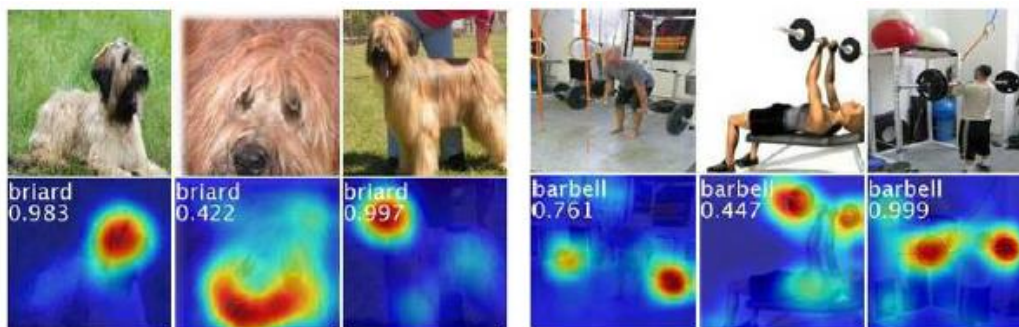
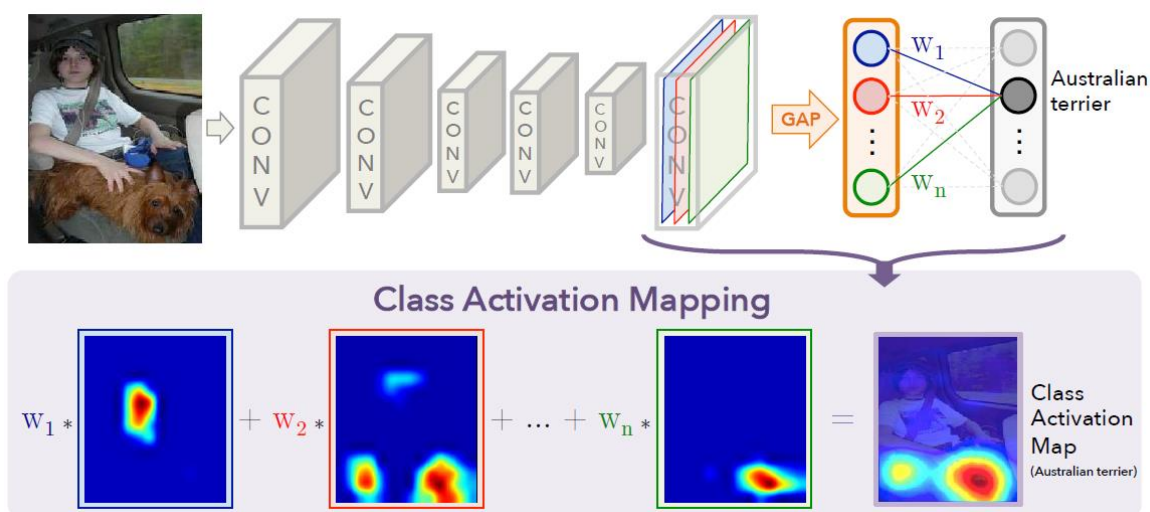


CNN یا class activation map با استفاده از لایه‌ی global average pooling در شبکه‌های CNN ساخته می‌شود. CAM برای یک دسته خاص، نواحی متمایز کننده تصویر استفاده شده توسط CNN برای شناسایی آن دسته را نشان می‌دهد. به عنوان مثال در شکل زیر CAMها نواحی متمایز کننده تصویر که برای طبقه‌بندی تصویر استفاده شده یعنی، سر حیوان برای سگ برید و صفحات در هالتر، را برجسته می‌کنند.



از معماری مشابه Network in Network و گوگل نت استفاده می‌کنیم. شبکه عمدتاً از لایه‌های کانولوشنال تشکیل شده است و درست قبل از لایه خروجی نهایی (در مورد طبقه بندی softmax)، global average pooling را روی نقشه‌های ویژگی‌های کانولوشن انجام می‌دهیم و از آن‌ها به عنوان ویژگی برای یک لایه کاملاً متصل که خروجی مورد نظر را تولید می‌کند استفاده می‌کنیم. با توجه به این ساختار اتصال ساده، می‌توانیم اهمیت نواحی تصویر را با بازتاب دادن وزن‌های لایه خروجی بر روی نقشه‌های ویژگی کانولوشنال شناسایی کنیم، تکنیکی که آن را نگاشت فعال سازی کلاس (class activation mapping) می‌نامیم.



