

1. از لحاظ معماری :

Diffusion model : مدل دیفوزن از لحاظ معماری از دو بخش forward process و reverse process تشکیل شده است که بخش reverse آن می تواند از یک نویز تصادفی به یک عکس جدید برسد.

GAN : GAN از دو شبکه generator و discriminator تشکیل شده است. شبکه discriminator باید داده های واقعی را از داده های فیک تشخیص دهد و شبکه generator باید از یک توزیع $p(z)$ عکس هایی تولید کند که شبکه discriminator را فریب دهد.

VAE : اتوانکودر variational از دو بخش encoder و decoder تشکیل شده است. این شبکه بطوری کلی یک ورودی عکس را به خودش در خروجی تبدیل می کند. اما در فضای latent خود می انالیز و واریانس توزیع داده ها را به ما می دهد که می توانیم از آن سمپل بگیریم و با بخش decoder عکس جدید تولید کنیم.

روش آموزش :

همه سه مدل بصورت supervised آموزش داده می شوند.

Diffusion : دیفوزن بتدریج به ورودی نویز اضافه می کند تا مرحله ای که به یک نویز خاص برسد و در فرایند reverse ، نویز را به عکس تبدیل می کند.

GAN : ابتدا discriminator روی داده های واقعی با لیبِل 1 آموزش داده می شود سپس generator با استفاده از ورودی با توزیع $p(z)$ سعی می کند عکس تولید کند و شبکه discriminator روی داده ها خروجی generator (که فیک هستند) با لیبِل 0 آموزش داده می شود و خروجی 0 از discriminator به generator داده می شود و باعث می شود generator عکس هایی تولید کند که شبکه discriminator را فریب دهد.

VAE : در این شبکه ، ورودی و خروجی هر دو یک عکس داده می شود و این شبکه باید یاد بگیرد عکس را تبدیل به feature map کرده و دوباره از feature map ، reconstruct کند.

کیفیت و تنوع نمونه ها :

GAN می تواند عکس های با کیفیت بالا و واقعی تولید کند و همچنین عکس های متنوع ارائه کند. اما ممکن است دچار mode collapse شود و چند حالت از توزیع داده را تولید کند و بقیه را نادیده بگیرد. VAE می تواند نمونه های متنوع تولید کند اما خروجی هایی ممکن است نویزی و blurry باشد و نمی تواند داده های پیچیده با فضای latent ایجاد بالا تولید کند.

Diffusion می تواند نمونه های واقعی و با جزئیات بالا تولید کند. چون تعداد زیادی step برای تولید نمونه و همچنین شبکه عصبی قدرتمندی برای reverse دارد. عکس های تولید شده توسط Diffusion از تنوع و کیفیت بیشتری نسبت به VAE و GAN برخوردار است.

کاربرد :

از کاربردهای مدل های generative می توان ساخت عکس ، افزودن data (augmentation) تشخیص anomaly را نام برد.

تولید و ساخت عکس توسط diffusion و VAE تنوع بیشتری دارد. اما GAN دارای تنوع کمتر اما نمونه های با کیفیت می باشد.

از GAN و diffusion در شبکه image-image translation استفاده می شود که برای تبدیل یک فرمت عکس به فرمت دیگر است مثل style transfer یا domain transform. همچنین چون مدل قادر به تولید عکس های جدید دارند می توانند به عنوان data aug استفاده شوند.

برای تشخیص outlier یا anomaly data هم می توان از VAE استفاده کرد که این کار را با محاسبه reconstruction loss انجام می دهد یا از GAN و Diffusion استفاده کرد که این کار را به ترتیب با اندازه گیری امتیاز discriminator و امتیاز reverse diffusion انجام می دهند.

2. مدل های Diffusion در شکل های مختلفی مثل segmentation, denoising, image-image translation در تصاویر پزشکی استفاده می شوند.

دلایل استفاده از آن ها عبارتند از: (مزیت ها)

- توانایی حذف کردن رتبه های پیچیده با ابعاد بالا مثل تصاویر پزشکی 3D
- با داده های پرچسب دار و همچنین غیر پرچسب دار آموزش داده می شوند و چون در حوزه پزشکی داده های با پرچسب محدود است مناسب است.
- توانایی تولید عکس هایی با کیفیت و جزئیات بالا و همچنین حفظ distribution داده ها را دارند.
- می توانند از دانش های قبلی مثل segmentation mask و اطلاعات کلیکی نیز استفاده کنند و با آن فرایند تولید را انجام دهد.

