

- مشکل گرادیان محو شونده یک مسئله رایج است که در شبکه های عصبی عمیق رخ می دهد. این زمانی اتفاق می افتد که شیب های تابع $loss$ کوچک می شوند، زیرا در طول آموزش در لایه های شبکه منتشر می شوند، و در نهایت منجر به این می شود که وزن لایه های قبلی اصلاً به روز نمی شود. این می تواند باعث شود که این لایه ها اساساً "محو شوند" یعنی همانطور که لایه های بیشتری با استفاده از توابع فعال سازی خاص به شبکه های عصبی اضافه می شوند، گرادیان $loss$ function به صفر نزدیک می شود و آموزش شبکه را سخت می کند. این به این دلیل است که برخی از توابع فعال سازی، مانند تابع سیگموئید، یک فضای ورودی بزرگ را به یک فضای ورودی کوچک بین ۰ و ۱ نظیر می کنند. بنابراین، تغییر زیاد در ورودی تابع سیگموئید باعث تغییر کوچکی در خروجی می شود. از این رو، مشتق کوچک می شود. برای شبکه های کم عمق با تنها چند لایه که از این فعال سازی ها استفاده می کنند، این مشکل بزرگی نیست. با این حال، زمانی که از لایه های بیشتری استفاده می شود، می تواند باعث شود که گرادیان برای اینکه بتواند آموزش به طور موثر کار کند بسیار کوچک باشد. گرادیان های شبکه های عصبی با استفاده از پس انتشار یافت می شوند. به زبان ساده، پس انتشار مشتقات شبکه را با حرکت لایه به لایه از لایه نهایی به لایه اولیه پیدا می کند. بر اساس قانون زنجیره، مشتقات هر لایه در شبکه ضرب می شوند (از لایه نهایی تا لایه اولیه) تا مشتقات لایه های اولیه را محاسبه کنند. با این حال، وقتی n لایه پنهان از یک فعال سازی مانند تابع سیگموئید استفاده می کنند، n مشتق کوچک با هم ضرب می شوند. بنابراین، با انتشار به لایه های اولیه، گرادیان به صورت تصاعدی کاهش می یابد. یک گرادیان کوچک به این معنی است که وزن ها و سوگیری های لایه های اولیه به طور موثر در هر $train$ به روز نمی شوند. از آنجایی که این لایه های اولیه اغلب برای شناسایی عناصر اصلی داده های ورودی بسیار مهم هستند، می تواند منجر به عدم دقت کلی کل شبکه شود.

راه حل ها:

ساده ترین راه حل استفاده از توابع فعال سازی دیگر است، مانند $ReLU$ ، که مشتق کوچکی ایجاد نمی کند. لایه های عادی سازی دسته ای $batch$ normalization layers نیز می توانند مشکل را حل کنند. همانطور که قبلاً گفته شد، مشکل زمانی ایجاد می شود که یک فضای ورودی بزرگ به یک فضای کوچک نگاشت می شود و باعث ناپدید شدن مشتقات می شود. تکنیک های منظم سازی $regularization$ techniques مانند انصراف $dropout$ یا منظم سازی $L1/L2$ به نوبه خود می تواند به بهبود عملکرد شبکه و کاهش گرادیان های محوشونده کمک کند. راه حل دیگر استفاده از تکنیک های اولیه سازی وزن است که می تواند به اطمینان حاصل شود که گرادیان ها ناپدید نمی شوند. یکی از روش های رایج، مقداردهی اولیه $Xavier$ نام دارد که وزن اولیه هر نورون را به گونه ای تنظیم می کند که به جلوگیری از کوچک شدن گرادیان ها کمک می کند.

