

# گزارش تحلیل رگرسیون خطی

هیوا ابوالهادی زاده ۴۰۰۴۰۵۰۰۴  
درس یادگیری ماشین

۶ آذر ۱۴۰۳

## ۱ مقدمه

در این گزارش به بررسی و پیاده‌سازی روش Linear Regression پرداخته شده است. مسائل بررسی شده شامل دو بخش زیر می‌باشند:

□ رگرسیون خطی با یک متغیر

□ رگرسیون خطی با چند متغیر

هدف از این تحلیل، پیش‌بینی مقادیر خروجی و ارزیابی دقت مدل‌ها با استفاده از روش‌های مختلف مانند حل بسته (Closed-form Solution)، روش نزول گرادیان (Gradient Descent) و تنظیم L2 است. نتایج شامل پیش‌بینی‌ها، نمودارها، و مقایسه مدل‌ها ارائه می‌شود.

## ۲ بخش اول: رگرسیون خطی با یک متغیر

### ۱.۲ فرمول تابع هزینه

برای رگرسیون خطی، تابع هزینه به صورت زیر تعریف می شود:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 \quad (1)$$

که در آن:

□  $h_{\theta}(x_i) = \theta_0 + \theta_1 x_i$ : مقدار پیش بینی شده.

□  $m$ : تعداد نمونه ها.

□  $y_i$ : مقدار واقعی خروجی.

### ۲.۲ روش های تخمین پارامترها

برای تخمین مقادیر  $\theta_0$  و  $\theta_1$ ، سه روش زیر استفاده شده است:

۱. حل بسته با روش MSE: با حل مستقیم معادلات، پارامترها محاسبه شدند.

۲. نزول گرادیان تصادفی (Stochastic GD): به ازای هر نمونه ورودی، پارامترها با استفاده از شیب محاسبه شده به روز شدند.

۳. نزول گرادیان دسته ای (Batch GD): به روزرسانی پارامترها با استفاده از میانگین گرادیان کل داده ها.

### ۳.۲ نمودارها

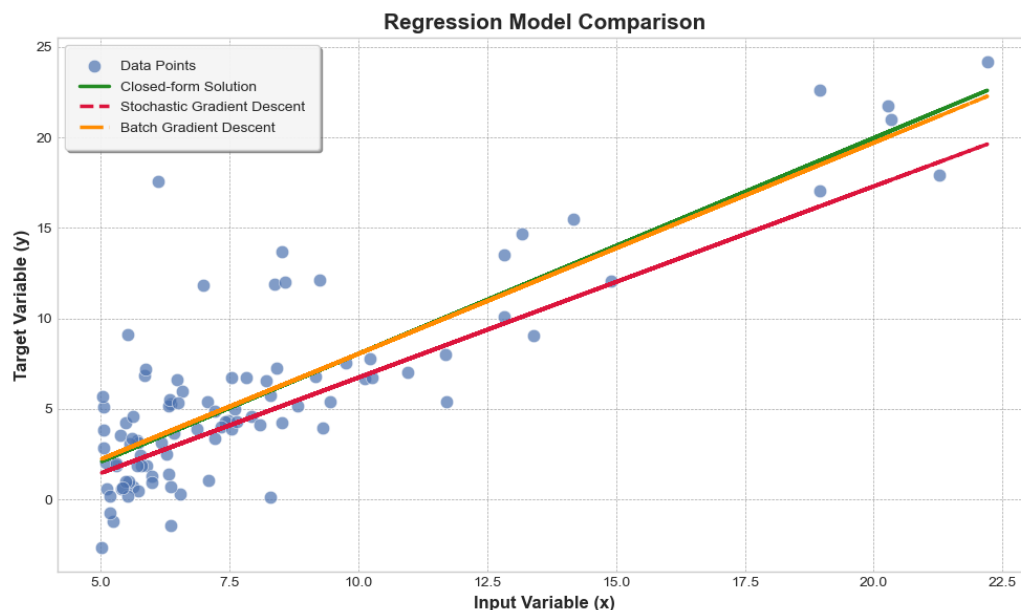
□ نمودار داده های ورودی و خروجی (Scatter Plot) رسم شده است.