همانطور که در شکل فوق نشان داده شده است، global average pooling میانگین مکانی نقشه ویژگی هر واحد را در آخرین لایه کانولوشنی به دست می دهد. از مجموع وزن دار این مقادیر برای تولید خروجی نهایی استفاده می شود. به طور مشابه، ما مجموع وزن دار از نقشه های ویژگی آخرین لایه کانولوشن را محاسبه می کنیم تا class activation mapهای خود را به دست آوریم. این تکنیک برای softmax توضیح داده می شود اما برای رگرسیون و سایر توابع ضرر نیز به همین شکل هست. برای یک تصویر $f_k(x,y)$ خروجی فعالسازی واحد k ام در آخرین لایه کانولوشنی است. پس از آن نتیجه انجام global average pooling بر روی واحد k ام برابر است با:

$$F^k = \sum_{x,y} f_k(x,y)$$

بنابراین برای کلاس c ورودی تابع softmax به صورت زیر است:

$$S_c = \sum_k w_k^c F_k$$

که w_k^c وزن متناظر با واحد k ام برای کلاس c است که اهمیت F_k را برای کلاس c نشان می دهد. در نهایت خروجی softmax برای کلاس c به صورت زیر محاسبه می شود:

$$P_c = (\exp(S_c)) / (\sum_c \exp(S_c))$$

با جایگذاری $F^k = \sum_{x,y} f_k(x,y)$ در فرمول S_c داریم:

$$S_c = \sum_k w_k^c \sum_{x,y} f_k(x,y) = \sum_{x,y} \sum_k w_k^c f_k(x,y)$$

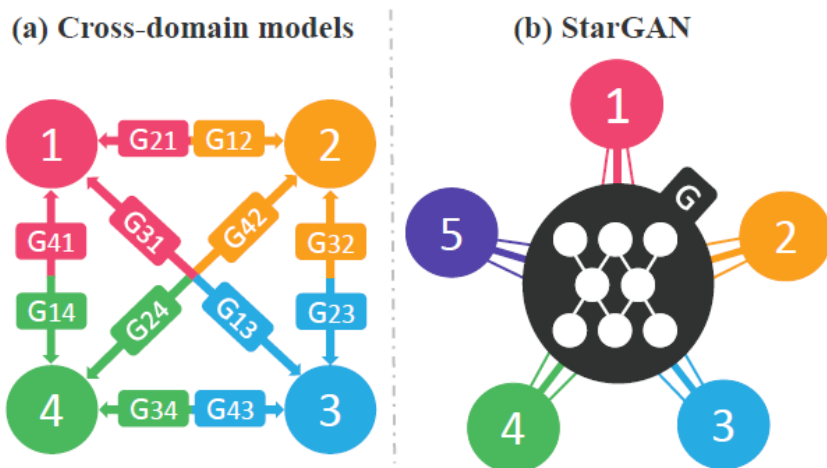
حال class activation map برای هر کلاس با توجه به مختصات فضایی آن به صورت زیر تعریف می شود:

$$M_c(x, y) = \sum_k w_k^c f_k(x, y).$$

در نتیجه $S_c = \sum_{x,y} M_c(x, y)$ و از این رو $M_c(x, y)$ مستقیماً اهمیت فعال سازی در گرید فضایی (x, y) را نشان می دهد که منجر به طبقه بندی یک تصویر به کلاس c می شود.

