①

الف ) ساده ترین حالت، محاسبه کانولوشن به صورت مستقیم با ابعاد $N \times N$ ، هزینه زمانی $O(N^2 \times N^2)$

برابر با $O(N^4)$ خواهد داشت.

ب ) اگر بخواهم دو چندجمله ای را با یکدیگر ضرب کنیم، می توانیم آنها را نقاطی یعنی در چندجمله ها را با درجه ای بالاتر نسبت به بیشترین ضرب این دو چندجمله محاسبه شده به این کار کانولوشن گفته می شود. دیگر این است که برای این محاسبه ابتدا چندجمله ها را به فضای دیگری می بریم و سپس ضرب را به صورت elementwise بین مولفه ها هر لیست انجام کنیم. این فضا فضای فوریه است که می توان به صورت رو به شمارش داد $FFT(f) \circ FFT(g)$ سپس بلافاصله در این حاصل ضرب فضای این اعمال می کنیم که با عکس تبدیل فوریه امکان پذیر است. بنابراین برای ضرب چندجمله ای های $f$ و چندجمله ای $g$ می توانیم از روش روبرو استفاده نمایم $FFT^{-1}(FFT(f) \circ FFT(g))$ این روش هزینه زمانی $O(n \log n)$ دارد.

چندجمله ای با سه مؤلفه بالایی و سه درجه داشته باشد. حال می توانیم از این روش برای کانولوشن دو بعدی هم استفاده نمایم. برای این منظور می خواهیم عمل می کنیم که دو ماتریس را از حالت دو بعدی به یک بعدی تبدیل می کنیم ، Flatten آرایه های یک بعدی شده را کانولوشن را انجام می دهیم و نتیجه را به حالت دو بعدی برمی گردانیم. در ابتدا باید هر دو ماتریس را با صفر padding کنیم تا ماتریس که از هم بزرگ به تصویر قرار بگیرند. بنابراین اگر مثل $H \times W$ باشد کافیست

s.a.m

در روش بعدی ما می‌توانیم از ایده‌ی این استفاده کنیم که تعداد درایه‌های Convolution یعنی N+W-1 است، پس padding حدودا

flatten می‌کنیم، با FFT نمونه برداری را انجام می‌دهیم، سپس عمل ضرب دو تایی را انجام می‌دهیم، کرنل k×k را داریم، k

... ... ... ... ... ... ... ... ... ... ... ... ... ... ... ...

با cut off کردن $\frac{k-1}{2}$ ... ... ... ... ... ... $N = N+k-1-2(\frac{k-1}{2})$

... ... ... ... $O(MN(logN+logM)) = O(MNlogMN)$ خواهد بود. حالت اینکه

تصویر ما N×N ... ... ... ... $O(N^2logN)$ خواهد بود، روش نوشتن شبه این را باید

نوشته است، این روش به ماسک می‌آید.

الف) فرض کنید هزینه با یک regularization term به شکل زیر باشد:

$$\min \quad \mathcal{J}(w) + \lambda \|w\|_2^2 = \min \quad \frac{1}{2}\|Xw - Y\|_2 + \lambda \|w\|_2^2$$

اگر نمونه بیشتری بیاید جواب خوبی بدست می‌آید. اما می‌خواهیم یک minimizer، این دانش را به عنوان prior (دانش قبلی) می‌شود در اختیارش بگذاریم اینکه این پاسخ با نرم کمتر پاسخ بهتری هست.

فضای ما این است. حال ما می‌خواهیم بنویسیم:

$$w = X^T \beta + v$$
$$s.t. \quad v^T x_i = 0 \quad \text{for all } i$$

$$\Rightarrow \min_{\beta, v} \quad \frac{1}{2}\sum_{i=1}^{n} loss(\beta^T X x_i, y_i) + \lambda \|X^T\beta\|_2^2 + \lambda \|v\|_2^2$$

اگر برحسب $v$ مینیمایز کنیم بهترین جواب این ما پاسخ بهینه $v=0$ است. نتیجه پاسخ به صورت زیر در می‌آید.

$$w = X^T \beta$$

نشان دادیم شده پاسخ در ابعاد داده‌ها ورودی است، و گرچسا با مساوی n است. به این تکنیک (یا

representer theorem گفته می‌شود.