□ خطوط مدل های تخمین زده شده توسط هر سه روش روی نمودار رسم شده اند.

این نمودار مقایسه ای بین سه روش مختلف برای رگرسیون خطی با یک متغیر را نشان می دهد. محور افقی نشان دهنده متغیر ورودی (x) و محور عمودی متغیر خروجی (y) است. نقاط آبی رنگ نشان دهنده داده های واقعی (نمونه ها) هستند و سه خط رنگی مربوط به مدل های تخمین زده شده است:

۱. خط سبز (حل بسته Closed-form Solution):

□ این خط نشان دهنده مدل بهینه ای است که با حل مستقیم معادلات به دست آمده است.



شکل ۱: نمودار ۱

□ همان‌طور که مشاهده می‌شود، این خط دقیقاً از مرکز داده‌ها عبور کرده و به نظر می‌رسد به بهترین شکل ممکن خطای میانگین مربعات (MSE) را حداقل کرده است.

## ۲. خط قرمز (نزول گرادیان تصادفی Stochastic Gradient Descent):

- این خط از الگوریتم نزول گرادیان تصادفی به دست آمده است.
- گرچه این روش عملکرد خوبی دارد، اما به وضوح کمی کمتر دقیق‌تر از روش حل بسته عمل کرده و خطای بیشتری را نسبت به آن نشان می‌دهد.
- علت این اختلاف می‌تواند به دلیل نوسانات ناشی از انتخاب تصادفی داده‌ها در هر تکرار باشد.

## ۳. خط نارنجی (نزول گرادیان دسته‌ای Batch Gradient Descent):

- این خط با استفاده از الگوریتم نزول گرادیان دسته‌ای (استفاده از کل داده‌ها در هر تکرار) ایجاد شده است.
- تطابق این خط با روش حل بسته بسیار زیاد است و تفاوت محسوسی مشاهده نمی‌شود، زیرا این روش گرادیان میانگین تمام داده‌ها را محاسبه می‌کند که دقت بیشتری دارد.

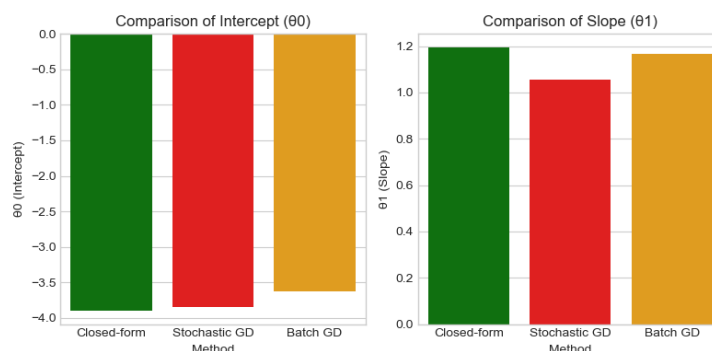
## ۳ تحلیل کلی

- روش حل بسته (Closed-form Solution): دقت بسیار بالایی دارد و بهترین خط را نسبت به داده‌ها ارائه داده است.
- نزول گرادیان دسته‌ای (Batch Gradient Descent): عملکرد مشابه روش حل بسته دارد اما ممکن است در داده‌های بسیار بزرگ به زمان بیشتری نیاز داشته باشد.
- نزول گرادیان تصادفی (Stochastic Gradient Descent): سرعت بالایی دارد، اما دقت آن نسبت به روش‌های دیگر کمتر است و معمولاً با نوسانات همراه است.

### ۱.۳ مقایسه پارامترها

مقدار  $\theta$  های تخمین زده شده توسط هر روش در یک نمودار به صورت مقایسه‌ای نشان داده شده است. این نمودار مقایسه مقادیر  $\theta_0$  (عرض از مبدأ) و  $\theta_1$  (شیب) را در سه روش Closed-form Solution، Stochastic Gradient Descent و Batch Gradient Descent نشان می‌دهد.

- مقادیر  $\theta_0$  و  $\theta_1$  در روش Closed-form Solution و Batch Gradient Descent بسیار نزدیک به یکدیگر هستند.
- روش Stochastic Gradient Descent نسبت به دو روش دیگر اختلاف اندکی دارد که به دلیل ماهیت تصادفی آن است.
- این نتایج نشان می‌دهد که روش‌های Closed-form Solution و Batch Gradient Descent دقت بیشتری دارند، در حالی که روش Stochastic Gradient Descent سرعت را به دقت ترجیح می‌دهد.