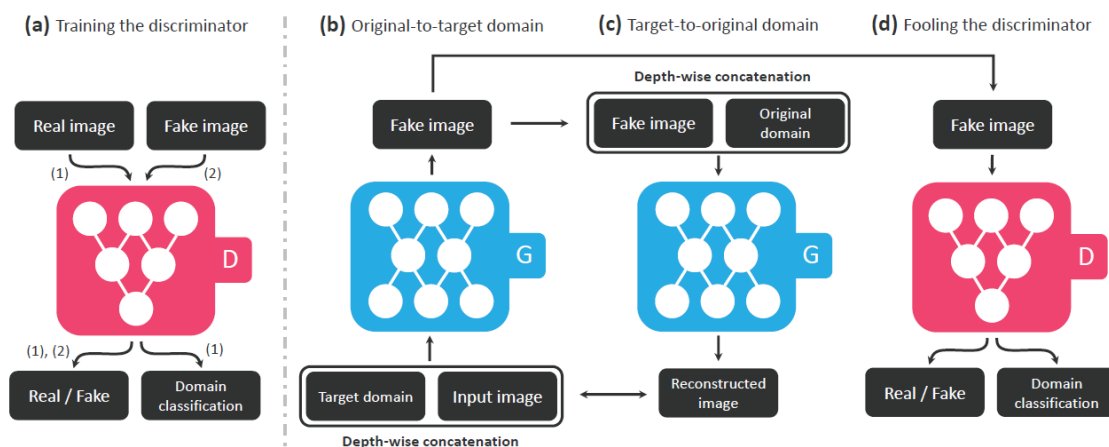
وظیفه ترجمه تصویر به تصویر این است که یک جنبه خاص از یک تصویر داده شده را به دیگری تغییر دهد، به عنوان مثال، تغییر حالت چهره یک فرد از خنده به اخم. این کار پس از معرفی شبکه‌های مولد رقابتی (GANs)، با نتایجی از تغییر رنگ مو، بازسازی عکس‌ها از نقشه‌های لبه و تغییر فصول تصاویر مناظر، پیشرفت‌های قابل توجهی را تجربه کرده است. با داشتن داده‌های آموزشی از دو دامنه متفاوت، این مدل‌ها یاد می‌گیرند که تصاویر را از یک دامنه به دامنه دیگر ترجمه کنند. ما اصطلاحات *attribute* را به عنوان یک ویژگی معنادار ذاتی در یک تصویر مانند رنگ مو، جنسیت یا سن، و مقدار ویژگی را به عنوان مقدار خاصی از یک ویژگی، به عنوان مثال، سیاه/بور/قهوه ای برای رنگ مو یا مرد/زن برای جنسیت نشان می‌دهیم. ما همچنین دامنه را به عنوان مجموعه ای از تصاویر با ارزش مشخصه یکسان نشان می‌دهیم. به عنوان مثال، تصاویر زنان می‌تواند یک حوزه را نشان دهد در حالی که تصاویر مردان نشان دهنده حوزه دیگری است. در ترجمه تصویر به تصویر چند دامنه تصاویر را با توجه به ویژگی‌های چندین دامنه تغییر می‌دهیم. مدل‌های موجود در چنین کارهای ترجمه تصویر چند دامنه‌ای کم‌بازده (inefficient) و ناکارآمد (ineffective) هستند. کم‌بازده بودن آن‌ها از این واقعیت ناشی می‌شود که برای یادگیری همه نگاشت‌ها در میان k دامنه، $k(k-1)$ مولد باید آموزش داده شوند. همانطور که از شکل زیر پیداست دوازده شبکه مولد مجزا باید برای ترجمه تصاویر بین چهار دامنه مختلف آموزش داده شوند. در همین حال، آنها ناکارآمد هستند یعنی اگرچه ویژگی‌های سراسری وجود دارد که می‌توان از تصاویر همه دامنه‌ها یاد گرفت، مانند اشکال چهره، اما هر مولد نمی‌تواند به طور کامل از کل داده‌های آموزشی استفاده کند و فقط می‌تواند از دو دامنه از k دامنه یاد بگیرد. عدم استفاده کامل از داده‌های آموزشی احتمالاً کیفیت تصاویر تولید شده را محدود می‌کند. علاوه بر این، آنها قادر به آموزش مشترک دامنه‌ها از مجموعه داده‌های مختلف نیستند، زیرا هر مجموعه داده به صورت جزئی برچسب گذاری شده است. به عنوان راه‌حلی برای چنین مشکلاتی، StarGAN را پیشنهاد می‌کنیم، یک شبکه مولد رقابتی که قادر به یادگیری نگاشت میان دامنه‌های متعدد است. همانطور که در شکل زیر قسمت (b) نشان داده شده است، این مدل داده‌های آموزشی چندین دامنه را دریافت می‌کند و نگاشت بین تمام دامنه‌های موجود را تنها با استفاده از یک مولد یاد می‌گیرد. StarGAN به جای یادگیری یک ترجمه ثابت (مثلاً موهای سیاه به بلوند)، هم تصویر و هم اطلاعات دامنه را به عنوان ورودی می‌گیرد و یاد می‌گیرد که به طور انعطاف پذیر تصویر ورودی را به دامنه مربوطه ترجمه کند. در طول آموزش، به‌طور تصادفی یک برچسب دامنه هدف را تولید می‌کنیم و به مدل آموزش می‌دهیم تا به‌طور انعطاف‌پذیر یک تصویر ورودی را به دامنه هدف ترجمه کند. با انجام این کار، می‌توانیم برچسب دامنه را کنترل کرده و تصویر را در مرحله آزمون به هر دامنه دلخواه ترجمه کنیم. همچنین یک رویکرد ساده اما مؤثر را معرفی می‌شود که آموزش مشترک بین دامنه‌های مجموعه داده‌های مختلف را با

