

1.1.a

در مدل SimCLR از لاس contrastive استفاده می‌شود. طبق این objective زوج‌های مثبت باید به هم نزدیک شوند و زوج‌های منفی نیز باید از هم دور شوند. اهمیت در کردن زوج‌های منفی از هم این است که در صورت عدم وجود آن مدل بهینه‌های محلی بدی می‌تواند پیدا کند. مثلاً می‌تواند تمام ورودی‌ها را بزرگ تقسیم کند و همچنان لاس منفی بکند.

مدل MoCo نیز از همین لاس استفاده می‌کند و بنابراین در محاسبه لاس در سدن نمونه‌های منفی را دارد و این زوج‌های منفی مهم هستند. اما برای نگه‌داری این زوج‌های منفی از یک منف استفاده می‌کند تا مشکل batch size بزرگ را نداشته باشد. (نمونه‌های داخل منف را توسط شبکه‌ی momentum enc ایجاد می‌کند)

در مدل BYOL اما از لاس contrastive استفاده نمی‌شود و فقط نمونه‌های مثبت به هم نزدیک می‌شوند. BYOL از 2 شبکه جدا به نام‌های online و target استفاده می‌کند. شبکه‌ی online یکی از augment های ورودی را می‌گیرد و شبکه‌ی target یکی دیگر از augment ها را. لاس به این صورت تعریف می‌شود که خروجی دو شبکه نزدیک هم باشد یعنی به بیان دیگر شبکه online بتواند شبکه target را پیش‌بینی کند.

نزدیک آمدن شدن شبکه‌ی target توسط گزاردان نیست و moving avg شبکه online است که باعث می‌شود بهینه‌های محلی ذکر شده پیش‌بینی نیاید و نیازی به زوج‌های منفی نباشد.

1.1.b

آموزش مدل با batch size بزرگ وقتی از SGD/Momentum عادی استفاده شود unstable است. برای stability بیشتر از بهینه‌ساز LARS استفاده می‌شود. همچنین برای batch size های کوچک‌تر نشان داده شده که اسکال کردن learning rate به اندازه چند batch size کمک کننده است.

همچنین از global batch norm استفاده می‌شود. به این معنا که batch norm را به صورت global در سطح تمام device ها محاسبه می‌کنیم نه per device . این کار باعث می‌شود مدل برای یادگیری بازتابی‌ها نتواند shortcut بزند.

1.2.a

اگر پارامترهای مدل teacher را θ_t و پارامترهای مدل student را θ_s بنامیم، مدل teacher از طریق گزاردان آپدیت می‌شود و به صورت moving average از روی شبکه student به دست می‌آید یعنی $\theta_t = m\theta_t + (1-m)\theta_s$ که در آن m بسیار نزدیک 1 است.

تجسیری که می‌توان برای شبکه teacher ارائه داد آن است که می‌توان آن را به صورت ensemble چندین مدل دانست. در حقیقت این شبکه ensemble درزن‌های مختلفی از student است و از هر کدام آن‌ها بهتر است ← می‌تواند شبکه student را حلیت کند.

1.2.b

به ازای هر تصویر دو کراپ global و چندین کراپ local داریم. تنها global view ها به شبکه teacher ورودی داده می‌شوند و باعث ایجاد $\text{local to global correspondence}$ در آن می‌شوند. به این صورت شبکه teacher سعی دارد امپدینگ‌های سازگاری برای این کراپ‌ها به دست آورد و بنابراین به اشتراک‌های این دو کراپ global توجه بیشتری می‌کند و از آن‌ها که object اصلی در هر دو کراپ global وجود دارد شبکه نیز به آن توجه بیشتری می‌کند تا در امپدینگ لحاظ شود.

1.2.c

به طور کلی از روش های مختلفی برای جلوگیری از collapse استفاده می شود مانند BN, contrastive loss, اما در این روش استفاده از تکنیک های centering و sharpening بر روی خروجی های شبکه teacher برای جلوگیری از collapse کافی است. تکنیک centering مانع آن می شود که یک بُعد سایر ابعاد را dominate کند اما مدل را ترجیح می دهد که همه ابعاد خروجی آن شبیه uniform شوند. تکنیک sharpening نیز برعکس آن را انجام می دهد. بنابراین استفاده از هر دو می تواند اثراتشان را با هم ترکیب می کند و باعث می شود مدل collapse نکند. تکنیک centering معادل امانه کردن یک بایاس به شبکه teacher است و sharpening نیز با انتخاب معاداری کوچک برای γ در نرمالیزیشن softmax شبکه teacher انجام می شود.

1.2.d

efficient implementation

• fast and memory efficient attention

نسخه مدرنی از Flash Attention پیاده سازی شده که از قبل بهتر است و همچنین هاسیو پارامترهایی نیز برای استفاده بهینه عنوان شده است:

embedding dim = 1536 , number of heads = 24 vit-g backbone has 1.1B parameters

• sequence packing

یک تکنیک NLP است که در آن sequence های حاصل از همه کراپ های مختلف با هم concat می شوند و به مدل ورودی داده می شوند. این برخلاف DINO است که مجبور بود برای هر کراپ این کار را به صورت جدا انجام دهد. همچنین برای اینکه کراپ های مختلف روی هم تداخل نداشته باشند از ماسک های block diagonal استفاده می شود و بهبود سرعت قابل توجهی حاصل می شود.

• efficient stochastic depth

حساب dropped residuals که در نسخه قبلی داشتیم اسلایپ می شود به جای آنکه محاسبه شده و نتایج mask شوند. این کار حافظه و توان پردازشی را save می کند.

• fully sharded data parallel (FSDP)

پارامترهای 4 مدل teacher, student, optimizer first moments, optimizer second moments بر روی 4 GPU تقسیم شده است. در مدل قبلی روی یک دستگاه بودند که باعث می شد 16GB بر GPU نیاز باشد. این کار توسط پیاده سازی FSDP در pytorch انجام شده.

• Model distillation

به جای آموزش مدل های کوچکتر from scratch، آنها را از distill شده مدل بزرگ vit-g به دست می آوریم.

1.2.e

• در BYOL از لاس MSE استفاده می‌شود اما در DINO از cross entropy استفاده می‌شود.
• فضا در DINO قبل از اعمال لاس فرضی با sharpen, center می‌شود.

• در DINO از backbone های مبتنی بر Vit استفاده شده اما در BYOL از ResNet

• در DINO از تکنیک multi crop استفاده شده.

2.1.a

برای permutation equivariant بودن باید وزن همسایه ها به صورت shared باشد. این طوری مدل نمی تواند یاد بگیرد که به برخی همسایه ها وزن بیشتری نسبت دهد. در این بخش نیز وزن همسایه ها همگی w_2 است پس معتبر است

2.1.b

این حالت معتبر نیست زیرا وزن های متفاوتی برای همسایه ها در نظر گرفته شده

2.1.c

این حالت معتبر است زیرا وزن تمام همسایه ها w_2 است. ضمناً می دانیم که ماکسیمم گرفتن مولفه ای نیز ترتیبی ایجاد نمی کند زیرا خود max نسبت به ترتیب آرگومان ها invariant است

2.2

ابتدا فرم ماتریسی ارسال پیام را می نویسیم :

$$H^{(t+1)} = \sigma(A H^{(t)} W' + H^{(t)} W)$$

عبارت اول ایندیک همسایه ها را جمع می کند و عبارت دوم ایندیک خود خود. در حقیقت برای عبارت دوم داریم:

$$H W = \begin{bmatrix} h_1^T W \\ \vdots \\ h_n^T W \end{bmatrix}$$

یعنی ایندیک هر رد در ماتریس W ضرب می شود و در سطر ماتریس جدید قرار می گیرد

حال یک جایگشت روی رد های گراف اعمال می کنیم و آنرا به GNN ورودی می دهیم. برای این کار باید سطر و ستون های ماتریس مجاورت (A) جایگشت پیدا کنند و همچنین سطر های ماتریس $H^{(t)}$ نیز همان جایگشت را داشته باشند پس به جای A باید $P A P^T$ را قرار دهیم و به جای $H^{(t)}$ باید $P H^{(t)}$ را قرار دهیم:

$$\begin{aligned} & \sigma \left(P A P^T P H^{(t)} W' + P H^{(t)} W \right) \\ &= P \sigma \left(A H^{(t)} W' + H^{(t)} W \right) = P H^{(t+1)} \end{aligned}$$

