



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

دانشکده برق

**موضوع: Mobile FaceNet**

درس: هوش محاسباتی

استاد: جناب آقای دکتر طالبی

نگارش: سجاد قدیری

محمد برآبادی

مارال مرداد

محیا حقگو

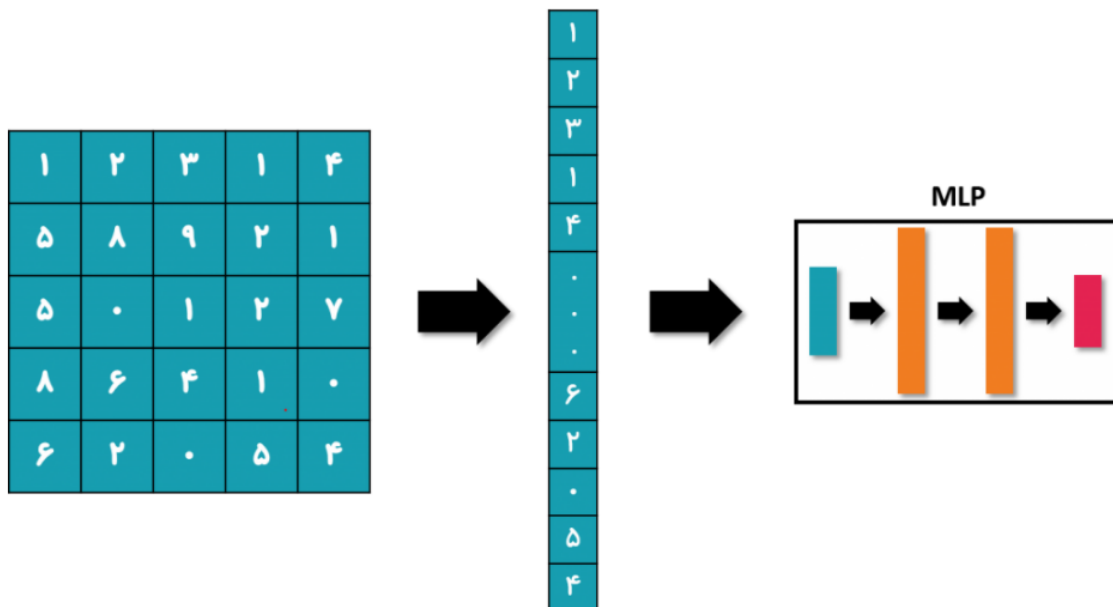
شبکه عصبی کانولوشن<sup>۱</sup>:

شبکه عصبی کانولوشن بسیار شبیه MLP است. این شبکه ها از نورون هایی تشکیل شده اند که دارای وزن های قابل یادگیری و بایاس هستند. به هر نورون یک سری ورودی داده می شود و با یک ضرب داخلی و عبور از تابع فعال ساز خروجی را می سازد.

وزن ها در هر سری یادگیری بر اساس مینیمم سازی یک تابع هزینه آپدیت می شوند. پس تا اینجا شبکه MLP و CNN بسیار شبیه به هم بودند. تفاوت آن ها در ورودی است.

اگر در شبکه MLP یک تصویر  $100 \times 100$  ورودی بدهیم در کل 10000 پیکسل دارد، شبکه برای دریافت ورودی این عکس را به یک آرایه یک بعدی 10000 تایی تبدیل می کند.

برای ساخت ورودی شبکه MLP نیاز است 10000 نورون در لایه ورودی قرار دهیم. که این تعداد نورون بسیار زیاد باعث زمان بر بودن فرآیند یادگیری و اتفاق افتادن بیش برآزش<sup>۲</sup> در زمان یادگیری می شود همچنین پارامترهای شبکه به شدت زیاد می شوند.



<sup>1</sup> CNN  
<sup>2</sup> overfitting

شبکه CNN برخلاف MLP ساختار ورودی را تغییر نمی‌دهد. ورودی شبکه CNN می‌تواند ماتریس های یک بعدی مانند سیگنال، دو بعدی مانند تصویر و طیف صوت، سه بعدی مانند ویدئو و تصاویر حجمی و داده‌های چهاربعدی مانند تصاویر حجمی همراه با زمان باشند. این شبکه به ارتباط بین پیکسل‌های همسایه اهمیت می‌دهد.

