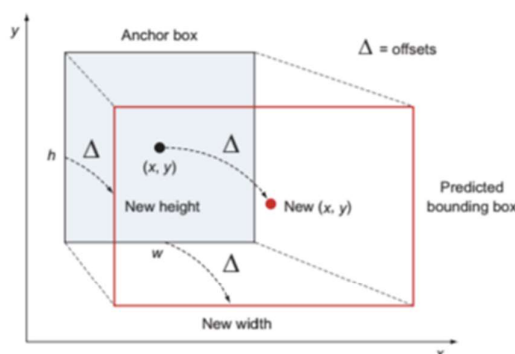


تمرین سری سوم

* ۱۰ نمره از ۱۱۰ نمره تمرین امتیازی است *

مسئله ۱. (۴)

توابع خطای Fast-RCNN و Faster-RCNN را موشکافانه مورد بررسی قرار دهید. همچنین ضمن بیان ارتباط مفهوم RPN با تصویر زیر^۱، مشخص کنید برای کدام یک مطرح شد و به چه چیزی اشاره دارد و موجب چه پیشرفتی شده است؟



مسئله ۲. (۴)

با در نظر گرفتن شارنوری و الگوریتم‌های Lucas & Kanade و Horn & Schunck به موارد زیر پاسخ دهید.

- الف) تفاوت اصلی بین دو الگوریتم را توضیح دهید. ضمن بیان دلیل مقاومت در برابر نویز در روش Lucas & Kanade نسبت به Horn & Schunck؛ پیش فرض‌های هر یک را نیز با یکدیگر مقایسه کنید.
- ب) نقش پارامتر رگولاریزاسیون α در روش Horn-Schunck چیست؟ تغییر آن چگونه بر میدان شار نوری حاصل تاثیر می‌گذارد؟ (ارائه مثال توصیه می‌شود)

ج) مشکل aperture در optical flow و چگونگی تاثیر آن بر motion estimation را توضیح دهید

¹ From course slides

مسئله ۳. (۶) *

با در نظر گرفتن مدل Flownet به موارد زیر پاسخ دهید.

الف) این مدل دو واریانت اصلی دارد؛ FlownetS و FlownetC. این دو نوع را با تمرکز بر تفاوت های معماری و مزایای خاصی که هر نوع برای تخمین شارنوری ارائه می دهند، مقایسه کنید.

ب) اجزای اصلی معماری FlownetC² مانند ساختار encoder، decoder، warping operations و لایه های کوریلیشنی و ارتباط این اجزا را مورد بررسی قرار دهید. در مورد چالش های مطرح شده توسط این مدل در تخمین شار نوری، به ویژه برای جابجایی های بزرگ و مناطق مسدود بحث کنید.



مسئله ۴. (۶)

ضمن توضیح کامل و با جزئیات تفسیر مفهوم خودتوجهی^۳ در زمینه ترنسفرمرهای بینایی؛ یک گونه از آن را با مشخصات زیر در نظر بگیرید و با توجه به آن تعداد پارامتری قابل آموزش مدل و تعداد پچ هایی که از تصویر ورودی تولید می شوند را بدست آورید. همچنین مزایا و محدودیت های ترنسفرهای بینایی را نسبت به شبکه های عصبی کانولوشنی (CNNs) در وظایف بینایی بررسی کنید.

- اندازه تصویر ورودی: ۲۲۴ در ۲۲۴ پیکسل
- اندازه پچ: ۱۶ در ۱۶ پیکسل
- تعداد سرهای توجه: ۸
- بعد تعبیه: ۵۱۲
- تعداد لایه ها: ۱۲

² From course slides

³ Self-attention

سوالات کامپیوتری (۹۰)

هدف از انجام این تمرین، بکارگیری روش‌های تشخیص اشیاء، تحلیل حرکت و مقدمه‌ای بر مبدل‌های بینایی^۴ و همچنین پوشش مطالب درس تا انتهای اسلاید ۷ است. لطفاً برای اینکه قالب تمامی پاسخ‌ها یکدست باشد، مراحل زیر را به ترتیب دنبال کنید.

