ۰,۹ انتخاب می‌شود.

$$E[g^2](t) = \beta E[g^2](t-1) + (1-\beta)g^2(t)$$

بنابراین در روش AdaDelta علاوه بر محاسبه Learning Rate برای هر پارامتر، از Momentum نیز استفاده می‌شود. همچنین از مشکل Vanishing که در Adagrad مطرح شد.

• Adam

در روش AdaDelta از نرخ یادگیری متفاوت برای آپدیت کردن هر پارامتر استفاده می‌شود. در روش Adam (Adaptive Moment Estimation) علاوه بر Learning Rate، از ترم Momentum متفاوت نیز استفاده می‌شود.

$$\widehat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

بنابراین الگوریتم Adam از رابطه زیر برای آپدیت کردن پارامترها استفاده می‌کند.

$$\theta_{t+1} = \theta_t - \frac{\alpha}{\sqrt{\hat{v}_t} + \varepsilon} \cdot \hat{m}_t$$

مقادیر رایج برای β_1 ، β_2 و α به ترتیب ۰.۹، ۰.۹۹۹ و ۰.۰۱ است. درواقع روش Adam منجر به همگرایی سریع‌تر نسبت به سایر الگوریتم‌های بهینه‌سازی می‌شود. همچنین با مشکلات سایر الگوریتم‌ها همچون همگرایی نرخ یادگیری به صفر و کاهش سرعت همگرایی تابع خطا، مواجهه نمی‌شود.

به‌طور کلی الگوریتم Adam بهتر از سایر الگوریتم‌های Adaptive عمل می‌کند. حال اگر داده‌های ورودی به اصطلاح sparse باشند، روش‌هایی مانند SGD و Momentum ضعیف عمل می‌کنند و باید از روش‌های Adaptive استفاده کرد. برای دستیابی به همگرایی سریع‌تر در مدل‌های عمیق و پیچیده، الگوریتم Adam بهتر از سایرین عمل می‌کند.

۲.۱ تکنیک‌هایی برای بهبود CNN

• تکنیک Batch Normalization

دید شده که نرمالایز کردن ورودی به دلیل هم‌مقیاس کردن، حساسیت یادگیری شبکه نسبت به مقیاس‌های متفاوت را کاهش می‌دهد. از طرفی هر Feature Map یک ورودی برای لایه‌های بعدی است. بنابراین بهتر است قبل از هر Hidden Layer، داده‌ها نرمالایز شوند. درواقع Batch Normalization باعث می‌شود توزیع داده‌های Train و Test نزدیک به هم شوند. بدین‌صورت که در هر بعد، مقادیر مربوط به هر sample را با یکدیگر جمع کرده و میانگین و سپس واریانس هر بعد را محاسبه می‌کند. بنابراین داده‌های هر بعد منهای میانگین و تقسیم بر واریانس مربوط به بعد خود می‌شوند. البته از یک‌ترم نیز استفاده شده است تا درصورت صفر شدن واریانس، داده‌ها به بین‌هایت میل نکنند. با این کار ورودی‌ها دارای میانگین صفر و واریانس یک خواهند شد. سپس با استفاده از پارامترهای (و) بر روی داده‌های جدید عمل Scale و Shift را پیاده می‌کند. پارامترهای (و) با استفاده از روش‌های Gradient Search، Learn می‌شوند. درواقع Batch Normalization سعی می‌کند

پراکندگی هر بعد را، مستقل از بعدهای دیگر، سفید بکند تا قدرت Partitioning را بهبود بخشد و با نحوه پیاده‌سازی Batch Normalization بر روی داده‌ها آمده است.

Input: Values of x over a mini-batch: $B = \{x_1 \dots x_m\}$;

Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(x_i)\}$

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

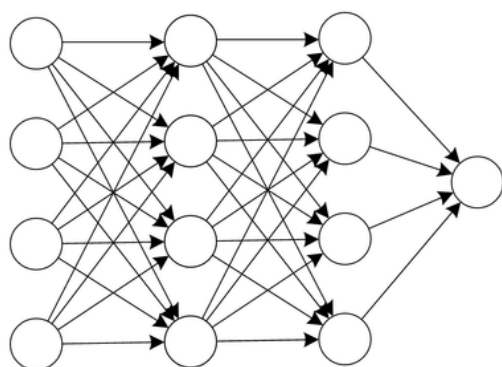
$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i)$$

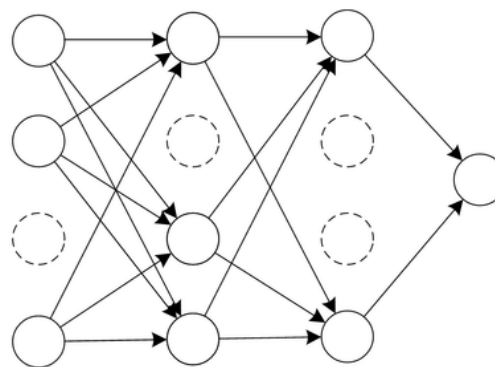
• تکنیک Drop Out

تکنیک Drop Out، در هر ایپاک به صورت تصادفی درصدی از نورون‌های هر لایه را deactivate می‌کند و سعی دارد با نورون‌های باقی‌مانده، شبکه را در مرحله call و recall آموزش دهد. در ایپاک بعد، ماسک تصادفی برای انتخاب نورون‌های فعال و غیرفعال عوض می‌شود اما درصد آن‌ها ثابت است. همانند Regularization term. هربار تابع به شکل تصادفی نادقیق‌تر می‌شود. البته اینکار منجر به کاهش دقت شبکه نمی‌شود؛ درواقع Drop Out با حذف برخی نورون‌ها، نویز تصادفی به شبکه اضافه می‌کند تا شبکه از Domain of Attraction های بد خارج شود و حتی با partially کار کردن برخی نورون‌ها، شبکه از Domain of Attraction خوب خارج نشود. البته زمانی می‌توان برخی نورون‌ها را غیرفعال کرد که به شبکه redundancy افزوده شده باشد تا با استفاده از Drop Out فضای پارامتری با دقت یادگیری و قدرت تعمیم بهتر ایجاد شود. لازم به ذکر است با خاموش کردن برخی نورون‌ها لازم است وزن نورون‌های فعال افزایش یابد تا قدرت فلو اطلاعاتی کاهش نیابد. بنابراین Drop Out قدرت یادگیری شبکه را بهتر کرده و شبکه را نسبت به خاموش شدن برخی نورون‌ها مقاوم می‌سازد. همچنین پس از پایان روند یادگیری، شبکه از قدرت تعمیم بهتری برخوردار خواهد بود. عملکرد Drop

Out در شکل زیر نشان داده شده است. شکل سمت چپ، شبکه را به صورت Fully Connected نشان می‌دهد. پس از اعمال Drop Out شبکه به صورت Partially Wiring شده در شکل سمت راست قابل مشاهده است.



(a) Standard Neural Network



(b) Network after Dropout

شکل ۵ - اعمال تکنیک Drop Out و نحوه سیم‌بندی

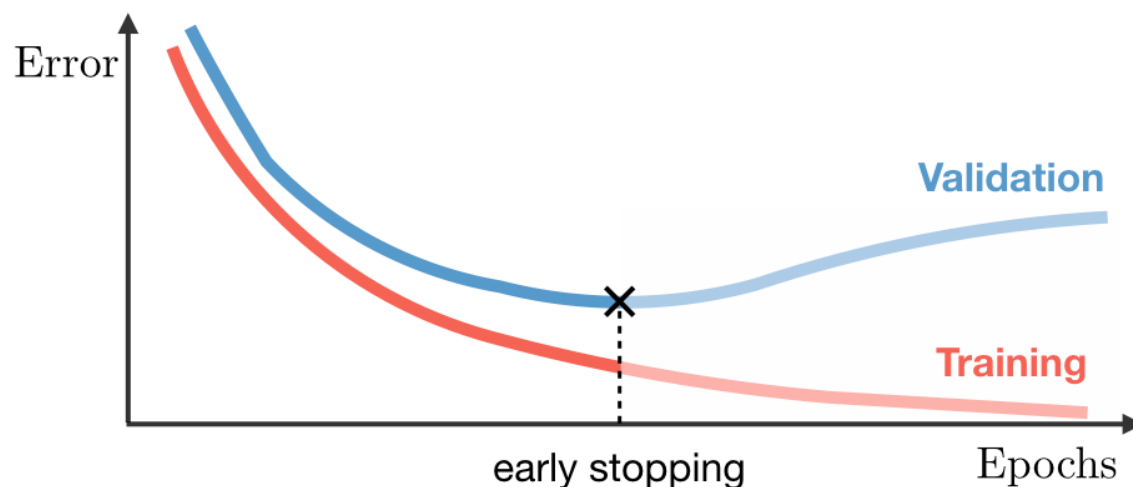
درواقع Drop Out یک Hyperparameter است و برای اعمال آن لازم است برخی نکات در نظر گرفته شود که در ادامه آمده است:

- برای لایه ورودی نباید از Drop Out استفاده کرد یا حداقل ۹۰ درصد ورودی باید به شبکه اعمال گردد.
- ممکن است یک لایه دارای ابعاد بزرگی نسبت به سایر لایه‌ها باشد. در این حالت شبکه نسبت به لایه مذکور بیشترین احتمال Overfitting را دارد. با استفاده از Drop Out می‌توان یک سری نوروں‌های مربوط به آن لایه را به صورت رندوم حذف کرد. البته در لایه‌های دیگر نیز می‌توان از Drop Out استفاده کرد اما لایه‌ای که احتمال Overfitting در آن بیش‌تر است، نوروں‌های بیش‌تری در آن حذف می‌شود.
- در ابتدا لازم است بدون Drop Out تابع Loss محاسبه گردد و اگر تابع به صورت Monotonically Decreasing بود، Drop Out بر روی لایه‌های مدنظر اعمال گردد.
- برای داده‌های Test نباید از تکنیک Drop Out استفاده کرد.

