

(1) سوال

(1.1)

ازوما خیر، چون بر اساس توزیع داده ممکن است داده های بیشتر هم که بدھیم توزیع درستی نداشته باشند. یک مثال ساده مثلا اگر در تشخیص سگ و گربه تعداد زیادی سگ داشته باشیم هر چقدر هم که دیتا زیاد باشد مدل به سمت تشخیص سگ اورفیت میکند و در تشخیص گربه عاجز است. دلیل دیگر اصلا پیچیدگی خود مدل است که اگر خیلی پیچیده باشد نسبت به دیتای ورودی آن ، هر چقدر هم دیتا بیاوریم اورفیت میکند. در ضمن اگر آموزش زیاد طولانی شود هم همچنان مشکل اورفیت ممکن است به وجود بیاید.

(1.2)

درست است، زیرا طبق تعریف اورفیت این مشکل در زمانی پیدا میشود که مدل روی دیتایی که ترین شده خیلی دقیق شده باشد ولی وقتی به سراغ دیتای دنیای واقعی می رویم (مثلا اعتبارسنجی یا تست) مدل دیگر خوب عمل نمی کند و قابلیت generalization ندارد.

(1.3)

درست است زیرا L1 باعث میشود برخی از وزن ها دقیقاً به صفر برسند و در نتیجه برخی ویژگی ها حذف شوند، که این به معنای کاهش ابعاد مدل و انتخاب ویژگی است. این خاصیت باعث میشود مدل به صورت اسپارس (sparse) شود و فقط ویژگی های مهمتر باقی بمانند. اما L2 وزن ها را به سمت صفر کوچک میکند اما هیچ کدام را دقیقاً صفر نمی کند، بنابراین همه ویژگی ها در مدل باقی میمانند ولی با وزن های کوچکتر. این باعث کاهش پیچیدگی مدل و جلوگیری از بیش برازش میشود اما حذف ویژگی انجام نمی دهد.

(1.4)

غلط است. زیرا در زمان استنتاج اگر بخواهیم در اپ اوت بزنیم ممکن است نتایج مختلفی بگیریم که اصلا مطلوب نیست، بلکه فقط در زمان آموزش برای دلایلی مثل اینکه وزن همه نورن ها بهتر آپدیت شوند و فیچر های بهتری پیدا شوند و حتی گاهی مدل زیاد پیچیده نشود استفاده می شود.

(1.5)

غلط. یک روش regularization هست منتها برای جلوگیری از overfit است.

(1.6)

درست است. زیرا همه داده ها را باهم یا به صورت دسته ای نمی بینیم و لزوما داریم هر دفعه یک نقطه را در نظر میگیریم که همان روش گردایان نزولی تصادفی می شود.

(1.7)

غلط. نرمال سازی اتفاقا در activation function هایی مثل tanh یا sigmoid به آموزش سریع تر کمک میکند.

(1.8)

منظور سوال را دقیق متوجه نشدم ولی پاسخ آن را توضیح میدهم. برای تقسیم بندی داده اگر داده ها خیلی زیاد باشند درصد کمی از آن برای تست می رود ولی اگر داده تعداد معمولی داشته باشد حدود 20 درصد آن مثلا برای تست میروند. از طرفی دیگر برای

مجموعه‌ی validation یا همان dev روشی داریم به نام cross validation که هی زیر مجموعه‌های مختلفی استفاده می‌کنیم تا مطمئن شویم روی قسمتی از داده متمرکز نشده ایم و مذل واقعاً نتیجه خوبی داشته باشیم.

(1.9) غلط است. به گرادیان‌های اخیر وزن بیشتری داده می‌شود و وزن گرادیان‌های قدیمی‌تر به صورت نمایی کاهش می‌یابد.

(2) سوال

(2.1) چند دلیل دارد که بعضی از آن‌ها به این شرح اند: خروجی آن vanishing gradient است. در حول نقطه‌ی صفر مشکل آن نسبت به sigmoid کمتر است (در بقیه نقاط لزومنه) بنابراین برای داده نرمال میتواند بهتر باشد. همگرایی سریع تری دارد. مشتق آن حداقل ۱ میباشد در صورتی که در در sigmoid برابر $\frac{1}{4}$ است.

(2.2) زیرا در هر مرحله آموختش به صورت تصادفی برخی نورون‌ها را خاموش می‌کند، بنابراین ساختار شبکه در هر اپوک کمی مقاومت است. به دلیل حذف تصادفی نورون‌ها، مقدار هزینه محاسبه شده در هر اپوک شامل نویز بیشتری است و ممکن است کاهش یا افزایش‌های موقتی نشان دهد، درواقعتابع هزینه فرمول مشخصی ندارد در هر مرحله (حداقل از چشم‌ما).

(2.3) زیرا به داده‌های جدید وزن بیشتری داده می‌شود و داده‌های قدیمی‌تر به صورت نمایی وزن کمتری دارند. این باعث می‌شود تغییرات ناگهانی و نویز‌های کوتاه‌مدت که در داده‌های اخیر دیده می‌شوند، کمتر تأثیرگذار باشند و روند کلی داده‌ها بهتر نمایان شود، این وزن‌دهی نمایی باعث می‌شود که فیلتر به عنوان یک فیلتر پایین‌گذر (Low-pass filter) عمل کند و نوسانات سریع و نویزهای با فرکانس بالا را کاهش دهد.

(2.4) میتواند از exploding gradient و نیز vanishing gradient جلوگیری کند. حتی میتواند به همگرایی سریعتر منجر شود اگر وزن‌ها درست انتخاب شوند. همچنین با ایجاد عدم تقارن بین نورون‌ها، مقداردهی تصادفی و مناسب وزن‌ها باعث می‌شود نورون‌ها در لایه‌های مختلف متقاومت عمل کنند و یادگیری بهتر انجام شود.

(3) سوال

$$v_i = \beta * v_{i-1} + (1 - \beta) * \theta_i$$

$$v_1 = 0.9 * 22 + 0.1 * 20 = 21.8$$

$$v_2 = 0.9 * 24 + 0.1 * 21.8 = 23.78$$

$$v_3 = 0.9 * 20 + 0.1 * 23.78 = 20.378$$

(4) سوال

$$v_4 = 0.9 * 19 + 0.1 * 20.378 = 19.1378$$

$$v_5 = 0.9 * 21 + 0.1 * 19.1378 = 20.81378$$

حساسیت به این صورت است که به صرارت تخمینی ما داریم در هر مرحله حدود $1 - \beta$ داده را در نظر میگیریم و در نتیجه با افزایش آن به سمت 1 درواقع تعداد بیشتری از داده های اخیر را در نظر میگیریم و برای همین منحنی نرم تر و به نویز کم حساسیت تر میشود و بر عکس اگر بتا به سمت صفر برود خیلی به نویز حساس میشود و تعداد داده های اخیر کمتری را در نظر میگیرد. بنابراین مثلا بتای 0.5 به نویز حساس تر از 0.9 است.