



دانشگاه صنعتی امیرکبیر
(پلی‌تکنیک تهران)
دانشکده ریاضی و علوم کامپیوتر

پروژه پنجم درس داده کاوی محاسباتی
علوم کامپیوتر، هوش مصنوعی و محاسبات نرم

فشرده سازی مدل های عمیق با تجزیه رتبه پایین و ارتقای کیفیت صنعتی

نگارش
هومن ذوالفقاری

استاد راهنما
دکتر مهدی قطعی

دی 1404



به نام خدا

تاریخ:

تعهدنامه اصالت اثر

اینجانب هومن ذوالفقاری متعهد می‌شوم که مطالب مندرج در این پایان نامه حاصل کار پژوهشی اینجانب تحت نظارت و راهنمایی اساتید دانشگاه صنعتی امیرکبیر بوده و به دستاوردهای دیگران که در این پژوهش از آنها استفاده شده است مطابق مقررات و روال متعارف ارجاع و در فهرست منابع و مآخذ ذکر گردیده است. این پایان نامه قبلاً برای احراز هیچ مدرک هم‌سطح یا بالاتر ارائه نگردیده است.

در صورت اثبات تخلف در هر زمان، مدرک تحصیلی صادر شده توسط دانشگاه از درجه اعتبار ساقط بوده و دانشگاه حق پیگیری قانونی خواهد داشت.

کلیه نتایج و حقوق حاصل از این پایان نامه متعلق به دانشگاه صنعتی امیرکبیر می‌باشد. هرگونه استفاده از نتایج علمی و عملی، واگذاری اطلاعات به دیگران یا چاپ و تکثیر، نسخه‌برداری، ترجمه و اقتباس از این پایان نامه بدون موافقت کتبی دانشگاه صنعتی امیرکبیر ممنوع است.

نقل مطالب با ذکر مآخذ بلامانع است.

در صفحه تعهدنامه اصالت اثر، در قسمت بالا سمت چپ، تاریخ دفاع خود را جایگزین تاریخ نوشته شده کنید.

همچنین در صفحه تعهدنامه اصالت اثر، در خط اول، نام و نام خانوادگی خود را به صورت کامل با نام و نام خانوادگی نمونه، جایگزین کنید. در انتهای متن تعهد، در قسمت امضا نیز باید نام و نام خانوادگی کامل خود را وارد نمایید.

هومن ذوالفقاری

امضا

چکیده

مدل‌های یادگیری عمیق علی‌رغم دقت بالا، معمولاً دارای تعداد زیادی پارامتر و هزینه محاسباتی قابل توجهی هستند که استفاده از آن‌ها را در محیط‌های عملیاتی و صنعتی محدود می‌کند. در این پروژه، از تکنیک‌های پیشرفته جبر خطی، به‌ویژه تجزیه مقادیر منفرد (SVD)، برای فشرده‌سازی لایه‌های یک شبکه عصبی عمیق استفاده شده است. بدون تغییر در معماری کلی مدل، لایه‌های حجیم به تقریب‌های رتبه پایین تبدیل شده‌اند. نتایج نشان می‌دهد که با کاهش قابل توجه تعداد پارامترها، افت دقت ناچیز بوده و با استفاده از تنظیم دقیق (Fine-tuning)، حتی بهبود دقت نیز حاصل می‌شود. همچنین، تأثیر فشرده‌سازی بر سرعت استنتاج در محیط CPU بررسی شده است.

Contents

| | |
|----|--|
| أ | چکیده |
| 1 | فصل اول مقدمه |
| 3 | فصل دوم تحلیل طیفی وزن های مدل |
| 6 | فصل سوم پیاده سازی عملی فشرده سازی |
| 8 | فصل چهارم ارزیابی دقت و سرعت |
| 10 | فصل پنجم سوالات تحلیلی |
| 13 | فصل ششم نتیجه گیری |
| 15 | منابع و مراجع |

فصل اول

مقدمه

مقدمه

لینک کد در گیت‌هاب:

<https://github.com/hoomanzolfaghari84/Data-Mining-Course.git>

با گسترش کاربردهای صنعتی یادگیری عمیق، نیاز به مدل‌هایی با مصرف حافظه کمتر و سرعت استنتاج بالاتر بیش از پیش احساس می‌شود. بسیاری از شبکه‌های عصبی عمیق دارای افزونگی پارامتری هستند و وزن‌های آن‌ها را می‌توان با تقریب‌های کم‌رتبه نمایش داد. تجزیه مقادیر منفرد (SVD) یکی از ابزارهای قدرتمند جبر خطی است که امکان تحلیل طیفی وزن‌ها و فشرده‌سازی مؤثر لایه‌های خطی را فراهم می‌کند.

هدف این پروژه، بررسی عملی امکان فشرده‌سازی یک مدل از پیش آموزش‌دیده بدون تغییر معماری کلی آن، و ارزیابی اثر این فشرده‌سازی بر دقت، تعداد پارامترها و سرعت استنتاج است.

فصل دوم

تحلیل طیفی وزن های مدل

تحلیل طیفی وزن های مدل

مدل و داده

در این پروژه، یک شبکه عصبی چندلایه (MLP) آموزش دیده بر روی دیتاست Fashion-MNIST مورد استفاده قرار گرفت. مدل شامل 4 لایه خطی اصلی است که یکی از لایه های میانی (512→512) و لایه نهایی (512→10) برای تحلیل انتخاب شدند.

تحلیل مقادیر منفرد

با اعمال تجزیه SVD روی ماتریس وزن ها، طیف مقادیر منفرد و واریانس تجمعی آن ها محاسبه شد.

نتایج به دست آمده:

• لایه میانی:

○ رتبه لازم برای حفظ ۹۵٪ واریانس : 275

○ درصد مقادیر منفرد استفاده شده: 53.71%

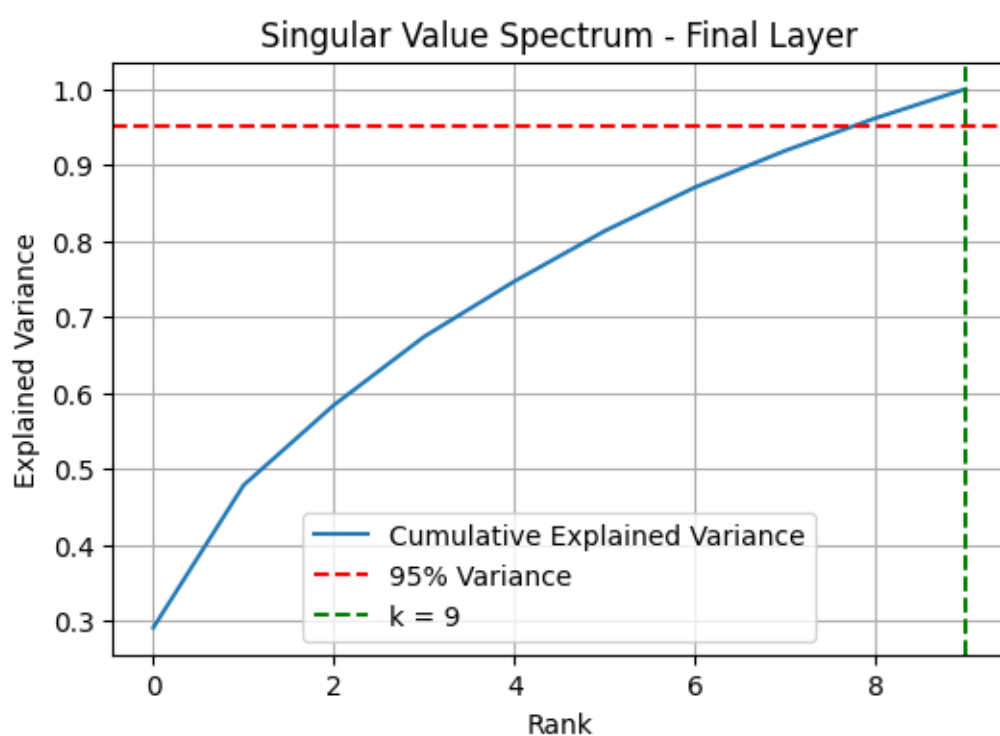
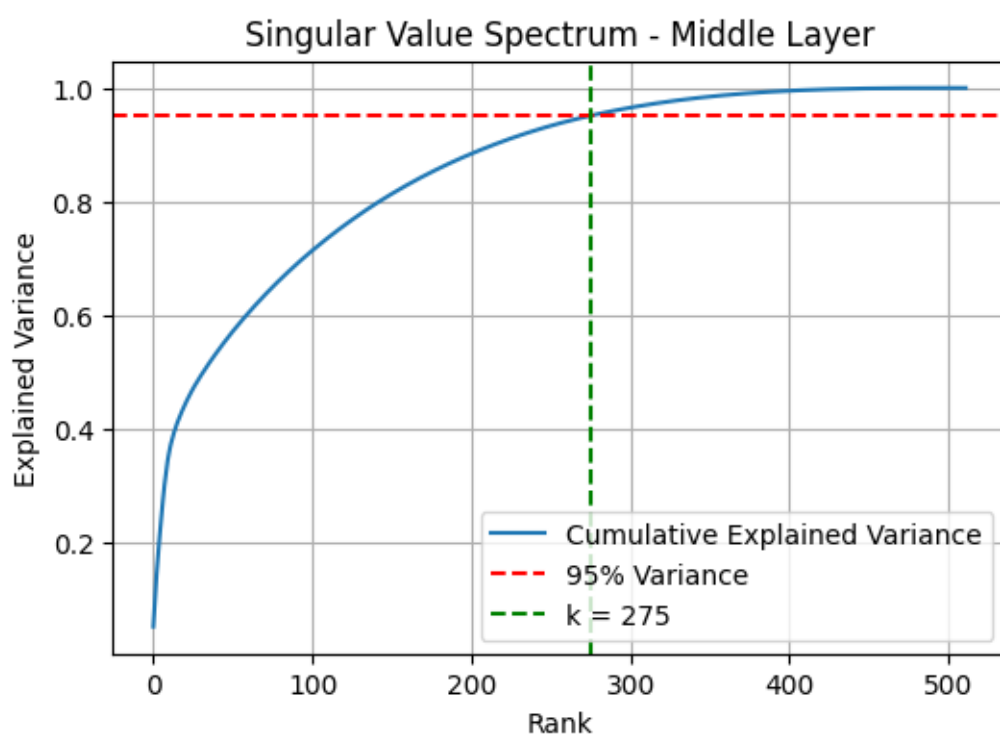
• لایه نهایی:

○ رتبه لازم برای حفظ ۹۵٪ واریانس : 9

○ درصد مقادیر منفرد استفاده شده: 90.00%

تحلیل

این نتایج نشان می دهند که لایه های میانی دارای افزونگی ساختاری بالاتری هستند و گزینه مناسبتری برای فشرده سازی محسوب می شوند، در حالی که لایه نهایی به دلیل ابعاد کوچکتر، ظرفیت فشرده سازی محدودی دارد.



فصل سوم

پیاده سازی عملی فشرده سازی

پیاده سازی عملی فشرده سازی

روش فشرده سازی

لایه خطی انتخاب شده با استفاده از SVD به دو لایه متوالی با رتبه k جایگزین شد. این کار بدون تغییر عملکرد ریاضی کلی لایه و تنها با یک تقریب کم رتبه انجام گرفت.

سناریوهای آزمایشی

مدل در سه حالت بررسی شد:

- بدون فشرده سازی
- فشرده سازی با ۵۰٪ رتبه
- فشرده سازی با ۸۰٪ رتبه

دقت اولیه مدل

| دقت (%) | حالت مدل |
|---------|-----------------|
| 86.92 | بدون فشرده سازی |
| 86.93 | فشرده سازی ۵۰٪ |
| 86.98 | فشرده سازی ۸۰٪ |

نتایج نشان می دهد که حتی بدون Fine-tuning ، فشرده سازی منجر به افت دقت قابل توجهی نشده است. (در اینجا بهبود یافته، گویا مدل اندکی اورفیت کرده است)

مقایسه تعداد پارامترها

| تعداد پارامتر | حالت مدل |
|---------------|-----------------|
| 932362 | بدون فشرده سازی |
| 932362 | فشرده سازی ۵۰٪ |
| 774666 | فشرده سازی ۸۰٪ |

کاهش محسوس تعداد پارامترها، کارایی حافظه مدل را بهبود می دهد.