افزودن یک بردار ماسک به برچسب دامنه امکان پذیر می سازد. روش پیشنهادی تضمین می کند که مدل می تواند برچسب های ناشناخته را نادیده بگیرد و بر روی برچسب ارائه شده توسط یک مجموعه داده خاص تمرکز کند.



همانطور که از شکل فوق پیداست برای مدیریت چندین دامنه، CycleGAN که یک مدل بین دامنه ای است باید برای هر جفت دامنه تصویر ساخته شود، در حالی که StarGAN قادر به یادگیری نگاشت بین چندین دامنه با استفاده از یک مولد واحد است. این باعث می شود تعداد پارامترهایی که StarGAN برای ترجمه تصویر احتیاج دارد بسیار کمتر از CycleGAN باشد.

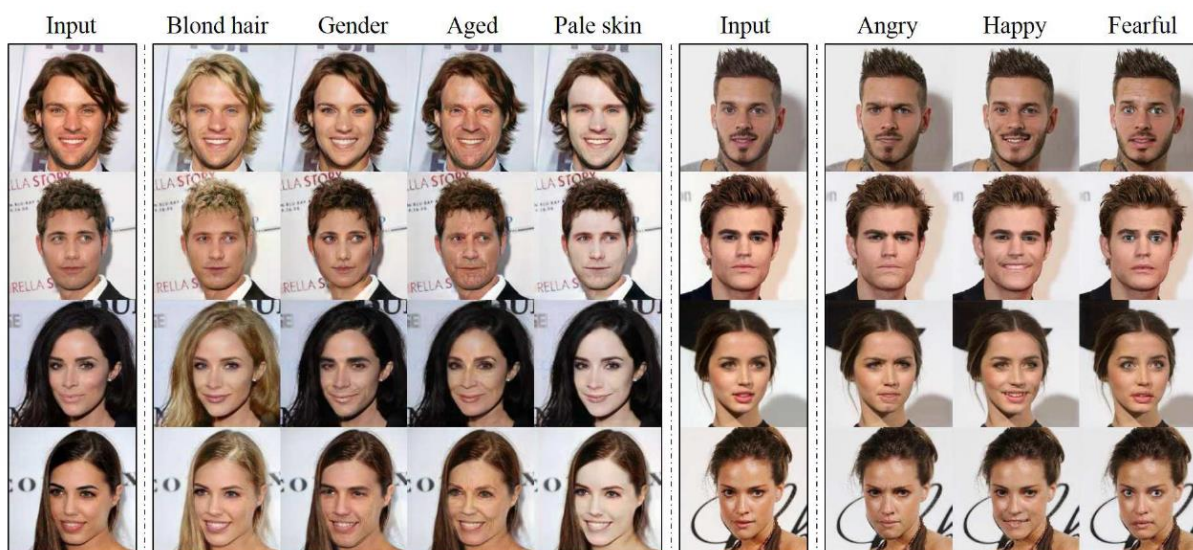
همچنین معماری StarGAN به شکل زیر است:



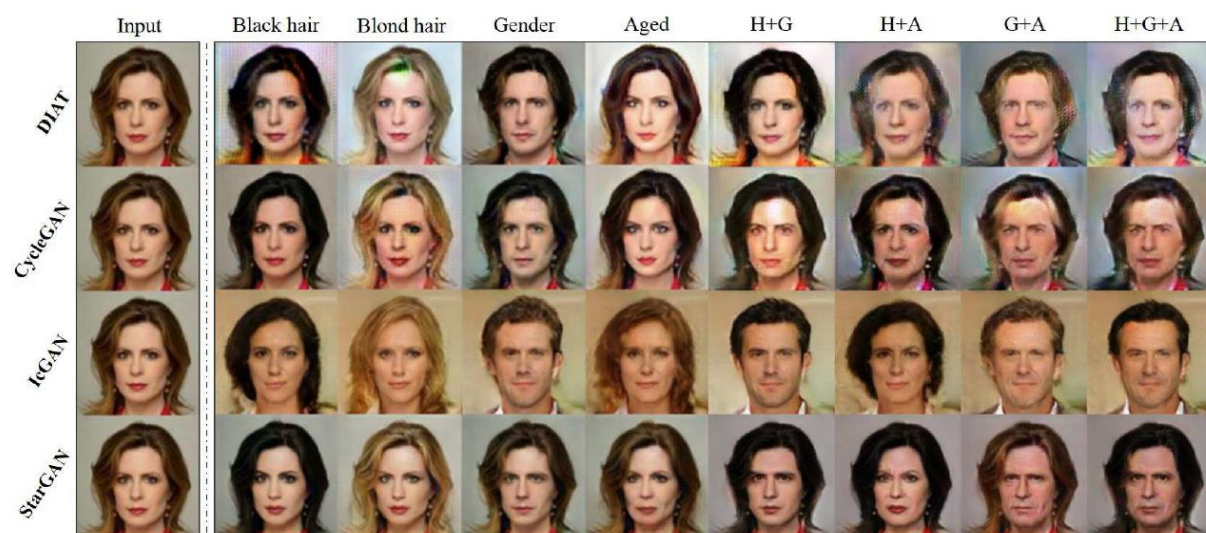
(a) D یاد می گیرد که بین تصاویر واقعی و جعلی تمایز قائل شود و تصاویر واقعی را به حوزه مربوطه خود طبقه بندی کند. (b) G هم تصویر و هم برچسب دامنه هدف را به عنوان ورودی می گیرد و یک تصویر جعلی

تولید می‌کند. برچسب دامنه هدف به صورت مکانی تکرار شده و با تصویر ورودی الحاق می‌شود. (c) G سعی می‌کند تصویر اصلی را از تصویر جعلی با توجه به برچسب دامنه اصلی بازسازی کند. (d) G سعی می‌کند تصاویر غیر قابل تشخیص از تصاویر واقعی و طبقه‌بندی به عنوان دامنه هدف توسط D ایجاد کند.

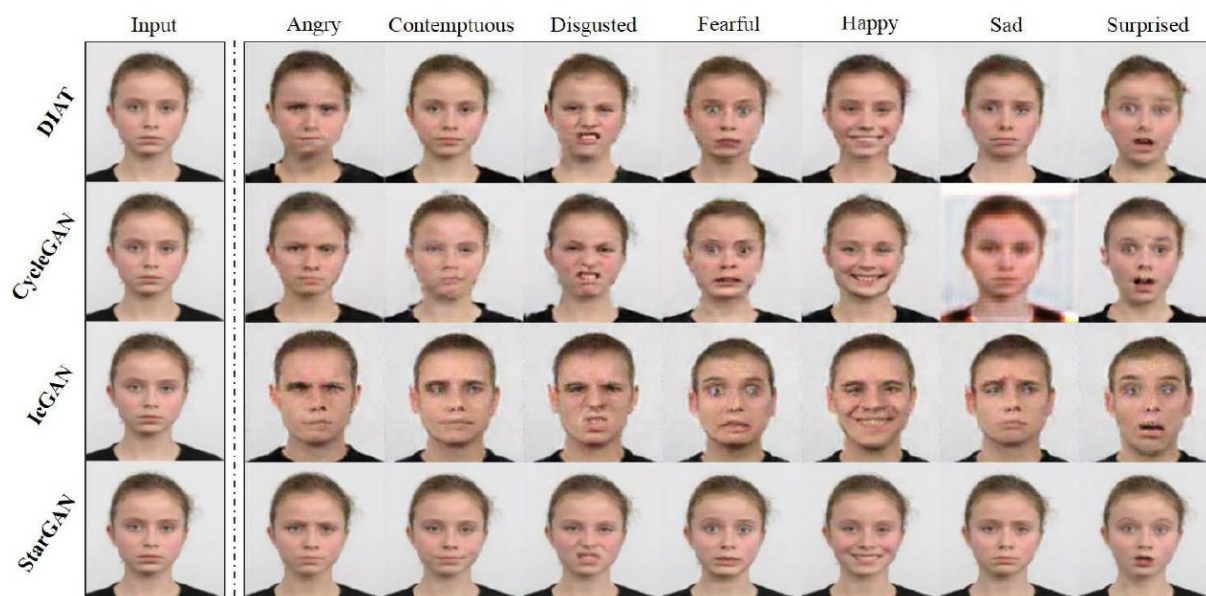
مقاله آزمایشات خود را بر روی دو مجموعه داده CelebA و RaFD انجام می‌دهد که در اولی رنگ مو، جنسیت و سن افراد تغییر می‌کند اما در دومی حالت‌های هیجانی مختلف چهره فرد ضبط شده است. شکل زیر نتایج ترجمه تصویر به تصویر چند دامنه‌ای در مجموعه داده CelebA از طریق انتقال دانش آموخته شده از مجموعه داده RaFD را نشان می‌دهد. ستون‌های اول و ششم تصاویر ورودی را نشان می‌دهند در حالی که ستون‌های باقیمانده تصاویر تولید شده توسط StarGAN هستند. توجه داشته باشید که تصاویر توسط یک شبکه مولد واحد تولید می‌شوند.



