

## پیش گزارش آزمایش دوم (پایه سازی پرسپترون)

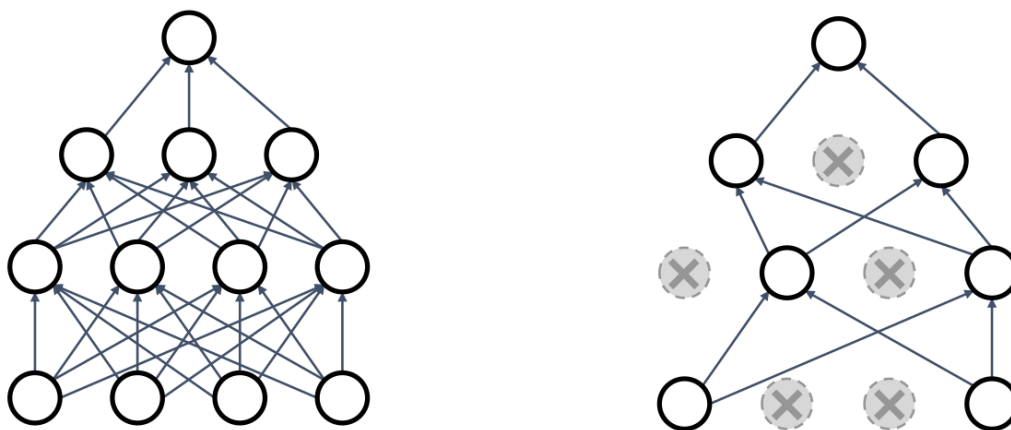
امیرحسین احمدی آشتیانی ۹۹۲۳۵۰۱

محمد رضا امیری ۹۹۲۶۰۴۰

محمد مهدی نوروزی ۹۹۲۳۰۸۵

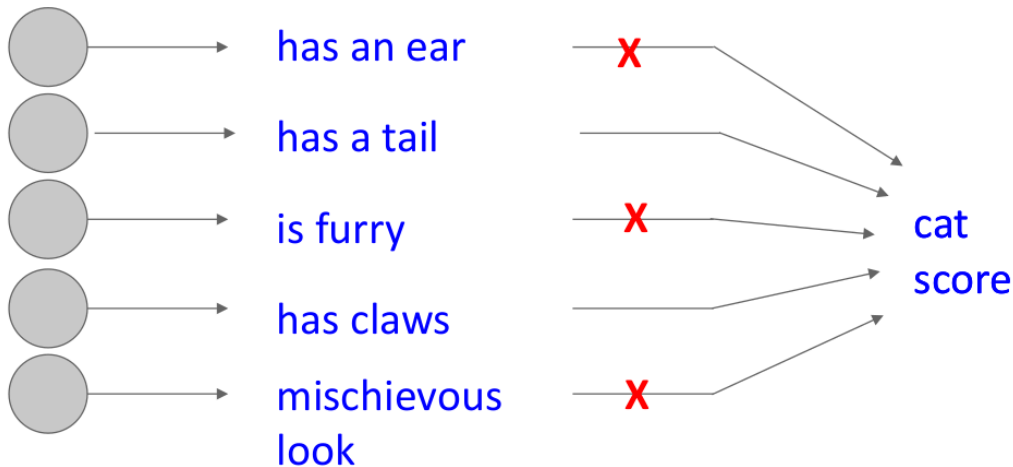
(1) عملکرد dropout را توضیح دهید.

این الگوریتم را در روش های regularization قرار دارد. هدف regularization جلوگیری از overfit است. ساده ترین regularization همان اضافه کردن مجموع ضرایب به loss function است که به نام L1 و L2 شناخته می شود. پس dropout هم سعی میکند overfit را خنثی کند. در این روش، در forward pass به صورت رندم برخی از نورون ها را صفر میکنیم. احتمال صفر کردن نورون یک هایپر پارامتر است.



اما این روش تاثیر های دیگری نیز دارد که به برخی از آنها اشاره میکنیم:

- شبکه را مجبور می کند که یک نمایش اضافی داشته باشد. برای مثال وقتی میخواهیم یک گره را تشخیص دهیم، حتما نباید تمام ویژگی ها را ببیند و با کنار هم قرار دادن تمام آنها، گره بودن را تشخیص دهد. شکل زیر این مورد را نشان میدهد:



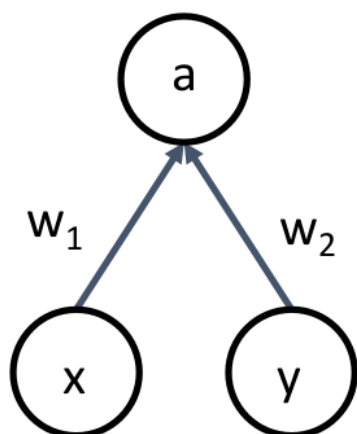
- یک تفسیر دیگر این است: در دنیای کلاسیفیکیشن، از ensemble برای تقویت مدل کلاسیفایر استفاده میشود. یعنی جواب کلاس بندی چندین مدل مختلف را باهم ترکیب کرده و بدین ترتیب مدل بسیار کارکرد بهتری خواهد داشت. Dropout نیز همین کار را میکند؛ چندین مدل مختلف (با نوروں های حذف شده مختلف) را کنار هم قرار میدهد و بعد تصمیم گیری میکند (در زمان تست). اما سوال مهم این است که چطوری در زمان تست عمل میکند؟ پاسخ این است که dropout خروجی را random میکند چون نوروں ها رندم هستند:

$$\text{Output (label)} \quad y = f_W(\text{Input (image)} \quad x, z) \quad \text{Random mask}$$