شکل ۲: مقایسه ضرایب  $\theta_0$  و  $\theta_1$  در روش‌های مختلف

## تحلیل نمودار هزینه: SGD در مقابل BGD

این نمودار هزینه  $J(\theta)$  را بر اساس تعداد دوره‌ها (Epochs) برای دو روش بهینه‌سازی نشان می‌دهد:

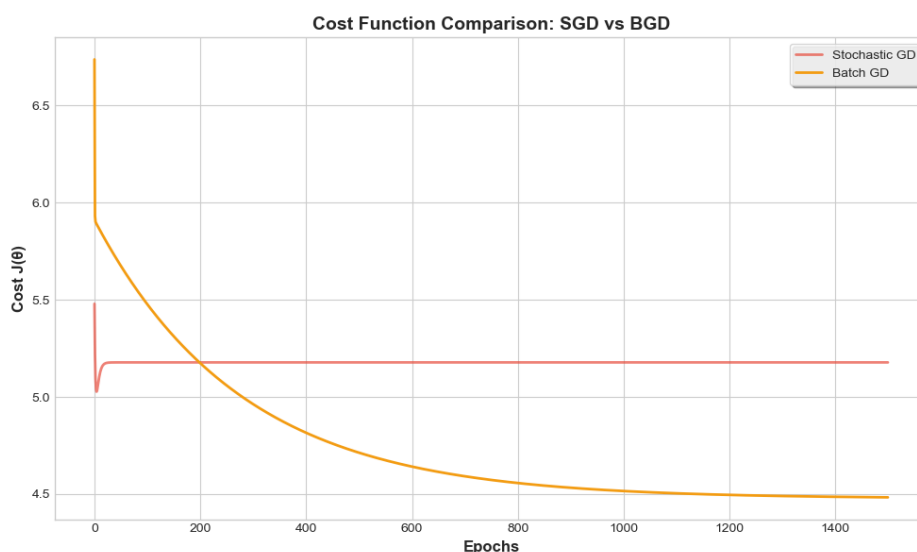
□ **Batch Gradient Descent (BGD)**: کاهش یکنواخت و پایدار هزینه.

□ **Stochastic Gradient Descent (SGD)**: نوسانات بیشتر در کاهش هزینه و رسیدن به یک مقدار ثابت بزرگ‌تر.

دلایل این رفتار:

□ در SGD، وزن‌ها برای هر نمونه داده به‌روزرسانی می‌شوند که منجر به نوسانات بیشتر می‌گردد.

□ در BGD، به‌روزرسانی وزن‌ها پس از مشاهده کل داده‌ها انجام می‌شود که نتایج پایدارتر و دقیق‌تری دارد.



شکل ۳:

## نتیجه گیری

با توجه به نمودار، روش Batch Gradient Descent به دلیل کاهش یکنواخت تر و پایدارتر هزینه و دستیابی به مقدار بهینه ی پایین تر، انتخاب بهتری نسبت به Stochastic Gradient Descent در این مسئله است.

## ۴ بخش دوم: رگرسیون خطی با چند متغیر

### ۱.۴ داده ها و پیش پردازش

داده ها شامل ویژگی هایی مانند جنسیت، منطقه، سیگار کشیدن و غیره می باشند. برای پیش پردازش:

□ متغیرهای جنسیت و سیگار کشیدن با استفاده از Integer Encoding رمز گذاری شدند.

□ متغیر منطقه با استفاده از One-Hot Encoding (OHE) رمز گذاری شد.

### ۲.۴ چرا روش One-Hot Encoding بهتر است؟

روش One-Hot Encoding برای کد گذاری ویژگی های منطقه ای مناسب تر است، زیرا برخلاف Integer Encoding، از ایجاد روابط عددی اشتباه بین مناطق جلوگیری کرده و استقلال هر منطقه را حفظ می کند.

### ۳.۴ روش های تخمین

۱. حل بسته با MSE: پارامترهای مدل با فرض  $BMI^2$  به عنوان ویژگی پایه محاسبه شدند.

