

درس یادگیری ماشین

گزارش تکلیف Linear Regression

استاد درس:

دکتر افتخاری

نگارش:

امیرحسین ابوالحسنی

شماره دانشجویی: ۴۰۰۴۰۵۰۳

پاییز 1403

فهرست مطالب

۱	مقدمه	۲
۲	رگرسیون خطی تک متغیره	۲
۱.۲	تابع هزینه رگرسیون خطی چیست؟	۲
۲.۲	آموزش رگرسیون خطی	۲
۱.۲.۲	راه حل بسته	۲
۲.۲.۲	گرادیان کاهشی تصادفی	۳
۳.۲.۲	گرادیان کاهشی دسته‌ای	۳
۳.۲	مصور سازی مدل با داده‌ها	۳
۴.۲	مقایسه پیش‌بینی مدل‌ها	۳
۵.۲	مقایسه مدل‌های آموزش دیده	۴
۶.۲	بررسی رفتار تابع هزینه	۴
۷.۲	کدام روش بهینه سازی ترجیح داده می‌شود؟	۵
۳	رگرسیون چند متغیره	۵
۱.۳	پیش پردازش داده‌ها	۵
۱.۱.۳	انکود داده‌ها	۵
۲.۱.۳	استانداردسازی	۵
۲.۳	پارامترها با راه حل بسته	۵
۳.۳	آموزش مدل	۶
۱.۳.۳	گرادیان کاهشی تصادفی	۶
۲.۳.۳	گرادیان کاهشی دسته‌ای	۶
۳.۳.۳	مقایسه روند کاهش تابع هزینه	۶
۴.۳.۳	اضافه کردن L_2 به راه حل بسته	۶

۱ مقدمه^۱

رگرسیون خطی یکی از ساده‌ترین و پرکاربردترین روش‌ها در یادگیری نظارت‌شده برای مدل‌سازی رابطه بین یک متغیر وابسته و یک یا چند متغیر مستقل است. با فرض وجود یک رابطه خطی بین ویژگی‌های ورودی و متغیر هدف، رگرسیون خطی تلاش می‌کند بهترین خط را پیدا کند که خطای پیش‌بینی را به حداقل برساند. هنگامی که تنها یک متغیر مستقل وجود دارد، این روش به عنوان رگرسیون خطی ساده شناخته می‌شود. با این حال، مشکلات دنیای واقعی اغلب شامل چندین عامل مؤثر بر متغیر هدف است. در چنین مواردی از رگرسیون خطی چندگانه استفاده می‌شود، جایی که مدل چندین متغیر مستقل را برای پیش‌بینی متغیر وابسته در نظر می‌گیرد. به صورت ریاضی، این رابطه به صورت زیر مدل‌سازی می‌شود:

$$y = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n + \epsilon$$

در اینجا y مقدار پیش‌بینی‌شده، w_0 عرض از مبدأ، w_1, w_2, \dots, w_n ضرایب متغیرهای مستقل x_1, x_2, \dots, x_n و ϵ خطای مدل است. در عمل، رگرسیون خطی می‌تواند با استفاده از تکنیک‌های مختلفی حل شود، مانند حل بسته (closed-form solution) که از کمینه‌سازی میانگین مربعات خطا (MSE) مشتق شده است یا الگوریتم‌های بهینه‌سازی مانند گرادیان نزولی. علاوه بر این، تکنیک‌های تنظیم (regularization) مانند تنظیم L_2 (رگرسیون ریدج) می‌توانند برای جلوگیری از بیش‌برازش با جریمه کردن مقادیر بزرگ ضرایب اعمال شوند. این روش‌ها به‌ویژه در مقابله با چندخطی بودن یا داده‌های پرنویز، پایداری مدل را تضمین می‌کنند.

این گزارش به بررسی رگرسیون خطی ساده و چندگانه می‌پردازد و بر حل بسته در حالت معمولی و با تنظیم L_2 ، بهینه‌سازی گرادیان نزولی در دو حالت گرادیان نزولی تصادفی^۲ و دسته‌ای^۳ تأکید دارد. همچنین عملکرد مدل در روش‌های مختلف با استفاده از داده‌های آموزش و آزمایش مقایسه شده است.