ج) میتوانیم این مسئله رابصورت زیر به نویسیم:

$$\min_{w} \|w\|_2^2 \quad \text{subject to} \quad y = Xw$$

* میتوانیم مسئله را با دخیل کردن محدودیت بصورت یک مسئله ی جدید با استفاده از ضریب لاگرانژ در قالب $\lim \lambda \to \infty$ بنویسیم ، و دلیل اینکه $\lambda$ به سمت بینهایت می رود آن است که خطا (error) صفر شود.

$$\min_{w} \left( \|w\|_2^2 + \lambda \|y - Xw\|_2^2 \right)$$

اینک این مسئله Convex است . بنابراین از آن مشتق گرفته و برابر صفر قرار می دهیم:

$$2w - 2\lambda X^T (y - Xw) = 0 \implies w(I + \lambda X^T X) = \lambda X^T y$$

$$\xrightarrow{\lambda \to \infty} w^* = (X^T X)^{-1} X^T y \xrightarrow{w = X^T \alpha} w^* = (X^T X)^{-1} X^T X X^T \alpha = X^T \alpha = w$$

* پس ثابت شد که اگر $w = X^T \alpha$ باشد کمینه ی تابع برقرار است برای $y = Xw$ بنابراین $w^* = X^T \alpha$ است و ثابت شد.

SGD به شکل آسان است .

قضیه اول میگه اگر شما M تا لایه مخفی که داری ترکیب خطی linear regression

هستش، میتونه... چندتا لایه مخفی ترکیب خطی هستن، که نهایتا ترکیب خطی میشن

دوم: اگر... با... M تا لایه مخفی... ترکیب خطی میشن.

$$f^1 = mx + n$$
$$f^2 = px + q$$ } activation functions

$$a^1 = f^1(z^1) = f^1(w^1x^1 + b^1) = m[w^1x^1 + b^1] + n = mw^1x^1 + mb^1 + n$$

$$a^2 = f^2(z^2) = f^2(w^2a^1 + b^2) = f^2(w^2(mw^1x^1 + mb^1 + n) + b^2)$$

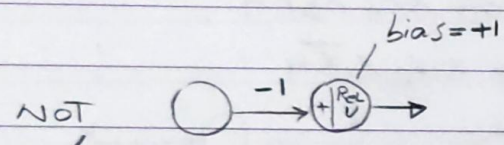$$= f^2(w^2mw^1x^1 + mb^1w^2 + w^2n + b^2) = p[w^2mw^1x^1 + mb^1w^2 + w^2n + b^2] + q$$

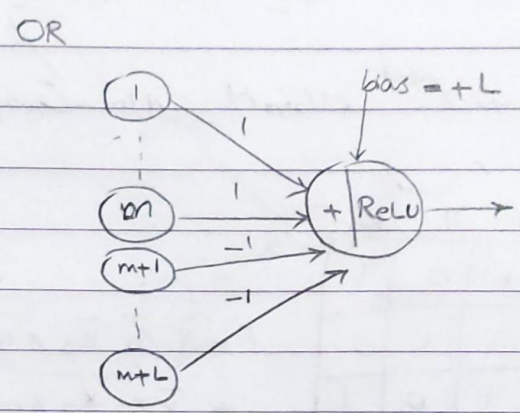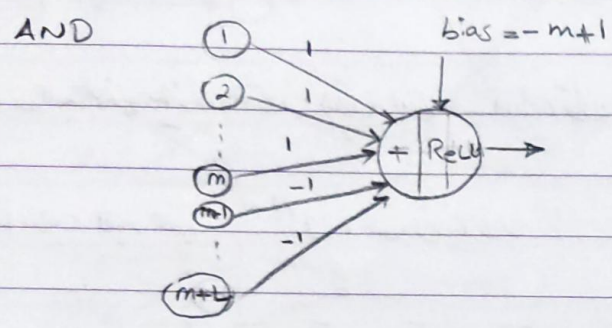$$= \underbrace{pw^2mw^1x^1}_{W} + \underbrace{pmb^1w^2 + pw^2n + pb^2 + q}_{B} = WX^1 + B$$

هم یک ترکیب خطی هستش

و بدون لایه مخفی میشه آخرش

الف) با استفاده شبکه perceptron می‌توانیم گیت‌های منطقی AND ، OR ، NOT را بسازیم و با سری

داشتن این گیت‌ها می‌توانیم هر فرمولی به فرم SOP بسازیم و با شبکه یک لایه مخفی رابطه را نشان دهیم.