- انفجار گرادیان مشکلی است که در یادگیری عمیق زمانی رخ می دهد که شیب ها در طول $back$ propagation بیش از حد بزرگ می شوند. این می تواند باعث شود که وزن شبکه در هر تکرار به شدت به روز شود که می تواند منجر به بی ثباتی و عملکرد ضعیف شود. در شبکه های عمیق یا شبکه های عصبی مکرر، گرادیان های خطا می توانند در طول به روزرسانی جمع شوند و منجر به گرادیان های بسیار بزرگ شوند. اینها به نوبه خود منجر به به روزرسانی های بزرگ برای وزن شبکه و به نوبه خود یک شبکه ناپایدار می شود. در نهایت، مقادیر وزن ها می توانند آنقدر بزرگ شوند که سرریز شده و به مقادیر NaN منجر شوند. انفجار از طریق رشد نمایی با ضرب مکرر گرادیان در لایه های شبکه که مقادیر بزرگتر از ۱.۰ دارند، رخ می دهد. در شبکه های عصبی مکرر، شیب های انفجاری می تواند منجر به شبکه ای ناپایدار شود که قادر به یادگیری از داده های آموزشی نیست و در بهترین حالت شبکه ای که نمی تواند از طریق توالی های ورودی طولانی اطلاعات یاد بگیرد.

راه حل ها:

استفاده از برش گرادیان $gradient$ clipping که شامل تعیین حداکثر آستانه برای گرادیان ها است. اگر گرادیان ها از این آستانه فراتر رفت، آنگاه آنها را کاهش می دهند تا باعث به روز شدن بیش از حد وزن ها نشوند. روش دیگر استفاده از نرمال سازی دسته ای $batch$ normalization است که فعال سازی هر لایه را قبل از انتقال به لایه بعدی نرمالایز می کند. این می تواند به کاهش بزرگی گرادیان ها و بهبود پایداری شبکه کمک کند. همچنین می توان از الگوریتم های بهینه سازی مختلف مانند $Adam$ یا $Adagrad$ استفاده کرد، که برای مدیریت موثرتر انفجارهای گرادیان نسبت به روش های بهینه سازی سنتی طراحی شده اند. یکی از این روش ها تنظیم وزن $weight$ regularization است که شامل اضافه کردن یک عبارت جریمه به تابع $loss$ است که مانع از وزنه های بزرگ می شود. $L1$ (absolute weights) or an $L2$ (squared weights) penalty، این می تواند به جلوگیری از برازش بیش از حد شبکه کمک کند و شیب ها را تحت کنترل نگه دارد. در

شبکه‌های عصبی عمیق، گرادیان‌های انفجاری ممکن است با طراحی مجدد شبکه برای داشتن لایه‌های کمتر مورد بررسی قرار گیرند. همچنین ممکن است استفاده از اندازه دسته کوچکتر در حین آموزش شبکه مزایایی داشته باشد. در شبکه‌های عصبی مکرر، به‌روزرسانی در گام‌های کمتر زمانی قبلی در طول آموزش، که به آن **truncated Backpropagation** در طول زمان گفته می‌شود، ممکن است مشکل گرادیان انفجاری را کاهش دهد. شیب‌های انفجاری را می‌توان با استفاده از واحدهای حافظه بلند مدت (LSTM) و احتمالاً ساختارهای نورونی نوع دروازه‌ای مرتبط کاهش داد. اتخاذ واحدهای حافظه LSTM بهترین روش جدید برای شبکه‌های عصبی مکرر برای پیش‌بینی توالی است.

-۲-

اگر مقدار عددی دقیق یک ویژگی مهم نیست و تأثیر آن فقط به محدوده‌ای که در آن قرار دارد بستگی دارد، می‌توانیم از مقیاس‌گذاری ویژگی برای تغییر مقادیر ویژگی استفاده کنیم. مقیاس‌بندی ویژگی تکنیکی است که برای استاندارد کردن دامنه ویژگی‌ها در یک مجموعه داده استفاده می‌شود، به طوری که اندازه‌های مشابهی داشته باشند.