شبکه CNN از لایه‌های مختلفی تشکیل شده است که عبارتند از:

- لایه ورودی (Input layer)
- لایه کانولوشن (Convolutional layer)
- لایه غیرخطی (Non-linear activation function)
- لایه پولینگ (Pooling layer)
- لایه فولی کانکتد (Fully connected layer)

معروف ترین ورودی شبکه CNN تصویر می‌باشد.

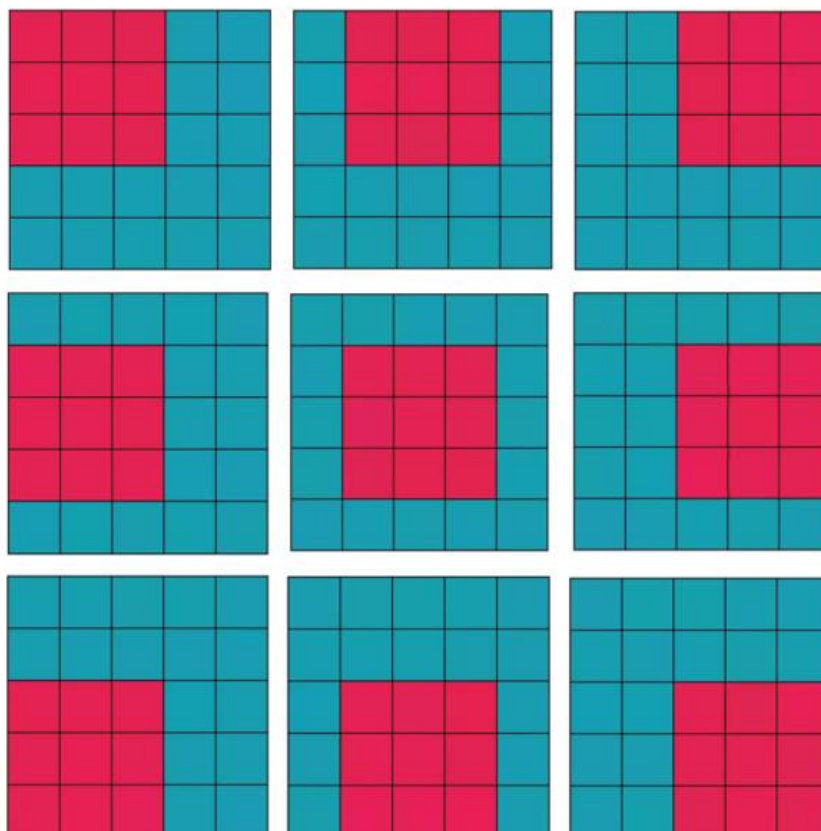
تصاویر ورودی می‌توانند خاکستری یا رنگی باشند.

تصویر رنگی از سه صفحه تشکیل شده است. این صفحات عبارتند از صفحه قرمز (R)، صفحه سبز (G) و صفحه آبی (B) به همین علت به این تصاویر، تصاویر RGB نیز گفته می‌شود.

هر یک از این صفحات یک ماتریس دوبعدی از اعداد بین 0 تا 255 می‌باشند. از ترکیب این سه صفحه یک تصویر رنگی حاصل می‌شود.

برای توضیح عملکرد کانولوشن می‌توان گفت، عملگر کانولوشن، کرنل یا فیلتر کانولوشنی را برمی‌دارد و روی تصویر یا ماتریس ورودی می‌لغزاند. فیلتر ابتدا هر سطر را ستون به ستون طی می‌کند و بعد یک سطر پایین می‌آید و دوباره ستون به ستون جلو می‌رود و این فرآیند تا آخر ادامه دارد.

برای مثال روند حرکت فیلتر بر تصویر ورودی به صورت زیر می‌باشد:



هر فیلتر کانولوشنی، شامل مجموعه‌ای عدد است. با قرار گرفتن فیلتر روی هر بخش از تصویر، اعداد فیلتر درایه به درایه در پیکسل تصویر متنظر ضرب می‌شوند و در نهایت همه اعداد با هم جمع می‌شوند.

