



یادگیری ماشین

تمرین سری ششم

عمادالدین رستمیان

403206479

پاسخ سوال ۱

در یادگیری خودنظراتی، وظیفه پیش‌متن یک وظیفه کمکی است تا به مدل کمک کند که ویژگی‌های موثری را از داده‌های بدون برچسب یاد بگیرد. در واقع هدف این است که مدل بتواند با استفاده از برچسب‌های تولید شده به صورت خودکار (بدون نیاز به دخالت انسان) این وظیفه را انجام دهد و در نتیجه آن، مدل representation‌هایی را یاد بگیرد که در انجام وظایف بعدی نیز عملکرد خوبی داشته باشد. وظایف پیش‌متن معمولاً ساده هستند، اما مدل را مجبور می‌کنند که الگوها یا ساختارهای معنی‌داری را از داده‌های ورودی استخراج کند.

الف) پیش‌بینی چرخش: در این وظیفه پیش‌متن تصاویر به صورت رندوم و با زوایای مشخص (مانند ۹۰، ۱۸۰ یا ۲۷۰ درجه) چرخیده می‌شوند و مدل باید زاویه چرخش را برای هر تصویر پیش‌بینی کند. آموزش مدل برای این وظیفه باعث می‌شود تا مدل جهت‌گیری‌های فضایی و ساختار کلی اجسام درون تصاویر را یاد بگیرد.

ب) رنگ‌آمیزی: با تغییر رنگ پیکسل‌های تصویر به رنگ‌های دیگر یا به grayscale مدل برای پیش‌بینی درست رنگ پیکسل‌ها یا نواحی مختلف تصویر آموزش می‌دهیم. این وظیفه پیش‌متن مدل را مجبور به یادگیری توزیع رنگ‌های اجسام، روابط معنایی بین آن‌ها و بافت‌های مختلف تصویر می‌کند.

ج) حل پازل: در این تسک تصویر به patch‌های مختلف تقسیم می‌شود. این پچ‌ها به صورت رندوم shuffle می‌شوند و مدل برای بازسازی تصویر اصلی آموزش داده می‌شود. ویژگی‌های استخراجی از این تسک شامل روابط فضایی و معنایی در تصاویر است.

پاسخ سوال ۲

حل پازل

الف) تصاویر ماهواره ای قرارگیری های مشخص و پیچیده ای را در بردارند مانند خیابان ها، ساختمان ها، مزارع و غیره. با آموزش مدل بر روی وظیفه پیش‌من حل پازل، می‌توانیم ویژگی های فضایی و ارتباط بین بخش های مختلف این تصاویر را استخراج کنیم که برای وظایف بعدی بسیار مفید واقع می‌شوند.

ب) ابتدا تصاویر ماهواره ای را به پنج های کوچک مربعی تقسیم می‌کنیم. سپس آن ها را شافل کرده و به مدل می‌دهیم تا بتواند بر اساس این پنج ها تصویر اصلی را بازسازی کند.

ج) با چرخش تصاویر ماهواره ای معمولاً تغییر معنی داری می‌سر نمی‌شود و در نتیجه مدل ویژگی های بخصوصی را یاد نمی‌گیرد. از طرفی تصاویر ماهواره ای معمولاً grayscale هستند و حتی اگر رنگی نیز باشند ویژگی های مفیدی را که می‌توان از حل پازل استخراج کرد را حاصل نمی‌کنند.

پاسخ سوال ۳

(الف)

$$N = \frac{224 \times 224}{16 \times 16} = 196$$
$$16 \times 16 \xrightarrow{\text{flatten}} 256$$
$$\text{embedding} = W \times p_{\text{flat}} + b$$
$$W \in \mathbb{R}^{128 \times 256}$$
$$b \in \mathbb{R}^{128}$$

$$\text{embedding}[i] \leftarrow \text{embedding}[i] + \text{PE}[I], i = 1, 2, \dots, 196 \quad (\text{ب})$$

بدون جاسازی موقعیتی مدل نمی‌تواند اطلاعاتی در مورد روابط فضایی بین پچ‌ها را یاد بگیرد. با افزودن این جاسازی به مدل این قابلیت را میدهیم که بتواند روابط نسبی یا مطلق بین پچ‌ها را یاد بگیرد تا بتواند ویژگی‌های فضایی مهم را استخراج کند تا در تسك‌های مختلف پردازش تصویر مفید واقع شود.

