



سامانه‌های یادگیری ماشین توزیع شده (پاییز ۱۴۰۴)

تمرین کامپیوتري ۲

مهلت ارسال: ۱۴۰۴ / ۰۹ / ۰۵

استاد درس: دکتر دوستی

دستیار طراح: علی رمضانی

قوانين و ملاحظات

نحوه ارسال تمرین:

- تمامی فایل‌ها باید در یک فایل فشرده با نام DMLS-CA2-StudentID ارسال شوند.
- کدهای مربوط به هر بخش را بر اساس جدول انتهای همین فایل ذخیره کرده و همراه گزارش ارسال کنید.
- تمامی کدهای ارسال شده باید امکان اجرای مجدد داشته باشند. اگر تنظیمات خاصی برای اجرا نیاز است، آن را ذکر کنید.
- کدهای ارسال شده باید توسط خودتان اجرا شده باشند و نتایج اجرا در فایل ارسالی مشخص باشد.

رعایت اصول آکادمیک و صداقت علمی:

- این تمرین باید به صورت **فردی** انجام شود. هرگونه همکاری یا نوشتن تمرین به صورت گروهی ممنوع است.
- در صورت مشاهده تشابه در پاسخ، تمامی افراد درگیر نمره صفر دریافت خواهند کرد و موضوع به استاد گزارش خواهد شد.

استفاده از ابزارهای هوش مصنوعی:

استفاده از ابزارهایی مانند Copilot، Gemini، ChatGPT و موارد مشابه مجاز است، اما تحت شرایط زیر:

- نحوه استفاده از این ابزارها را در گزارش خود توضیح دهید (ابزارهای استفاده شده، کاربردهای مشخص و موارد مرتبط).
- تمامی پرامپت‌ها و لینک‌های استفاده شده را در انتهای گزارش قرار دهید.
- عدم ارائه این اطلاعات به منزله سرقت علمی محسوب شده و منجر به نمره صفر خواهد شد.

مهلت ارسال و جریمه تأخیر:

- امکان ارسال تمرین **با تأخیر تا ۲ روز** وجود دارد. برای **هر روز تأخیر ۱۰ درصد جریمه** اعمال می‌شود.
- تأخير به صورت ساعتی محاسبه شده و پس از دو روز تأخیر، تمرین پذیرفته نخواهد شد.

ارزیابی حضوری:

- ارزیابی تمرین به صورت حضوری انجام خواهد شد.
- محل ارزیابی: آزمایشگاه NLP، طبقه منفی یک، دانشکده مهندسی برق و کامپیوتر ساختمان شماره ۲.

لطفاً پیش از شروع کار بر روی تمرین، به نکات زیر توجه فرمایید.

- برای سؤال اول می‌توانید از Kaggle یا Google Colab نیز استفاده کنید و نیازی به استفاده از سرور ندارید.
- برای سؤال دوم حتماً باید از سرور سبلان استفاده کنید.
- ویدئوهای PyTorch DDP و CUDA Programming را قبل از انجام تمرین مشاهده نمایید.
- برای راحتی در توسعه و تست کد، از ماشین مجازی لینوکس خود استفاده نمایید تا ترافیک کلاستر (بهخصوص در ساعت آخر مهلت تمرین) افزایش نیابد. پس از اطمینان از عملکرد کد، می‌توانید آن را روی کلاستر اجرا کنید.
- سوالات خود را در گروه تلگرام درس مطرح نمایید. به هیچ وجه کد یا پاسخ سوالات را در گروه به اشتراک نگذارید.
- هدف از انجام این تمرین، آشنایی بیشتر با CUDA Programming و PyTorch DDP و بررسی پارامترهای مختلف بر تغییر سرعت و حافظه مصرفی در آموزش مدل‌های شبکه‌ی عصبی است.
- در سؤال اول این تمرین با دیتاست MNIST و در سؤال دوم با دیتاست STL-10 کار می‌کنیم. دیتاست دوم در محل /storage/dmls/stl10_data در سرور ذخیره شده است.
- برای در دسترس بودن کتابخانه‌های لازم، لازم است ابتدا محیط مجازی زیر را فعال کنید:

```
source /home/dmls/pytorch/bin/activate
```

- هر جایی که از شما خواسته شده، باید کدها را در گزارش نیز توضیح دهید.
- می‌توانید از طریق ایمیل زیر در ارتباط باشید:

ali.ramezani.96@ut.ac.ir

- برای انجام این تمرین به سرور زیر متصل شوید:

Sabalani 1: 172.18.32.194

۵۰ نمره

سؤال ۱. آشنایی با CUDA Kernel و پیاده‌سازی Kernel Fused ReLU + Convolution

در این تمرین قرار است با مفاهیم پایه و کاربردی CUDA Kernel و نقش آن در تسریع عملیات یادگیری عمیق آشنا شوید. هدف، نوشتمن یک کرنل برای ترکیب دو عملیات Convolution و ReLU است، ثبت این کرنل به عنوان یک لایه در PyTorch، و بررسی عملی سرعت و دقیقت آن نسبت به پیاده‌سازی‌های استاندارد است.
اهداف این تمرین:

- توانایی نوشتمن و درک عملکرد CUDA Kernel اختصاصی
- امکان استفاده از این کرنل‌ها در شبکه‌های عصبی سفارشی
- بنچمارک و تحلیل سرعت و کارایی کرنل اختصاصی در مقایسه با cuDNN و PyTorch

- تشخیص اینکه چه زمانی پیاده‌سازی کرنل اختصاصی منطقی و مؤثر است

شما برای این تمرین یک ویدئوی آموزشی و یک نوتبوک راهنمایی دریافت کرده‌اید؛ بسیاری از بخش‌های کدنویسی تکمیل شده و تنها بخش‌هایی کلیدی را باید کامل کنید تا بتوانید روی درک و تحلیل تمرکز بیشتری داشته باشید.

۱.۱ مفاهیم اولیه (۵ نمره)

۱. CUDA Kernel چیست و چه کاری انجام می‌دهد؟

۲. cuDNN چیست؟ چرا PyTorch و سایر فریمورک‌ها از این کتابخانه استفاده می‌کنند؟ مزیت‌های cuDNN چیست و چه مشکلی را حل می‌کند؟ نسبت cuBLAS و cuDNN چطور است؟

۲.۱ پیاده‌سازی کرنل Fused Convolution + ReLU (۱۵ نمره)

۱. یک کرنل پایتونی بنویسید که دو عملیات ReLU و Convolution را همزمان برای یک تصویر تک کاناله انجام دهد. (کافیست توابع conv_relu_py و conv_relu_kernel_py را کامل کنید).

ویژگی‌های کانولوشن مد نظر:

```
C_in = 1, C_out = 1
kernel_size = 3x3
padding = 1, stride = 1
x: [N, 1, H, W]
w: [1, 1, 3, 3]
b: [1]
out: [N, 1, H, W]
```

۲. کرنل بالا را به زبان CUDA/C++ بنویسید. قسمت‌های خالی رشتہ‌ی cuda_src را کامل کنید. دقت کنید ایندکس‌دهی، Padding و سایر جزئیات باید دقیقاً مطابق نسخه‌ی پایتونی باشد.

۳.۱ ثبت کرنل جدید به عنوان یک لایه در پایتورچ (۵ نمره)

در این بخش ما این کرنل را به عنوان یک لایه جدید در پایتورچ ثبت کردیم. شما کافی است تابع کلاس FusedConvReLUFn را کامل توضیح دهید.

۴.۱ بنچمارک کردن لایه جدید ایجاد شده (۵ نمره)

با استفاده از داده‌ی MNIST، دو شبکه کانولوشنی با ۲ و ۱۰ لایه convolution تعریف کنید و در هر کدام سناریوهای زیر را تست کنید:

۱. شبکه با لایه اصلی روی GPU

۲. شبکه با لایه جدید (کرنل Fused) روی GPU

۳. شبکه با لایه اصلی روی GPU با غیرفعال بودن cuDNN

۴. شبکه با لایه اصلی روی CPU

در هر سناریو، زمان اجرای ۵ اپیک (Epoch) را برای کل فرآیند، همچنین فقط بخش forward و فقط بخش backward در طور جداگانه محاسبه و گزارش کنید. دقت نهایی شبکه (Test Accuracy) رانیز برای هر مورد بنویسید.

۵.۱ تحلیل و بررسی (۲۰ نمره)

۱. نکات جالب توجه، اختلاف‌های زمانی و دقت را توضیح دهید.
۲. در چه موقعی پیاده‌سازی کرنل اختصاصی (Custom Kernel) می‌تواند مفید باشد و چه موقع بهتر است به پیاده‌سازی‌های داخلی cuDNN و PyTorch اعتقاد کنیم؟

سؤال ۲. بررسی تغییر سرعت آموزش در شرایط استفاده از چند GPU

۳۰ نمره

در این سؤال، قصد بررسی تغییر سرعت آموزش در شرایط استفاده از چند GPU را داریم.