می دانیم که ماتریس های permutation ، unitary هستند پس $P^T = P^{-1}$

بنابراین اثبات شد که اگر ماتریس را جایگشت دهیم در GNN ورودی بدیم معادل آن است که ماتریس اولیه را به GNN بدیم نتیجه این را جایگشت دهیم ← permutation equivariant

3.1

آشنایی فضاهای فرکانس نوع خاصی از آشنایی فضاهای image agnostic است. تفاوت آن با حالت عادی آن است که یک آشنایی کوچک به صورت ثابت به تمام ورودی ها اعمال می شود و نیازی نیست که برای هر ورودی، به صورت جدا آشنایی حساب شود. بنابراین این نوع آشنایی که magnitude کوچکی دارد و توسط چشم (انسان قابل رؤیت نیست به تعداد ورودی اعمال شده باعث می شود مدل دسته بندی، کلاسی استباهی را به عنوان کلاسی آن تصویر اعلام کند.

3.2

یافتن آشنایی ها از منظر تأمین امنیت کاربران مدل و به طور کلی robust کردن مدل بسیار حائز اهمیت است. به خصوص این نوع آشنایی به دلیل universal بودن آن می تواند بسیار خطرناک باشد و همان طور که در مقاله اشاره شده است به یازدهی وسیعی از مدل ها نیز apply می شود. فرض کنید یک خودروی خودران می خواهد تا بلوهای پیش روی خود را شناسایی کند. با این روش می توان یک آشنایی فراگیر و ثابت یافت و به هر تصویر ورودی مدل اضافه کرد. به این صورت مدل باعث ایجاد مداخلات قابل ملاحظه ای می شود. همچنین از دید تکنیکال نیز می توان گفت وجود چنین آشنایی هایی فرآیند نشان دهنده وجود correlation های مخفی می باشد در مرتبه تقسیم گیری دسته بندیها در ابعاد بالا است و اینکه جهت های در فضای ورودی وجود دارند که تنها با حرکت در آن رانندگی (صرف نظر از ورودی) می توان دسته بندی را فریب داد.

3.3

we want to find v that satisfies :

$$\|v\|_p \leq \epsilon$$

$$\mathbb{P}_{x \sim \mu}(\hat{k}(x+v) \neq \hat{k}(x)) \geq 1-\delta \rightarrow \text{fooling rate}$$

یعنی احتمال فریب خوراندن مدل k بر روی داده هایی که از توزیع μ می آیند از $1-\delta$ که آنرا تعیین می کنیم، بیشتر باشد. معیار به آنکه اندازه آشنایی نیز از ϵ کمتر باشد.

initialize $v \leftarrow 0$

while $\text{Err}(X_v) \leq 1-\delta$:

for $x_i \in X$:

if $\hat{k}(x_i+v) = \hat{k}(x_i)$:

$$\Delta v_i \leftarrow \arg \min_r \|r\|_2 \quad \text{s.t.} \quad \hat{k}(x_i+v+r) \neq \hat{k}(x_i)$$

$$v \leftarrow P_{p,\epsilon}(v + \Delta v_i)$$

return v

$$\begin{aligned} X_v &\triangleq \{x_1+v \dots x_m+v\} \\ \text{Err}(X_v) &\triangleq \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\hat{k}(x_i+v) \neq \hat{k}(x_i)} \\ P_{p,\epsilon}(v) &\triangleq \arg \min_{v'} \|v-v'\|_2 \quad \text{s.t.} \quad \|v'\|_p \leq \epsilon \end{aligned}$$

به این منظور الگوریتم زیر را خواهیم داشت:

به این معنی که روی دیتاست X حرکت می کنیم و به آرای هر x_i که مدل با ورودی x_i+v فریب نخورده بود، Δv_i ای بردست می آوریم که با اضافه کردن آن به v ، مدل قریب بخورد. برای اضافی قید $\|v\|_p \leq \epsilon$ بردار $v+\Delta v_i$ را بر روی دایره p به شعاع ϵ پروژکت می کنیم و v را با نتیجه آن به روز رسانی می کنیم.

⊗ تابع g که در صورت سوال به آن اشاره شده، حداقل احتمال موفقیت حمله را نشان می دهد و معادل $1-\delta$ می باشد. دیتاست را نیز به جای D ، X می گوییم.

4.1

در 3 مدل سعی دارند که multi modal learning انجام دهند و بتوانند امپدینگ های مشترک متن و تصویر ایجاد کنند.