همچنین در شکل زیر می‌توان نتایج StarGAN در تغییر رنگ مو، جنسیت و سن افراد یا ترکیبی از آن‌ها را با سایر مدل‌ها مقایسه کرد.



در شکل زیر نیز می‌توان نتیجه StarGAN در تغییر حالت‌های هیجانی چهره افراد را با سایر مدل‌ها مقایسه کرد.



تفاوت اصلی cycleGAN و starGAN بالاتر با هیالات زرد مشخص شده‌است.

(۳)

توضیحات مربوط به چگونگی کارکرد کد مربوط به DCGAN در فایل نوتبوک مربوطه نوشته شده است.

شبهه کد اصلی برای لاس به صورت زیر است:

$$L_{\text{genAB}} = E_{\text{realA} \sim P_{\text{data}}(\text{realA})} [1 - \log(\text{discB}(\text{genAB}(\text{realA})))] \\ + \lambda * E_{\text{realA} \sim P_{\text{data}}(\text{realA})} [\|\text{genBA}(\text{genAB}(\text{realA})) - \text{realA}\|_1]$$

اما با توجه به اینکه این فرمول در زمان آپدیت گرادیان باعث می شود زمانی که به نقطه بهینه نزدیک شویم تغییرات شدید در آپدیت گرادیان به وجود بیاید و بهینه سازی را با مشکل روبرو کند بنابراین بهتر است به شبهه کد زیر تغییر کند:

$$L_{\text{genAB}} = - E_{\text{realA} \sim P_{\text{data}}(\text{realA})} [\log(\text{discB}(\text{genAB}(\text{realA})))] \\ + \lambda * E_{\text{realA} \sim P_{\text{data}}(\text{realA})} [\|\text{genBA}(\text{genAB}(\text{realA})) - \text{realA}\|_1]$$

mode collapse: معمولاً می خواهید GAN شما خروجی های متنوعی تولید کند. برای مثال، برای هر ورودی تصادفی به مولد چهره خود، یک چهره متفاوت می خواهید. با این حال، اگر یک مولد خروجی قابل قبولی تولید کند، مولد ممکن است یاد بگیرد که فقط آن خروجی را تولید کند. در واقع، مولد همیشه در تلاش است تا خروجی ای را بیلد که برای ممیز (discriminator) قابل قبول ترین به نظر می رسد. اگر مولد بارها و بارها شروع به تولید همان خروجی (یا مجموعه کوچکی از خروجی ها) کند، بهترین استراتژی ممیز این است که یاد بگیرد همیشه آن خروجی را رد کند. اما اگر نسل بعدی ممیز در یک مینیمم محلی گیر کند و بهترین استراتژی را پیدا نکند، برای تکرار بعدی مولد یافتن معقول ترین خروجی برای ممیز فعلی بسیار آسان است. هر تکرار از مولد برای یک ممیز خاص بیش از حد بهینه (over-optimize) می شود و ممیز هرگز نمی تواند راه خود را از تله بیاموزد. در نتیجه مولدها از طریق مجموعه کوچکی از انواع خروجی می چرخند. این شکل از خرابی GAN، mode collapse نامیده می شود. بنابراین mode collapse زمانی است که GAN تعداد کمی از تصاویر را با موارد تکراری (حالت های) زیاد تولید می کند. این زمانی اتفاق می افتد که مولد قادر به یادگیری یک نمایش ویژگی غنی نباشد، زیرا یاد می گیرد خروجی های مشابه را به چندین ورودی مختلف مرتبط کند. برای بررسی mode collapse، تصاویر تولید شده را بررسی کنید. اگر تنوع کمی در خروجی وجود داشته باشد و برخی از آنها تقریباً یکسان باشند، احتمال mode collapse وجود دارد. اگر مشاهده کردید که این اتفاق می افتد، سعی کنید توانایی مولد برای ایجاد خروجی های متنوع تر را با موارد زیر افزایش دهید:

افزایش ابعاد داده های ورودی به مولد

افزایش تعداد فیلترهای مولد برای ایجاد تنوع گسترده‌تری از ویژگی‌ها
آسیب رساندن به ممیز با دادن برجسب‌های نادرست تصادفی به تصاویر واقعی

<https://ch.mathworks.com/help/deeplearning/ug/monitor-gan-training-progress-and-identify-common-failure-modes.html>

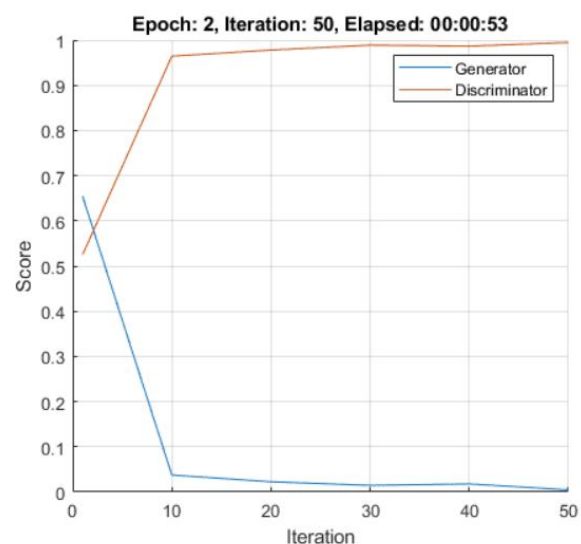
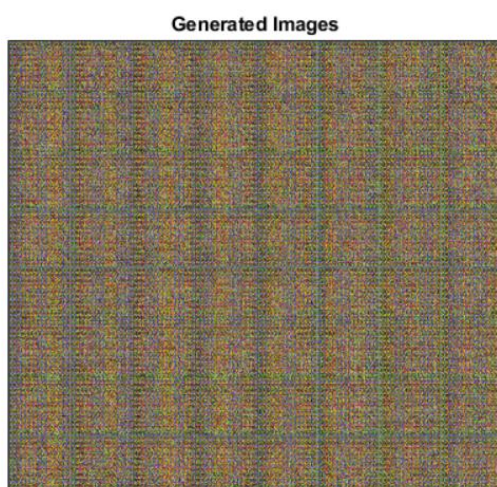
<https://wandb.ai/authors/DCGAN-ndb-test/reports/Measuring-Mode-Collapse-in-GANs--VmlldzoxNzg5MDk>

convergence failure

convergence failure زمانی اتفاق می‌افتد که مولد و ممیز در طول فرایند آموزش به تعادل نرسند.

ممیز غلبه می‌کند:

این سناریو زمانی اتفاق می‌افتد که score مولد به صفر یا نزدیک به صفر برسد و score برای ممیز به یک یا نزدیک به یک برسد. نمودار زیر نمونه‌ای از غلبه ممیز بر مولد را نشان می‌دهد. توجه داشته باشید که score مولد به صفر نزدیک می‌شود و بازیابی نمی‌شود. در این حالت، ممیز بیشتر تصاویر را به درستی طبقه‌بندی می‌کند. به نوبه خود، مولد نمی‌تواند هیچ تصویری تولید کند که ممیز را فریب دهد و در نتیجه در یادگیری شکست می‌خورد.



