

به نام خدا

مبانی بینایی کامپیوتر

دکتر محمدی

تمرین شش

علی عطاریان - ۹۹۵۲۱۴۵۱

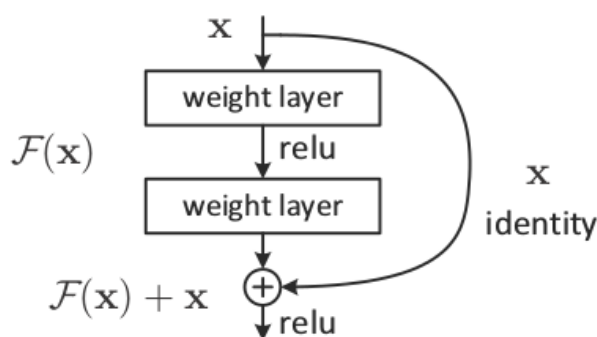
سوال (۱) الف)

هنگام استفاده از روش gradient learning با backpropagation، برای بهبود وزن‌ها به ترتیب از لایه آخر تا لایه اولیه، با قاعده زنجیره‌ای، مشتق گرفته می‌شود. در n لایه n مشتق در یکدیگر ضرب می‌شوند، اگر این مشتق‌ها بزرگ باشد به exploding gradient و اگر مشتق‌ها کوچک باشند به vanishing gradient منجر می‌شود. در حالت انفجار گرادیان مدل بسیار ناپایدار می‌شود و یادگیری انجام نمی‌گیرد. همچنین در موارد حاد، به مشکل overflow برمی‌خوریم که وزن‌ها مقدار NaN می‌گیرند. در حالت ناپدیدشدن گرادیان چون مقدار مشتق‌ها بسیار کوچک است تغییر ملموسی در وزن‌ها بوجود نمی‌آید و باز هم مدل یاد نمی‌گیرد. در بدترین حالت مشتق صفر منجر به گرادیانت صفر می‌شود و وزن‌ها صفر می‌شوند.

چند راه حل برای این مشکل مطرح است: با کاهش تعداد لایه‌ها تعداد ضرب‌ها کم می‌شود و احتمال هر دو مشکل کاهش می‌یابد اما پیچیدگی شبکه کم می‌شود. همچنین با gradient clipping حد بالا و پایین برای گرادیان تعریف می‌کنیم. علاوه بر این دو روش با مقداردهی اولیه صحیح وزن‌ها (استفاده از Xavier یا He) می‌توان از این مشکل جلوگیری کرد.

سوال (۱) ب)

شبکه‌های قبل ResNet با مشکل vanishing gradient مواجه بودند. ResNet با معرفی Residual blocks و روش skip connections این مشکل را حل کرد. البته باید توجه داشت استفاده از residual network تنها برای حل مشکل VGP نیست و یادگیری مدل را نیز بهبود می‌بخشد. یک residual block در واقع تشکیل شده از چندین لایه کانولوشنی و یک shortcut connection است. خروجی هر بلاک مجموع خروجی لایه‌های کانولوشنی و ورودی است.



در این حالت شبکه بر روی خروجی هر بلاک train می‌شود و این یعنی بسیاری از لایه ها skip می‌شوند و تعداد ضربها کم می‌شود. در واقع اضافه کردن ورودی به خروجی نهایی اجازه می‌دهد هنگام backpropagation گرادیانها مستقیما از طریق shortcut connection (identity path) پروپگیت شوند. با این معماری می‌توانیم شبکه‌هایی با صدها لایه بسازیم.

سوال (۲) الف)

پارامترها: در هر کدام از لایه‌های $1*1$ Conv، برای هر کدام از ۳ کانال ۳۲ ضریب داریم که با احتساب bias می‌شود $128 = 32 * (1 + 3)$. در مورد $3*3$ Conv می‌توان گفت که ورودی همه آنها ۳۲ کاناله است و هر کانولوشن ۹ ضریب دارد که با احتساب bias می‌شود:

$$9248 = 32 * (1 + 9 * 32)$$

لایه آخر نیز یک ورودی ۹۶ کاناله دارد (با فرض کانکت کردن ۳ تا خروجی ۳۲ کاناله) که با در نظر گرفتن بایاس می‌شود $24832 = 256 * (1 + 1 * 96)$ پس جمعا داریم:

$$128 + 128 + 128 + 9248 + 9248 + 9248 + 24832 = 52960$$

RF: در هر ۳ لایه $1*1$ Conv هر پیکسل معادل یک پیکسل است پس میدان تاثیر ۱ است. در دوتا لایه ابتدایی $3*3$ Conv پنجره $3*3$ داریم و هر پیکسل ورودی معادل ۱ پیکسل است پس میدان تاثیر $3*3$ است. در لایه $3*3$ Conv بعدی پنجره $3*3$ است و هر پیکسل ورودی نماینده ۹ پیکسل تصویر اصلی است پس میدان تاثیر $5*5$ است. در لایه آخر نیز پنجره $1*1$ است و از بین ورودی‌های این لایه نیز بیشترین میدان تاثیر $5*5$ است پس میدان تاثیر این لایه همان $5*5$ است. برای کل شبکه نیز بیشترین RF لایه‌ها را در نظر می‌گیریم که $5*5$ است.

سوال (۲) ب) A)

پارامتر: مشابه توضیحات بخش اول برای لایه اول داریم: $448 = 16 * (1 + 3 * 9)$ و برای لایه دوم داریم $4640 = 32 * (1 + 16 * 9)$ پس $448 + 4640 = 5088$

RF: در لایه اول میدان تاثیر $3*3$ و در لایه دوم $5*5$ است و در کل $5*5$ است.

سوال ۲) ب (B)

پارامتر: برای لایه اول داریم $(3*9+1)*(n-2)*(n-2)*16$

و برای لایه دوم داریم که $(16*9+1)*(n-4)*(n-4)*32$

RF: مانند مورد قبلی لایه اول $3*3$ و لایه دوم $5*5$ و کل شبکه $5*5$ است.