**AND**

bias = -m+1

$$ \text{①} \xrightarrow{1} $$
$$ \text{②} \xrightarrow{1} $$
$$ \text{⋮} \xrightarrow{1} $$
$$ \text{ⓜ} \xrightarrow{-1} \quad +\, |ReLU \rightarrow $$
$$ \text{ⓜ₊₁} $$
$$ \text{ⓜ₊ₗ} \xrightarrow{-1} $$

**OR**

bias = +L

$$ \text{①} $$
$$ \text{⋮} \xrightarrow{1} $$
$$ \text{ⓜ} \xrightarrow{1} \quad +\, |ReLU \rightarrow $$
$$ \text{ⓜ₊₁} \xrightarrow{-1} $$
$$ \text{⋮} \xrightarrow{-1} $$
$$ \text{ⓜ₊ₗ} $$

**NOT**

bias = +1

$$ \bigcirc \xrightarrow{-1} +\,|ReLU \rightarrow $$

حالا هر تابع که داشته باشیم می‌توانیم با تبدیل کردن به فرم SOP و با استفاده از AND، OR مدل شبکه نوشته شده را بسازیم

برای نشان دادن این سعی از شبکه رابطه دو ظرفیتی رابطه $X Y + \bar{Y} Z$ می‌خواهیم بسازیم. به طور مثال این می می‌توانیم آن را بسازیم

زیر نشان می‌دهیم:

$$ X \xrightarrow{1} $$
$$ \xrightarrow{1} \; \text{(-1)} $$
$$ Y \xrightarrow{-1} \quad \text{(0)} \rightarrow X Y + \bar{Y} Z $$
$$ Z \xrightarrow{1} \; \text{(+1)} $$

input
layer

۱- تابع XOR با ۴ ورودی است و می‌خواهیم شبکه عصبی طراحی کنیم. می‌دانیم چون

OR تعداد ورودی زیاد است ولی حداقل یک نورون در لایه مخفی است و در تابع AND را داریم. یعنی درمیان نورون‌های هر لایه حداقل یک نورون $2^{(4-1)}$ خط است و می‌توان گفت

ور در لایه مخفی نیاز به نورون دارند. این مسئله برای آن معین داریم انتخاب کنیم ( N ورودی) $2^{N-1}$ خط صفحه می‌خواهد

بنابراین لایه مخفی $2^{N-1}$ ، نورون نیاز دارد

| $X_1X_2$ \ $X_3X_4$ | 00 | 01 | 11 | 10 |
|---|---|---|---|---|
| 00 | 0 | 1 | 0 | 1 |
| 01 | 1 | 0 | 1 | 0 |
| 11 | 0 | 1 | 0 | 1 |
| 10 | 1 | 0 | 1 | 0 |

$$F = (\overline{X_1}\,\overline{X_2}\,\overline{X_3}\,X_4) + \overline{X_1}\,\overline{X_2}\,X_3\,\overline{X_4} + \overline{X_1}\,X_2\,\overline{X_3}\,\overline{X_4}$$
$$+ \overline{X_1}\,X_2\,X_3\,X_4 + X_1\,X_2\,\overline{X_3}\,X_4 + X_1\,X_2\,X_3\,\overline{X_4}$$
$$+ X_1\,\overline{X_2}\,\overline{X_3}\,\overline{X_4} + (X_1\,\overline{X_2}\,X_3\,X_4)$$

$$F = X_1 \oplus X_2 \oplus \cdots \oplus X_N = \{((( X_1 \oplus X_2) \oplus (X_3 \oplus X_4)) \oplus ((X_5 \oplus X_6) \oplus (X_7 \oplus X_8)))$$

$$\oplus \cdots$$



* علامت XOR، بهره‌وری ۶ سالاته، عمق مدار پیاده‌سازی است $2 \log_2 N$ و تعداد دروازه مورد نیاز $3(N-1)$ طبقه (است) برابر پیاده‌سازی است.

عمق مسیر + تعداد دروازه‌ها است.

$$\text{Sigmoid} \qquad \frac{1}{1+e^{-x}}$$

مزایا: ۱) مشتق پذیر است و در هر نقطه می‌توانیم مشتق را به دست آوریم.

۲) این تابع monotonic است، نسبت اعداد بزرگ (کوچک) را حفظ می‌کند.

۳) گرادیان smooth دارد و باعث می‌شود تغییرات ناگهانی نداشته باشیم.

۴) مقدار خروجی آن بین ۰ و ۱ است، بنابراین در مواقعی که بخواهیم احتمال محاسبه کنیم، بسیار کارآمد است.

۵) مقدار تابع صفات‌های بالای ۲ و پایین‌تر از ۲- به ترتیب بسیار به ۱ و ۰ نزدیک می‌شود و به ما اجازه تخمین روشنی را می‌دهد.

معایب: ۱) در نقاط ابتدایی و انتهای این تابع مقدار گرادیان بسیار کوچک می‌شود باعث می‌شود فرآیند یادگیری بسیار کُند شود ولی اگر مقدار گرادیان صفر شود کلاً یادگیری متوقف داشت و در پیشه حل نکرده باشیم. این مسئله vanishing gradient نامیده می‌شود.

۲) خروجی این تابع zero-centered نیست.