فصل چهارم

ارزیابی دقت و سرعت

ارزیابی دقت و سرعت

تنظیم دقیق (Fine-tuning)

مدل‌های فشرده‌شده به مدت ۳ اپیاک Fine-tune شدند. پس از تنظیم دقیق، دقت‌ها به شکل زیر بهبود یافت:

| دقت بعد از FT (%) | حالت مدل |
|-------------------|-----------------|
| 86.91 | بدون فشرده‌سازی |
| 88.75 | فشرده‌سازی ۵۰٪ |
| 88.66 | فشرده‌سازی ۸۰٪ |

این افزایش دقت نشان می‌دهد که فشرده‌سازی می‌تواند به‌عنوان نوعی منظم سازی (Regularization) عمل کند.

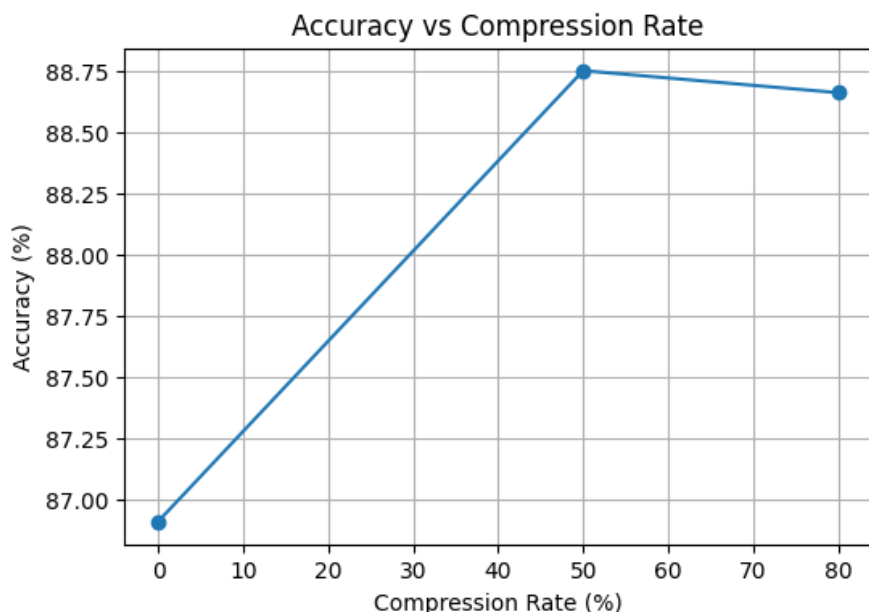
سرعت استنتاج روی CPU

| زمان استنتاج (ms) | حالت مدل |
|-------------------|-----------------|
| 0.29 | بدون فشرده‌سازی |
| 0.30 | فشرده‌سازی ۵۰٪ |
| 0.26 | فشرده‌سازی ۸۰٪ |

تفاوت زمان‌ها ناچیز است که نشان می‌دهد افزایش تعداد لایه‌ها اثر منفی قابل توجهی بر سرعت CPU در این مقیاس نداشته است.

نمودار دقت بر حسب نرخ فشرده‌سازی

نمودار رسم‌شده نشان می‌دهد که تا نرخ‌های بالای فشرده‌سازی، دقت مدل حفظ شده و حتی پس از Fine-tuning افزایش یافته است.



فصل پنجم

سوالات تحلیلی

سوالات تحلیلی

تأثیر تأخیر ناشی از افزایش تعداد لایه‌ها

در محیط‌هایی مانند GPU، افزایش تعداد لایه‌های کوچک ممکن است به دلیل parallelism تأثیر کمی داشته باشد، اما در CPU و سیستم‌های embedded، تأخیر ناشی از فراخوانی لایه‌ها می‌تواند غالب شود. بنابراین فشردن فشرده‌سازی همیشه به افزایش سرعت منجر نمی‌شود و به معماری سخت‌افزار وابسته است.