1. چون داده های ریاست ما یک manifold ای از فضای \mathbb{R} را اشغال می کنند و اگر به داده ها نویز اضافه نکنیم تا بفک کنند و سبتر از manifold را اشغال کنند، گرایان هایمان در نقطه initial، بی معنی (non sense) خواهد بود و به سمت توزیع درست داده ها، هدایت نمی شویم. پس نیاز به فرایند ریفرم داریم تا در هر مرحله یک نویز که به ریاست اضافه کنیم تا $\nabla_{\alpha} P(\alpha)$ مقدار معنی داری پیدا کرده و ما را به توزیع اصلی نزدیک کند.

$$P(\alpha, t) = \sum_{j=1}^n w_j \frac{1}{\sqrt{2\pi (\sigma_j^2 + \sigma^2 t)}} \exp \left\{ -\frac{[\alpha - \mu_j]^2}{2(\sigma_j^2 + \sigma^2 t)} \right\} \quad .2$$

$$\nabla_{\alpha} P(\alpha, t) = \sum_{j=1}^n w_j \frac{1}{\sqrt{2\pi (\sigma_j^2 + \sigma^2 t)}} \cdot e^{\left[-\frac{(\alpha - \mu_j)^2}{2(\sigma_j^2 + \sigma^2 t)} \right]} \cdot \frac{-(\alpha - \mu_j)}{\sigma_j^2 + \sigma^2 t}$$

$$\begin{aligned} \nabla_{\alpha} \log P(\alpha, t) &= \frac{\nabla_{\alpha} P(\alpha, t)}{P(\alpha, t)} \\ \nabla_{\alpha} \log P(\alpha, t) &= \frac{\sum_{j=1}^n w_j \frac{1}{\sqrt{2\pi (\sigma_j^2 + \sigma^2 t)}} e^{\left[-\frac{(\alpha - \mu_j)^2}{2(\sigma_j^2 + \sigma^2 t)} \right]} \cdot \frac{-(\alpha - \mu_j)}{\sigma_j^2 + \sigma^2 t}}{\sum_{j=1}^n w_j \frac{1}{\sqrt{2\pi (\sigma_j^2 + \sigma^2 t)}} e^{\left[-\frac{(\alpha - \mu_j)^2}{2(\sigma_j^2 + \sigma^2 t)} \right]}} \\ t=0 \Rightarrow \nabla_{\alpha} \log P(\alpha, 0) &= \frac{\sum_{j=1}^n \frac{w_j}{\sigma_j \sqrt{2\pi}} e^{\left[-\frac{(\alpha - \mu_j)^2}{2\sigma_j^2} \right]} \cdot \frac{-(\alpha - \mu_j)}{\sigma_j^2}}{\sum_{j=1}^n \frac{w_j}{\sigma_j \sqrt{2\pi}} e^{\left[-\frac{(\alpha - \mu_j)^2}{2\sigma_j^2} \right]}} \end{aligned}$$

ادامہ سوال ②

جنس 2. بہ ازای $t=0$ می بینیم کہ ترم t^2 در تابع Score از بین می رود و
تأثیرگذار نیست. اما بہ ازای t دلفواہ غیر صفر با t Scale ترم t^2
را در تابع Score دایم و تأثیرگذار است.

1. به طور مثال اگر روی شبکه تشخیص وضعیت سرطان کلیه کار می کنیم، می توانیم 100 نمونه استیاه کلاس بندی شده را تحلیل کنیم. مثلاً اگر تحلیل بصورت زیر باشد:

	benign	malignant	noisy	
1				
2	✓	✓		
3	✓	✓	✓	
⋮				
100		✓		
	7%	63%	6%	
				total dev error 15%
				benign ~ 1%
				malignant ~ 9.5%
				noisy ~ 0.9%

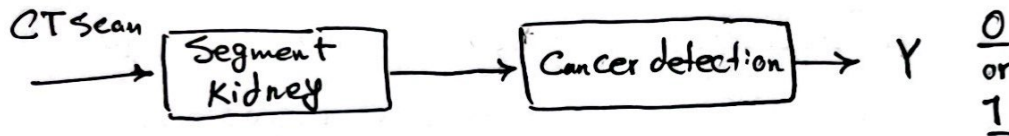
خب در اینصورت وقت گذاشتن برای تهیه داده های غیر نویزی یا denoise کردن آن ها با یک شبکه عصبی دیگر با توجه به اینکه 0.9% بهبود در error (تست dev خواهد داشت، بهبوده است.

بهتر است علت error زیاد روی کلاس malignant را بررسی کنیم و با افزایش داده های مربوط به این کلاس یا تکنیک در 1035، بهبود در خطای این کلاس و بهبود قابل توجه (حدود 9.5%) در خطای کل (تست dev) بوجود بیآوریم.