۱.۲ طراحی و آموزش روی یک GPU (۱۰ نمره)

یک شبکه‌ی کانولوشنال طراحی کنید که بتواند به دقت بالای ۵۰ درصد برای طبقه‌بندی داده‌های تست برسد. این شبکه را بر روی یک GPU آموزش دهید. از بهینه‌ساز Adam و اندازه‌ی Batch برابر با ۳۲ استفاده نمایید.

۲.۲ پیاده‌سازی موازی روی دو GPU (۱۰ نمره)

کد بخش قبل را به گونه‌ای تغییر دهید که بتوان آن را با کمک دستور python و با استفاده از PyTorch Multiprocessing در روی دو GPU اجرا کرد.

۳.۲ مقایسه و تحلیل (۱۰ نمره)

در هر آزمایش (تک GPU و دو GPU)، دقت مدل بر روی داده‌های تست، زمان اجرای آموزش و میزان حافظه‌ی مصرفی GPU را به دست آورده، با یکدیگر مقایسه کرده و تحلیل نمایید.

سؤال ۳. بررسی تأثیر تغییر اندازه‌ی Batch

۱۰ نمره

در این سؤال، قصد بررسی تأثیر تغییر اندازه‌ی Batch در سرعت آموزش را داریم. بخش دوم سؤال دوم (آموزش بر روی دو GPU) را با اندازه‌های Batch برابر با ۱۶، ۳۲، ۶۴ و ۱۲۸ اجرا نمایید. در هر آزمایش، زمان اجرای آموزش، حافظه‌ی مصرفی GPU و دقت نهایی مدل را به دست آورده، جدول تغییر زمان آموزش، دقت مدل و حافظه‌ی مصرفی GPU را بر حسب تغییر اندازه‌ی Batch رسم نمایید و نتایج را تحلیل کنید.

سؤال ۴. بررسی‌های مختلف Backend در communication DDP PyTorch

۱۰ نمره

در این سؤال، قصد بررسی استفاده از Backends مختلف در communication DDP PyTorch را داریم. سؤال را با بهینه‌ساز Adam و با اندازه‌های Batch برابر با ۳۲ و ۱۲۸ و با استفاده از backends nccl و gloo اجرا کنید و نتایج را از نظر حافظه‌ی مصرفی GPU و زمان آموزش مقایسه نمایید.

۱. گزارش پروژه را در قالب فایل PDF تهیه کنید.

۲. تمام فایل‌های مربوط به تمرین را طبق ساختار جدول زیر در پوشه‌های مشخص شده قرار داده و در انتهای با فرمت zip در سامانه elearn بارگذاری کنید.

نام فایل	بخش	سؤال
CUDA.ipynb	بخش ۲	۱
Classifier.py	بخش ۱	۲
Classifier_mp.py	بخش ۲	

سؤال	بخش	نمره	توضیح
۱ (۵۰ نمره)	۱.۱ (بخش الف)	۲	تعريف و توضیح CUDA Kernel
۲ (۳۰ نمره)	۲.۱ (بخش ب)	۳	توضیح cuBLAS، cuDNN، مزایا و ارتباط با Fused Convolution + ReLU
۳ (۱۰ نمره)	۳.۱ (بخش ج)	۱۵	پیاده‌سازی کرنل در پایتون و CUDA
۴ (۱۰ نمره)	۴.۱ (بخش د)	۵	توضیح دقیق کلاس FusedConvReLUFn و تعامل با autograd
۵ (۱۰ نمره)	۵.۱ (بخش ه)	۵	پیاده‌سازی و اجرای آزمایش‌های MNIST در سناریوهای مختلف
۶ (۱۰ نمره)	۶.۱ (بخش ی)	۲۰	تحلیل اختلاف‌های زمانی و دقیق، نمودارها و بحث در مورد کرنل اختصاصی
۷ (۱۰ نمره)	۱.۲ (بخش الف)	۱۰	طراحی شبکه و آموزش روی تک GPU با دقت بالاتر از ۵۰%
۸ (۱۰ نمره)	۲.۲ (بخش ب)	۱۰	پیاده‌سازی آموزش توزیع شده روی دو GPU با DDP
۹ (۱۰ نمره)	۳.۲ (بخش ج)	۱۰	مقایسه سرعت، دقت و حافظه مصرفی میان تک GPU و دو GPU
۱۰ (۱۰ نمره)	سوال سوم	۱۰	اجرای آزمایش‌های Batch Size مختلف، رسم جدول/نمودار و تحلیل نتایج
۱۱ (۱۰ نمره)	سوال چهارم	۱۰	مقایسه Backend‌های nccl و gloo از نظر زمان و حافظه و تحلیل نهایی