(ج) با افزودن توکن [CLS] به ابتدای دنباله، تعداد توکن‌ها به ۱۹۷ میرسد که بعد هر توکن ۱۲۸ می‌باشد. در ترنسفورمر تمامی توکن‌ها با توکن‌های دیگر توسط مکانیز self-attention تعامل می‌کنند. تعامل توکن [CLS] با سایر توکن‌ها، باعث می‌شود تا این توکن تمامی اطلاعات مربوط به همه توکن‌ها را جمع آوری کند. این جمع آوری باعث می‌شود تا مدل یک دید کلی یا گلوبال به پچ‌ها داشته باشد، گویا کل تصویر اصلی را دیده است و ویژگی‌های آن را استخراج کرده است.

پاسخ سوال ۴

الف) ابتدا تصویر از انکودر ViT میگذرد تا یک embedding از تصویر بدست آید. از طرف دیگر متن مورد نظر از یک انکودر ترنسفورمر عبور میکند تا آن نیز بدست آید. سپس با استفاده از cosine similarity تشابه این دو embedding محاسبه میشود. مقدار حاصل نشان میدهد که چه میزان تصویر با متن همخوانی دارد.

برای تصویر «سیب قرمز» متن «یک سیب قرمز آبدار روی میز» بیشترین امتیاز را خواهد داشت. کلمات «سیب»، «قرمز» به صورت مستقیم در متن آمده اند که باعث نزدیکی بردار embedding تصویر و متن خواهد بود.

ب) این رفتار نشان میدهد که مدل به رنگ و شکل اجسام وزن بیشتری در امبدینگ میدهد. زیرا با وجود اینکه کلمه «سیب» به طور مستقیم در متن دوم آورده شده، مدل به سرخ بودن و گرد بودن اولویت بیشتری داده. همچنین نشان میدهد که واژه «سیب» نزدیکی معنایی بین تصویر و متن را کاهش میدهد که نشان دهنده بایاس مدل به رنگ، نسبت به روابط معنایی است.

پاسخ سوال ۵

:Attention Pooling (الف)

- میانگین‌گیری از وزن‌های attention برای محاسبه feature map
- وزن‌های attention به صورت پویا یادگرفته می‌شوند و به نواحی مختلف feature map متصل می‌شوند که باعث می‌شود مدل بخش‌های مهم تر تصویر را یاد بگیرد.
- خروجی context-aware است و از جمع میانگین فیچر مپ‌ها بدست می‌آید.
- میتواند برای تسک‌های مختلف خروجی‌های متفاوتی بدهد.

:Global Average Pooling

- میانگین هر فیچر مپ را بر روی تمامی ابعاد محاسبه می‌کند. تمامی اطلاعات فضایی برای هر ویژگی در یک مقدار خلاصه می‌شود.
- خروجی یک میانگین یکنواخت از تمامی ویژگی‌های است.
- از دست رفتن اطلاعات را به همراه دارد که از توانایی مدل برای تمرکز بر روی نواحی مهم می‌کاهد.

(ب) ماتریس لیبل یک ماتریس $N \times N$ با درایه‌های ۰ (نشان دهنده جفت‌های منفی) و ۱ (نشان دهنده جفت‌های مثبت) است. اگر فرض کنیم که تصاویر لیبل مشترک ندارند، تنها درایه‌های قطر اصلی مقادیر ۱ خواهند داشت. پس ماتریس N تا 1 و $N^2 - N$ صفر دارد.

(ج) بر اساس مقاله، عملکرد zero-shot مدل کلیپ در تسک‌هایی ضعیف است که :

۱. نیاز به دانش تخصصی دارند مانند تصاویر پزشکی، تصاویر ماهواره‌ای و ...
۲. نیاز به نتیجه‌گیری و استنتاج‌های پیچیده و انتزاعی دارند مانند شمردن اشیا، تشخیص ویژگی‌های جزئی و ...
۳. نیاز به اطلاعات سطح بالاتری دارند مانند علائم رانندگی

دلیل ضعف کلیپ در موارد ذکر شده عبارتند از:

- ۱- نیاز به تعداد دیتای بالا و متنوع
- ۲- تکیه مدل صرفا بر الگوهای بصری مانند رنگ، شکل، بافت
- ۳- ناتوانی زبان طبیعی در توصیف برخی تسک‌ها