2. اگر (تست dev) مان به اندازه کافی بزرگ باشد می توانیم devset را به دو مجموعه eyeball devset و blackbox devset تقسیم کنیم. eyeball مجموعه کوچکی از devset است که اجازه داریم به آن با چشم نگاه کنیم و error analysis را روی آن انجام دهیم.

تقسیم devset به این دو مجموعه به این دلیل است که ماسچی می کنیم که مدل خودمان را روی eyeball، overfit کنیم و بعد از انجام error analysis و بهبود مدل اگر بخواهیم مدل را دوباره ارزیابی کنیم احتیاج به داده های labeled جدید یا مجموعه blackbox داریم. روی مجموعه blackbox در این زمان کار بردی است اما فقط اجازه داریم روی این مجموعه ارزیابی کرده، میزان خطا را محاسبه کنیم و با توجه به آن هایپر پارامترها را set کنیم ولی اجازه دیدن داده های این مجموعه را نداریم.

3. بطور مثال مانند بخش اول می‌خواهیم داده‌های CT Scan را به دو کلاس دارای توده کلیه و فاقد توده کلیه تقسیم بندی کنیم. برای این کار یک سیستم طراحی می‌کنیم که دارای دو بخش است. بخش اول محدوده کلیه را در تصاویر CT سکانت می‌کند. بخش دوم محدوده کلیه را گرفته و وجود توده را Classify می‌کند.



بروسه روش error analysis by part می‌توان هر خطا در خروجی (Y) را به یک بخش از سیستم که مقصر بوده، نسبت بدهیم.

بطور مثال برای یک داده خروجی 0 (فاقد سرطان) تشخیص داده شده در حالی که سرطان داشته است. باید برای این مورد بررسی کنیم خروجی شبکه Segmentation چه بوده است. اگر شبکه Segmentation، محدوده کلیه را درست شناسایی نکرده باشد و عضو دیگری را Segment کرده باشد پس شبکه Cancer detection درست کار کرده و خطا مربوط به بخش Segmentation است.

همینطور برای خطاهای مختلف سیستم این تحلیل را انجام می‌دهیم تا ببینیم کدام بخش از سیستم تأثیر بیشتری در خروجی اشتباه داشته است و سپس در اولویت بالاتر باید به تصحیح آن بخش از سیستم بپردازیم.

$$J_{naive}(\theta) = \frac{1}{2} \int d\mathbf{x} dt \, p(\mathbf{x}, t) \left[S_{\theta}(\mathbf{x}, t) - \nabla_{\mathbf{x}} \log p(\mathbf{x}, t) \right]^2$$

base on bayes theorem:

$$p(\mathbf{x}, t) = p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) p(\mathbf{x}^{(0)}, 0)$$

and also:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}, t) = \frac{\nabla_{\mathbf{x}} p(\mathbf{x}, t)}{p(\mathbf{x}, t)} = \frac{\nabla_{\mathbf{x}} [p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) \cdot p(\mathbf{x}^{(0)}, 0)]}{p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) p(\mathbf{x}^{(0)}, 0)}$$

$$\Rightarrow \nabla_{\mathbf{x}} \log p(\mathbf{x}, t) = \frac{[\nabla_{\mathbf{x}} p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0)] \cdot p(\mathbf{x}^{(0)}, 0)}{p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) p(\mathbf{x}^{(0)}, 0)}$$

$$\Rightarrow \nabla_{\mathbf{x}} \log p(\mathbf{x}, t) = \frac{\nabla_{\mathbf{x}} p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0)}{p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0)}$$

$$\Rightarrow \nabla_{\mathbf{x}} \log p(\mathbf{x}, t) = \nabla_{\mathbf{x}} \log p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) \left[\frac{\partial}{\partial \mathbf{x}} p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) \cdot p(\mathbf{x}^{(0)}, 0) + \left[\frac{\partial}{\partial \mathbf{x}} p(\mathbf{x}^{(0)}, 0) \right] \cdot p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) \right]$$

صفر

So:

$$J_{naive}(\theta) = \frac{1}{2} \int d\mathbf{x} dt \, p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) p(\mathbf{x}^{(0)}, 0) \left[S_{\theta}(\mathbf{x}, t) - \nabla_{\mathbf{x}} \log p(\mathbf{x}, t | \mathbf{x}^{(0)}, 0) \right]^2$$

$$\Rightarrow J_{naive}(\theta) = J_{mod}(\theta)$$

پس زمانی که دو تابع با هم برابر باشند،
global min یکسانی دارند.