۲ رگرسیون خطی تک متغیره

۱.۲ تابع هزینه رگرسیون خطی چیست؟

تابع هزینه رگرسیون خطی با نام میانگین مربع خطاها یا Mean Square Error (MSE) شناخته می‌شود.

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y}^i)^2$$

۲.۲ آموزش رگرسیون خطی

۱.۲.۲ راه حل بسته

در حل بسته، باید رابطه زیر حل شود تا θ پیدا شود.

$$\theta = (X^T X)^{-1} X^T \vec{y}$$

البته ممکن است $X^T X$ معکوس پذیر نباشد (که در تکلیف همین مورد برای داده‌های Train اتفاق می‌افتد) که در این صورت می‌توان از Moore-Penrose pseudo-inverse در رابطه راه حل بسته استفاده کرد:

$$\theta = (X^T X)^+ X^T \vec{y}$$

پاسخ به دست آمده:

$$y = \theta_0 + \theta_1 x$$

$$\theta_0 = -3.8957, \theta_1 = 1.1930$$

^۱مقدمه با مدل GPT 3 نوشته شده است.

^۲Stochastic Gradient Descent

^۳Batch Gradient Descent

۲.۲.۲ گرادیان کاهشی تصادفی

در این روش از الگوریتم گرادیان کاهشی تصادفی برای بهینه سازی تابع هزینه استفاده می شود. پاسخ به دست آمده:

$$y = \theta_0 + \theta_1 x$$

$$\theta_0 = -3.8481, \theta_1 = 1.0570$$

۳.۲.۲ گرادیان کاهشی دسته ای

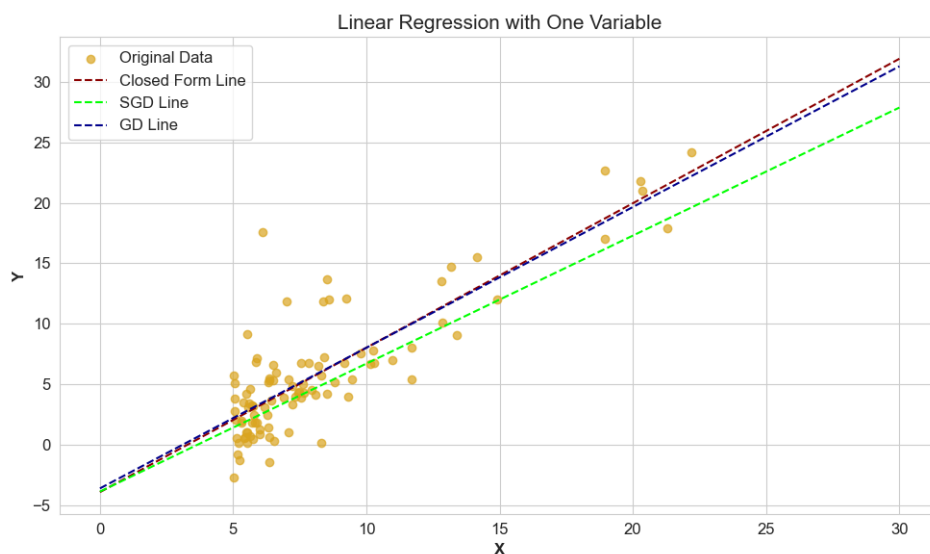
در این روش از الگوریتم گرادیان کاهشی دسته ای برای بهینه سازی تابع هزینه استفاده می شود. پاسخ به دست آمده:

$$y = \theta_0 + \theta_1 x$$

$$\theta_0 = -3.5858, \theta_1 = 1.1619$$

۳.۲ مصور سازی مدل با داده ها

هر سه پاسخ به دست آمده در کنار داده ها در شکل ۱ به تصویر کشیده شده است. همانطور که دیده می شود گرادیان کاهشی دسته ای بهتر از تصادفی عمل کرده است.