روش‌های مختلفی برای مقیاس‌بندی ویژگی وجود دارد، از جمله نرمال‌سازی و استانداردسازی.

عادی‌سازی شامل مقیاس‌بندی مقادیر یک ویژگی در محدوده‌ای بین ۰ و ۱ با استفاده از فرمول است:

$$X' = (X - X_{\min}) / (X_{\max} - X_{\min})$$

که در آن X مقدار جدید مقیاس شده، X مقدار اصلی، X_{\min} حداقل مقدار ویژگی، و X_{\max} حداکثر مقدار ویژگی است.

استانداردسازی شامل مقیاس‌بندی مقادیر یک ویژگی به گونه‌ای است که میانگین ۰ و انحراف معیار ۱ با استفاده از فرمول:

$$X' = (X - \text{mean}(X)) / \text{std}(X)$$

که در آن X مقدار جدید مقیاس شده، X مقدار اصلی، $\text{mean}(X)$ میانگین ویژگی، و $\text{std}(X)$ انحراف استاندارد ویژگی است.

نرمال‌سازی و استانداردسازی هر دو می‌توانند در تغییر مقادیر ویژگی مؤثر باشند به طوری که اندازه‌های مشابهی داشته باشند و راحت‌تر قابل مقایسه و تجزیه و تحلیل باشند. انتخاب روش مورد استفاده به نیازهای خاص مسئله و ماهیت مجموعه داده بستگی دارد.

-۳-

مشکل مجموعه داده‌های نامتعادل، این است که مدل‌های یادگیری ماشینی می‌توانند نسبت به طبقه اکثریت (در این مورد، افراد سالم) تعصب داشته باشند. این به این دلیل است که مدل در معرض داده‌های طبقه اکثریت قرار می‌گیرد و ممکن است تشخیص صحیح طبقه اقلیت را نیاموزد. این می‌تواند منجر به عملکرد ضعیف در هنگام پیش‌بینی طبقه اقلیت شود، که اغلب طبقه مورد علاقه است (در این مورد، افراد مبتلا به بیماری).

برای رفع این مشکل می‌توان از چندین تکنیک استفاده کرد:

کم نمونه برداری: حذف تصادفی برخی از نمونه‌های کلاس اکثریت به طوری که مجموعه داده متعادل شود. این ممکن است منجر به از دست دادن اطلاعات شود، اما می‌تواند به مدل کمک کند تا طبقه اقلیت را تشخیص دهد.

Oversampling: ایجاد نمونه‌های مصنوعی از کلاس اقلیت برای افزایش نمایش آن در مجموعه داده. این را می‌توان با استفاده از تکنیک‌هایی مانند **SMOTE** (تکنیک نمونه برداری بیش از حد اقلیت مصنوعی) یا **ADASYN** (نمونه‌گیری مصنوعی تطبیقی) انجام داد. نمونه برداری بیش از حد می‌تواند از نظر محاسباتی گران باشد، اما می‌تواند به مدل کمک کند تا طبقه اقلیت را بهتر تشخیص دهد.

یادگیری حساس به هزینه: تغییر تابع ضرر الگوریتم یادگیری ماشین برای در نظر گرفتن عدم تعادل در مجموعه داده. این را می‌توان با اختصاص وزن‌های بالاتر به نمونه‌های کلاس اقلیت انجام داد.

روش‌های مجموعه: ترکیب چندین مدل یادگیری ماشین آموزش‌دیده بر روی نسخه‌های مختلف مجموعه داده (مثلاً با تکنیک‌های مختلف کم‌نمونه‌سازی یا بیش‌نمونه‌سازی) برای بهبود عملکرد کلی.

توجه به این نکته مهم است که انتخاب تکنیک به مجموعه داده‌ها و مسئله خاص در دست بستگی دارد، و ممکن است برای یافتن بهترین راه حل نیاز به آزمایش با رویکردهای مختلف باشد.