بسته به اندازه فیلتر، اندازه ماتریس خروجی تغییر می‌کند. اگر ماتریس ورودی و فیلتر مربعی باشند و ابعاد آن‌ها به ترتیب  $n$  و  $k$  باشند، بعد ماتریس خروجی ( $m$ ) به صورت زیر می‌باشد:

$$m = n - (k - 1)$$

اگر بخواهیم بعد ماتریس خروجی برابر با بعد ماتریس ورودی باشد، می‌توان از روش لایه گذاری<sup>۳</sup> استفاده کرد. یک راه ساده و رایج لایه گذاری، اضافه کردن سطر و ستون صفر به صورت متقارن به دور ماتریس ورودی است. به این روش لایه گذاری صفر<sup>۴</sup> گفته می‌شود.

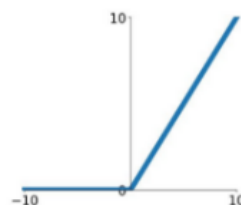
فیلتر به دنبال پیدا کردن نواحی مشابه خود در تصویر است. پس کانولوشن منجر به یافتن الگوهای خاص در تصویر با توجه به فیلتر می‌شود. اعداد موجود در فیلتر از طریق آموزش به دست می‌آیند.

هر صفحه از تصویر رنگی به طور مجزا برای خودش یک فیلتر دارد. تصویر رنگی سه صفحه دارد بنابراین فیلتر هم باید سه صفحه داشته باشد. سه صفحه از فیلتر به طور موازی با هم روی سه صفحه تصویر حرکت می‌کنند و در نهایت نتیجه ضرب سه صفحه با هم جمع می‌شوند و خروجی نهایی یک تصویر تک صفحه است. هر یک از فیلترها به تنهایی برای تشخیص یک الگوی خاص کاربرد دارد و برای تشخیص یک مجموعه الگو باید از چندین فیلتر استفاده کرد. تعداد فیلترها معمولاً نمایی از 2 (بین 32 تا 4096) می‌باشد. استفاده از فیلترهای بیشتر موجب قدرتمندتر شدن شبکه عصبی می‌شود.

بنابراین هر لایه کانولوشن در شبکه عصبی شامل مجموعه ای فیلترهاست و خروجی از کانولوشن فیلترها و ورودی به دست می‌آید. به خروجی لایه کانولوشنی Feature map گفته می‌شوند.

مشابه با سایر شبکه های عصبی، شبکه عصبی کانولوشن هم از تابع تحریک غیرخطی بعد از لایه کانولوشن استفاده می‌کند. از تابع ReLU به عنوان تابع غیرخطی استفاده می‌شود که به صورت زیر می‌باشد:

$$\text{ReLU} \\ \max(0, x)$$



هدف لایه پولینگ کاهش اندازه مکانی Feature map به دست آمده از لایه کانولوشنی است.

لایه پولینگ پارامتر قابل آموزش ندارد. عملکرد پولینگ مشابه عملکرد کانولوشن است. یک پنجره از پیش تعریف شده که روی تصویر حرکت می‌کند و در هر پنجره ماکزیمم را انتخاب می‌کند و بقیه را دور می‌ریزد.

این سه مرحله چندبار تکرار می‌شوند تا به ارور مطلوبی برسیم.

<sup>3</sup> padding  
<sup>4</sup> Zero padding

آخرین لایه های شبکه عصبی کانولوشن، لایه فولی کانکتد است که کاربرد آن طبقه بندی (classification) می باشد. مجموعه ویژگی های استخراج شده با استفاده از لایه های کانولوشنی یک بردار هستند که در نهایت این بردار ویژگی ها به فولی کانکتد داده می شود تا کلاس درست را شناسایی کند.

### معماری های متفاوت شبکه های CNN:

معماری های قدیمی تر به سادگی از لایه های کانولوشن انباشته تشکیل شده می شدند ولی معماری های مدرن راه های جدید و نوآورانه ای را برای ساخت لایه های کانولوشن به گونه ای که امکان یادگیری کارآمدتر شود ارائه داده اند. تقریباً همه این معماری ها بر اساس یک واحد تکرارپذیر هستند که در سراسر شبکه استفاده می شود.

