

به نام خدا



دانشگاه تهران  
پردیس دانشکده‌های فنی  
دانشکده برق و کامپیوتر



یادگیری عمیق با کاربرد در  
بینایی ماشین و پردازش صوت  
پروژه‌ی امتیازی  
گروه ۴

نام و نام خانوادگی:

سارا جاهد آزاد ۸۱۰۶۹۹۱۴۹

محمد مهدی مهمانچی ۸۱۰۱۹۹۲۸۷

عرفان میرزایی ۸۱۰۱۹۹۲۸۹

مرداد ماه ۱۴۰۰

## فهرست

سوال ۱.....	۵
۱. مقدمه و هدف مقاله.....	۵
۲. پیاده سازی تولید تصاویر احساسات چهره.....	۵
۲,۲. معماری شبکه.....	۷
۲,۲,۱. تولیدکننده ی عکس.....	۸
۲,۲,۲. طبقه بند.....	۱۰
۲,۳. تابع هزینه.....	۱۲
2.4. آموزش شبکه.....	۱۲
۳. پیاده سازی تولید تصاویر اعداد دستنویس.....	۱۵
۳,۱. آماده کردن مجموعه دادگان:.....	۱۵
۳,۲. معماری شبکه.....	۱۶
۳,۳. آموزش شبکه.....	۱۹
۳,۴. نتایج.....	۲۰
۳,۵. نتیجه گیری:.....	۲۱
۴. افزودن دادگان.....	۲۲
۴,۱. حالت اول آموزش:.....	۲۲
۴,۳. نتایج.....	۲۳
۴,۴. نتیجه گیری.....	۲۴
پیوست ۱: روند اجرای برنامه.....	۲۵
مراجع.....	۲۶

## فهرست اشکال

- شکل ۱- تصویر پس از انجام برش و تشخیص چهره..... ۶
- شکل ۲- تصویر موجود در مجموعه دادگان بدون پیش پردازش..... ۶
- شکل ۳- معماری شبکه تولیدکننده..... ۸
- شکل ۴- معماری شبکه طبقه بند..... ۱۰
- شکل ۵- نمودار تابع هزینه در حین آموزش بر دادگان KDEF..... ۱۳
- شکل ۶- نمودار دقت طبقه بند احساسات در حین آموزش بر دادگان KDEF..... ۱۳
- شکل ۷- تصاویر تولید شده به همراه تصاویر اصلی از مجموعه دادگان آموزش KDEF..... ۱۴
- شکل ۸- تصاویر تولید شده (سطح-بالا) به همراه تصاویر اصلی از مجموعه دادگان آزمون KDEF..... ۱۴
- شکل ۹- تصاویر تولید شده (سطح-پایین) به همراه تصاویر اصلی از مجموعه دادگان آزمون KDEF..... ۱۴
- شکل ۱۰- نمودار هزینه در شبکه‌ی مربوط به اعداد دست‌نویس..... ۲۰
- شکل ۱۱- نمودار دقت در شبکه‌ی مربوط به اعداد دست‌نویس..... ۲۰
- شکل ۱۲- نمونه‌هایی از خروجی تولیدکننده شبکه در فاز آموزشی..... ۲۱
- شکل ۱۳- نمونه‌هایی از خروجی تولیدکننده شبکه در فاز آزمون..... ۲۱
- شکل ۱۴- تصاویر افزوده، دارای احساسات متعجب، منزجر، عادی و خوشحال..... ۲۳
- شکل ۱۵- نمودار هزینه طبقه بندی احساسات بدون افزودن داده..... ۲۳
- شکل ۱۶- نمودار هزینه طبقه بندی احساسات پس از افزودن داده..... ۲۳

## فهرست جداول

- جدول ۱- ساختار شبکه‌ی تولیدکننده..... ۹
- جدول ۲- ساختار شبکه‌ی طبقه‌بند..... ۱۰
- جدول ۳- ابرپارامترهای مورد استفاده در آموزش شبکه..... ۱۲
- جدول ۴- توضیح برچسب دادگان qmnist..... ۱۵
- جدول ۵- بخش تولیدکننده مربوط به شبکه‌ی تصاویر اعداد دست‌نویس..... ۱۷
- جدول ۶- بخش طبقه‌بند، مربوط به شبکه‌ی تصاویر اعداد دست‌نویس..... ۱۸
- جدول ۷- ابرپارامترهای استفاده‌شده در آموزش شبکه‌ی اعداد دست‌نویس..... ۱۹

## چکیده

در این پروژه، مقاله‌ی [1] پیاده‌سازی‌شد. شبکه‌ی معرفی‌شده در این مقاله می‌تواند تصاویر جدیدی را با ویژگی‌های<sup>۱</sup> دلخواه تولیدکند. این مقاله برای این امر، از یک ساختار دو مرحله‌ای بهره‌برده‌است. در قسمت اول ابتدا شبکه‌ی تولیدکننده<sup>۲</sup> تصاویری را تولید می‌نماید و سپس در ادامه یک شبکه‌ی طبقه‌بندی<sup>۳</sup> به طبقه‌بندی ویژگی‌های موجود در تصویر ایجاد شده می‌پردازد. ساختار مورد استفاده در این مقاله شبیه به شبکه‌های مولد تخصصی<sup>۴</sup> است با این تفاوت که در این مقاله، تولیدکننده و طبقه‌بند به همراه یکدیگر کار می‌کنند و نه برخلاف یکدیگر.

این شبکه برای مجموعه دادگان<sup>۵</sup> تصاویر چهره (KDEF) پیاده‌سازی‌شد. در مرحله‌ی اول تصاویری با احساسات متفاوت و هم چنین با چرخش‌های متفاوت در بخش شبکه‌ی تولیدکننده تولیدشدند، و در ادامه در قسمت شبکه‌ی طبقه‌بند، هویت فرد، احساس او و جهت چرخش تصویر او مورد طبقه‌بندی قرار گرفت. این شبکه برای مجموعه دادگان اعداد دست‌نویس (QMNIIST)، نیز پیاده‌سازی‌شد. در این حالت ویژگی‌های مورد طبقه‌بندی، فرد نویسنده‌ی آن عدد، خود عدد نوشته‌شده، رنگ عدد و جهت چرخش آن بودند.

در انتها از شبکه‌ی پیاده‌سازی‌شده برای مجموعه‌ی دادگان KDEF، برای تولید تصاویر جدید به منظور افزودن دادگان<sup>۶</sup> استفاده‌شد، و این داده‌ها برای آموزش شبکه‌ای با دقت بیشتر به کار گرفته‌شدند.