• تکنیک Early Stopping

در نمودار Loss Function برحسب تعداد Epoch، اگر خطای مربوط به داده‌های Training رسم شود، مشاهده می‌شود که به صورت Monotonically Decreasing است. حال این سؤال مطرح می‌شود که تا کجا فرایند Learning شبکه و Update کردن پارامترها ادامه پیدا کند؟ و کدام Loss مناسب خواهد بود؟ در این قسمت از داده‌های Validation کمک گرفته می‌شود. بدین‌صورت که نمودار خطا مربوط به داده‌های

Validation دارای یک نقطه مینیمم است و پس از آن نقطه خطا داده‌های Validation افزایش یافته و شبکه قدرت تعمیم خود را از دست می‌دهد و به اصطلاح Overfit می‌شود. واضح است که در این نقطه تابع J یا همان خطای مربوط به داده‌های Training به‌طور کامل کاهش نیافته است و به آن Early Stopping می‌گویند. یکی از فواید داده‌های Validation آن است که در توقف فرایند یادگیری شبکه کمک می‌کنند.



شکل ۶ - نمودار خطا داده‌های Training و Validation بر حسب تعداد اپیاک

• تکنیک Norm Penalty

ممکن است یک شبکه عصبی در فرایند یادگیری از دقت مناسب برخوردار باشد اما برای داده‌های Validation یا Test درست عمل نکند و به اصطلاح Overfit شود. به عنوان مثال اگر یک شبکه دارای Dataset کوچک باشد، معمولاً با مشکل Overfitting همراه است. بدین صورت که شبکه همه ویژگی‌های داده‌های Training حتی اطلاعات مربوط به نویز و پترن‌های کم تکرار را حفظ می‌کند. در واقع شبکه داده‌های Training را یاد نگرفته بلکه در حافظه خود ذخیره کرده است. یکی از راه‌حل‌های افزایش قدرت تعمیم شبکه برای داده‌های جدید، استفاده از Regularization است. به گونه‌ای که در Loss Function، یک Regularization Term نیز با تابع جمع می‌شود. تابع خطا بدون در نظر گرفتن Regularization، تنها خروجی را به Target نزدیک می‌کند؛ اما با استفاده از Regularization Term، سعی می‌شود تا اندازه وزن‌ها نیز کوچک بماند و جواب‌هایی که در کرانه‌های فضا قرار دارند، در نظر گرفته نشوند.

رویکرد Regularization براساس محدود کردن گنجایش مدل‌های شبکه عصبی است که با افزودن Norm Penalty به تابع خطا ممکن می‌شوند. ترم (یک Hyperparameter است که میزان استفاده از

Regularization را نشان می‌دهد. بدیهی است که $\alpha=0$ هیچ‌ترم رگولاریزیشن به تابع خطا اضافه نمی‌کند. لازم به ذکر است بایاس‌ها در ترم رگولاریزیشن دخالت داده نمی‌شوند.

$$J(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

حال Norm Penalty می‌تواند انواع مختلفی داشته باشد. استفاده از L1 Norm مجموع وزن‌های شبکه را با تابع خطا جمع می‌کند. رایج‌ترین شیوه پیاده‌سازی Regularization استفاده از L2 Norm است که مجموع توان دوم وزن‌ها را با تابع خطا جمع می‌کند.

۲. طبقه‌بندی دیتاست Fashion MNIST

هدف از این بخش، طبقه‌بندی داده‌ها مربوط به Fashion MNIST با استفاده از Convolutional Neural Network و مقایسه شبکه طراحی‌شده با شبکه‌های کلاسیک است. دیتاست Fashion MNIST، دارای ۷۰۰۰۰ تصویر Grayscale در ابعاد ۲۸ (۲۸) است. همانند دیتاست MNIST، ۶۰۰۰۰ نمونه به‌عنوان داده‌های Training و ۱۰۰۰۰ نمونه باقی‌مانده به‌عنوان داده‌های Test به کار گرفته می‌شوند. خروجی‌های شبکه، ده دسته بوده که شامل ده نمونه لباس است. در ادامه ۲۵ نمونه از داده‌های Fashion MNIST به همراه برچسب مربوطه آمده است.



شکل ۷ - ۲۵ نمونه از داده‌های Fashion MNIST

۱,۲ طراحی یک شبکه CNN

شبکه عصبی مدنظر برای طبقه‌بندی داده‌های ورودی دیتاست با استفاده از Framework کراس طراحی شده است. دیتاست Fashion MNIST با استفاده از Keras در برنامه فراخوانی شده است. سپس داده‌های ورودی و خروجی نرمالایز شده‌اند؛ پیش‌تر گفته شد که ابعاد تصاویر (۲۸) ۲۸ پیکسل شامل مقادیر ۰ تا ۲۵۵ هستند. بنابراین به منظور مستقل از بعد شدن شبکه، پیکسل تصاویر Training و Test به عدد ۲۵۵ تقسیم شده‌اند. حال شبکه مدنظر به صورت زیر طراحی گردید.

- لایه Conv2D اول

داده‌های Training به عنوان ورودی شبکه به یک لایه Conv شامل ۶۴ فیلتر با ابعاد (۳) ۳ داده شده است. بدین صورت که یک Patch (۳) ۳ در تصاویر ورودی sweep شده و خانه‌های متناظر patch و فیلتر ضرب داخلی شده و سپس با یکدیگر و با بایاس جمع می‌شوند. در نهایت تابع فعال‌ساز ReLU بر آن‌ها پیاده شده و یک Feature Map ساخته می‌شود. به دلیل استفاده از ۶۴ فیلتر، Feature Map ساخته می‌شود. از آنجا که ابعاد ورودی و فیلتر به ترتیب (۲۸) ۲۸ و (۳) ۳ بوده، ابعاد Feature Map های ساخته شده، (۲۶) ۲۶ خواهد شد. بنابراین در این لایه از Padding استفاده نشده و Stride=1 در نظر گرفته شده است. لازم به ذکر است مقادیر فیلتر اعم از وزن‌ها و بایاس در هنگام استفاده از Optimizer در راستای مینیمم‌سازی تابع خطا مشخص می‌شوند. این پارامترها برای هر فیلتر شامل ۹ وزن و یک بایاس بوده که برای ۶۴ فیلتر، ۶۴۰ پارامتر می‌شوند.

- لایه Maxpooling اول

در لایه قبل ۶۴ Feature Map با ابعاد (۲۶) ۲۶ ساخته شد. حال با استفاده از یک لایه Maxpooling عمل Scaling بر روی ورودی‌های جدید پیاده می‌شود. اندازه Maxpooling ۲*۲ است که مقدار ماکزیمم ۴ خانه مجاور را انتخاب کرده و به Feature Map انتقال می‌دهد. عمل Maxpooling در اینجا منجر به نصف شدن ابعاد تصاویر می‌گردد. بنابراین Feature Map های جدید دارای ابعاد ۱۳*۱۳ خواهند بود. بدیهی است که Maxpooling تغییری در تعداد ورودی‌های به وجود نمی‌آورد و همچنین پارامتری برای Learn شدن ندارد.

- لایه Conv2D دوم

درواقع خروجی Maxpooling به صورت ورودی این لایه است؛ پس ورودی این لایه به صورت ۶۴ Feature Map با ابعاد ۱۳*۱۳ خواهد بود. در این قسمت نیز مجدد از ۶۴ فیلتر ۳*۳ استفاده می‌شود و تابع فعال‌ساز ReLU بر آن‌ها اعمال می‌گردد. بنابراین خروجی این لایه ۶۴ Feature Map با ابعاد ۱۱*۱۱ خواهد بود.

- لایه MaxPooling دوم

در این بخش نیز از یک لایه MaxPooling با ابعاد 2×2 استفاده شده است که ماکزیمم مقدار ۴ پیکسل کنار هم را به Feature Map انتقال می‌دهد. لازم به ذکر است از آنجایی که ابعاد ورودی این بخش 11×11 بوده و ابعاد Maxpooling 2×2 است، در هر ردیف خانه‌های آخر در Maxpooling شرکت داده نمی‌شوند. بنابراین خروجی به صورت ۶۴ Feature Map با ابعاد 5×5 خواهد بود. بدیهی است در این بخش همانند Maxpooling اول، پارامتری برای Learn شدن وجود ندارد.

- لایه Flatten

از لحاظ هندسی، شبکه جداپذیر شده است و می‌توان همانند شبکه‌های MLP با آن برخورد کرد. بنابراین با افزودن دو لایه Dense می‌توان طبقه‌بندی را تکمیل نمود. حال با استفاده از لایه Flatten، ویژگی‌های استخراج‌شده در لایه قبل به فرم یک بردار درمی‌آیند. خروجی لایه قبل، ۶۴ Feature Map با ابعاد 5×5 بود. بنابراین به صورت ۱۶۰۰ درایه برای یک بردار هستند. درواقع ویژگی‌های باقیمانده شامل Frequent Pattern ها بوده و nullity و disturtion از آن‌ها حذف گردیده است.

- لایه Dense اول

در این لایه ۱۶۰۰ ویژگی استخراج‌شده به ۱۲۸ نورون به صورت Fully Connected داده می‌شود. بدیهی است که ورودی در وزن مربوطه ضرب داخلی و با بایاس جمع می‌شود و سپس تابع ReLU بر آن‌ها پیاده می‌گردد. بنابراین تعداد پارامترها ۲۰۴۹۲۸ خواهد بود که شامل $128 + 204800$ است.