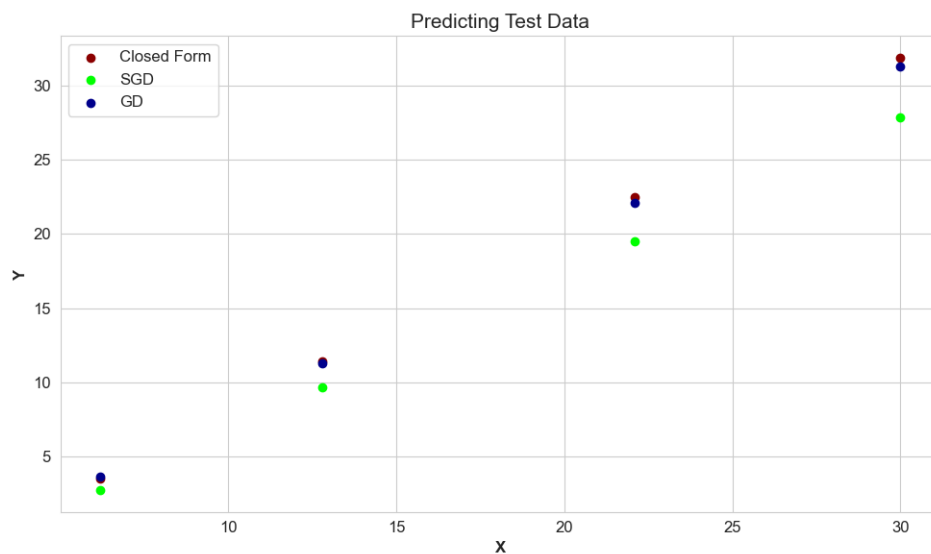
شکل ۱: خطوط به دست آمده از سه روش مختلف

۴.۲ مقایسه پیش بینی مدل ها

پس از آموزش هر مدل، برای تست کردن، به آن ها داده های

$$X = 6.2, 12.8, 22.1, 30$$

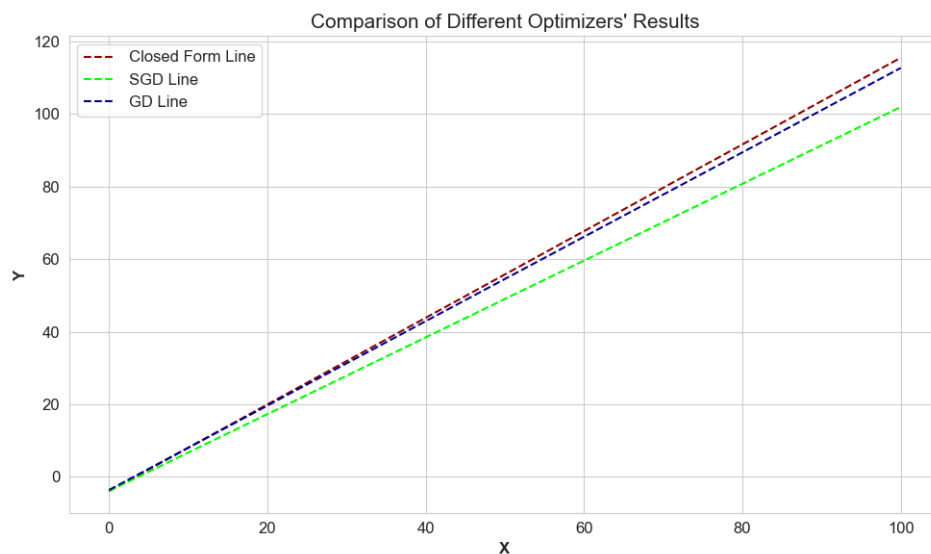
داده می شود. نزدیکی خروجی هر مدل، به مدلی که از راه حل بسته به دست آمده است می تواند مشخص کند آن شیوه بهینه سازی چقدر خوب توانسته پارامترهایی را پیدا کند که تابع هزینه را کمینه می کنند.