- (۱) یک دایرکتوری با عنوان CV-CHW3-[Student ID] بسازید.
- (۲) برای هر سوال، یک نوتبوک با عنوان Q[number].ipynb در آن ایجاد کنید.
- (۳) برای سوالات موردی ابتدای هر مورد، سلولی مارک‌داون ایجاد کنید.
- (۴) مورد سوال را در قالب Q[number]-[part] در ابتدای این سلول بنویسید. (مثلاً Q4-A)
- (۵) هر جا سوالات نیاز به پاسخ تئوری داشت؛ مارک‌داون ایجاد کنید و توضیحات لازم را بنویسید.
- (۶) همگی ضمایم بصورت یکجا در این [لینک](#)^۵ قابل دسترسی هستند.

مسئله اول کامپیوتری (۱۵)

در این سوال، هدف ما طراحی یک مدل تشخیص شیء ساده برای شناسایی گربه‌ها در تصاویر با الهام از معماری‌های تدریس‌شده در کلاس؛ با طی کردن مراحل ذکر شده و کامل کردن نوتبوک پیوست شده است. در ابتدا داده‌ها را پیش‌پردازش کنید. (به این صورت که تصاویر و حاشیه‌نگاری‌ها، را بخوانید، ابعاد تصاویر را به اندازه مشخصی (۴۱۶ در ۴۱۶ پیکسل) تغییر دهید و نرمال‌سازی مقادیر پیکسل‌ها و مختصات جعبه‌های محدودکننده^۶ است.) در ادامه، معماری مدل را با ساخت یک شبکه هرمی ویژگی^۷ برای استخراج ویژگی‌ها از مراحل مختلف ResNet50 به همراه یک متد Detection Head شامل دو شاخه برای رگرسیون جعبه محدودکننده^۸ و پیش‌بینی کلاس طراحی کنید. سپس، با یکپارچه‌سازی شبکه پایه^۹، شبکه هرمی ویژگی و detection head، مدل نهایی را بسازید. در پایان، مدل را با استفاده از مجموعه داده آموزش دهید و دقت مدل و نمونه نتایج تشخیص را نمایش دهید. (تمامی ضمایم در Q1.zip قابل مشاهده است)

⁴ Vision Transformers

⁵ <https://drive.google.com/drive/folders/1aTEUFIJ43ZRwY1Pw1qs4bCp073KC0n6I?usp=sharing>

⁶ Annotation

⁷ Bounding box

⁸ FPN

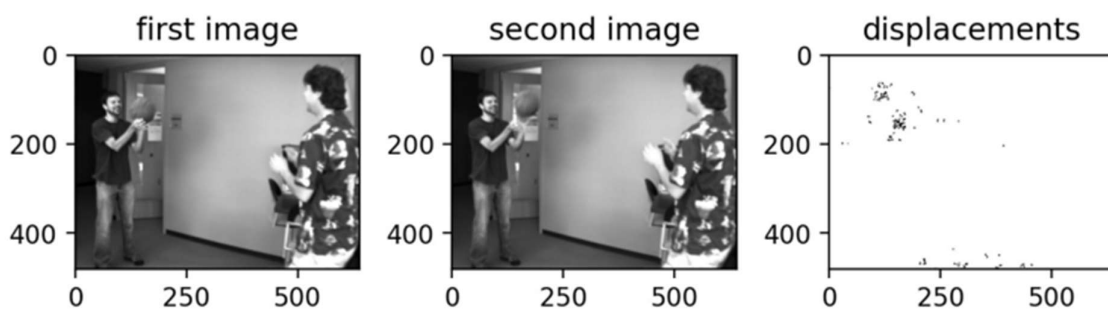
⁹ Regression bounding box

¹⁰ Backbone

مسئله دوم کامپیوتری (۱۰)