---

<sup>۱</sup>.Features

<sup>۲</sup>.Generator

<sup>۳</sup>.Classifier

<sup>۴</sup>.Generative Adversarial Network(GAN)

<sup>۵</sup>.Dataset

<sup>۶</sup> Data Augmentation

## ۱. مقدمه و هدف مقاله

در این پروژه به پیاده‌سازی مقاله‌ی Generative Cooperative Net for Image Generation and Data Augmentation پرداخته شده است. ساختن یک مدل خوب برای تولید تصاویر، یکی از مسئله‌های مطرح در بینایی ماشین است. به طور کلی در این مقاله، یک روش تولیدی جدید مطرح شده که می‌تواند مفاهیم سطح بالا را یاد بگیرد و تصاویر جدید با کیفیت بالا که ویژگی‌های<sup>۱</sup> دلخواه ما را دارند را تولید کند. آزمایش‌های مقاله روی مجموعه دادگان<sup>۲</sup> اعداد دستنویس (QMNIIST) و احساسات چهره (KDEF) انجام شده است که برای این کار، شبکه‌ای متشکل از دو بخش تولیدکننده<sup>۳</sup> و طبقه‌بند<sup>۴</sup> استفاده شده است. همچنین مقاله از شبکه‌ی ایجاد شده به منظور افزودن داده<sup>۵</sup> نیز بهره برده است. در این پروژه ما هر سه مرحله ذکر شده را پیاده‌سازی کردیم و به مقایسه و بررسی نتایج بدست آمده با نتایج ذکر شده در مقاله پرداخته‌ایم. در ادامه به توضیح جزئیات پیاده‌سازی هر یک از مراحل می‌پردازیم؛ بدین منظور ابتدا به شرح تولید تصاویر احساسات چهره (KDEF) می‌پردازیم، سپس به آزمایش روی مجموعه دادگان اعداد دستنویس (QMNIIST) پرداخته می‌شود و در انتها به طبقه‌بندی احساسات موجود در تصاویر صورت به کمک عکس‌های افزوده شده<sup>۶</sup> توسط شبکه مرحله قبل می‌پردازیم.

## ۲. پیاده‌سازی تولید تصاویر احساسات چهره

مقاله‌ی مربوط به این پروژه، به پیاده‌سازی کاربرد شبکه‌ی خود برای تولید تصاویر از چهره افراد در حالت‌های گوناگون شامل احساسات متفاوت و هم چنین چرخش‌های متفاوت تصویر پرداخته است. بدین منظور این مقاله، از مجموعه‌ی داده‌ی KDEF بهره‌برده است که در ادامه به آماده‌سازی داده‌ها برای به کارگیری در شبکه پرداخته می‌شود.

### ۲.۱. آماده‌کردن مجموعه داده:

در مقاله‌ی این پروژه از مجموعه داده KDEF استفاده شده است. این مجموعه شامل تصاویر چهره ۷۰ فرد (۳۵ زن و ۳۵ مرد) است که از پنج زاویه مختلف تصویر برداری شده است. در مقاله فقط از زاویه مستقیم (تمام

<sup>۱</sup>.Features

<sup>۲</sup>.Dataset

<sup>۳</sup>.Generator

<sup>۴</sup>.Classifier

<sup>۵</sup>.Data Augmentation

<sup>۶</sup>.Augmented Images

رخ) استفاده شده است. همچنین هر فرد ۷ حالت گوناگون چهره به خود گرفته است که در مقاله از چهار حالت عادی<sup>۱</sup>، خوشحال<sup>۲</sup>، منزجر<sup>۳</sup> و متعجب<sup>۴</sup> استفاده شده است. در ضمن از هر فرد در دو نوبت عکس گرفته شده است که ما فقط از عکس های نوبت اول استفاده کرده ایم. نام هر تصویر یک رشته به طول ۷ یا ۸ می باشد که شامل کلیدی اطلاعات لازم برای آن تصویر می باشد. برای نمونه نام یک تصویر به صورت AF06SUS.JPG می باشد که به کمک آن اطلاعات زیر درباره تصویر به دست می آید:

- حرف اول نشان دهنده نوبت تصویر برداری است: A یا B
- حرف دوم نشان دهنده جنسیت است: F یا M
- حرف های سوم و چهارم نشان دهنده شماره فرد است: از 00 تا ۳۵
- حرف های پنجم و ششم مشخص کننده حالت چهره است: AF, AN, DI, HA, NE, SA, SU
- حرف های هفتم و هشتم (در صورت وجود) زاویه تصویر برداری را نشان می دهد: FL, HA,S,HR,FR

در این بخش مطابق روشی که مقاله ذکر کرده است با استفاده از پکیج opencv و الگوریتم CascadeClassifier آن در هر تصویر چهره فرد را شناسایی می کنیم و عکس را برش<sup>۵</sup> می زنیم تا فقط شامل صورت باشد. در نهایت ابعاد این تصویر را به اندازه ۱۵۸ در ۱۵۸ در می آوریم. در شکل زیر یک نمونه تصاویر پیش از و پس از اعمال پیش پردازش آورده شده است.



شکل ۲- تصویر موجود در مجموعه دادگان بدون پیش پردازش



شکل ۱- تصویر پس از انجام برش و تشخیص چهره

<sup>۱</sup>.Neutral  
<sup>۲</sup>.Happy  
<sup>۳</sup>.Disgusted  
<sup>۴</sup>.Surprise  
<sup>۵</sup>.Crop

در ادامه هر یک از تصاویر بدست آمده را در سه زاویه ۹۰ و ۱۸۰ و ۲۷۰ می چرخانیم<sup>۱</sup>. در نهایت تصویر اصلی و سه چرخش آن را نسبت به محور تقارن چهره بازتاب<sup>۲</sup> می کنیم؛ بدین ترتیب از هر تصویر اصلی ۷ تصویر جدید به دست آورده ایم و مجموعه داده را غنی تر کرده ایم.

در مقاله دو دسته نتایج به عنوان خروجی آورده شده است که این دو دسته نیازمند دو نوع دادگان آزمون می باشند. به همین علت دادگان آزمون به دو صورت زیر ساخته شده اند:

چهار نفر را به صورت تصادفی انتخاب می کنیم، برای هر کدام یک حالت چهره و یک زاویه چرخش را باز هم به صورت تصادفی انتخاب می کنیم، تصویر مربوط به این فرد و حالت چهره و زاویه چرخش ( همراه با بازتاب) را در مجموعه تست قرار می دهیم. (مجموعاً ۸ تصویر)

دو نفر را به صورت تصادفی انتخاب می کنیم، برای هر کدام یک حالت چهره را باز هم به صورت تصادفی انتخاب می کنیم، تصویر مربوط به این فرد و حالت چهره را به ازای تمام زوایای چرخش ( همراه با بازتاب) در مجموعه تست قرار می دهیم. (مجموعاً ۱۶ تصویر)

بدین صورت مجموعه دادگان آزمون شامل ۲۴ تصویر ساخته می شود.

## ۲.۲ معماری شبکه

شبکه‌ی این مقاله از دو بخش عمده تشکیل شده است:

۱- تولیدکننده‌ی عکس

۲- طبقه‌بند

که در ادامه هر یک از این بخش ها به تفصیل توضیح داده می شوند.