۲. افزایش تدریجی داده ها: خطای تست (MSE) به ازای افزایش تدریجی تعداد نمونه های آموزشی از ۱۰ تا ۱۰۰۰ نمونه محاسبه و گزارش شده است.

۳. نزول گرادیان دسته ای و تصادفی: دو روش Batch GD و Stochastic GD برای حل مسئله استفاده شده اند.

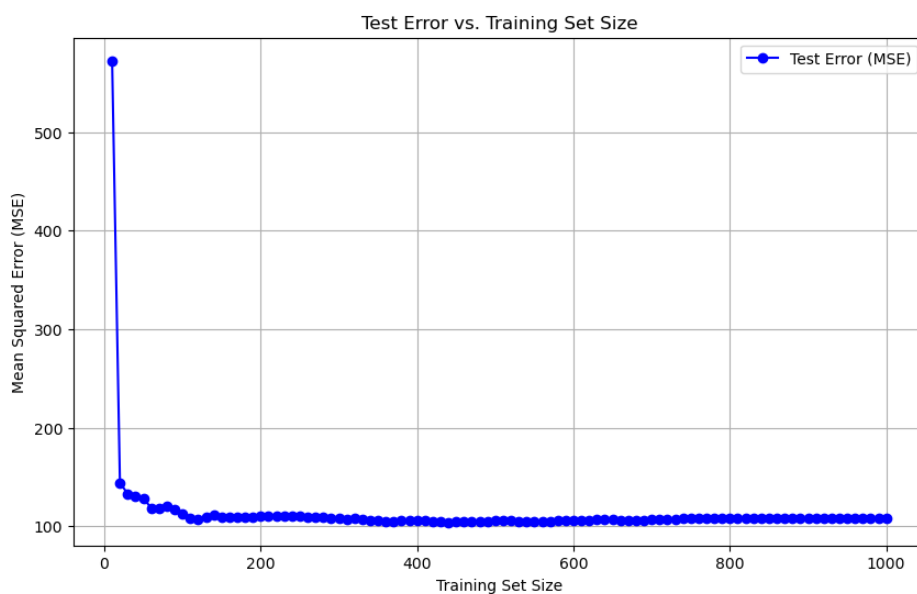
نمودار حل با روش بسته نشان می دهد که:

□ با افزایش اندازه مجموعه داده های آموزشی، مقدار خطای آزمون (MSE) به طور چشمگیری کاهش می یابد.

□ پس از یک حد مشخص از اندازه داده، خطا به مقدار ثابتی همگرا می شود.

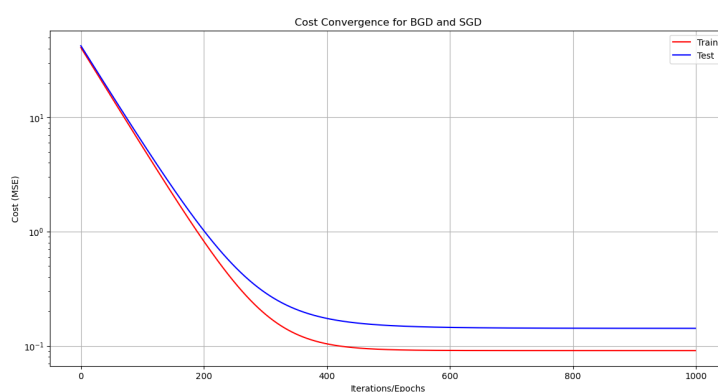
این رفتار نشان دهنده ی این است که استفاده از داده های آموزشی بیشتر می تواند دقت مدل را تا یک حد مشخص بهبود دهد، اما پس از آن، افزایش داده تأثیر قابل توجهی نخواهد داشت.

نمودار ۵ نشان می دهد که:



شکل ۴:

- خطای آموزش (Train Error) و خطای آزمون (Test Error) هر دو با افزایش تعداد تکرارها کاهش می‌یابند.
- خطاها به مقدار ثابتی همگرا می‌شوند که نشان‌دهنده بهینه‌سازی مناسب مدل است.
- فاصله‌ی بین خطای آموزش و آزمون کم است که نشان‌دهنده عدم وجود overfitting قابل توجه است.