۴-

روش های مختلفی برای این منظور وجود دارد.

One-hot/dummy encoding

داده های طبقه بندی شده به صورت بردارهای صفر و یک نمایش داده می شوند. این کار با استفاده از یک متغیر ساختگی مجزا برای هر دسته انجام می شود و در صورتی که مشاهده متعلق به آن دسته باشد، مقدار متغیر ساختگی را ۱ و در غیر این صورت ۰ قرار می دهیم. برای داده های NOMINAL بسیار خوب و محبوب است.

Label / Ordinal encoding

این نوع رمزنگاری برای داده هایی که نسبت به هم اولویت دارند به کار میرود (همان مثال تهران ۱۰) و برای هر دسته یک عدد در نظر گرفته میشود.

سایر روش ها::

Frequency / count encoding-Binary Encoding-Feature Hashing-Target Encoding...

با عدد دادن به شهر ها عملاً داریم برتری ناخواسته ایجاد میکنیم و مدل را به اشتباه می اندازیم. چرا که مدل اعداد را بر اساس مقدارشان در نظر میگیرد نه به عنوان ایندکس. پس اگر م به تهران عدد ۱۰ بدهیم و به اردبیل عدد ۱ بدهیم و به همدان عدد ۳۲ بدهیم یعنی ارزش خانه ها در همدان بیشتر از تهران و بیشتر از اردبیل است که این درست نیست و حتی اگر بخواهیم با مقداردهی این کار را انجام دهیم باید بر اساس ملاک هایی مثل کلان شهر بودن، پایتخت بودن و... این اعداد داده شود. یعنی اگر همه موارد دو خانه یکسان باشند و تنها شهرشان فرق کند قیمت کدام بیشتر خواهد بود و بر این اساس شماره بدهیم.

۵-

جلوگیری از بیش‌برازش (overfitting) :

۱ Hold-out (data) .

به جای استفاده از تمام داده‌هایمان برای آموزش، می‌توانیم به سادگی مجموعه داده‌های خود را به دو مجموعه تقسیم کنیم: آموزش و تست. یک نسبت تقسیم معمول ۸۰ درصد برای آموزش و ۲۰ درصد برای تست است. ما مدل خود را تا زمانی آموزش می دهیم که نه تنها در مجموعه آموزشی بلکه برای مجموعه تست نیز عملکرد خوبی داشته باشد. این نشان دهنده قابلیت تعمیم خوب است زیرا مجموعه آزمایشی داده های دیده نشده را نشان می دهد که برای آموزش استفاده نشده اند. با این حال، این رویکرد به یک مجموعه داده به اندازه کافی بزرگ برای آموزش حتی پس از تقسیم نیاز دارد.

۲ Cross-validation (data) .

ما می‌توانیم مجموعه داده‌های خود را به k گروه تقسیم کنیم . (k-fold cross-validation) در این روش اجازه می دهیم یکی از گروه ها مجموعه تست باشد و بقیه به عنوان مجموعه آموزشی، و این روند را تا زمانی که هر گروه جداگانه به عنوان مجموعه تست استفاده شود) به عنوان مثال، k تکرار (تکرار شود. برخلاف Hold-out ، این روش اجازه می‌دهد تا در نهایت از تمام داده‌ها برای آموزش استفاده شود، اما همچنین از نظر محاسباتی گران‌تر از Hold-out است.

۳ Data augmentation (data) .

یک مجموعه داده بزرگتر باعث کاهش بیش از حد بیش‌برازش می‌شود. اگر نمی‌توانیم داده‌های بیشتری جمع‌آوری کنیم و محدود به داده‌هایی هستیم که در مجموعه داده فعلی خود داریم، می‌توانیم افزایش داده‌ها (data augmentation) را برای افزایش مصنوعی اندازه مجموعه داده‌مان اعمال کنیم. به عنوان مثال، اگر در حال آموزش برای یک کار طبقه‌بندی تصویر هستیم، می‌توانیم تغییر شکل‌های تصویری مختلفی را در مجموعه داده‌های تصویر خود انجام دهیم (به عنوان مثال، چرخش، تغییر مقیاس، شیفت).