روند اولیه برای توسعه مدل های CNN عمیق تر و پیچیده تر کردن شبکه برای رسیدن به دقت بالاتر بود. مانند VGGNet, GoogleNet, ResNet و DenseNet. اگرچه این روند بهبود مدل ها باعث افزایش سایز آن ها می شد. این افزایش سایز باعث افزایش محاسبات و حافظه مورد نیاز برای ذخیره سازی آنها می شود. برای حل این مشکل راندمان پایین ناشی از مدل های پیچیده محققین شبکه های CNN سبک وزن را پیشنهاد دادند. معماری های خاص و کارآمدی را برای ساخت CNN های سبک وزن پیشنهاد کردند.

SqueezeNet یک شبکه سبک وزن نسبتاً اولیه است که در ICLR2017<sup>5</sup> پیشنهاد شده است. که می تواند به دقت مشابه مدل AlexNet روی داده های ImageNet با کاهش 50 برابری پارامترها برسد. ماژول آتش استفاده شده در SqueezeNet عامل اصلی کاهش تعداد پارامترها است و از squeeze convolution module و ماژول گسترش تشکیل شده است. به فشردن تعداد کانال های ورودی کمک می کند. و فیلترهای  $3 \times 3$  ماژول گسترش را با فیلترهای  $1 \times 1$  جایگزین می کند. ShuffleNetV1 از عملیات کانولوشن گروهی نقطه ای و shuffle کردن کانال برای کاهش پارامترها و محاسبات و در عین حال حفظ دقت استفاده می کند. ShuffleNetV2 برای افزایش سرعت آموزش دیگر از تعداد زیادی کانولوشن گروهی استفاده نمی کند اما عملکرد شبکه را از طریق عملیات تقسیم کانال بهبود می بخشد. MobileNetV1 یک عملیات کانولوشن جدید را پیشنهاد کرد. اسم این عملیات کانولوشن، کانولوشن عمقی تفکیک پذیر است<sup>6</sup> که کانولوشن استاندارد را به دو بخش کانولوشن عمقی و کانولوشن نقطه ای<sup>7</sup> تقسیم می کند و محاسبات اضافی را تا حد زیادی کاهش

<sup>5</sup> International Conference on Learning Representations

<sup>6</sup> Depthwise separable convolutions

<sup>7</sup> pointwise

می‌دهد این عملیات کانولوشن در ادامه به صورت کامل توضیح داده شده است. MobileNetV2 برای بهبود مدل MobileNetV1 باقی‌مانده‌های معکوس<sup>۸</sup> و گلوگاه خطی<sup>۹</sup> معرفی کرد.

برای رسیدگی به مشکلات ناشی از مدل‌های پیچیده محققان چندین معماری شبکه خاص را برای تسک تشخیص چهره طراحی کردند.

MobileFaceNet یک مدل بر پایه MobileNetV2 است که برای تشخیص چهره استفاده می‌شود و به صورت تئوری به دقت 99.55٪ بر روی دادگان LFW<sup>۱۰</sup> می‌رسد.

## معماری Mobile facenet :

خلاصه‌ای از معماری Mobile facenet در جدول زیر آورده شده است:

Input	Operator	t	c	n	s
$112 \times 112 \times 3$	conv $3 \times 3$	-	64	1	2
$56 \times 56 \times 64$	depthwise conv $3 \times 3$	-	64	1	1
$56 \times 56 \times 64$	bottleneck	2	64	5	2
$28 \times 28 \times 64$	bottleneck	4	128	1	2
$14 \times 14 \times 128$	bottleneck	2	128	6	1
$14 \times 14 \times 128$	bottleneck	4	128	1	2
$7 \times 7 \times 128$	bottleneck	2	128	2	1
$7 \times 7 \times 128$	conv $1 \times 1$	-	512	1	1
$7 \times 7 \times 512$	linear GDConv $7 \times 7$	-	512	1	1
$1 \times 1 \times 512$	linear conv $1 \times 1$	-	128	1	1

که پارامترهای s , n , c , t به ترتیب ضریب گسترش، تعداد کانال خروجی، تعداد تکرار و طول گام می‌باشد. توضیح چند سطر از جدول در زیر آورده شده است:

سطر اول جدول کانولوشن استاندارد ورودی سه کاناله با ابعاد  $112 \times 112$  و کرنل سه کاناله با ابعاد  $3 \times 3$  می‌باشد. تعداد تکرار و گام و کرنل‌ها به ترتیب 1 و 2 و 64 می‌باشد. در نتیجه خروجی 64 کاناله است.