- لایه Dense دوم

در این لایه Feature های قبلی، ۲۰۴۹۲۸، با اعمال تابع فعال‌ساز Softmax به ۱۰ نورون داده می‌شود تا ۱۰ کلاس مدنظر خروجی شبکه گردند.

درنهایت شبکه طراحی شده به صورت زیر معرفی می‌گردد. به‌طور کلی الگوریتم Optimization باید ۲۴۳۷۸۶ پارامتر را برای شبکه آموزش دهد. Optimizer استفاده شده Adam است و از تابع خطا Sparse Categorical Cross-entropy استفاده شده است.

Model: "sequential_4"

Layer (type)	Output Shape	Param #
conv2d_7 (Conv2D)	(None, 26, 26, 64)	640
max_pooling2d_7 (MaxPooling2D)	(None, 13, 13, 64)	0
conv2d_8 (Conv2D)	(None, 11, 11, 64)	36928
max_pooling2d_8 (MaxPooling2D)	(None, 5, 5, 64)	0
flatten_4 (Flatten)	(None, 1600)	0
dense_7 (Dense)	(None, 128)	204928
dense_8 (Dense)	(None, 10)	1290
Total params: 243,786		
Trainable params: 243,786		
Non-trainable params: 0		

شکل ۸ - خلاصه‌ای از شبکه تعریف شده

در نهایت با استفاده از Batch-size=32 و تعداد ۵۰ اپیک، شبکه آموزش داده شد و به دقت ۹۹,۴۳ درصد برای داده‌های Training رسید. زمان آموزش شبکه در حدود ۳۶۸ ثانیه طول کشید. همچنین شبکه برای داده‌های Test به دقت ۹۰,۸۶ درصد رسید.

۲,۲ مقایسه شبکه طراحی شده با شبکه‌های عصبی ساده

شبکه Multi Layer Perceptron با تک‌لایه مخفی توانایی تقریب توابع convex را با دقت کافی دارد. درحالی که برای مسائل non-convex باید تعداد نورون‌های شبکه تک‌لایه را زیاد کرد. روش بهتر افزودن یک لایه مخفی دیگر است. حال شبکه MLP با دو لایه مخفی، می‌توان توابع non-convex را تقریب زد. اگر ورودی‌ها correlation، distortion و disturbance نداشته باشند، دو لایه مخفی کفایت می‌کند و افزودن لایه‌های دیگر redundancy ایجاد می‌کند.

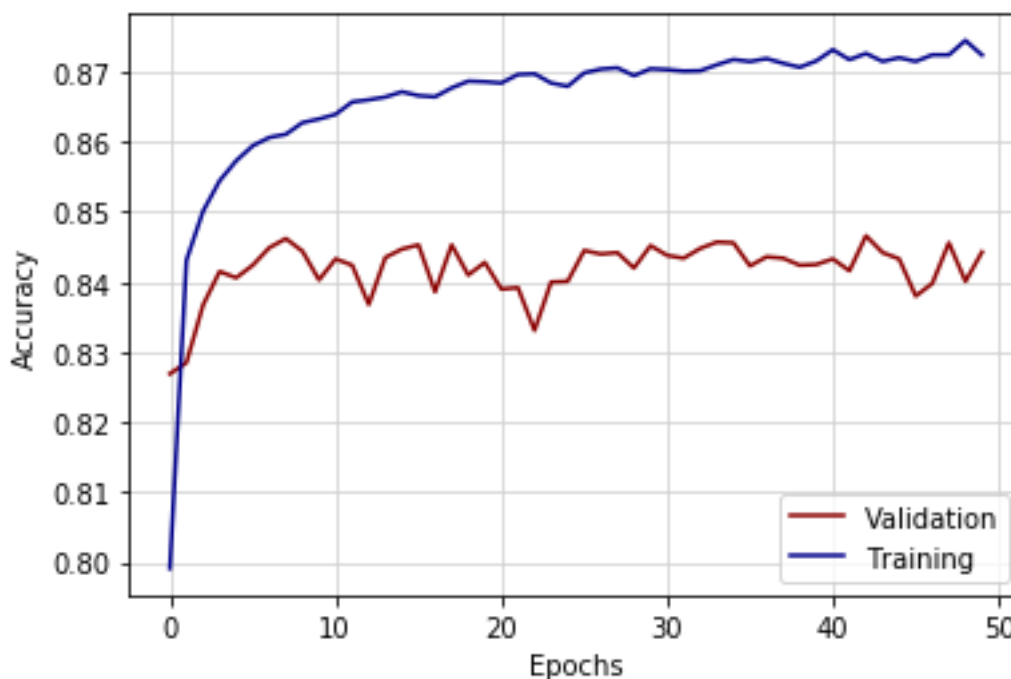
در این قسمت شبکه CNN طراحی شده با سه شبکه عصبی ساده با تعداد ۰، ۱ و ۲ لایه مخفی مقایسه گردید. در ابتدا معماری هر شبکه توضیح داده می‌شود و سپس چهار شبکه مذکور مقایسه می‌گردند.

منظور از صفر لایه مخفی، آن است که ورودی‌های شبکه به خروجی داده شوند. بنابراین از همان لایه خروجی شبکه CNN استفاده شده است با این تفاوت که ورودی‌ها نیز به آن داده شده است.

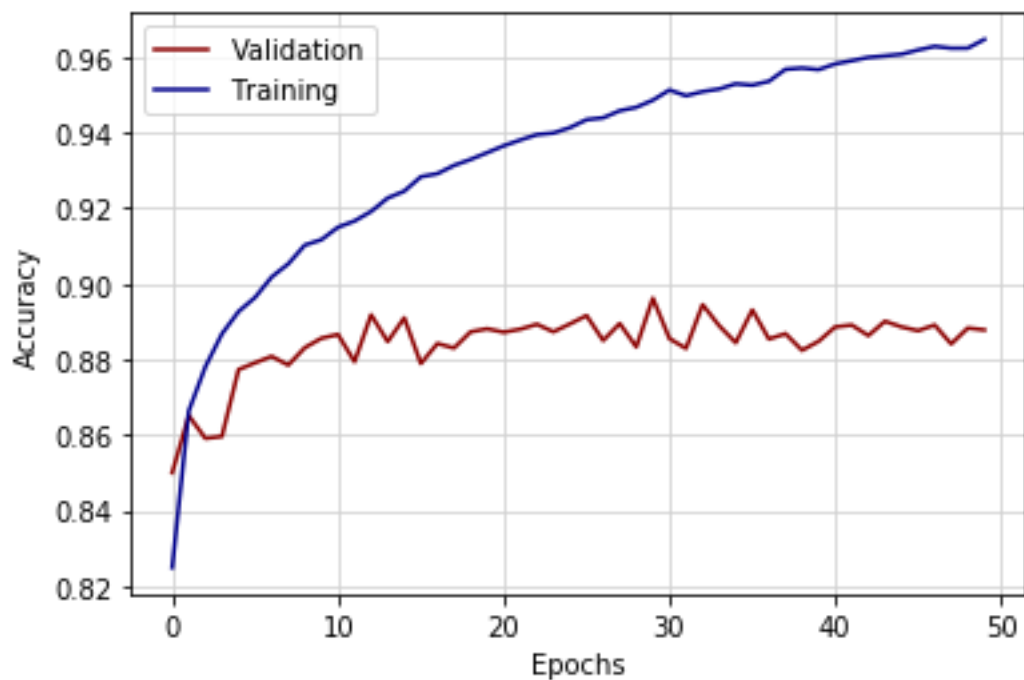
برای طراحی شبکه عصبی با یک لایه مخفی، ورودی‌ها به یک لایه با ۱۲۸ نورون داده شده و از تابع فعال‌ساز ReLU استفاده شده است. در نهایت خروجی لایه مخفی به یک لایه خروجی با ۱۰ نورون و تابع Softmax داده شده است.

در شبکه عصبی با دو لایه مخفی، دو لایه Fully Connected با تعداد نورون‌های ۲۵۶ و ۳۲ و با تابع ReLU طراحی شده است.

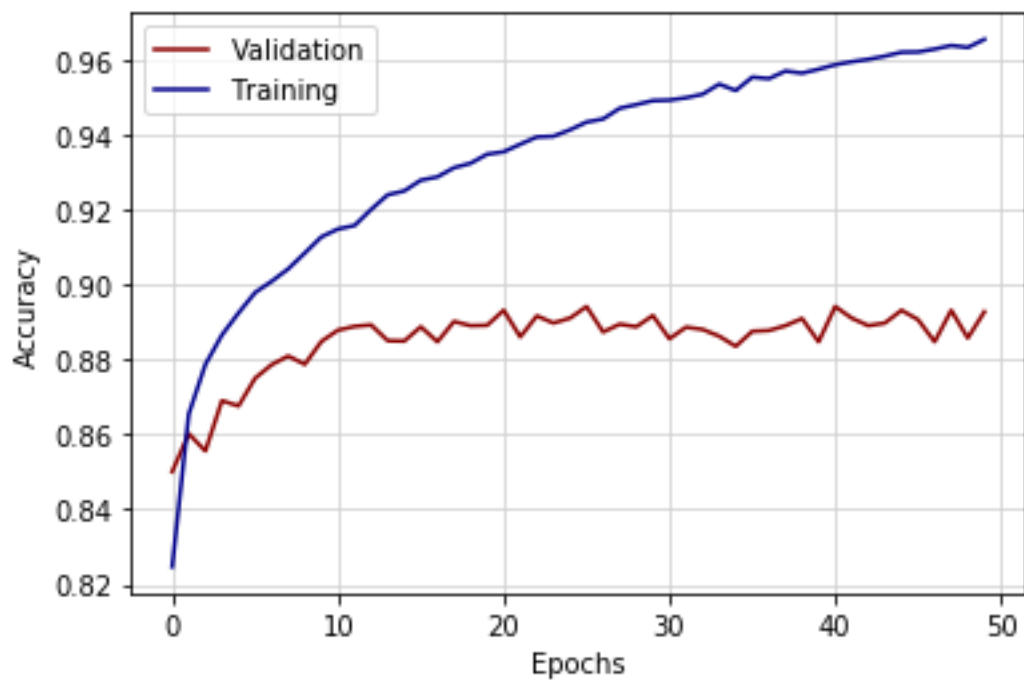
برای سه شبکه ساده، از اپتیمایزر Adam و تابع خطا Cross-entropy استفاده شده است. همچنین تعداد ایپاک ۵۰ و Batch-size=32 در نظر گرفته شده است. در نهایت نمودار دقت داده‌های Train و Validation برحسب تعداد ایپاک برای هر شبکه رسم گردید که در ادامه آمده است.