برای اینکه تشخیص دهیم که به طور متوسط خروجی چه مقداری خواهد بود باید expectation آنرا محاسبه کنیم.

$$y = f(x) = E_z[f(x, z)] = \int p(z) f(x, z) dz$$

برای محاسبه بیاید یک نوروں را در نظر بگیریم:



در زمان تست داریم:

$$E[a] = w_1x + w_2y$$

اما در زمان ترین داریم:

$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$

پس در زمان تست نیازی به حذف کردن نرون ها نیست، فقط خروجی را محاسبه و در احتمال dropout ضرب کنیم.

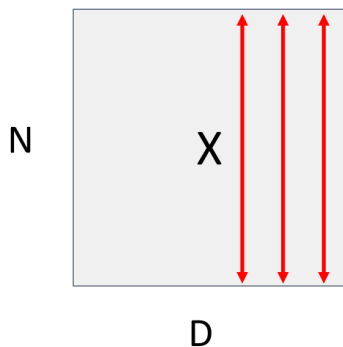
(2) روش bach normalization را توضیح دهید.

ایده اصلی این است که خروجی هر لایه رو نرمال کنیم به طوری که میانگین صفر و واریانس ۱ داشته باشد. این کار به optimaz کردن کمک میکند. اما چطوری؟ در ادامه بررسی میکنیم. به روش زیر خروجی یک لایه را نرمال میکنیم:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

این یک تابع مشتق پذیر است و میتوان در لایه های شبکه عصبی قرار داد و عملیات backprop را روی آنها انجام داد. از اسم این روش مشخص است که با تعدادی داده کار داریم (batch). میانگین و انحراف معیار به صورت زیر محاسبه خواهد شد:

**Input:**  $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j} \quad \text{Per-channel mean, shape is D}$$

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 \quad \text{Per-channel std, shape is D}$$

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}} \quad \text{Normalized x, Shape is N x D}$$

از محاسبات بالا متوجه میشویم که این نرمال کردن روی یک داده نیست. باید فیچر ها نرمال شوند. خب حالا سوال پیش می آید که اگر میانگین صفر و واریانس واحد شرایط مطلوبی نباشد باید چه می شود؟ پاسخ این است که برای حل کردن این مشکل دو پارامتر جدید برای لایه batch norm تعریف میشود: scale و shift.

اگر گاما که همان پارامتر scale است برابر انحراف معیار شود و بتا که همان shift است برابر میانگین شود، لایه batch norm یک لایه همانی خواهد شد.

مشکلی که وجود دارد این است که اگر بخواهیم در زمان تست این کار را انجام دهیم، خروجی بر اساس ورودی تغییر خواهد کرد، یعنی یک لایه ای درون شبکه عصبی وجود دارد که از قبل تعیین نشده و در زمان تست مشخص میشود که این مطلوب نیست. برای حل این مشکل میانگین و

واریانس را از مرحله ترین استفاده میکنند. یعنی running average میانگین و انحراف معیار در مرحله ترین، در مرحله تست استفاده میشود.

اگر از بچ نرمالیزیشن استفاده کنیم، میتوانیم لرنینگ ریت را افزایش دهیم بدون اینکه نگران شویم که واگرا خواهیم شد.

(3) الگوریتم پس انتشار خطا در حالت کلی آپدیت کردن وزن ها است با توجه به خطایی که در خروجی داشتیم. پس باید مشتق خطا نسبت به تمام وزن هارا داشته باشیم. یک الگوریتم ساختارمند که از مشتق زنجیره ای حاصل شده است، گراف محاسبات است. پس اگر گراف محاسبات را محاسبه کنیم، فقط آپدیت کردن باقی می ماند. پس در زیر گراف محاسبات را توضیح می دهیم. همانند شکل زیر ابتدا باید forward انجام شود:

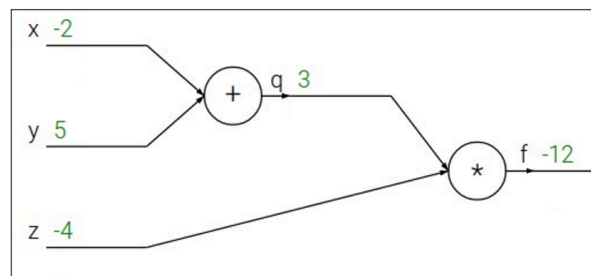
Backpropagation:  
Simple Example

$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

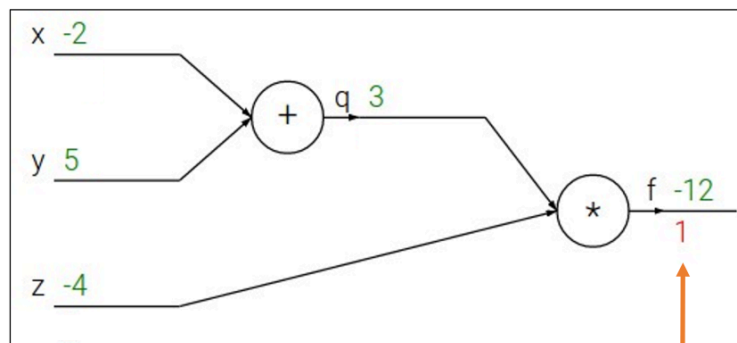


حال مقادیر زیر را می خواهیم:

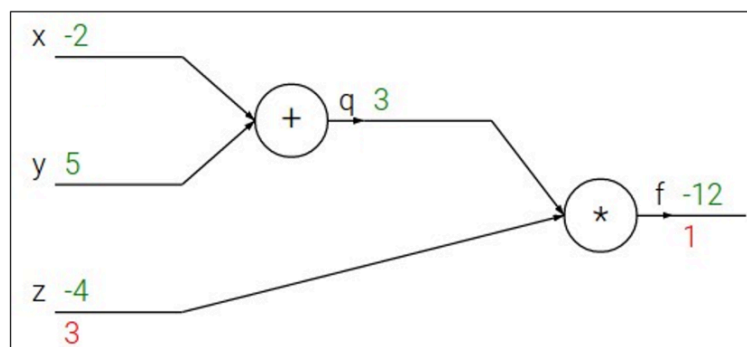
2. Backward pass: Compute derivatives

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

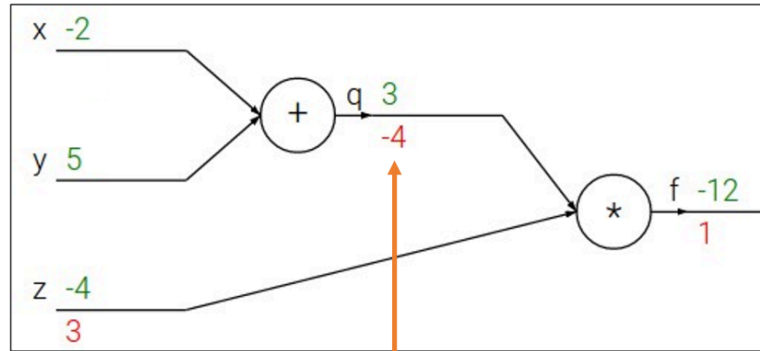
تصاویر زیر محاسبات را نشان می دهد:



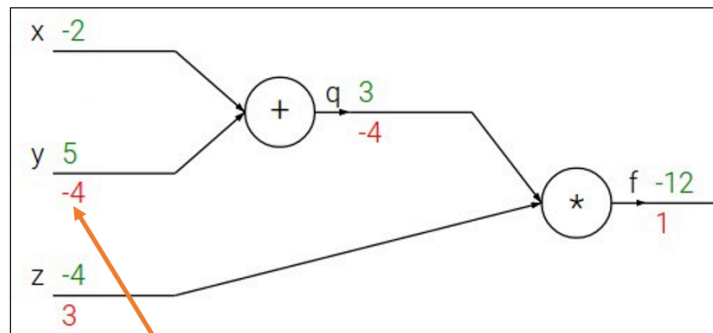
$$\frac{\partial f}{\partial f}$$



$$\frac{\partial f}{\partial z}$$



$$\frac{\partial f}{\partial q} = z$$



**Chain Rule**

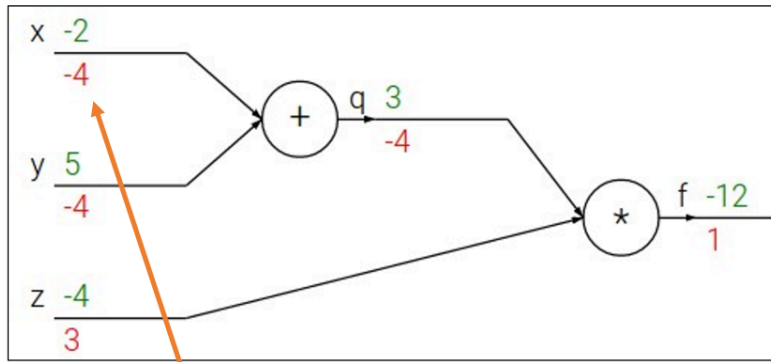
$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \frac{\partial f}{\partial q}$$

$$\frac{\partial q}{\partial y} = 1$$

Downstream  
Gradient

Local  
Gradient

Upstream  
Gradient



**Chain Rule**

$$\frac{\partial f}{\partial x} = \frac{\partial q}{\partial x} \frac{\partial f}{\partial q}$$

$$\frac{\partial q}{\partial x} = 1$$

Downstream Gradient      Local Gradient      Upstream Gradient