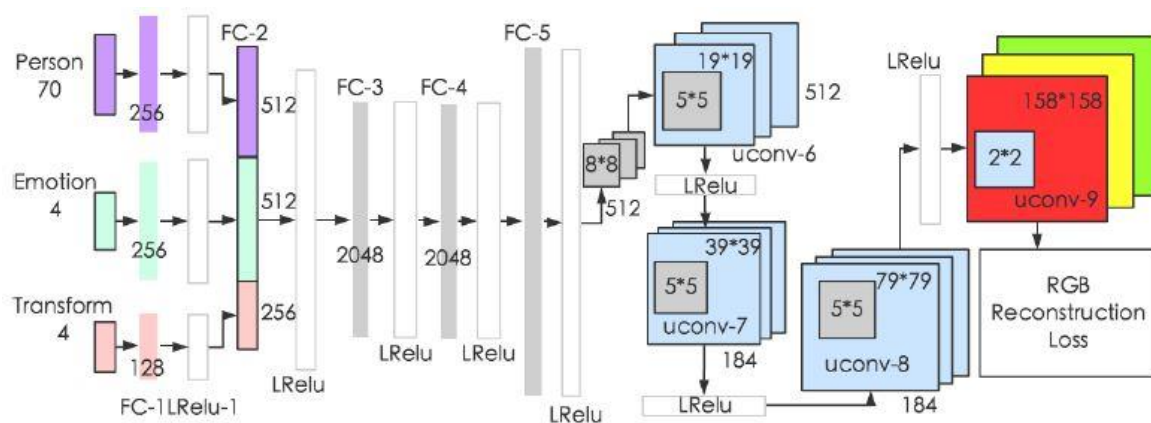
---

<sup>۱</sup>. Rotate

<sup>۲</sup>. Mirror

## ۲.۲.۱. تولیدکننده‌ی عکس

قسمت نخست شبکه برای تولید تصاویر به کار می‌رود. این بخش سه نوع ورودی دارد؛ بردار مربوط به هویت شخص، بردار مربوط به احساس شخص و بردار مربوط به چرخش تصویر فرد. این سه بردار به صورت وان‌هات<sup>۱</sup> از سه ورودی وارد شبکه می‌گردند، از لایه‌های تمام-متصل<sup>۲</sup> مربوط به خود عبور می‌کنند و با هم پیوست<sup>۳</sup> می‌شوند، سپس پس از عبور از چند لایه‌ی تمام-متصل دیگر و پس از آن چند لایه‌ی پیچشی<sup>۴</sup> دیگر، در نهایت یک خروجی به ابعاد  $158 \times 158$  با سه کانال به دست خواهد آمد که با تصویر اصلی مقایسه می‌گردد (یکی از بخش‌های تابع هزینه در این قسمت محاسبه می‌شود؛ یعنی مجموع خطای مربعات خروجی تولیدکننده با عکس‌های اصلی آموزش). معماری بخش تولیدکننده طبق معماری پیشنهادی در مقاله‌ی [4] بوده است. معماری دقیق این بخش در ادامه مشاهده می‌گردد.



شکل ۳- معماری شبکه تولیدکننده

<sup>۱</sup> One Hot

<sup>۲</sup> Fully-Connected

<sup>۳</sup> Concatenate

<sup>۴</sup> Convolutional Layers

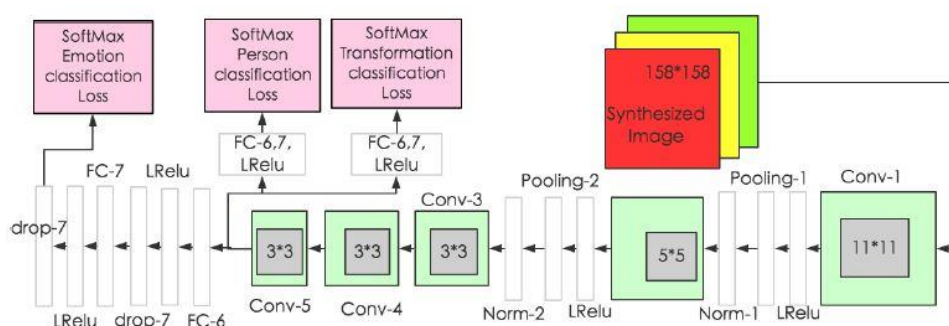


جدول ۱- ساختار شبکه‌ی تولیدکننده

Layer name	Input Feature Size	Kernel Size	Padding	Stride	Output Feature Size
Person FC 1	70	–	–	–	256
Leaky ReLU (negative slope=0.1)					
Person FC 2	256	–	–	–	512
Emotion FC 1	4	–	–	–	256
Leaky ReLU (negative slope=0.1)					
Emotion FC 2	256	–	–	–	512
Transform FC 1	4	–	–	–	128
Leaky ReLU (negative slope=0.1)					
Transform FC 2	128	–	–	–	256
Then		these three outputs		being concatenated	
Leaky ReLU (negative slope=0.1)					
FC 3	1280	–	–	–	2048
Leaky ReLU (negative slope=0.1)					
FC 4	2048	–	–	–	2048
Leaky ReLU (negative slope=0.1)					
FC 5	2048	–	–	–	16384
Leaky ReLU (negative slope=0.1)					
Reshape to 256*8*8					
deConv 6	256	5	0	2	512
Leaky ReLU (negative slope=0.1)					
deConv 7	512	5	1	2	184
Leaky ReLU (negative slope=0.1)					
deConv 8	184	5	1	2	184
Leaky ReLU (negative slope=0.1)					
deConv 9	184	2	0	2	3

## ۲,۲,۲. طبقه‌بند

خروجی تولیدکننده‌ی عکس به عنوان ورودی وارد طبقه‌بند می‌شود. قسمت پیش‌بینی<sup>۱</sup> این شبکه از معماری Alexnet[5] الهام گرفته‌شده است؛ با این تفاوت که گام<sup>۲</sup> در لایه‌ی اول به ۲ کاهش یافته (زیرا دادگان استفاده شده در این مقاله ابعاد کوچک‌تری نسبت به مجموعه دادگانی دارد که Alexnet[5] در اصل با آن آموزش پیدا کرده است). همچنین لایه‌های local response normalization به قسمت پیش‌بینی این مقاله اضافه شده‌اند که در مقاله‌ی اصلی Alexnet نیستند. طی بررسی‌های انجام‌شده، جایگزینی لایه‌های batch normalization با LRN پاسخ بهتری به دست می‌دهد. سپس سه قسمت متشکل از لایه‌های تمام-متصل، هر یک برای طبقه‌بندی اشخاص، احساسات و چرخش به انتهای بخش پیش‌بینی شبکه اضافه می‌گردند. در ادامه، معماری بخش طبقه‌بند با جزئیات بیشتر آورده شده است.