توضیح تأثیر افزایش محیط تعداد لایه‌ها

هر عملیات محاسباتی پشت سر هم انجام می‌شود. افزایش بسیار محسوس CPU تک‌هسته‌ای لایه‌ها با زمان استنتاج (تقریباً خطی) رشد می‌کند

پردازش موازی روی هزاران هسته انجام می‌شود، ولی همچنان لایه‌های بزرگ (مثل Conv با فیلترهای زیاد) باعث کمتر محسوس GPU تأخیر می‌شوند.

حتی چند لایه اضافی به دلیل محدودیت حافظه و توان پردازشی می‌تواند زمان استنتاج را به طور قابل توجهی شدید افزایش دهد.

کاهش رتبه یا فشردن فشرده‌سازی لایه‌ها می‌تواند تأخیر را کاهش دهد، چون تعداد پارامتر و FLOPs کم می‌شود.

نقش عدد وضعیت در انتخاب لایه‌های مستعد فشرده‌سازی

عدد وضعیت (Condition Number) معیاری از حساسیت ماتریس نسبت به اغتشاشات است. لایه‌هایی با عدد وضعیت بزرگ و افت سریع مقادیر منفرد، معمولاً افزونگی بالایی دارند و گزینه‌های مناسبی برای فشردن فشرده‌سازی کم‌رتبه محسوب می‌شوند.

عدد وضعیت (Condition Number) برای یک ماتریس W به صورت زیر تعریف می‌شود:

$$\kappa(W) = \frac{\sigma_{\max}}{\sigma_{\min}}$$

σ_{\max} و σ_{\min} بزرگ‌ترین و کوچک‌ترین مقادیر منفرد هستند. عدد وضعیت نشان می‌دهد که ماتریس تا چه حد «نزدیک به تنزل رتبه» است و چقدر حساس به نویز/خطا است.

لایه با عدد وضعیت کوچک:

- مقادیر منفرد بزرگ و کوچک به صورت یکنواخت پخش شده‌اند.
- تقریب رتبه پایین باعث از دست رفتن اطلاعات مهم می‌شود → کمتر مناسب فشردسازی.

لایه با عدد وضعیت بزرگ:

- ماتریس دارای چند مقدار منفرد بزرگ و بقیه نزدیک صفر است.
- این لایه مستعد فشردسازی رتبه پایین است، چون می‌توان بیشتر مقادیر کوچک را حذف کرد بدون کاهش دقت مدل.

فصل ششم نتیجه‌گیری

نتیجه‌گیری

در این پروژه نشان داده شد که با استفاده از تجزیه رتبه پایین مبتنی بر SVD می‌توان مدل‌های یادگیری عمیق را بدون تغییر معماری کلی، به‌طور مؤثری فشرده‌سازی کرد. نتایج تجربی حاکی از کاهش تعداد پارامترها، حفظ یا حتی بهبود دقت پس از Fine-tuning ، و عدم افزایش محسوس زمان استنتاج هستند. این رویکرد، راهکاری عملی و صنعتی برای بهینه‌سازی مدل‌های عمیق در محیط‌های با محدودیت منابع محسوب می‌شود.

منابع و مراجع

- [1] این پروژه نیاز کمی به منابع داشت و بیشتر از دانش خودم استفاده کردم
- [2] <https://docs.pytorch.org/>
- [3] https://en.wikipedia.org/wiki/Low-rank_approximation
- [4] https://openaccess.thecvf.com/content_CVPR_2020/papers/Idelbayev_Low-Rank_Compression_of_Neural_Nets_Learning_the_Rank_of_Each_CVPR_2020_paper.pdf
- [5] منابع اصلی درس



**Amirkabir University of Technology
(Tehran Polytechnic)**

... Department ...

MSc or PhD Thesis

Title of Thesis

**By
Name**

**Supervisor
Dr.**

**Advisor
Dr.**

Month & Year