۳) از نظر محاسباتی سنگین است.

$$\frac{e^x - e^{-x}}{e^x + e^{-x}} \qquad tanh$$

مزایا: این تابع شبیه سیگموید است sigmoid و تفاوت‌هایی با آن دارد. این تابع
خروج آن در بازه (1و1-) است، تفاوت دیگر این است که این تابع zero centered است. بنابراین شیب تندتری است
دارد بنابراین گرادیان تندتر قادیر نشست میره و بنابراین یادگیری سریع تر نسبت به sigmoid است.

معایب: این تابع همچنان مثل sigmoid دچار مشکل vanishing gradient می‌شود زیرا در لبه‌های دو طرف گرادیان رو به صفر است

معایب دیگر:

۲) این تابع از نظر محاسباتی تابع سنگین است.

$$\begin{cases} 0 & x \leq 0 \\ x & x > 0 \end{cases} \qquad ReLU$$

مزایا: این تابع از نظر محاسباتی بسیار ساده تر از ۲ تابع قبلی است و بدلیل همین سادگی سرعت همگرایی خودش را دارد و ... بیشتر ... به سیگموید رسیده است. به همین دلیل سریع تر عمل می‌کند و همین دلیل که شبکه‌های عمیق بیشتر به این می‌رسد.
۲) اگر مجموع وزن ورودی‌ها مثبت و بزرگ تر از صفر باشد تابع ReLU از نقطه نظر عمل نشست می‌کند

خوب است.

معایب: مشکلی به نام dead relu دارد وقتی ورودی منفی می‌شود مقدار گرادیان صفر شده و آموزش صورت نمی‌گیرد ولی

مشکل vanishing gradient ندارد.

s.a.m

اگر ۳ عدد موجود باشد با سیستم سه نرون می‌توانیم بیشترین عدد را از بین آنها تشخیص بدهیم :

اگر بنویسیم :

$\left.\begin{array}{l} a \geqslant b \\ a \geqslant c \end{array}\right\} \Rightarrow \quad a+a \geqslant b+c \Rightarrow 2a-b-c \geqslant 0 = 2a-b-c+0.5 \geqslant 0.5 \geqslant 0$

بنابراین می‌توانیم شبکه‌ای مانند روبرو برای شناسایی اولویت شبکه ایجاد بکنیم :



$y_1 = Relu(2a-b-c+0.5)$

$y_2 = relu(2b-a-e+0.5)$

$y_3 = relu(2c-a-b+0.5)$

weight decay $\Rightarrow w = (1-\lambda)w - \alpha \Delta C_0$

$\Delta C_0$ دلیل ضریب $\lambda$ باعث نمی‌جریم گامی را از مجموعه regularization term استفاده نمی‌کنیم.

L2 regularization $\Rightarrow C = C_0 + \frac{\lambda}{2}\|w\|_2^2$

$\frac{\partial C}{\partial w} = \frac{\partial C_0}{\partial w} + 2\frac{\lambda w}{2}$  $\Delta C = \frac{\partial C}{\partial w}$  $w = w - \alpha \Delta C = w - \alpha(\Delta C_0 + \lambda w)$

$\Rightarrow w = w - \alpha \Delta C_0 - \alpha \lambda w = (1-\alpha\lambda)w - \alpha \Delta C_0$

طبیعتاً $\lambda' = \lambda \alpha$ از اسپلت پایین‌تر می‌شود این مسئله در موقع نرخ بزرگ‌تر می‌آید:

$w = (1 - \lambda')w - \alpha \Delta C_0$   $(\lambda' = \lambda \alpha)$ =

کن این طوری‌ها برای غلط" اینا به شرط اینکه از weight decay استفاده کنیم. اینا به SGD، $L_2$ Decay و weight Decay یا

مثلاً برای روش‌های دیگر مثل Adam این یکسان نیست.

Guided Backpropagation؛ این روش نشان می‌دهد که چطور ورودی روی خروجی شبکه اثر می‌گذارد. این روش در

detector ... برای این ... که ... شبکه ... نشان می‌دهد. اساس این روش back Propagation است.

قدمت این روش برای ... در ... می‌شود و ... استفاده می‌کنند.

Guided Grad CAM؛ این روش با ترکیب ... Class Activation Maps ... به ما نشان می‌دهد که کدام class

prediction ... این روش ... شبکه ... را نشان می‌دهد ... مثل این prediction

از global average pooling استفاده می‌کنند. روش دیگری به نام Grad CAM ، CAM

Class discriminative است ... که بیشتر ... می‌شود و نشان می‌دهد ... نسبت به ... این روش ترکیبی

از تلفیق است ... حمل ... Guided Grad CAM ... ترکیبی از Grad CAM و

Guided Back Prop. است. با این ترکیب که در ... نقش ... با تجمیع نشان می‌دهد و هم ...

در ... از back propagation ... ReLU ... استفاده می‌شود.