شکل ۴ - معماری شبکه طبقه‌بند

جدول ۲- ساختار شبکه‌ی طبقه‌بند

Layer	Input Feature Size	Kernel Size	Padding	Stride	Output Feature Size
Conv 1	3	11	2	2	64
Leaky ReLU (negative slope=0.2)					
Batch Normalization					
Maxpool layer 1	64	3	2	0	64
Conv 2	64	5	2	1	192
Leaky ReLU (negative slope=0.2)					
Batch Normalization					

<sup>۱</sup>. Convolutional

<sup>۲</sup>. Stride

Maxpool layer 2	192	3	2	0	192
Conv 3	192	3	1	1	384
Leaky ReLU (negative slope=0.2)					
Conv 4	384	3	1	1	256
Leaky ReLU (negative slope=0.2)					
Conv 5	256	3	1	1	256
Leaky ReLU (negative slope=0.2)					
Maxpool layer 5	256	3	2	0	256
Flattening Layer					
FC 6	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 7	4096	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 8	4096	-	-	-	70
FC 6	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 7	4096	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 8	4096	-	-	-	4
FC 6	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 7	4096	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 8	4096	-	-	-	4

### ۲.۳. تابع هزینه<sup>۱</sup>

تابع هزینه برای این شبکه به صورت زیر تعریف می‌گردد:

$$Loss: Euclidian \text{ pixel to pixel loss} + 10 * Crossentropy \text{ Loss}_{person} + 10 * Crossentropy \text{ Loss}_{emotion} + 10 * Crossentropy \text{ Loss}_{transformation}$$

تابع هزینه اقلیدسی، در واقع تابع هزینه‌ای است که بین تصاویر تولید شده در بخش تولیدکننده و تصاویر واقعی محاسبه می‌شود.

نمایش کلی تابع هزینه به صورت زیر خواهد بود:

$$\frac{1}{N} \left\{ \sum_{i=1}^N \left\| Pixel_r^i - Pixel_s(f_{i_1}, \dots, f_{i_M})^i \right\|^2 - \left[ \sum_{i=1}^N \sum_{f=1}^M \sum_{j=1}^{K_f} Weight_f \{y_f^i = j\} \log \frac{e^{\theta_{f_j}^T x_f^i}}{\sum_{l=1}^{K_f} e^{\theta_{f_l}^T x_f^i}} \right] \right\}$$

در عبارت بالا N تعداد تصاویر، M تعداد ویژگی‌های هر عکس و  $K_f$  بُعد کلاس برای یک ویژگی خاص f است. در فرمول بالا weight برابر با 10 قرار داده می‌شود.

### ۲.۴. آموزش شبکه

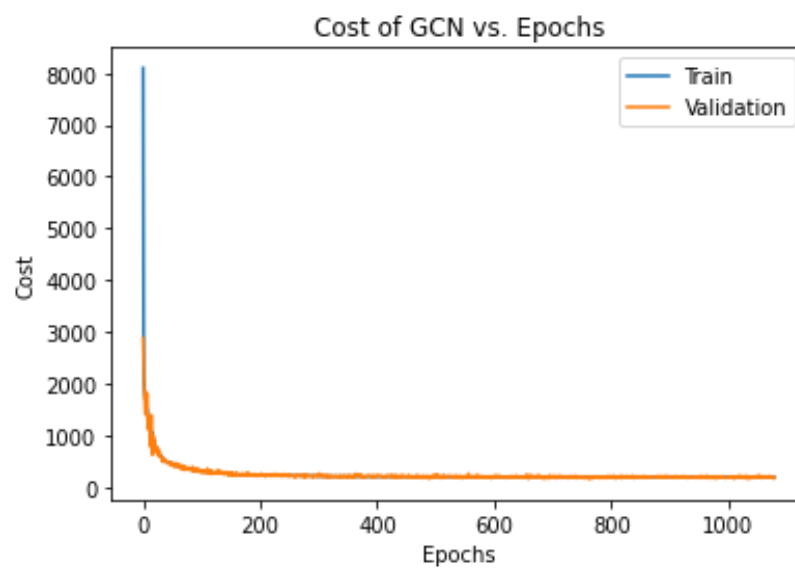
برای آموزش شبکه، از مجموعه ابرپارامترهای<sup>۲</sup> زیر استفاده کرده‌ایم:

جدول ۳- ابرپارامترهای مورد استفاده در آموزش شبکه

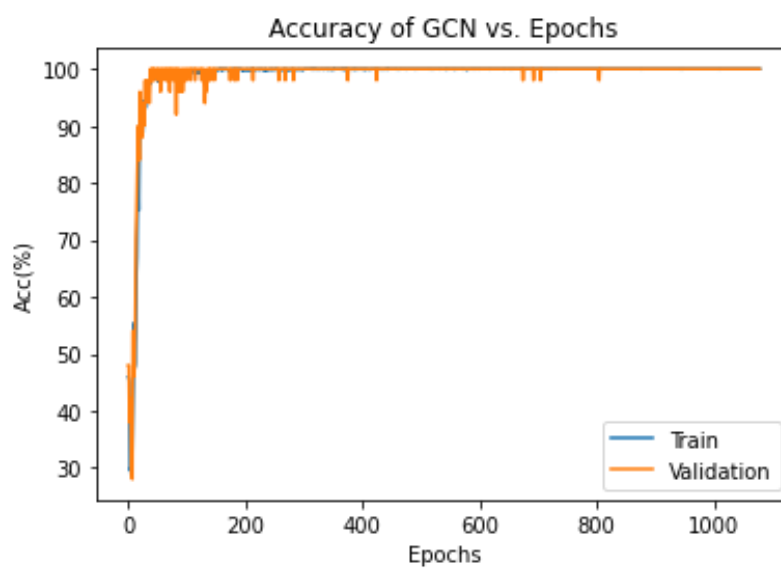
Optimizer	Adam
Learning rate	0.0002
Learning rate decay:	after every 150 epochs it is divided by 2
$\beta_1$	0.9
$\beta_2$	0.99
$\varepsilon$	$10^{-8}$
batch size	64
epochs	1100