اگر برای تکرارهای زیادی score از این مقادیر بازبایی نشد، بهتر است فرایند آموزش را متوقف کنید. اگر این اتفاق افتاد، سعی کنید عملکرد مولد و ممیز را با روش زیر متعادل کنید:

آسیب رساندن به ممیز با دادن برچسب های نادرست تصادفی به تصاویر واقعی

افزودن لایه های dropout ، به ممیز

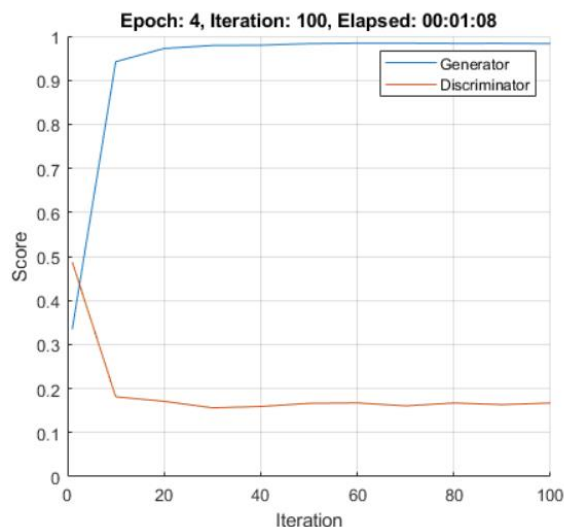
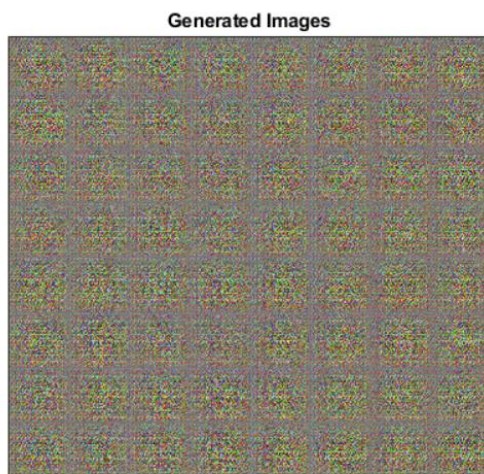
بهبود توانایی مولد برای ایجاد ویژگی های بیشتر با افزایش تعداد فیلترها در لایه های کانوولوشنی آن

آسیب رساندن به ممیز با کاهش تعداد فیلترهای آن

مولد غلبه می کند:

این سناریو زمانی اتفاق می افتد که score مولد به یک یا نزدیک به یک برسد.

نمودار زیر نمونه ای از غلبه مولد بر ممیز را نشان می دهد. توجه داشته باشید که score مولد برای چندین تکرار به یک می رسد. در این مورد، مولد تقریباً همیشه یاد می گیرد که چگونه ممیز را فریب دهد. وقتی این خیلی زود در فرایند آموزش اتفاق می افتد، احتمالاً مولد یک بازنمایی ویژگی بسیار ساده را یاد می گیرد که ممیز را به راحتی فریب می دهد. این بدان معناست که تصاویر تولید شده علیرغم داشتن امتیاز بالا می توانند بسیار ضعیف باشند. توجه داشته باشید که در این مثال، امتیاز ممیز خیلی به صفر نزدیک نمی شود، زیرا همچنان می تواند برخی از تصاویر واقعی را به درستی طبقه بندی کند.



اگر برای تکرارهای زیاد score از این مقادیر بازبایی نشد، بهتر است فرایند آموزش را متوقف کنید. اگر این اتفاق افتاد، سعی کنید عملکرد مولد و ممیز را با روش زیر متعادل کنید:

بهبود توانایی ممیز برای یادگیری ویژگی های بیشتر با افزایش تعداد فیلترها

آسیب رساندن به مولد با افزودن لایه های dropout

آسیب رساندن به مولد با کاهش تعداد فیلترهای آن

حالت زیر نیز یکی از موارد غلبه مولد است:

همانطور که مولد با آموزش بهبود می یابد، عملکرد ممیز بدتر می شود زیرا ممیز نمی تواند به راحتی تفاوت بین واقعی و جعلی را تشخیص دهد. اگر مولد به طور کامل موفق شود، پس ممیز دارای دقت ۵۰٪ است. در واقع انگار ممیز یک سکه می اندازد و با شیر یا خط پیش بینی خود را انجام می دهد. این پیشرفت یک مشکل برای همگرایی GAN ایجاد می کند: بازخورد ممیز با گذشت زمان کمتر معنادار می شود. اگر GAN از نقطه ای که ممیز بازخورد کاملاً تصادفی می دهد به آموزش ادامه دهد، مولد شروع به آموزش بر روی بازخورد ناخواسته می کند و کیفیت خودش ممکن است سقوط کند.

<https://ch.mathworks.com/help/deeplearning/ug/monitor-gan-training-progress-and-identify-common-failure-modes.html>

<https://developers.google.com/machine-learning/gan/training>