CLIP سعی می کند یک multi modal joint embedding به دست بیاورد که در آن تصاویر و متن های مشابه به هم نزدیک تر

باشند. SimVLM نیز در بخش enc خود سعی دارد امپدینگ مناسبی برای متن و تصویر ورودی به دست بیاورد تا در بخش dec بتواند

با توجه به آن امپدینگ مشترک، متن تولید کند. Coca نیز توسط لاس contrastive loss امپدینگ های حاصل از img enc و txt dec را نزدیک نگه می دارد

در 3 مدل قابلیت Zero shot learning را دارند می توان به عنوان مثال Zero shot classification انجام داد.

CLIP و Coca که هر دو قابلیت dual enc دارند و در درس رابع به آن صحبت شد. SimVLM نیز می تواند به طریق مشابه تمایز ای برای تبدیل کلاس ها به متن در نظر بگیرد و تصویر و آن متن را به بخش enc ورودی دهد. حال می توان از متن generate کرد، تعلق به آن کلاس را فهمید.

تفاوت ها

این 3 مدل در pretraining objective با هم تفاوت دارند. CLIP با contrastive pretraining سعی می کند pair های متن و تصویر مشابه را به هم نزدیک و از بقیه دور کند. SimVLM از prefix language model objective استفاده می کند که سعی می کند با توجه به امپدینگ های متن و تصویر ورودی، متن تولید کند. Coca اما این دو ایده را ترکیب می کند و علاوه بر اینکه با contrastive loss امپدینگ های متن و تصویر را align می کند، متن نیز برای آن تولید می کند.

از نظر کاربرد نیز این 3 مدل قادر به حل کردن تسک های متفاوتی هستند

	CLIP	SimVLM	Coca
txt to img retrieval	✓		✓
img to txt retrieval	✓		✓
img captioning		✓	✓
VQA		✓	✓

4.2.a

در بخش txt enc از transformer استفاده شده است. به طور دقیق تر می توان گفت که ما مدل ترنسفورمر آن

12 لایه با عرض 512 دارد و شامل 8 تا attention head است.

برای برداشت آوردن امپدینگ متن ابتدا از ترنسفورمر استفاده می شود و سپس یک لایه projection خطی روی آن اعمال می شود تا به فضای امپدینگ مشترک برسیم.

در بخش img enc از دو معماری استفاده می شود. یکی ResNet 50 که کمی تغییر یافته و دیگری VIT. بسته به کاربرد می توان از هر کدام استفاده کرد. در این بخش نیز پس از img enc یک لایه projection خطی روی آن اعمال می شود تا به فضای امپدینگ مشترک برسیم.

در نهایت در فضای مشترک لاس contrastive اعمال می شود

4.2.b

$$\frac{\partial s_{ij}}{\partial n_i} = \frac{\partial_j \|n_i\| \|d_j\| - \|d_j\| \frac{n_i}{\|n_i\|} (n_i \cdot d_j)}{\|n_i\|^2 \|d_j\|^2} = \frac{\|n_i\|^2 d_j - (n_i \cdot d_j) n_i}{\|n_i\|^3 \|d_j\|}$$

$$\frac{\partial L_1}{\partial n_i} = -\frac{1}{N} \frac{\sum_j e^{s_{ij}/\tau}}{e^{s_{ii}/\tau}} \left[\left(\frac{1}{\tau} e^{s_{ii}/\tau} \frac{\partial s_{ii}}{\partial n_i} \right) \sum_j e^{s_{ij}/\tau} - \left(\sum_j \frac{1}{\tau} e^{s_{ij}/\tau} \frac{\partial s_{ij}}{\partial n_i} \right) e^{s_{ii}/\tau} \right]$$

$$\times \frac{1}{\left(\sum_j e^{s_{ij}/\tau} \right)^2}$$

$$= \frac{(n_i \cdot d_i) n_i - \|n_i\|^2 d_i}{N \tau \|n_i\|^3 \|d_i\|} + \frac{\sum_j \left(e^{s_{ij}/\tau} \frac{\|n_i\|^2 d_j - (n_i \cdot d_j) n_i}{\|n_i\|^3 \|d_j\|} \right)}{N \tau \sum_j e^{s_{ij}/\tau}}$$

در برخی سوال ها با علی بابایی همکاری داشتم