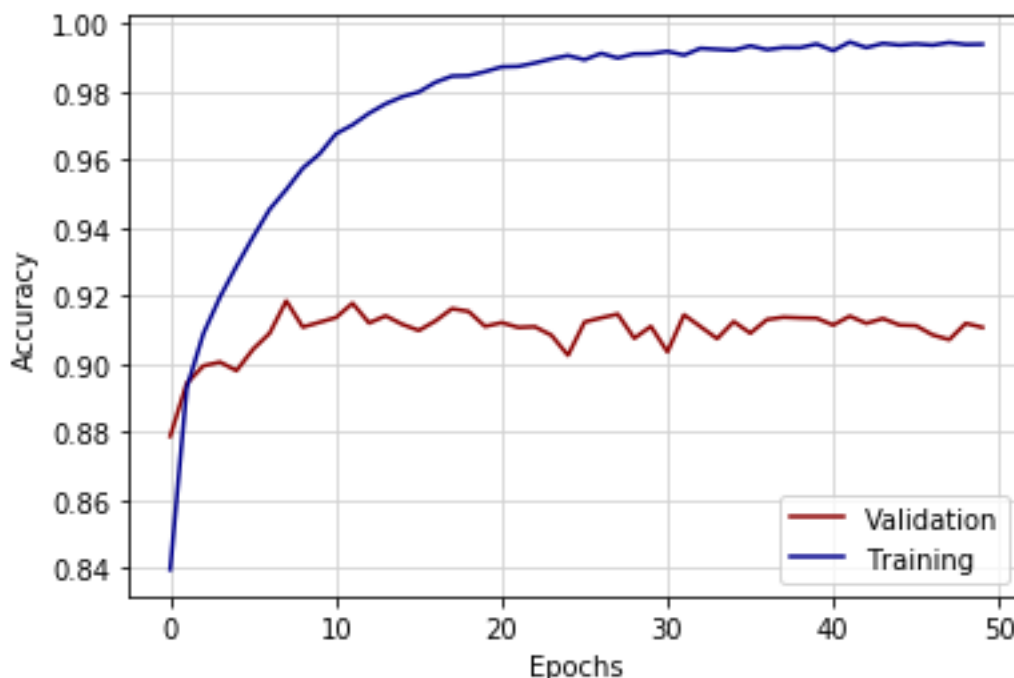
شکل ۹ - نمودار دقت برحسب ایپاک برای شبکه MLP بدون لایه مخفی



شکل ۱۰ - نمودار دقت برحسب اپیاک برای شبکه MLP با یک لایه مخفی



شکل ۱۱ - نمودار دقت برحسب اپیاک برای شبکه MLP با دو لایه مخفی



شکل ۱۲ - نمودار دقت برحسب اپاک برای شبکه CNN

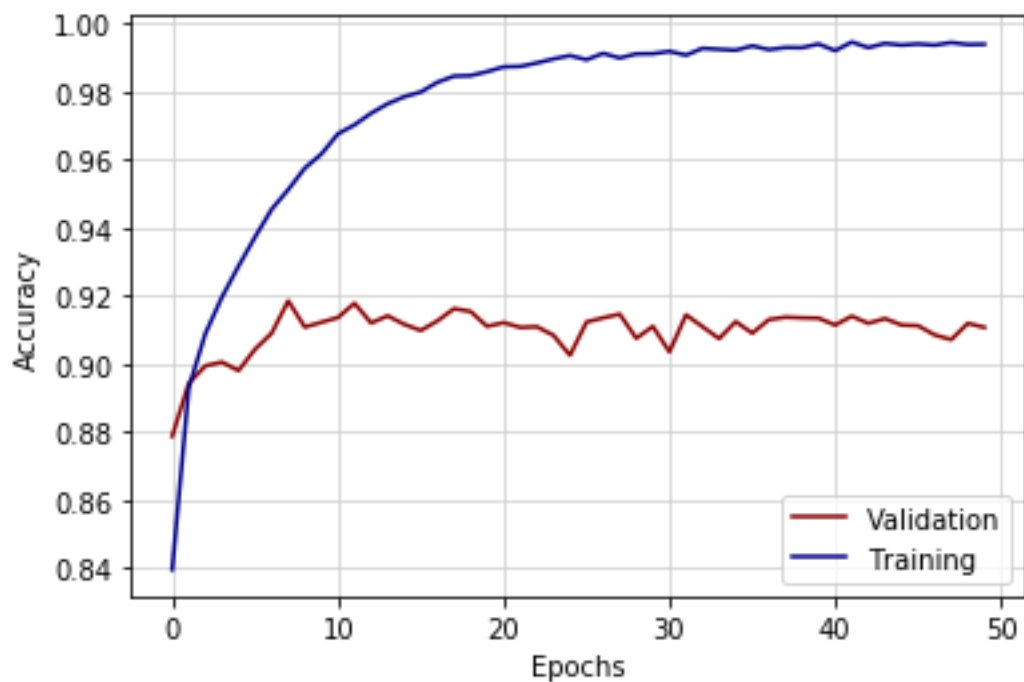
حال براساس نتایج نمودار دقت داده‌های Validation بر تعداد اپاک برای شبکه‌های مذکور، جدول زیر طراحی گردید. مشاهده می‌شود که دقت داده‌های validation در شبکه CNN بیش‌تر از بقیه است.

جدول ۱ - مقایسه شبکه‌های طراحی شده از منظر دقت داده‌های Validation

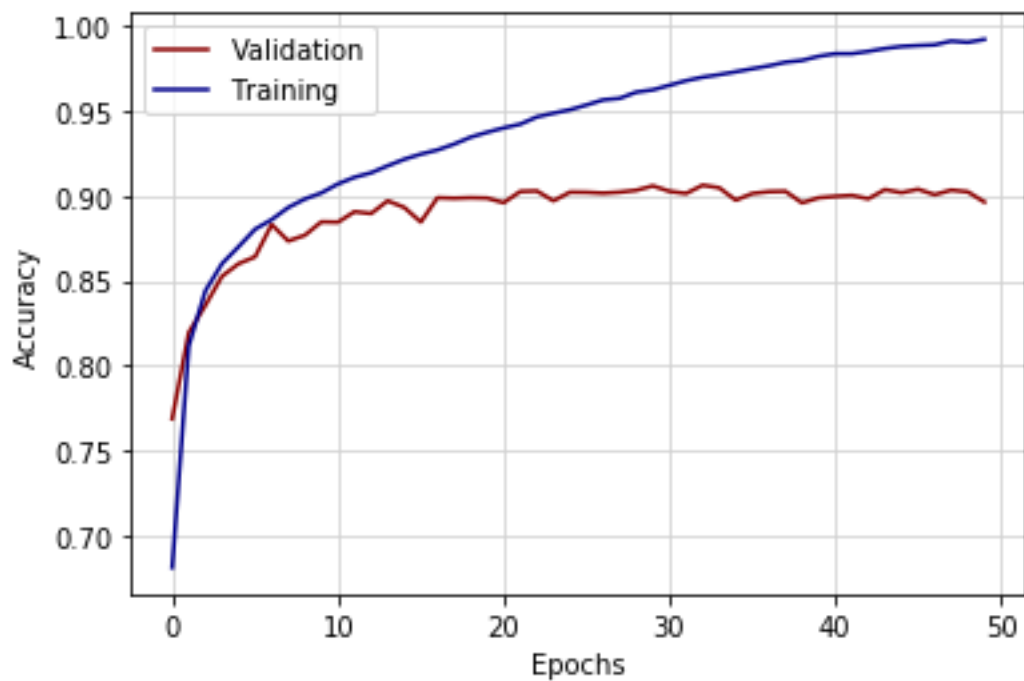
Network	Validation Accuracy	Run-time
MLP Without Hidden Layer	84.43%	374sec
MLP With One Hidden Layer	88.78%	376sec
MLP With Two Hidden Layer	89.27%	470sec
CNN	91.08%	410sec