سطر دوم جدول کانولوشن عمقی ورودی 64 کاناله با ابعاد  $56 \times 56$  و کرنل 64 کاناله با ابعاد  $3 \times 3$  می‌باشد. تعداد تکرار و گام و کرنل‌ها به ترتیب 1 و 1 و 64 می‌باشد. در نتیجه خروجی 64 کاناله است.

<sup>8</sup> Inverted residuals  
<sup>9</sup> Linear bottlenecks  
<sup>10</sup> Labeled Faces in the Wild

سطر سوم از لایه Bottleneck استفاده شده است. این لایه یک کانولوشن  $1 \times 1$  می باشد که ابعاد خروجی را به منظور جلوگیری از Overfitting و vanishing gradient کاهش می دهد.

### مقایسه Mobile facenet و سایر معماری های CNN :

Network	LFW	AgeDB-30	Params	Speed
MobileNetV1	98.63%	88.95%	3.2M	60ms
ShuffleNet (1×, g = 3)	98.70%	89.27%	<b>0.83M</b>	27ms
MobileNetV2	98.58%	88.81%	2.1M	49ms
MobileNetV2-GDConv	98.88%	90.67%	2.1M	50ms
<b>MobileFaceNet</b>	<b>99.28%</b>	<b>93.05%</b>	<b>0.99M</b>	24ms
MobileFaceNet (112 × 96)	99.18%	92.96%	0.99M	21ms
MobileFaceNet (96 × 96)	99.08%	92.63%	0.99M	<b>18ms</b>
MobileFaceNet-M	99.18%	92.67%	0.92M	24ms
MobileFaceNet-S	99.00%	92.48%	<b>0.84M</b>	23ms
MobileFaceNet (ReLU)	99.15%	92.83%	0.98M	23ms
MobileFaceNet (expansion factor ×2)	99.10%	92.81%	1.1M	27ms

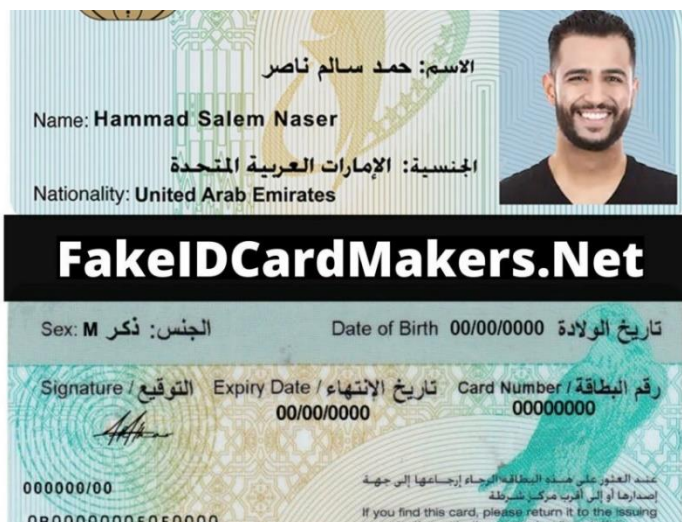
### کاربردها :

#### 1- Google Lens :

تشخیص مکان ها و بناهای تاریخی، متن ها، حیوانات و اشیاء مختلف ترجمه متن به زبان های مختلف از روی عکس گرفته شده توسط دوربین، ذخیره شماره تلفن و آدرس از روی عکس گرفته شده توسط دوربین از کاربردهای این نرم افزار می باشد.

برای نمونه با استفاده از این نرم افزار متن های موجود در این عکس از زبان عربی به انگلیسی ترجمه شده است.





## 2- TapTap See

به منظور تشخیص اشیاء و افراد موجود در عکس برای افراد کم بینا و نابینا می توان از این نرم افزار استفاده کرد. این نرم افزار متنی شامل اشیاء و افراد تشخیص داده شده در عکس را توسط صدا به فرد نابینا اعلام می کند.

برای نمونه با استفاده از این نرم افزار اشیاء موجود در عکس تشخیص داده شده است:

