

Gradient accumulation روشی مناسب دیتاهای بزرگ است که استفاده از حافظه را بهینه‌تر میکند. در واقع بجای اینکه بعد از عبور هر بچ دیتا، پارامترها آپدیت شوند، چند بچ عبور میکنند و گرادیان روی هم جمع میشود و در آخر آپدیت میشود. به این صورت، مموری موردنیاز برای ذخیره گرادیان ها کاهش می یابد. مثلا میتوانیم ۴ بچ را به مدل بدهیم و سپس مدل را آپدیت کنیم. وقتی دیتای ما بسیار بزرگ است، مقدار بچ توسط مقدار حافظه‌ای که داریم، محدود میشود. و همین باعث میشود سایز بچ را مجبور شویم کوچک انتخاب کنیم و مدل دیرتر همگرا میشود. با استفاده از تکنیک Gradient accumulation، هر بچ مدنظرمان را به چند مینی بچ تقسیم میکنیم که روی مموری جا بشوند، اما آپدیت را بعد از اینکه همه گرادیان‌های مینی بچ‌ها محاسبه شد، انجام میدهیم. در واقع effective batch size یعنی سایز بچ واقعی را افزایش میدهیم (با تقسیم به مینی بچ، یک بچ بزرگتر را شبیه سازی میکنیم)

```
# from https://kozodoi.me/blog/20210219/gradient-accumulation

# batch accumulation parameter
accum_iter = 4

# loop through enumerate batches
for batch_idx, (inputs, labels) in enumerate(data_loader):
    # extract inputs and labels
    inputs = inputs.to(device)
    labels = labels.to(device)

    # passes and weights update
    with torch.set_grad_enabled(True):
        # forward pass
        preds = model(inputs)
        loss = criterion(preds, labels)

    # normalize loss to account for batch accumulation
    loss = loss / accum_iter

    # backward pass
    loss.backward()

    # weights update
    if ((batch_idx + 1) % accum_iter == 0) or (batch_idx + 1 == len(data_loader)):
        optimizer.step()
        optimizer.zero_grad()
```

## Backpropagation in CNN

در CNN، چون هر نورون، تمام عکس (ورودی) را می بیند، شیوه مقایسه  $\Delta$  و سطح وزن متفاوت است.

$$x_j^l = \sum_m \sum_n w_{m,n}^{l-1} o_{i+m, j+n}^{l-1} + b_j^l \quad \text{مسیر forward}$$

فیلتر روی عکس با همان اندازه  $m \times n$

$x$  فیلتر می خروجی و ناوتل اندیس آن نورون در فیلتر می در لایه  $l-1$

$$E = \frac{1}{2} \sum_p (t_p - y_p)^2$$

وزن های فیلتر لایه ۱ را در نورون های لایه ۱ ضرب می کنیم تا به اندیس ناوتل ضرب می کنیم تا به اندیس ناوتل

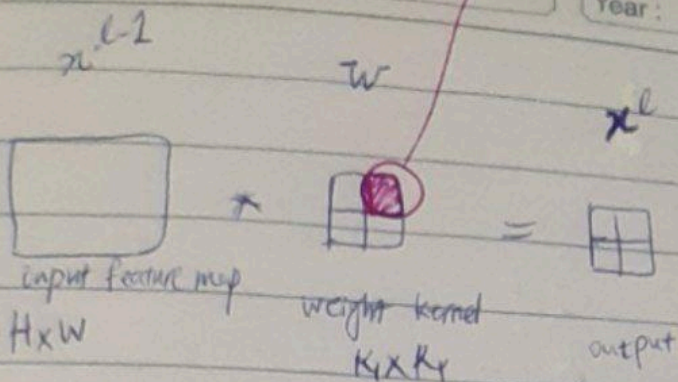
activation function (داتی)

Subject:

Year:

Month:

Date:



$$\frac{\partial E}{\partial w_{m',n'}^l} = \sum_{i=0}^{H-K_x} \sum_{j=0}^{W-K_y} \left( \frac{\partial E}{\partial x_{ij}^l} \times \frac{\partial x_{ij}^l}{\partial w_{m',n'}^l} \right) = \sum_{i=0}^{H-K_x} \sum_{j=0}^{W-K_y} \left( \delta_{ij}^l \frac{\partial x_{ij}^l}{\partial w_{m',n'}^l} \right)$$

iterate کردن روی  
output های لایه  
و وزن های لایه  
تأثیرات لایه های قبل را بر روی

$$\frac{\partial x_{ij}^l}{\partial w_{m',n'}^l} = \frac{\partial}{\partial w_{m',n'}^l} \left( \sum_m \sum_n w_{m,n}^l o_{i+m,j+n}^{l-1} + b^l \right)$$

$$= \frac{\partial}{\partial w_{m',n'}^l} \left( w_{0,0}^l o_{i+0,j+0}^{l-1} + \dots + w_{m',n'}^l o_{i+m',j+n'}^{l-1} + \dots + b^l \right)$$

$$= 0 \quad \text{for } m',n' \neq i,j$$

$$\rightarrow \frac{\partial E}{\partial w_{m,n}^l} = \sum_{i=0}^{H-K_x} \sum_{j=0}^{W-K_y} \delta_{ij}^l o_{i+m,j+n}^{l-1}$$





invariant شود. همچنین فیچرها را میتوانیم با مکس پولینگ سلکت کنیم (فیچرها مهمتر) و اورجی پولینگ نیز میانگینی از همه اطلاعات را حفظ میکند.

نقاط منفی: باعث از دست رفتن بخشی از دیتا میشود. باعث افزوده شدن تعدادی هایپراپارامتر مثل سایز پولینگ و موقعیتشان در شبکه میشود.

در یک عکس عادی، اگر فیلترهای پولینگ را اعمال کنیم، در عکسهایی با بگگراند سفید، مکس پولینگ میتواند اشیای تاریک را کلا محو کند. و مین پولینگ در عکسهایی با بگگراند مشکی، اشیای روشن را محو کند. اما در کل بیشتر از مکس پول استفاده میشود تا فیچرهای مهمتر از دست نروند.

استفاده بیش از حد از پولینگ نیز طبیعتاً منجر به از دست رفتن دیتا میشود که شاید برای تسک مهم باشد.

#### سوال ۴-

ترنسفر لرنینگ به ما امکان میدهد که از یک مدل از قبل آموزش داده شده استفاده کنیم تا تسک جدیدی را با فاین تون کردن مدل روی دیتاست جدید انجام دهیم. این کار باعث میشود بطور بهینه و بدون صرف انرژی زیاد، تسک را به انجام برسانیم. در واقع یک سری لایه اولیه با وزنهای شبکه‌ی pretrained فریز میشوند و فقط لایه های جدید اضافه شده، گرادیان را درون خود به عقب نشر میدهند.

از جمله این مدلها، alexnet برای ایمج کلسیفیکیشن، vgg که فیچرهای عکسها را برای ما اکسترکت میکند، resnet برای ایمج کلسیفیکیشن، ابجکت دیتکشن و سگمنتیشن، گوگل نت نیز برای ایمج کلسیفیکیشن و ابجکت دیتکشن.

#### سوال ۵-

-درست، وریفیکیشن یعنی یک چهره را با چهره سیو شده تطابق بدهیم تا بفهمیم آیا این فرد، شخص مورد نظر هست یا نه، یعنی خروجی یک باینری بله یا خیر است. اما رکوگنیشن یعنی تشخیص دهد که هویت این فرد کیست (از بین دیتابیس بزرگ).

-نادرست، چون این تعداد دیتا ممکن است باعث شود شبکه به خوبی یاد نگیرد که مثلا به نورهای مختلف invariant شود. بهتر است از هر فرد با لیبل مشخص، چندین عکس با زوایا و نور متفاوت داده شود.

-درست، چون لایه های اولیه معمولا فیچرهای ساده تر را اکسترکت میکنند و لایه های بعدی فیچرهای پیچیده تر هستند و این را در فرایند دیکانولوشن هم میتوانیم ببینیم که هرچه جلوتر میرویم کرنل ها اشیای پییده تری را دیتکت میکنند.