<sup>۱</sup>. Loss Function

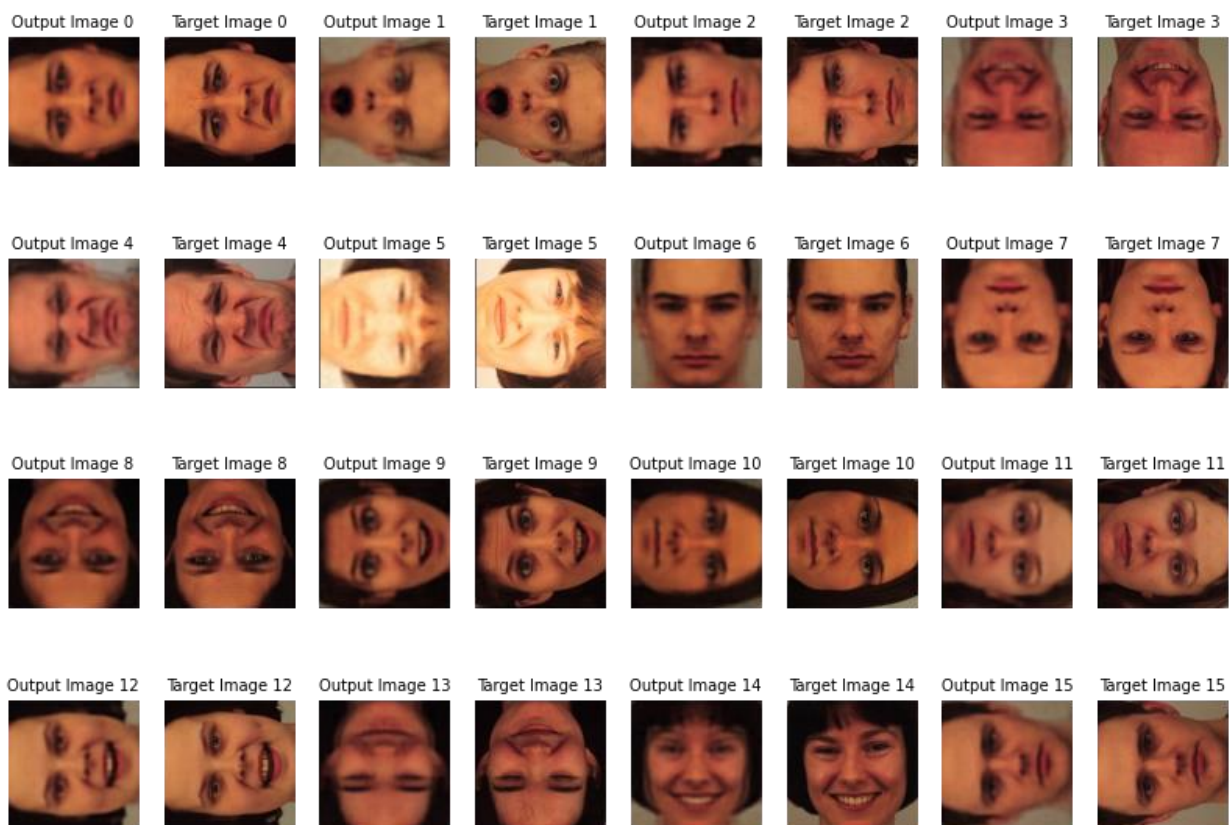
<sup>۲</sup>. Hyper-parameters



شکل ۵- نمودار تابع هزینه در حین آموزش بر دادگان KDEF



شکل ۶- نمودار دقت طبقه بند احساسات در حین آموزش بر دادگان KDEF



شکل ۷- تصاویر تولید شده به همراه تصاویر اصلی از مجموعه دادگان آموزش KDEF



شکل ۹- تصاویر تولید شده (سطح-پایین) به همراه تصاویر اصلی از مجموعه دادگان آزمون KDEF



شکل ۸- تصاویر تولید شده (سطح-بالا) به همراه تصاویر اصلی از مجموعه دادگان آزمون KDEF

## ۲,۶. نتیجه‌گیری:

همان طور که نتایج بالا نشان می‌دهد شبکه توانسته است ویژگی‌های سطح-پایین<sup>۱</sup> مانند چرخش تصاویر دیده شده و هم چنین ویژگی‌های سطح-بالا<sup>۲</sup> مانند احساسات موجود در چهره‌ی افراد را به خوبی یاد بگیرد. از مقایسه‌ی تصاویر ایجاد شده توسط شبکه بر روی دادگان آموزش و ارزیابی می‌توان تا حدی تفاوت عملکرد را حس کرد که این امر به خاطر کوچک بودن مجموعه دادگان این شبکه با توجه به بزرگ بودن شبکه تا حدی توجیه پذیر است. هم چنین همان طور که مشاهده می‌شود دقت طبقه بندی احساسات، بر روی مجموعه دادگان آموزش و تایید، تقریباً از همان ابتدای کار به ۱۰۰ درصد رسیده است.

## ۳. پیاده سازی تولید تصاویر اعداد دست‌نویس

مقاله‌ی مربوط به این پروژه، به پیاده‌سازی کاربرد شبکه‌ی خود برای اعداد دست‌نویس نیز پرداخته است. طبق ادعای این مقاله، مجموعه‌ی دادگان مورد استفاده در این پژوهش، `mnist` بوده است. باری با بررسی کدهای ارائه‌شده توسط نویسندگان [2]، این مسئله دریافت شد که این قسمت از پروژه با کمک دادگان `qmnist` پیاده‌سازی شده است، چرا که فاکتور دیگری تحت عنوان فرم<sup>۳</sup> یا استایل<sup>۴</sup> در برچسب‌های شبکه دخیل شده که دادگان `mnist` فاقد آن هستند. با بررسی بیشتر مشخص شد که دادگان به کار رفته در این تحقیق `qmnist` هستند. در ادامه به آماده‌سازی داده‌ها برای به کارگیری در شبکه پرداخته می‌شود.

### ۳,۱. آماده کردن مجموعه دادگان:

برای بارگذاری داده‌ها از روش به کار رفته در منبع [3] استفاده شده است.

دادگان این مجموعه شامل عکس‌هایی سیاه و سفید از اعداد دست‌نویس است که توسط افراد مختلفی نوشته شده و هر کدام از این نویسندگان با شماره‌ای در برچسب<sup>۵</sup> این تصاویر قابل شناسایی هستند. حالت کلی برچسب این تصاویر، از هشت بُعد تشکیل شده و به صورت زیر است:

جدول ۴- توضیح برچسب دادگان `qmnist`

Column	Description	Range
0	Character class	0 to 9
1	NIST HSF series	0, 1, or 4

<sup>۱</sup>. Low-level

<sup>۲</sup>. High-level

<sup>۳</sup> form

<sup>۴</sup> style

<sup>۵</sup>. Label

2	NIST writer ID	0-326 and 2100-2599
3	Digital index for this writer	0 to 149
4	NIST class code	30-39
5	Global NIST digit index	0 to 281769
6	Duplicate	0
7	Unused	0

شخص نویسنده، فرم یا استایل عدد دست‌نویس را تعیین می‌کند. ده نفر به صورت تصادفی از بین اشخاص نویسنده انتخاب شدند و تمام تصاویر دست‌نویس آنان با نامی به فرمت زیر در پوشه‌ای ذخیره شد:

[image number (0 to primary dataset size-1)]\_[writer new id (0 to 9)]\_[digit shown in the photo(0 to 9)].jpg

سپس تصاویر سیاه‌سفید و تک‌کاناله‌ی این پوشه به صورت عکس‌های رنگی سه‌کاناله‌ای در سه رنگ قرمز، آبی و سبز با فرمت نامی جدیدی ذخیره شدند. در انتها تصاویر رنگی به‌دست‌آمده از قبل، در چهار جهت چرخش صفر درجه، نود درجه، صد و هشتاد درجه و دویست و هفتاد درجه ذخیره می‌گردند. فرمت نامی تصاویر نهایی به صورت زیر است:

[image number]\_[writer id]\_[digit]\_[color id(0 to 2)]\_[rotation id(0 to 3)].jpg

برای هر یک از اشخاص در مجموعه دادگان ایجاد شده، یکی از اعداد انتخاب می‌گردد و یکی از رنگ‌ها و یا یکی از چرخش‌ها برای آن حذف می‌شوند. این دادگان حذف شده به عنوان مجموعه دادگان آزمون<sup>۱</sup> و باقی اعداد به عنوان مجموعه دادگان آموزش<sup>۲</sup> به کار می‌روند.