۳,۲ تأثیر توابع فعال‌ساز بر شبکه طراحی‌شده

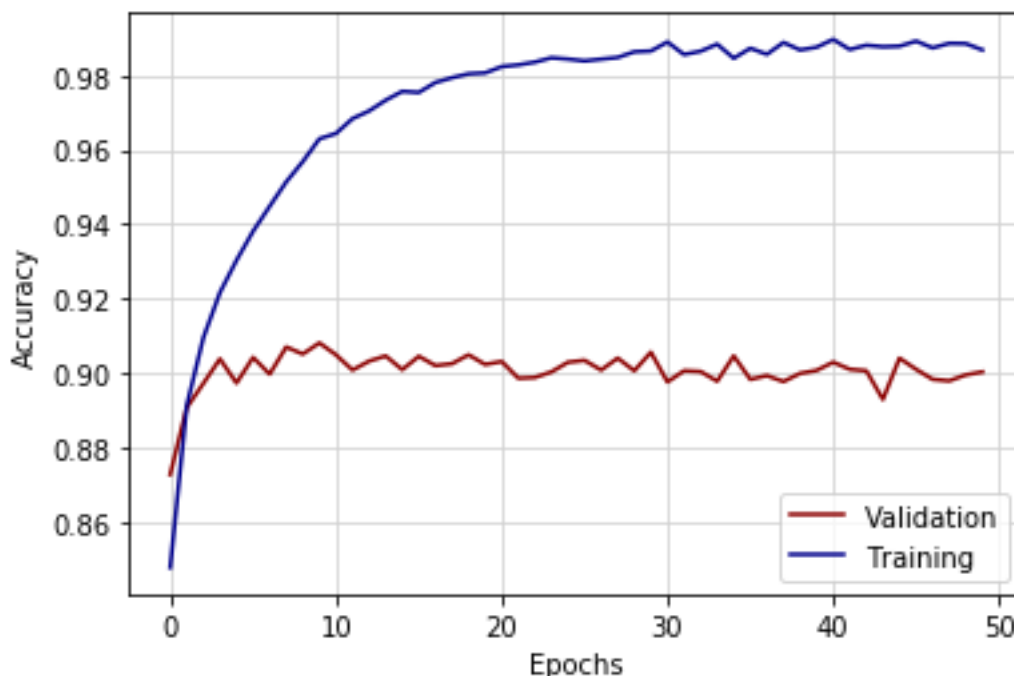
در این بخش، تأثیر توابع فعال‌ساز بررسی شده در بخش اول بر روی شبکه CNN بررسی می‌گردد. البته تابع فعال‌ساز در لایه خروجی در هر سه شبکه، با توابع ReLU, Sigmoid, Tanh و Softmax خواهد بود. پیش‌تر گفته شد که توابع ReLU و Tanh فقط برای لایه‌های میانی مناسب هستند و تابع Sigmoid اگر خروجی دو کلاس داشته باشد به کار می‌رود. پس اعمال توابع فعال‌ساز مختلف بر شبکه CNN نتایج زیر مشاهده گردید.



شکل ۱۳ - نمودار دقت برحسب اپاک برای شبکه CNN با تابع فعال‌ساز ReLU



شکل ۱۴ - نمودار دقت برحسب اپاک برای شبکه CNN با تابع فعال‌ساز Sigmoid



شکل ۱۵ - نمودار دقت برحسب اپیاک برای شبکه CNN با تابع فعال‌ساز Tanh

همان‌طور که در نمودارها مشخص است، دقت شبکه برای داده‌های Validation با استفاده از تابع فعال‌ساز ReLU بالاتر از دو حالت دیگر است. در قسمت قبل، مزایا و معایب توابع فعال‌ساز بیان شد. توابع Sigmoid و Tanh با مشکل Gradient Vanishing همراه هستند؛ درواقع با کوچک شدن گرادیان آن‌ها، به‌روزشدن پارامترها متوقف می‌گردد. حال که استفاده از تابع ReLU منجر می‌شود برخی از نورون‌ها در فرایند آموزش غیرفعال و شبکه سبک‌تر شود. براساس دقت داده‌های Validation در اپیاک پنجاهم یک جدول طراحی شده که در ادامه آمده است. همان‌طور که انتظار می‌رفت دقت شبکه با استفاده از تابع فعال‌ساز ReLU بیشتر از سایرین است. بین دو تابع Sigmoid و Tanh، تابع Tanh به دلیل داشتن گرادیان بزرگ‌تر دارای دقت بیشتری است.

جدول ۲ - مقایسه تأثیر تابع فعال‌ساز بر دقت داده‌های Validation

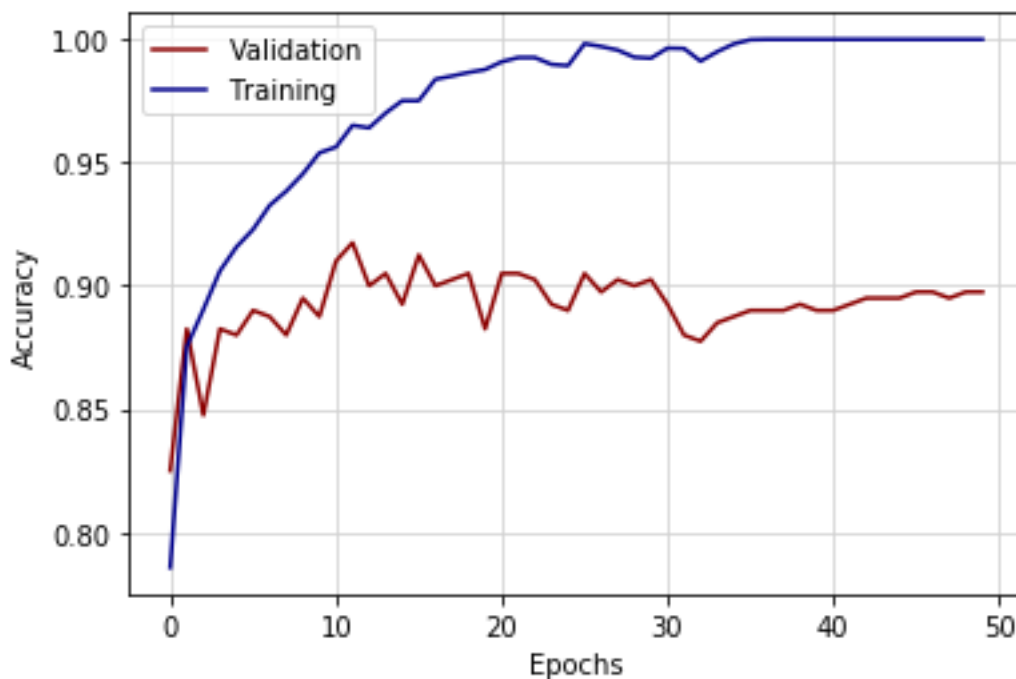
Activation Function	Validation Accuracy
ReLU	91.08%
Sigmoid	89.65%
Tanh	90.03%

۳,۳ بررسی تأثیر کاهش نمونه‌ها بر شبکه CNN طراحی شده

همان‌طور که گفته شد دیتاست Fashion MNIST شامل ۷۰۰۰۰ داده است که در ده کلاس طبقه‌بندی شده‌اند. در این حین ۶۰۰۰۰ داده مربوط به داده‌های Train و ۱۰۰۰۰ مربوط به داده‌های Test هستند. هدف از این بخش، کاهش نمونه‌های ورودی است به گونه‌ای که هر ۸۰۰ نمونه مربوط به یک کلاس باشند. بدیهی است دیتاست جدید شامل ۸۰۰۰ نمونه خواهد بود.

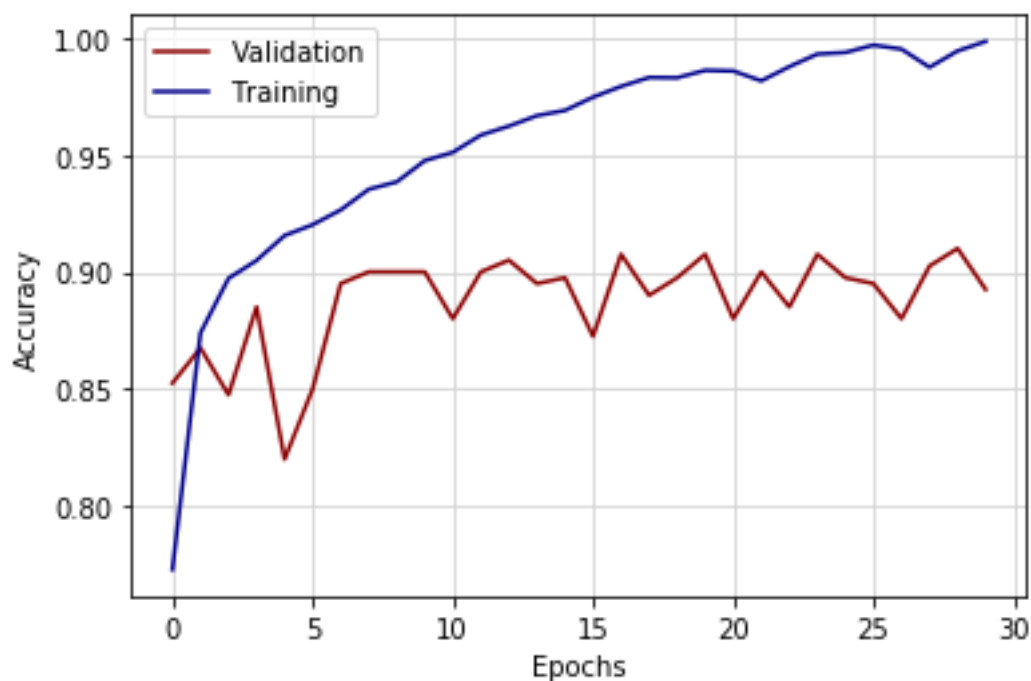
برای کاهش دیتاست شمارنده‌هایی برای هر کلاس در نظر گرفته شد تا تعداد داده‌ها دقیقاً ۸۰۰ باشد. بنابراین الگوریتم طراحی شده، هر نمونه را فراخوانی می‌کند. با استفاده از Label هر نمونه، با بررسی یک شرط آن را به کلاس متناظر اضافه می‌کند. شرط گفته‌شده بررسی می‌کند که اگر تعداد اعضای هر کلاس کمتر از ۸۰۰ بود، نمونه را به کلاس مربوطه اضافه می‌کند.

در نهایت دیتاست جدید به شبکه CNN طراحی شده اعمال گردید و با همان پارامترهای قبلی شبکه آموزش داده شد. لازم به ذکر است ۵ درصد از دیتاست جدید به عنوان داده‌های Validation در نظر گرفته شد. پس از آموزش شبکه در ۵۰ اپیک، دقت داده‌های Training ۱۰۰ درصد و دقت داده‌های Validation به ۸۹,۲۵ درصد رسید. همچنین نمودار دقت برحسب تعداد اپیک به صورت زیر مشاهده گردید.



شکل ۱۶ - نمودار دقت در ۵۰ اپیک برای دیتاست کاهش‌یافته

همان‌طور که در شکل مشخص است، دقت داده‌های Train از ایپاک ۳۰ به ۱۰۰ درصد میل می‌کنند و شبکه به اصطلاح Overfit شده است. بدین منظور بار دیگر شبکه با ایپاک ۳۰ آموزش داده‌شد. در این حالت دقت داده‌های Train، ۹۹٫۸۴ درصد و برای داده‌های Validation ۸۹٫۲۵ درصد شد. نمودار دقت برحسب تعداد ایپاک در ادامه آمده است.



شکل ۱۷ - نمودار دقت در ۳۰ ایپاک برای دیتاست کاهش‌یافته

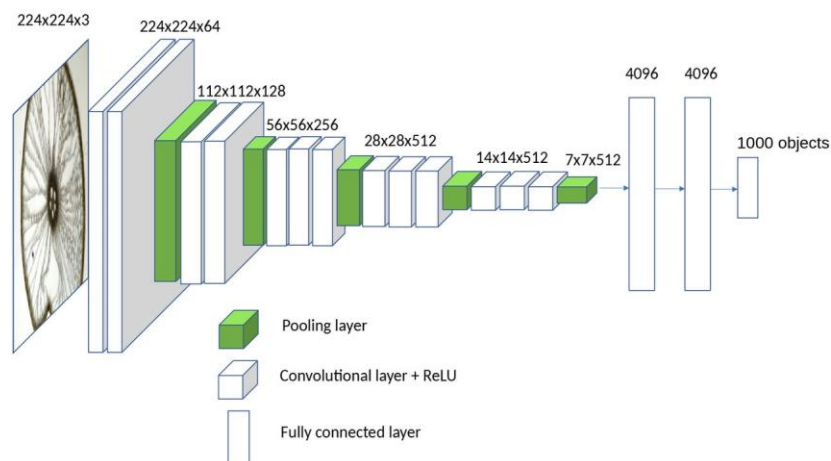
۳ بخش Transfer Learning

۱,۳ معماری VGGNet

در هنگام آموزش شبکه، ورودی لایه Conv اول، تصاویر سه کاناله RGB با سایز 224×224 هستند. تنها پیش پردازش انجام شده، تفریق میانگین هر پیکسل برای سه کانال، از مقادیر اصلی آنها است. سپس تصاویر به لایه‌های CONV اعمال می‌شود و با استفاده از فیلترها ویژگی‌های تصاویر استخراج می‌گردند. سایز استفاده شده برای فیلترها 3×3 است. البته در یکی از پیکربندی‌های از فیلترهای 1×1 نیز استفاده شده است. فیلترهای 1×1 به عنوان یک انتقال خطی از ورودی به Feature Map استفاده می‌شوند. درواقع هر فیلتر بر روی Patch متناظر از ورودی قرار گرفته و پس از محاسبه مقدار خانه مربوطه در Feature Map، بر روی ورودی Sweep می‌شود. در معماری VGG سوییچ به اندازه یک پیکسل است ($\text{stride}=1$). همچنین ابعاد تصاویر پس از اعمال فیلترها، حفظ می‌شود؛ بنابراین از Padding به اندازه ۱ پیکسل استفاده شده است. از لایه Max-pooling با ابعاد 2×2 و $\text{Stride}=2$ نیز استفاده شده است. پس از لایه‌های Conv و Pooling از سه لایه Fully-connected استفاده شده است. دو لایه Fully-connected اول ۴۰۹۶ نورون دارند. لایه آخر براساس استاندارد مسابقه ILSVRC که برای دسته‌بندی ۱۰۰۰ کلاس بوده، ۱۰۰۰ نورون داشته و تابع فعال ساز Softmax استفاده کرده است. سایر لایه‌ها، CONV، Max-pooling و Dense، از تابع فعال ساز ReLU بهره برده‌اند. تعداد کانال‌های استفاده شده در لایه‌های CONV کوچک است؛ به گونه‌ای که در اولین لایه از ۶۴ فیلتر استفاده شده و بعد از هر Max-pooling دوبار برابر گردیده تا به ۵۱۲ فیلتر رسیده است. پیکربندی‌های VGG عمق متفاوتی دارند؛ عمق پیکربندی‌ها بین ۱۱ تا ۱۹ لایه وزن دار است. به عنوان مثال در VGG16، ۱۳۴ میلیون پارامتر آموزش داده می‌شود.

در معماری VGG از الگوریتم بهینه‌سازی Mini-batch Gradient Descent با Momentum استفاده شده است؛ به گونه‌ای که مقادیر Batch-size و Momentum، به ترتیب ۲۵۶ و ۰,۹ انتخاب شده‌اند. برای رگولاریزیشن از Wight Decay ($L2$) استفاده شده است. همچنین دو لایه ابتدایی Fully-connected از Dropout با میزان ۵۰ درصد بهره برده‌اند.

در ابتدا از نرخ یادگیری ۰,۰۱ استفاده شده است؛ اما بدلیل موقوف شدن افزایش دقت داده‌های Validation، نرخ یادگیری ۰,۰۰۱ شده است. به طور کلی نرخ یادگیری سه بار کاهش یافته است. در نهایت فرایند یادگیری پس از ۳۷۰۰۰۰ تکرار (۷۴ اپیک) متوقف گردید.



شکل - معماری VGGNet

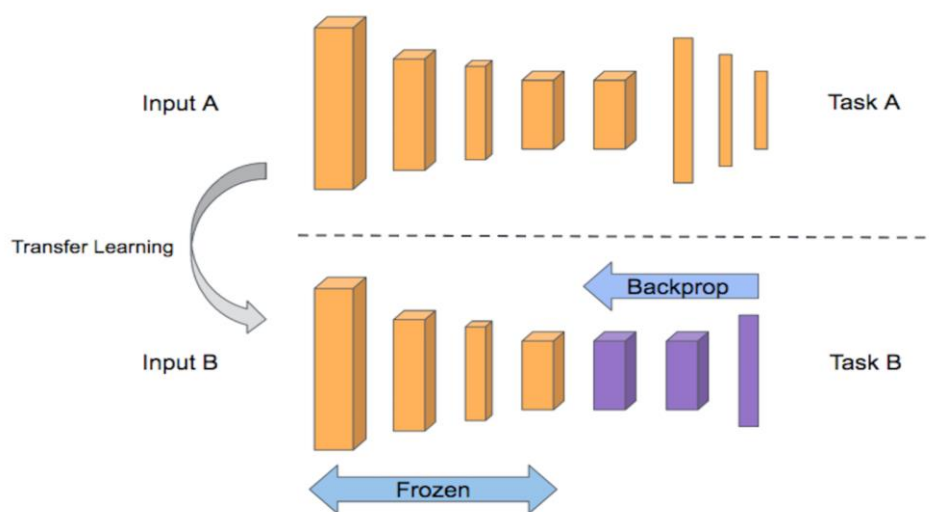
از مزایای VGG می‌توان به ساده و یک‌پارچگی آن اشاره کرد. درواقع در این نوع معماری از فیلترهای کوچک ولی با تعداد تکرار بیشتر استفاده شده است. معماری VGG از این نکته استفاده کرده است که دو فیلتر کانولوشنی 3×3 معادل یک فیلتر 5×5 و سه فیلتر 3×3 معادل یک فیلتر 7×7 است. بنابراین مزیت اصلی VGG، استفاده از دو 3×3 به جای یک 5×5 و سه فیلتر 3×3 به جای یک 7×7 است که افزایش غیرخطیگی را به دنبال دارد. زیرا هنگام استفاده از هر فیلتر یک تابع فعال‌ساز ReLU استفاده می‌شود. حال با افزایش تعداد فیلترها، از توابع ReLU بیش‌تری نیز استفاده می‌شود که منجر به افزایش Non-linearity شده و درواقع افزایش Non-linearity افزایش قدرت تعمیم شبکه و دقت داده‌های validation را به دنبال دارد.

معماری VGGNet با افزودن عمق به شبکه، توانست در سال ۲۰۱۴ با استفاده از بخشی از دیتا ImageNet (نه همه آن)، توانست به دقت بالاتری نسبت به رقبا خود دست پیدا کند. عمده کاربرد VGGNet در مسائل Classification است؛ خصوصاً هنگامی که دیتاست برای مسئله مدنظر کوچک است از شبکه VGGNet به‌عنوان Pre-train استفاده می‌شود. در بخش‌های آتی درمورد Transfer Learning و Pre-training صحبت می‌شود.

۲,۳ Transfer Learning

فرایند یادگیری شبکه‌های عصبی عمیق برای یک دیتاست بزرگ زمان‌بر بوده و ممکن است روزها یا هفته‌ها طول بکشد. یک راه میانبر استفاده مجدد از مدل‌های از پیش آموزش داده شده همانند مدل‌های تشخیص تصاویر ImageNet، برای دیتاست مدنظر است. درواقع Transfer Learning، فرایندی است که شبکه عصبی از یک

Task یاد گرفته و بر روی Task مشابه دیگر پیاده می‌شود؛ به گونه‌ای که با افزودن یک یا دو لایه، شبکه اولیه برای مسئله دوم قابل استفاده می‌گردد. از فواید Transfer Learning می‌توان به کاهش زمان آموزش اشاره کرد. همچنین از آنجایی که شبکه برای یک دیتاست اولیه آموزش داده شده، از قدرت تعمیم بالایی برخوردار است.



شکل ۱۸ - مفهوم Transfer Learning

یکی از دلایل استفاده از Transfer Learning آن است که دیتاست موجود کوچک است. بنابراین از سیستمی که با داده‌های زیاد Train شده استفاده و دیتاست کوچک را بر آن اعمال می‌شود. به عنوان مثال یک شبکه برای یک میلیون عکس گربه آموزش داده شده و برای دیتاست صدتایی رادیولوژی استفاده می‌شود. دلیل کارآمد بودن Transfer Learning آن است که بیشتر Low-level Feature ها همانند تشخیص لبه، منحنی، اشیا مثبت و ... در دیتاست اول، منجر به عملکرد بهتر در دیتاست دوم می‌شود. در واقع می‌توان از ویژگی‌های شبکه Pre-train، در شبکه مدنظر نیز بهره برد. البته نمی‌توان از شبکه‌ای استفاده کرد که داده‌های Pre-train کمتری نسبت به دیتاست مدنظر دارد.

از طرفی آموزش شبکه‌هایی همچون VGGNet، ResNet و ... با امکان دسترسی به GPU قدرتمند، در حدود چند هفته طول کشیده است. پس دلیل بعدی برای استفاده از Transfer Learning می‌تواند نبود امکانات و کم بودن زمان و هزینه باشد.

۳٫۳ پیاده‌سازی VGG16 با استفاده از Transfer Learning

در این بخش با استفاده از فریم‌ورک کراس، معماری VGG16 فراخوانی گردید؛ به گونه‌ای که برای قسمت‌های بعد از وزن‌های شبکه VGG استفاده می‌شود. در قسمت‌های بعدی با جزئیات بیش‌تر این قسمت توضیح داده می‌شود.

۴٫۳ لیست اشیا قابل‌شناسایی برای VGG16

برای دیدن لیست اشیایی که معماری VGG می‌تواند شناسایی کند، یا همان هزار دسته مربوط به ImageNet. از یک حقه استفاده شده است! به گونه‌ای که پس از فراخوانی VGG16 موجود در Keras، یک عکس گربه وارد گردید. به جای اختصاص دادن نزدیک‌ترین کلاس به تصویر، ۱۰۰۰ کلاس نزدیک مشاهده گردید. بدین‌صورت همه دسته‌های موجود در VGGNet مشاهده گردید. در ادامه اسامی این کلاس‌ها آمده است. همچنین یک فایل.txt. همراه کد این بخش آمده که اسامی اشیا قابل تشخیص معماری VGG در آن ذخیره گردیده‌است.

Egyptian_cat, tabby, tiger_cat, lynx, bow_tie, carton, radiator, vase, remote_control, ping-pong_ball, space_heater, window_screen, bath_towel, punching_bag, velvet, cup, mouse, pitcher, tub, feather_boa, lens_cap, sock, paper_towel, doormat, refrigerator, washer, washbasin, window_shade, lampshade, wall_clock, table_lamp, candle, crossword_puzzle, bucket, bathtub, goblet, quilt, bonnet, coffee_mug, plastic_bag, bookcase, cowboy_hat, iPod, water_bottle, toilet_tissue, fountain_pen, mailbox, Christmas_stocking, screen, shower_curtain, computer_keyboard, red_wine, pillow, swab, spotlight, jean, drumstick, fur_coat, wooden_spoon, Siamese_cat, plunger, cowboy_boot, rubber_eraser, hook, scale, hammer, paintbrush, television, pool_table, syringe, keyboard, hard_disc, mixing_bowl, monitor, beer_bottle, iron, printer, dishwasher, saltshaker, great_grey_owl, soap_dispenser, pajama, frying_pan, shower_cap, mailot, cellular_telephone, spatula, mitten, broom, purse, cloak, brassiere, cleaver, power_drill, hatchet, basenji, beer_glass, studio_couch, buckle, lighter, dumbbell, thimble, barometer, sombrero, hamper, letter_opener, screwdriver, switch, cauldron, laptop, Crock_Pot, rule, toilet_seat, miniskirt, perfume, wine_bottle, hair_spray, stole, water_jug, photocopier, microwave, wok, desktop_computer, bottlecap, wig, medicine_chest, radio, backpack, teapot, hummingbird, miniature_pinscher, projector, tray, loupe, jigsaw_puzzle, safety_pin, zucchini, horned_viper, meerkat, stethoscope, wool, cardigan, Doberman, analog_clock, modem, corkscrew, prairie_chicken, comic_book, kimonos, slide_rule, slug, alligator_lizard, totem_pole, mantis, Italian_greyhound, altar, cock, sewing_machine, lotion, marmoset, poncho, torch, sunglasses, toaster, centipede, loudspeaker, trench_coat, running_shoe, tricrator, dishrag, crab, Weimaraner, zebra, Chihuahua, piggy_bank, file, beaker, envelope, four-poster, violin, sunscreen, breastplate, coffeepot, kelpie, bubble, seat_belt, chime, miniature_schnauzer, maraca, grand_piano, motor_scooter, Eskimo_dog, pick, binder, microphone, orange, chain, pole, revolver, parallel_bars, bikini, cicada, barrel, fountain, nail, giant_schnauzer, ruffed_grouse, crayfish, shopping_basket, bolo_tie, silky_terrier, teddy, strainer, cocktail_shaker, balloon, gong, throne, vestment, reflex_camera, nipple, grasshopper, vizsla, mailbox, diaper, cassette, robin, stove, library, Angola, sundial, hay, French_bulldog, home_theater, butternut_squash, rain_barrel, pencil_box, dalmatian, scabbard, bearskin, barbell, tile_roof, digital_clock, quill, horizontal_bar, handkerchief, jaguar, hip, face_powder, car_mirror, German_short-haired_pointer, dial_telephone, joystick, rocking_chair, academic_gown, abaya, wombat, tripod, walking_stick, toy_terrier, kite, electric_guitar, baseball, Scotch_terrier, knee_pad, bull_mastiff, projectile, dining_table, Airedale, shield, pot, mongoose, gar, hamster, crutch, overskirt, birdhouse, green_lizard, snow_leopard, black-footed_ferret, Polaroid_camera, plate_rack, drum, lacewing, swimming_trunks, flute, burrito, cuirass, sunglass, shovel, espresso, acorn, frilled_lizard, bluetick, digital_watch, chain_mail, fire_screen, honeycomb, measuring_cup, traffic_light, desk, Siberian_husky, magnetic_compass, barrow, flagpole, Granny_Smith, ice_cream, Australian_terrier, French_horn, gamsack, banjo, barbershop, hoopskirt, soccer_ball, partridge, upright, grey_fox, junco, snail, electric_fan, wallaby, box_turtle, Norwegian_elkhound, long-horned_beetle, cricket, oscilloscope, coyote, terrapin, lycanoid, cradle, lab_coat, stopwatch, tobacco_shop, chocolate_sauce, bassinet, American_Staffordshire_terrier, titl, weasel, picklehaube, cocker_spaniel, dough, Great_Dane, cucumber, bittern, house_finch, ringlet, bassoon, chest, cinema, sandal, picket_fence, oboe, butcher_shop, basset, folding_chair, Brittany_spaniel, pretzel, waffle_iron, unicycle, whiptail, schipperke, church, swing, padlock, African_chameleon, admiral, sea_lion, redbone, combination_lock, Norfolk_terrier, pencil_sharpener, agama, neck_brace, safe, web_site, boxer, mink, African_grey, plate, hognose_snake, lemon, muffle, espresso_maker, chambered_nautilus, guacamole, restaurant, Bouvier_des_Flandres, clog, holster, barracouta, rugby_ball, mosque, Rotweiler, Mexican_hairless, crash_helmet, common_iguana, spider_web, maze, German_shepherd, obelisk, jersey, puffer, Cardigan, vine_snake, beagle, American_chameleon, hand-held_computer, fly, cornet, tricycle, milk_can, coho, hermit_crab, shoe_shop, Dungeness_crab, spiny_lobster, hen, barn_spider, tape_player, pomegranate, bathing_cap, night_snake, CD_player, wolf_spider, ocarina, grille, Kerry_blue_terrier, peacock, cello, reel, American_lobster, parking_meter, marimba, sax, crane, platypus, head_cabbage, moped, bell_pepper, Lakeland_terrier, carpenter's_kit, water_ouzel, whistle, racket, stretcher, black-and-white_coonhound, fig, toy_poodle, street_sign, car_wheel, starfish, chimpanzee, buckeye, ballplayer, paddle, entertainment_center, bloodhound, brass, pizza, Labrador_retriever, harp, carousel, nematode, hair_slide, loafer, prison, West_Highland_white_terrier, Border_terrier, macaque, polecat, guillotine, cassette_player, chiffonier, diamondback, rifle, electric_ray, missile, potter's_wheel, stupa, black_grouse, grocery_store, black_and_gold_garden_spider, cockroach, guenon, cab, jay, sulphur_butterfly, Rhodesian_rideback, trombone, lion, guinea_pig, vault, lorikeet, accordion, cauliflower, ibex, cheetah, sidewinder, jacamar, oil_filter, pug, bustard, basketball, Welsh_springer_spaniel, macaw, trilobite, cannon, otter, oxygen_mask, confectionery, bee, malinois, French_loaf, bulletproof_vest, albatross, china_cabinet, redshank, Gordon_setter, tick, Border_collie, Appenzeller, typewriter_keyboard, whippet, hotdog, banded_gecko, timber_wolf, Chesapeake_Bay_retriever, greenhouse, dome, quail, malamute, tree_frog, cabbage_butterfly, potpie, streetcar, chainlink_fence, Bernese_mountain_dog, menu, inststile, barber_chair, Madagascar_cat, bulbul, triumph_arch, space_bar, goldfish, king_penguin, whiskey_jug, vulture, mud_turtle, mashed_potato, vending_machine, magpie, skunk, jackfruit, Pembroke, indri, Old_English_sheepdog, carbonara, squirrel_monkey, red_fox, sea_urchin, boa_constrictor, dingo, bakery, cash_machine, kit_fox, spaghetti_squash, king_crab, lawn_mower, tiger_shark, shopping_cart, castle, English_springer, pirate, bagel, consomme, dragonfly, gazelle, palace, green_mamba, affenpinscher, hartbeest, binoculars, patio, red_wolf, rotisserie, scorpion, harvestman, koala, garden_spider, red-breasted_merganser, half-truck, leafhopper, red-backed_sandpiper, anemone_fish, brambling, Staffordshire_bullterrier, wire-haired_fox_terrier, football_helmet, orangutan, sea_slug, Irish_water_spaniel, thatch, schooner, rock_beauty, trifle, collie, acorn_squash, ladybug, Bedlington_terrier, siamang, space_shuttle, cardoon, container_ship, black_widow, yellow_lady's_slipper, flat-coated_retriever, badger, puck, toucan, English_setter, Pekinese, king_snake, worm_fence, drake, mosquito_net, goose, Greater_Swiss_Mountain_dog, flatworm, Shih-Tzu, Indian_cobra, hot_pot, stone_wall, mushroom, thunder_snake, pack_bench, golf_ball, cairn, liner, dock, golden_retriever, meat_loaf, Petri_dish, ski, soft-coated_wheaten_terrier, great_white_shark, Blenheim_spaniel, organ, pay_phone, isopod, hyena, gas_pump, sports_car, monastery, yawl, pier, rhinoceros_beetle, toyshop, Brabancon_ariffon, leatherback_turtle, Irish_terrier, hornbill, bullet_train, monarch, bald_eagle, sea_cucumber, Sussex_spaniel, assault_rifle, wing, African_crocodile, custard_apple, daisy, proboscis_monkey, cheseeburger, bell_cote, recreational_vehicle, tank, miniature_pool, greenhead, otterhound, hammerhead, limousine, papillon, Lhasa, bee_eater, Walker_hound, coral_fungus, rock_python, English_foxhound, parachute, artichoke, amphibian, trailer_truck, forklift, jirrikisha, viaduct, lionfish, water_snake, manhole_cover, impala, croquet_ball, black_swan, eel, garter_snake, Gila_monster, ringneck_snake, bullfrog, Scottish_deerhound, Dutch_oven, gorilla, curly-coated_retriever, Japanese_spaniel, jeep, Pomeranian, beaver, weevil, sturgeon, Saint_Bernard, Irish_setter, European_gallinule, seashore, earthstar, snorkel, freight_car, green_snake, groom, suspension_bridge, water_tower, pickup, Shetland_sheepdog, gyromitra, stinkhorn, tarantula, megalith, Great_Pyrnese, lesser_panda, ambulance, odometer, tailed_frog, lumbermill, ant, solar_dish, tractor, capuchin, minibus, canoe, American_alligator, Irish_wolfhound, common_newt, minivan, moving_van, brain_coral, Sealymen_terrier, gondola, American_black_bear, goldfinch, hippopotamus, school_bus, chain_saw, steel_drum, damselfly, coucal, snowplow, EntleBucher, tench, maypole, porcupine, axolotl, beach_wagon, American_egret, American_cot, Model_T, brown_bear, slot, volleyball, borzoi, mobile_home, Newfoundland, three-toed_sloth, Samoyed, clammer, warplane, oystercatcher, chow, garbage_truck, timararn, dowitcher, dhote, ox, racer, Maltese_dog, ground_beetle, ptarmigan, convertible, Tibetan_terrier, black_stork, wild_boar, beacon, sea_anemone, patas, geyser, bobbed_pane, langur, agaric, chiton, European_fire_salamander, rock_crab, bighorn, barn, ostrich, fire_engine, ruddy_turnstone, disk_brake, howler_monkey, crane, golfcart, cliff, eft, tow_truck, airship, catamaran, African_hunting_dog, standard_poodle, rapsessed_hog, sorrel, baboon, horse_cart, Saluki, pelican, tiger_beetle, radia_telecope, Limpkin, Arctic_fox, volcano, chickenade, firecat, marmot, kuvasz, hen-of-the-woods, flamingo, spoonbill, passenger_car, breakwater, police_van, komondor, aircraft_carrier, scoreboard, lifeboat, gibbon, white_wolf, planetarium, bolete, echidna, African_elephant, llama, Afghan_hound, keeshond, Indian_elephant, ram, loggerhead, coral_reef, lakeside, jellyfish, sloth_bear, grey_whale, killer_whale, indigo_bunting, airliner, Komodo_dragon, promontory, spotted_salamander, briard, stringray, alp, Leonberg, scuba_diver, harvester, little_blue_heron, sea_snake, submarine, oxcart, dogsled, leaf_beetle, Dandie_Dimont, spider_monkey, trolleybus, speedboat, paddlewheel, bicycle-built-for-two, dugong, drilling_platform, apiary, bathhouse, cliff_dwelling, mountain_bike, Tibetan_monastery, sandbar, electric_locomotive, sulphur-crested_cockatoo, valley, ice_bear, dam, colobus, steam_locomotive, plow, giant_panda, snowmobile, mountain_tent, bison, white_stork, warthog, steel_arch_bridge, water_buffalo, fiddler_crab, thresher, go-kart, Arabian_camel, furred_beetle, tusker, yurt,

شکل ۱۹ - لیست اشیا قابل‌شناسایی برای معماری VGG

۵,۳ فراخوانی یک تصویر و تشخیص اشیا

ابعاد تصویر گرفته شده ۴۰۳۲×۳۰۲۴ پیکسل در پیکسل است. در ابتدا ابعاد تصویر به ۲۴۴×۲۴۴ تبدیل می‌شود تا بتوان به‌عنوان ورودی به شبکه VGG16 اعمال کرد.



شکل ۲۰ - تصویر اعمال شده به شبکه VGG16

تصویر به فرمت PIL است و برای اعمال آن به شبکه عمیق مدنظر باید فرمت آن به Numpy Array تبدیل گردد. بدین منظور از تابع `img_to_array()` استفاده می‌شود. حال لازم است به شبکه تعداد نمونه‌ها اعمال گردد؛ از آنجایی که از یک تصویر استفاده شده است، فرم ورودی به شکل $(1, 224, 224, 3)$ درمی‌آید. سپس با استفاده از شبکه اپلیکیشن‌ها موجود در Keras، مدل VGG16 فراخوانی شده و خلاصه‌ای از آن مشاهده می‌گردد. همان‌طور که انتظار می‌رفت این شبکه در حدود ۱۴۰ میلیون پارامتر را به‌روزرسانی کرده‌است.

Model: "vgg16"

Layer (type)	Output Shape	Param #
input_9 (InputLayer)	(None, 224, 224, 3)	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0
flatten (Flatten)	(None, 25088)	0
fc1 (Dense)	(None, 4096)	102764544
fc2 (Dense)	(None, 4096)	16781312
predictions (Dense)	(None, 1000)	4097000

Total params: 138,357,544
 Trainable params: 138,357,544
 Non-trainable params: 0

شکل ۲۱ - خلاصه‌ای از مدل VGG16

درنهایت پیکسل‌ها با تابع preprocess_input پیش‌پردازش می‌شوند و به مدل VGG16 اعمال می‌گردند. همچنین با تابع decode_prediction سه نمونه از بیشترین احتمالات موجود در عکس، مشاهده می‌گردد که برای عکس مذکور، میز، ماوس و کوله‌پشتی تشخیص داده شده است. لازم به ذکر است اعداد جلو اسامی، درصد احتمال آن‌ها است.

```
[('desk', 73.25401902198792), ('mouse', 4.669767618179321), ('backpack', 2.305394969880581)]
```

شکل ۲۲ - تشخیص سه اشیا با بیش‌ترین احتمال

1. <https://www.coursera.org/specializations/deep-learning>
2. <https://www.coursera.org/specializations/tensorflow-in-practice?>
3. Optimizing Gradient Descent- <http://sebastianruder.com/optimizing-gradient-descent/>
4. Dean, J., Corrado, G. S., Monga, R., Chen, K., Devin, M., Le, Q. V, ... Ng, A. Y. (2012). Large Scale Distributed Deep Networks. NIPS 2012: Neural Information Processing Systems. <http://doi.org/10.1109/ICDAR.2011.95>
5. Qian, N. (1999). On the momentum term in gradient descent learning algorithms. Neural Networks: The Official Journal of the International Neural Network Society, 12(1), 145–151. [http://doi.org/10.1016/S0893-6080\(98\)00116-6](http://doi.org/10.1016/S0893-6080(98)00116-6)
6. Kingma, D. P., & Ba, J. L. (2015). Adam: a Method for Stochastic Optimization. International Conference on Learning Representations
7. Zaremba, W., & Sutskever, I. (2014). Learning to Execute, 1–25. Retrieved from <http://arxiv.org/abs/1410.4615>
8. Zhang, S., Choromanska, A., & LeCun, Y. (2015). Deep learning with Elastic Averaging SGD. Neural Information Processing Systems Conference (NIPS 2015).Retrieved from <http://arxiv.org/abs/1412.6651>
9. Darken, C., Chang, J., & Moody, J. (1992). Learning rate schedules for faster stochastic gradient search. Neural Networks for Signal Processing II Proceedings of the 1992 IEEE Workshop, (September). <http://doi.org/10.1109/NNSP.1992.253713>