شکل ۲: خروجی مدل در روش

۵.۲ مقایسه مدل‌های آموزش دیده

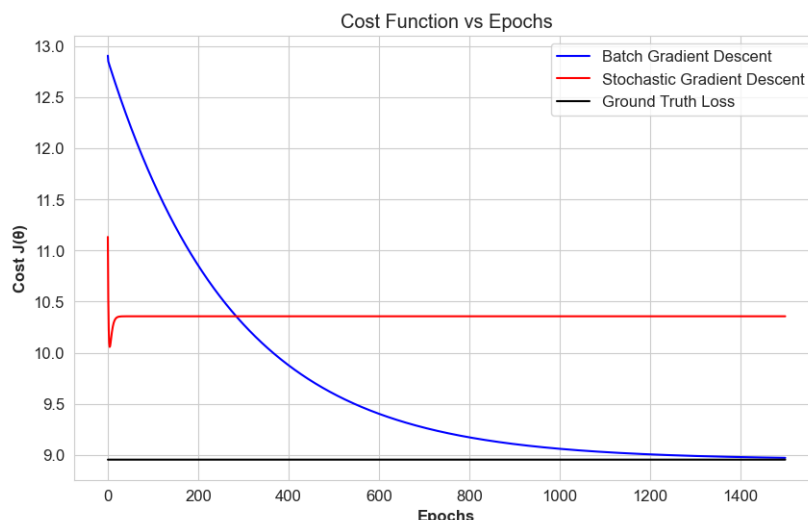
از آنجایی که پاسخ بسته بهترین پاسخی است که می‌توان به آن رسید، مقایسه هر خطی با پاسخ بسته می‌تواند نشان بدهد آن مدل چقدر خوب است. (شکل ۳)



شکل ۳: خروجی مدل در روش

۶.۲ بررسی رفتار تابع هزینه

از مهمترین روش‌های مانیتور کردن کیفیت یادگیری یک مدل، بررسی رفتار تابع هزینه در هنگام آموزش می‌باشد. همانطور که در شکل ۴ دیده می‌شود، گرادیان تصادفی بعد از ۲۰ تکرار دیگر نمی‌تواند تابع هزینه را مینیمم کند و احتمالاً در یک بهینه محلی گیرکرده است. از طرفی گرادیان کاهشی دسته‌ای توانسته به نزدیکی مقدار مینیمم کلی تابع هزینه برسد.



شکل ۴: نمودار مقدار تابع هزینه با توجه به هر تکرار در زمان یادگیری

۷.۲ کدام روش بهینه سازی ترجیح داده می شود؟

با توجه به بخش ۶.۲ و شکل ۴ می توان گفت از آنجایی که گرادیان کاهشی دسته ای نسبت به تصادفی بیشتر به کل داده ها دید دارد، بهتر می تواند مسیر خود را به سمت مینیمم تابع هزینه پیدا کند، مخصوصاً زمانی که داده ها کم هستند و میتوانند در رم باشند.

۳ رگرسیون چند متغیره

۱.۳ پیش پردازش داده ها

۱.۱.۳ انکود داده ها

بعضی از متغیرهای دیتاست کنگوریکال هستند و نیاز بود که اینها متغیرها انکود شوند. از One Hot Encoding (OHE) و Integer Encoding (IE) برای انکودینگ استفاده شد. متغیر *region* با استفاده از OHE انکود شده و برای ویژگی های *smoker*, *gender* از روش دوم انکودینگ استفاده شده است.

سوال: چرا برای این سه ویژگی از روش های متفاوت استفاده شد؟

زمانی از IE برای انکودینگ استفاده می کنیم که بدانییم ترتیب در مقادیر ویژگی مهم است. ممکن است این سوال پیش بیاید که پس چرا برای ویژگی *gender* از این روش استفاده می شود. می توان گفت عملاً تفاوتی نمی کند چون تنها دو مقدار برای این ویژگی وجود دارد. و اما زمانی که از روش OHE استفاده می کنیم دیگر ترتیب اثری ندارد و با هر ویژگی ایجاد شده که مقدار ۰ یا ۱ را دارد به طور عادلانه برخورد می شود، همانطور که از ماهیت ویژگی *region* مشخص است.

۲.۱.۳ استانداردسازی

در زمان انجام تکلیف، دیده شد که عدم استانداردسازی داده ها باعث انفجار گرادیان و پارامترها می شود. به همین علت استانداردسازی روی داده های آموزشی انجام شد و با همان پارامترها، داده های تست نیز استانداردسازی شدند.

۲.۳ پارامترها با راه حل بسته

در این بخش، چون نمی توان از $X^T X$ (که X ماتریس داده آموزشی باشد) معکوس گرفت، از Moore-Penrose pseudo-inverse استفاده می کنیم.

۳.۳ آموزش مدل

۱.۳.۳ گرادین کاهشی تصادفی

۲.۳.۳ گرادین کاهشی دسته‌ای

۳.۳.۳ مقایسه روند کاهش تابع هزینه

۴.۳ اضافه کردن L_2 به راه حل بسته