### ۳.۲. معماری شبکه

تصاویر مجموعه دادگان اعداد دست‌نویس به ابعاد  $28 \times 28$  هستند؛ به همین دلیل نمی‌توان از ساختار شبکه‌ی به کار رفته در قسمت ایجاد تصاویر استفاده نمود؛ زیرا که تصاویر چهره‌ی افراد ابعادی چند برابر تصاویر اعداد دست‌نویس دارند. همچنین، برچسب‌ها به چهار بخش فرد نویسنده(استایل)، شماره‌ی موجود در تصویر، رنگ و چرخش تقسیم می‌شوند. بدین ترتیب، نیاز به تغییراتی در ساختار شبکه احساس می‌گردد.

<sup>۱</sup>. Test Dataset

<sup>۲</sup> Train Dataset



جدول ۵- بخش تولیدکننده مربوط به شبکه‌ی تصاویر اعداد دست‌نویس

Layer name	Input Feature Size	Kernel Size	Padding	Stride	Output Feature Size
Writer id FC 1	10	–	–	–	256
Leaky ReLU (negative slope=0.1)					
Writer id FC 2	256	–	–	–	512
digit FC 1	10	–	–	–	256
Leaky ReLU (negative slope=0.1)					
digit FC 2	256	–	–	–	512
color FC1	3	–	–	–	32
Leaky ReLU (negative slope=0.1)					
color FC2	32	–	–	–	64
Transform FC 1	4	–	–	–	128
Leaky ReLU (negative slope=0.1)					
Transform FC 2	128	–	–	–	256
Then	these four		outputs	being concatenated.	
Leaky ReLU (negative slope=0.1)					
FC 3	1344	–	–	–	2048
Leaky ReLU (negative slope=0.1)					
FC 4	2048	–	–	–	2048
Leaky ReLU (negative slope=0.1)					
FC 5	2048	–	–	–	16384
Leaky ReLU (negative slope=0.1)					
Reshape to 256*8*8					
deConv 6	256	4	1	2	256
Leaky ReLU (negative slope=0.1)					
Conv 6	256	3	1	1	256
Leaky ReLU (negative slope=0.1)					
deConv 7	256	3	2	2	3
ReLU					

جدول ۶- بخش طبقه‌بند، مربوط به شبکه‌ی تصاویر اعداد دست‌نویس

Layer	Input Feature Size	Kernel Size	Padding	Stride	Output FeatureSize
Conv 1	3	2	0	2	128
Leaky ReLU (negative slope=0.2)					
Batch Normalization					
Conv 2	128	4	2	2	256
Leaky ReLU (negative slope=0.2)					
Batch Normalization					
Flattening Layer					
FC 3	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 4	4096	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 5	4096	-	-	-	10
FC 3	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 4	4096	-	-	-	1024
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 5	1024	-	-	-	10
FC 3	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 4	4096	-	-	-	1024
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 5	1024	-	-	-	3

FC 3	16384	-	-	-	4096
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 4	4096	-	-	-	1024
Leaky ReLU (negative slope=0.2)					
Dropout(0.5)					
FC 5	1024	-	-	-	4

### ۳,۳. آموزش شبکه

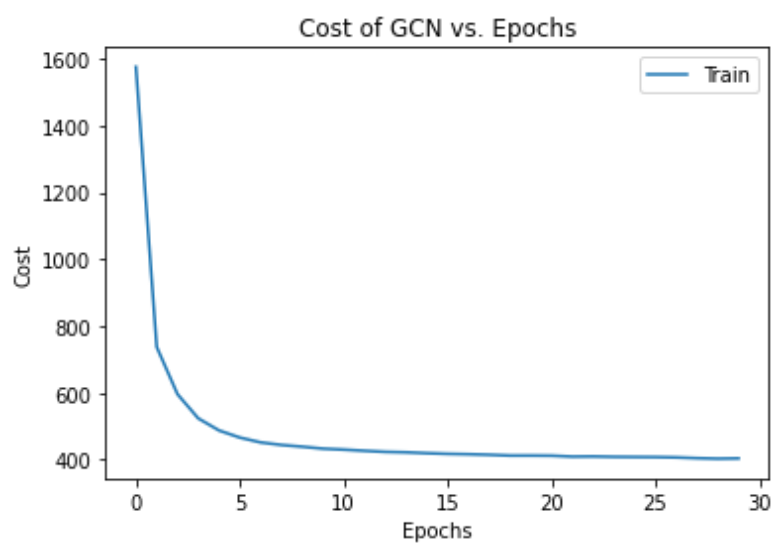
برای آموزش شبکه، از مجموعه ابرپارامترهای<sup>۱</sup> زیر استفاده کرده ایم:

جدول ۷- ابرپارامترهای استفاده شده در آموزش شبکه‌ی اعداد دسنویس

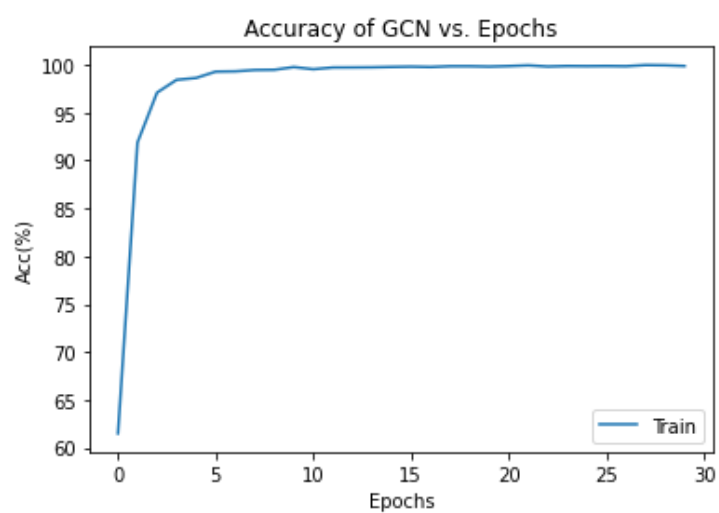
Optimizer	Adam
Learning rate	0.0002
Learning rate decay: after every 15 epochs	
$\beta_1$	0.9
$\beta_2$	0.995
$\varepsilon$	$10^{-8}$
batch size	64
epochs	30