در این مسئله، روش Lucas-Kanade برای تخمین شار نوری و تحلیل حرکت پیاده سازی خواهد شد. شار نوری حرکت اجسام بین فریم های متوالی است. از آنجایی که هر ویدیو دنباله ای از تصاویر است هر دو فریم متوالی را می توان به الگوریتم های محاسبه گر شار نوری داد که جابجایی^{۱۱} ویژگی ها را به صورت مکانی محاسبه می کند. در نظر داشته باشید فرض شده دو فریم متوالی گرفته شده از یک ویدیو با اختلاف زمانی ناچیز گرفته شده اند، پس می توان فرض کرد که جابه جایی هیچ جسمی زیاد نیست و روشنایی جسم تغییر چشم گیری نمی کند و همچنین تمام پیکسل های همسایه (شامل شی) در یک جهت حرکت می کنند. حال با استفاده از این مفروضات حرکت اجسام را در فریم های داده شده در تصاویر موجود در ضمایم محاسبه کنید.

(استفاده از کتابخانه های آماده برای Lucas-Kanade مجاز نیست و هدف پیاده سازی گام به گام تمام مراحل است.)

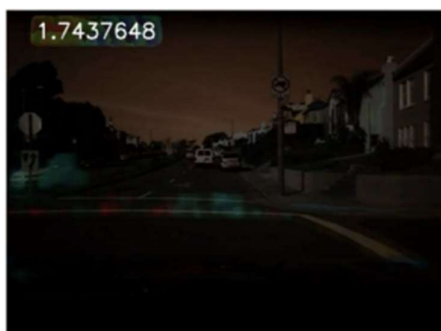


(تمامی ضمایم در Q2.zip قابل مشاهده است)

مسئله سوم کامپیوتری (۳۰)

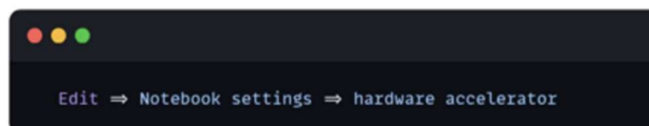
می خواهیم با استفاده از داده های دوربین ثبت وقایع یک خودرو در حال حرکت (رو به جلو)، سرعت آن را پیش بینی کنیم. این کار به علت شرایط نوری مختلف محیط چالش برانگیز است. یک رویکرد معمول استفاده از CNN ها با تصاویر و داده های برجسته گذاری شده است. با این حال، اگر ما از فریم های اصلی ویدئو به تنهایی بهره ببریم؛ مدل احتمالا دچار بیش برآزش خواهد شد. زیرا تنها تصاویر را با مقادیر سرعت به خاطر خواهد سپرد. بنابراین، باید از ویژگی های ورودی مناسب برای آموزش بهره ببریم. جریان نوری با استفاده از دو فریم متوالی ویدیویی در حالت خاکستری، ماتریسی با همان ابعاد تصویر ورودی ایجاد می کند. هر پیکسل این ماتریس نشان دهنده تغییر موقعیت و سرعت آن پیکسل نسبت به فریم قبلی است. جهت و بزرگی تغییرات در هر پیکسل از طریق کانال رنگی HSV ذخیره می شود. با آموزش شبکه با استفاده از جریان نوری، شبکه می آموزد که کدام

پیکسل‌های جریان نوری مهم هستند و وزن‌دهی مناسب را به آنها اختصاص می‌دهد. حال با استفاده از چهارچوب کاری PyTorch و معماری دلخواه، شبکه‌ای با استفاده از شار نوری آموزش دهید که سرعت حرکت خودرو را تخمین بزند و با معیارهای مناسب دقت مدل خود را گزارش دهید. همچنین در خروجی دو ویدئوی یکی با استفاده از ویدئوی اصلی و دیگری ویدئو ترکیبی جریان نوری و ویدئوی تست، هر دو با نمایش سرعت پیش‌بینی شده در بالای تصویر، تولید کرده و بصری‌سازی کنید. (برای درک بهتر دو فریم از موارد خواسته شده آورده شده است.) (ملاک ارزیابی این سوال بر روی دادگان آزمایش داده شده و همچنین رعایت موارد خواسته شده است.) (تمامی ضمایم در Q3.zip قابل مشاهده است)



مسئله چهارم کامپیوتری (۳۵) *

در این مسئله، هدف ما استفاده از ترنسفرمر بینایی برای دسته‌بندی تصاویر مجموعه داده CIFAR-10 است. ابتدا مجموعه داده CIFAR-10 را با تغییر اندازه تصاویر به ۲۲۴ در ۲۲۴ پیکسل و نرمال‌سازی مقادیر پیکسل‌ها پیش‌پردازش کنید. به دلیل سنگین بودن عملیات، ترجیحاً این نوت‌بوک را در گوگل کولب بارگذاری کرده و از مسیر زیر استفاده کنید تا شتاب‌دهنده سخت‌افزاری را بر روی GPU یا TPU قرار دهید و از سرعت پردازش بالاتری بهره‌مند شوید.



در مرحله بعد، یک مدل ترنسفرمر بینایی با استفاده از چهارچوب PyTorch پیاده‌سازی کنید. در فایل نوت‌بوک VisionTransformer که در ضمایم این سوال در دسترس است، تمامی مراحل با توضیحات مستندسازی شده در قابل مشاهده هستند و مقادیر هایپرپارامترها از پیش داده شده‌اند. وظیفه شما تکمیل ماژول‌های مشخص شده است تا مدل VisionTransformer بتواند به درستی عمل کند. حال برای این کار ترنسفرمر بینایی را روی

مجموعه داده پیش پردازش شده CIFAR-10 برای تعداد مشخصی مرحله آموزش دهید. خطا و روند بهینه سازی از قبل تعریف شده اند، و شما باید حلقه‌ی آموزش را پیاده سازی کنید. برای بهبود عملکرد ترنسفرمر بینایی می‌توانید هایپرپارامترهای مختلف را امتحان کنید. بدین منظور، در فایل نوت‌بوک یک سلول جدید ایجاد کرده و هایپرپارامترهای جدید را تعریف کنید و مدل را بر اساس آن ایجاد و حلقه‌ی آموزش را از ابتدا آغاز کنید. بهترین دقت نیز شامل نمره‌ی اضافی خواهد شد. پس از آموزش مدل، عملکرد آن را با محاسبه دقت روی یک مجموعه آزمایش جداگانه ارزیابی کنید. در نهایت، عملکرد ترنسفرمر بینایی را با مدل‌های پایه‌ای مانند ResNet یا VGG بر روی همان مجموعه داده مقایسه کنید. برای این منظور، مدل‌های از پیش آموزش داده شده موجود در کتابخانه PyTorch را در نوت‌بوک بارگذاری کرده و نتایج آن‌ها را نیز گزارش دهید. (تمامی ضمایم در Q4.zip قابل مشاهده اس

نکات:

- تحویل تکلیف در سامانه کوئرا و تا زمان مشخص شده خواهد بود.
 - تمارین تایپ شده شامل ۱۵ درصد نمره امتیازی (سوالات تئوری) می‌باشد.
 - مسئله سوم تئوری (الف)، چهارم کامپیوتری حاوی مواردی **امتیازی (*)** هستند
 - استفاده از [Mathcha](#) توصیه می‌شود.
 - فرمت فایل سوالات خود را حتماً به صورت زیر رعایت فرمایید.
- HW3/CHW3 - [Full Name] - [Student ID]
- در صورت مشاهده هرگونه تقلب، رونویسی و ... با افراد خاطی برخورد خواهد شد.