۴. Feature selection (data)

اگر فقط تعداد محدودی از نمونه‌های آموزشی داریم که هر کدام دارای تعداد زیادی ویژگی هستند، باید فقط مهم‌ترین ویژگی‌ها را برای آموزش انتخاب کنیم تا مدل ما نیازی به یادگیری برای این همه ویژگی نداشته باشد. ما می‌توانیم به سادگی ویژگی‌های مختلف را آزمایش کنیم، مدل‌های فردی را برای این ویژگی‌ها آموزش دهیم، و قابلیت‌های تعمیم را ارزیابی کنیم، یا از یکی از روش‌های مختلف انتخاب ویژگی استفاده کنیم.

۵. L1 / L2 regularization (learning algorithm)

رگولاریزیشن تکنیکی برای محدود کردن شبکه ما از مدلی است که بسیار پیچیده است. در رگولاریزیشن L1 یا L2، می‌توانیم یک عبارت جریمه را بر روی تابع هزینه اضافه کنیم تا ضرایب تخمینی را به سمت صفر برسانیم (و شدیدتر را نگیریم). رگولاریزیشن L2 به وزن‌ها اجازه می‌دهد تا به صفر کاهش یابند، اما نه به صفر، در حالی که L1 به وزن‌ها اجازه می‌دهد تا به صفر برسند.

۶. Remove layers / number of units per layer (model)

همانطور که در رگولاریزیشن L1 یا L2 ذکر شد، یک مدل بیش از حد پیچیده ممکن است دچار بیش‌برازش شود. بنابراین، ما می‌توانیم به طور مستقیم پیچیدگی مدل را با حذف لایه‌ها کاهش دهیم و اندازه مدل خود را کاهش دهیم. ما ممکن است پیچیدگی را با کاهش تعداد نورون‌ها در لایه‌های کاملاً متصل کاهش دهیم.

۷. Dropout (model)

با اعمال dropout، که نوعی منظم‌سازی است، در لایه‌های خود، زیر مجموعه‌ای از واحدهای شبکه خود را با احتمال مجموعه نادیده می‌گیریم. با استفاده از dropout، می‌توانیم یادگیری وابسته به هم را در بین واحدها کاهش دهیم، که ممکن است به بیش‌برازش منجر شده باشد. با این حال، با dropout، ما به آموزش بیشتری برای همگرایی مدل خود نیاز خواهیم داشت.

۸. Early stopping (model)

توقف زودهنگام یک تکنیک منظم‌سازی (رگولاریزیشن) برای شبکه‌های عصبی عمیق است که هنگامی که به‌روزرسانی‌های پارامتر در مجموعه اعتبارسنجی دیگر شروع به بهبود نمی‌کنند، آموزش را متوقف می‌کند. در اصل، ما بهترین پارامترهای فعلی را در طول آموزش ذخیره و به روز می‌کنیم، و زمانی که به روز رسانی پارامترها دیگر بهبودی حاصل نمی‌شود (پس از یک تعداد تکرار مشخص)، آموزش را متوقف می‌کنیم و از آخرین بهترین پارامترها استفاده می‌کنیم. این روش با محدود کردن روند بهینه‌سازی به حجم کمتری از فضای پارامتر، به عنوان یک تنظیم‌کننده عمل می‌کند.

جلوگیری از کم‌برازش (underfitting):

۱. افزایش پیچیدگی مدل

مدل ممکن است صرفاً به این دلیل که برای ثبت الگوها در داده‌ها به اندازه کافی پیچیده نیست، مناسب نیست. استفاده از یک مدل پیچیده تر، به عنوان مثال با تغییر از یک مدل خطی به یک مدل غیر خطی یا با افزودن لایه‌های پنهان به شبکه عصبی، اغلب کمک می‌کند.