<sup>۱</sup>.Hyper-parameters

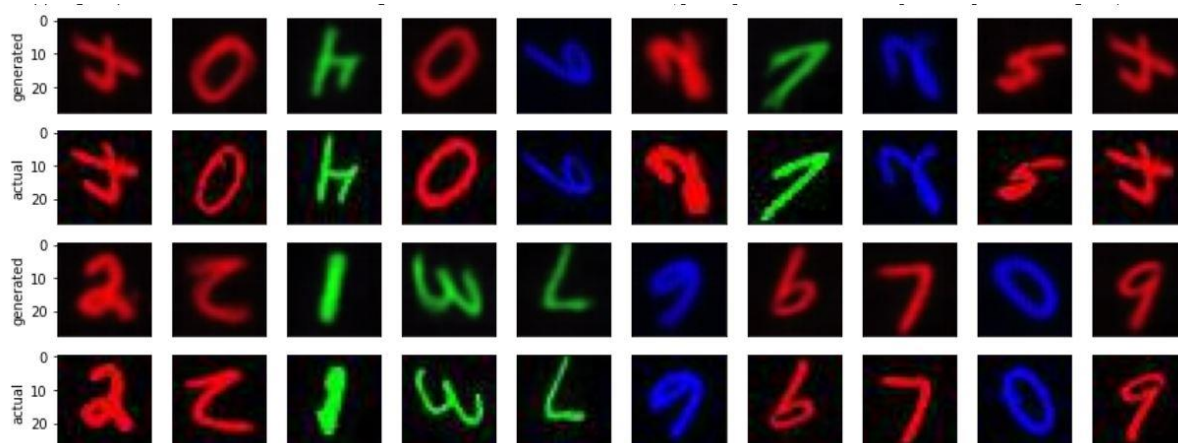
### ۳,۴. نتایج



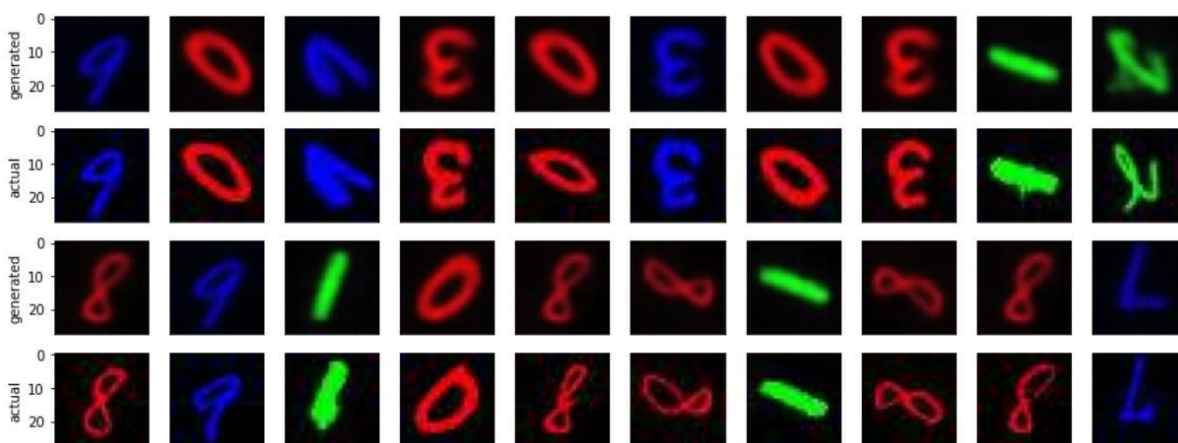
شکل ۱۰- نمودار هزینه در شبکه‌ی مربوط به اعداد دست‌نویس



شکل ۱۱- نمودار دقت در شبکه‌ی مربوط به اعداد دست‌نویس



شکل ۱۲-نمونه‌هایی از خروجی تولیدکننده شبکه در فاز آموزشی



شکل ۱۳-نمونه‌هایی از خروجی تولیدکننده شبکه در فاز آزمون

### ۳.۵. نتیجه‌گیری:

نتایج نشان می‌دهد که شبکه به صورت نسبی توانسته عملکرد خوبی کسب نماید. در تعدادی از اعداد تست، ممکن است مقداری تفاوت با استایل اصلی دیده شود؛ مثلاً امکان دارد قسمت پایینی تصویر عدد هشت در عکس اصلی، بزرگ‌تر از بخش بالایی آن باشد ولی در تصویر پیش‌بینی شده توسط شبکه این دو قسمت نسبتاً هم اندازه به نظر برسند. باید دقت کرد که پیچیدگی‌های تصاویر اعداد دست‌نویس و حتی ابعاد این تصاویر چندان نیست به همین دلیل نباید انتظار داشت که تصویر پیش‌بینی شده و اصلی «دقیقاً» مشابه هم باشند؛ چه بسا که یک فرد واحد در چند بار نوشتن یک عدد نمی‌تواند به صورت دقیق به یک شیوه بنویسد. بنابراین عملکرد فعلی شبکه به نظر قابل قبول است.

## ۴. افزودن دادگان<sup>۱</sup>

در این بخش در دو حالت زیر عمل طبقه بندی را روی احساس چهره افراد انجام می دهیم.

استفاده از تصاویر مجموعه دادگان KDEF

استفاده از تصاویر مجموعه دادگان KDEF به همراه تصاویر تولیدی در بخش قبلی

هدف این بخش این است که نشان دهیم استفاده از داده های تولیدی بخش قبلی می تواند دقت طبقه بندی را بالا ببرد.

### ۴.۱. حالت اول آموزش:

ابتدا شش نفر را انتخاب کرده و تمام تصاویر آنان را کنار می گذاریم. (در مجموع ۲۴ تصویر) تا مجموعه دادگان آزمون را تشکیل بدهند. ۶۴ نفر باقیمانده مجموعه دادگان آموزش را تشکیل می دهند (در مجموع ۲۵۶ تصویر). حال از یک شبکه عصبی عمیق برای طبقه بندی نوع احساس در تصاویر استفاده می کنیم. این شبکه عصبی دقیقاً همان شبکه ی طبقه بند در بخش ابتدایی است با این تفاوت که در خروجی آن به جای سه لایه سافت مکس<sup>۲</sup>، فقط یک لایه سافت مکس برای پیش بینی احساس وجود دارد.

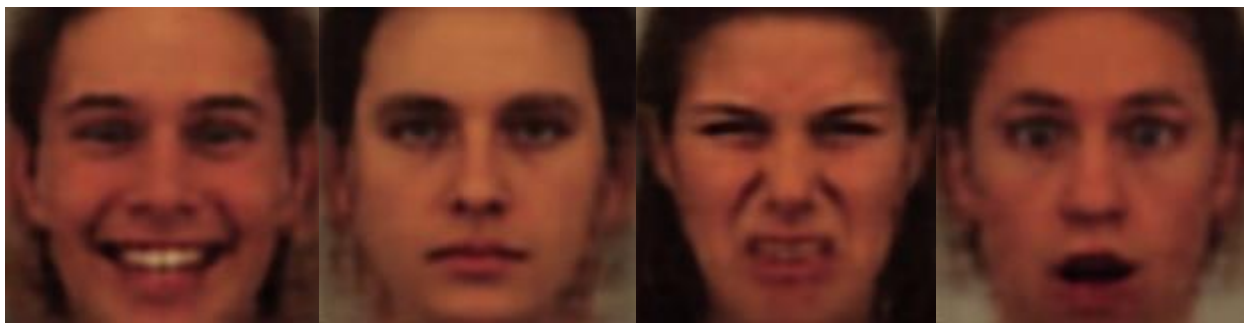
### ۴.۲. حالت دوم آموزش:

در این بخش علاوه بر ۲۵۶ تصویر بخش قبلی، با کمک شبکه تولیدکننده بخش قبل ۱۶۱۲۸ تصویر جدید تولید کرده و آن را به مجموعه دادگان آموزش اضافه می کنیم. شبکه تولیدکننده سه دسته ورودی دارد. ورودی اول یک بردار ۷۰ تایی برای مشخص کردن فرد است. ورودی های دوم و سوم نیز هر کدام یک بردار ۴ تایی برای مشخص کردن احساس و تبدیل هستند. برای تولید کردن نمونه جهت استفاده در طبقه بندی دو به دو افراد مختلف را که احساس و تبدیل یکسانی دارند را با هم ترکیب می کنیم. برای ترکیب کردن دو فرد، درایه های مربوط به آن دو فرد را در ورودی اول شبکه (بردار ۷۰ تایی) برابر ۰.۵ قرار داده و سایر درایه ها را صفر می گذاریم. دو ورودی دیگر شبکه نیز همان بردارهای one-hot احساس و تبدیل هستند. در شکل زیر چهار تصویر تولیدی به ازای چهار احساس متفاوت آورده شده است.

---

<sup>۱</sup> Data Augmentation

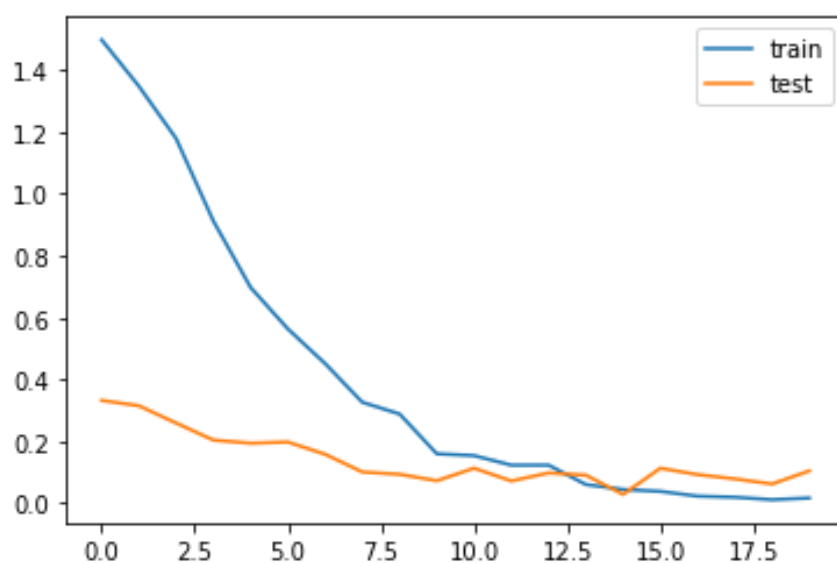
<sup>۲</sup> Softmax



شکل ۱۴- تصاویر افزوده، دارای احساسات متعجب، منزعج، عادی و خوشحال

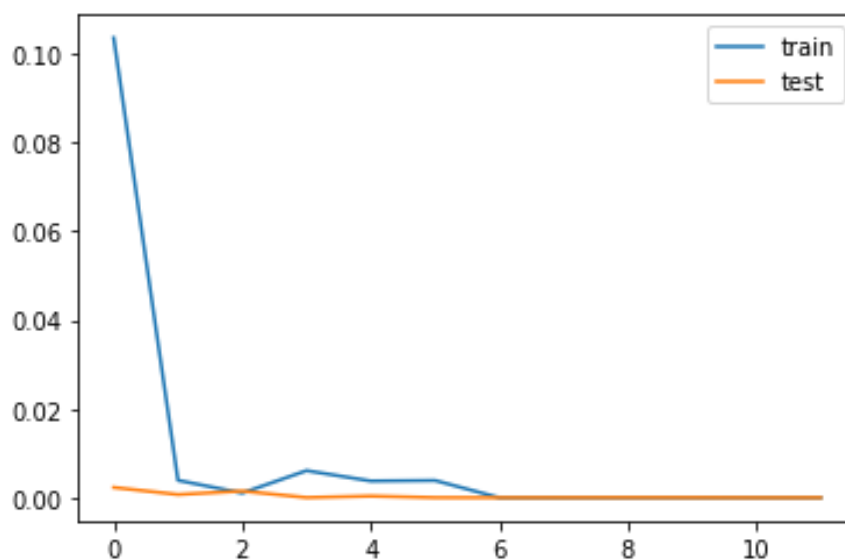
## ۴,۳. نتایج

در شکل زیر نمودار تابع هزینه برای داده های آموزش و آزمون حالت اول آورده شده است.



شکل ۱۵-نمودار هزینه طبقه بندی احساسات بدون افزودن داده

در شکل زیر نمودار تابع هزینه برای داده های آموزش و تست حالت دوم آورده شده است.



شکل ۱۶- نمودار هزینه طبقه بندی احساسات پس از افزودن داده

#### ۴,۴. نتیجه گیری

همان طور که از شکل ۱۵ پیداست، پس از ۱۴ ایپاک مدل دچار بیش‌برازش<sup>۱</sup> شده و لذا دقت آن بر روی مجموعه دادگان آزمون کاهش می یابد. بنابراین مقدار بهینه دقت بر روی مجموعه تست در ایپاک ۱۴ بدست می آید که برابر حدود ۹۲ درصد است. اما در شکل ۱۶، با افزایش دادگان آموزش به روش افزودن دادگان مشخص می‌شود که شبکه به خوبی آموزش دیده است. در این حالت دقت بر روی مجموعه تست برابر ۱۰۰ درصد می باشد.

---

<sup>۱</sup>.Overfit



## پیوست ۱: روند اجرای برنامه

برنامه‌ها در محیط گوگل کولب<sup>۱</sup> نوشته شدند. کدهای برنامه به صورتی نوشته شده است که نیازی به بارگذاری فایلی به صورت دستی در قسمتی نیست و می‌توان بدون نیاز به انجام عملیاتی دستی، تمام کدها را ران<sup>۲</sup> کرد. تنها نکته‌ای که شاید لازم باشد به آن اشاره نمود، بارگیری مجموعه دادگان تصاویر افراد از مسیر زیر و بارگذاری آن روی گوگل درایو<sup>۳</sup> است:

`!wget-P/content/drive/My\Drive/https://www.kdef.se/download/KDEF and AKDEF.zip`

تنها یک بار ران کردن این قسمت برای بارگذاری داده‌ها روی گوگل درایو کافی است. این قسمت در کد اصلی موجود است و تنها عمل مورد نیاز برای ران کردن کدها، بارگذاری خود فایل کد روی گوگل کولب و ران کردن آن است.

---

<sup>۱</sup> Google Colab

<sup>۲</sup> Run

<sup>۳</sup> Google Drive

- [1] Xu, Q., Qin, Z., & Wan, T. (2019, March). Generative cooperative net for image generation and data augmentation. In International Symposium on Integrated Uncertainty in Knowledge Modelling and Decision Making (pp. 284-294). Springer, Cham.
- [2] <https://github.com/Xharlie/MultiGen>
- [3] <https://www.programmersought.com/article/54863640507/>
- [4] Dosovitskiy, A., Tobias Springenberg, J., & Brox, T. (2015). Learning to generate chairs with convolutional neural networks. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 1538-1546).
- [5] Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. Advances in neural information processing systems, 25, 1097-1105.