۲. افزایش داده‌ها:

افزایش داده ها شامل ایجاد داده های آموزشی جدید با اعمال تبدیل به داده های موجود است. به عنوان مثال، می توانید تصاویر را به صورت افقی یا عمودی ورق بزنید، به داده ها نویز اضافه کنید، یا تصاویر را برای ایجاد نمونه های آموزشی جدید برش دهید. این می تواند به مدل کمک کند تا الگوهای پیچیده تر را یاد بگیرد و بهتر تعمیم دهد.

۳. تغییر معماری مدل:

به جای دستکاری مدل موجود، می توان از یک مدل متفاوت استفاده کرد. برای مثال، اگر از مدل رگرسیون خطی استفاده می شود، از یک شبکه عصبی یا درخت تصمیم استفاده کنیم.

۴. تنظیم Hyperparameter:

فرایپارامترها پارامترهایی هستند که قبل از آموزش مدل تنظیم می شوند. تنظیم این فرایپارامترها می تواند به بهبود عملکرد مدل کمک کند. برای مثال، می توانید نرخ یادگیری، قدرت منظم سازی یا تعداد دوره ها را تنظیم کنید.

۵. مهندسی ویژگی:

مهندسی ویژگی شامل انتخاب و تبدیل ویژگی های مرتبط در داده های شما برای بهبود توانایی مدل برای گرفتن الگوها است. می توانید ویژگی های جدیدی را اضافه کنید که ممکن است با مشکل مرتبط باشند، ویژگی های موجود را ترکیب کنید، یا ویژگی های نامربوط یا اضافی را حذف کنید.

هنگام آموزش یک مدل یادگیری ماشین، استفاده از مجموعه داده های آموزشی و اعتبارسنجی برای اطمینان از دقت و اثربخشی مدل ضروری است. مجموعه داده های آموزشی برای آموزش مدل استفاده می شود، به این معنی که الگوریتم الگوها و روابط را از این مجموعه داده یاد می گیرد. استفاده از مجموعه داده های آموزشی به اندازه کافی بزرگ و متنوع برای اطمینان از اینکه مدل به اندازه کافی آموزش داده شده است و می تواند به داده های جدید و دیده نشده تعمیم یابد، ضروری است. از سوی دیگر، مجموعه داده های اعتبارسنجی برای ارزیابی عملکرد مدل در طول آموزش استفاده می شود. این مجموعه داده برای نظارت بر دقت مدل در طول فرآیند آموزش و اطمینان از عدم تناسب بیش از حد مدل با مجموعه داده های آموزشی استفاده می شود. تطبیق بیش از حد زمانی اتفاق می افتد که مدل بیش از حد پیچیده باشد و مجموعه داده های آموزشی را به خاطر بسپارد، به جای یادگیری الگوهای کلی که می توانند روی داده های جدید اعمال شوند. هنگام تقسیم داده های عملیاتی، از مجموعه داده های آموزشی برای آموزش مدل استفاده می شود، در حالی که مجموعه داده های اعتبارسنجی برای ارزیابی دقت مدل در طول آموزش استفاده می شود. عملکرد مدل در مجموعه داده های اعتبارسنجی را می توان برای تنظیم پارامترهای مدل و جلوگیری از برازش بیش از حد استفاده کرد. هنگامی که مدل آموزش دید و اعتبارسنجی شد، می توان از آن برای پیش بینی داده های جدید و دیده نشده در مجموعه داده های عملیاتی استفاده کرد. بنابراین، استفاده از هر دو مجموعه داده آموزشی و اعتبارسنجی هنگام تقسیم داده های عملیاتی ضروری است تا اطمینان حاصل شود که مدل به اندازه کافی آموزش داده شده است و بر روی داده های جدید و دیده نشده به خوبی عمل می کند.