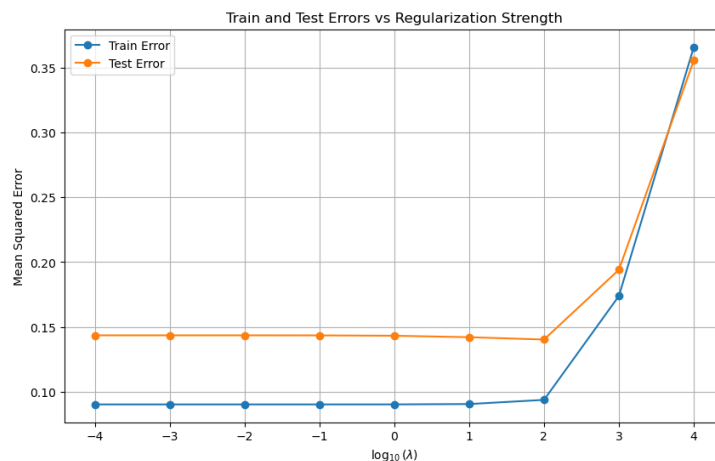


شکل ۵: نمودار ۵

## ۴.۴ تنظیم L2

یک منظم کننده L2 به تابع هزینه اضافه شده است و مقدار بهینه  $\lambda$  در بازه  $10^{-4}$  تا  $10^4$  محاسبه شده است. نمودار خطای آموزشی و تست نسبت به مقادیر مختلف  $\lambda$  رسم شده است. نمودار ۶ نشان می‌دهد:

- با افزایش قدرت منظم سازی ( $\lambda$ )، خطای آموزش به طور پیوسته کاهش می‌یابد.
- خطای آزمون ابتدا کاهش و سپس افزایش می‌یابد، که نشان‌دهنده وقوع بیش‌برازش (Overfitting) است.
- مقدار بهینه  $\lambda$  حدود  $\log_{10}(\lambda) = 0$  است، جایی که خطای آزمون کمینه می‌شود و مدل تعادل مناسبی بین پیچیدگی و تعمیم‌پذیری دارد.
- بعد از  $\log_{10}(\lambda) = 0$ ، خطای آزمون افزایش یافته و مدل به سمت کم‌برازش (Underfitting) می‌رود.



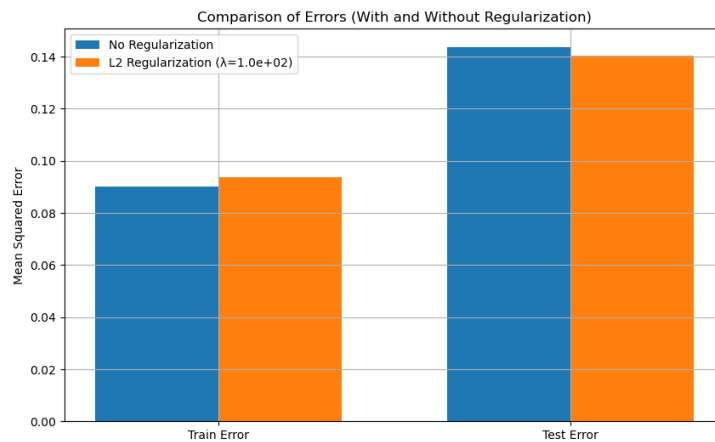
شکل ۶: اثر L2

هدف انتخاب مقدار  $\lambda$  بهینه است که خطای آزمون را به حداقل برساند.

## ۵.۴ مقایسه مدل‌ها با و بدون تنظیم

مدل تنظیم شده خطای کمتری در داده‌های تست داشته و تقریباً از Overfitting جلوگیری کرده است.





شکل ۷: اثر L2

## ۵ نتیجه گیری

□ روش Closed-form Solution سریع و دقیق است، اما برای داده‌های بزرگ محدودیت‌هایی دارد.

□ روش‌های Gradient Descent برای داده‌های بزرگ مقیاس پذیرتر هستند. روش Stochastic GD سرعت بالاتری دارد، درحالی‌که روش Batch GD دقت بیشتری ارائه می‌دهد.

□ افزودن تنظیم L2 باعث کاهش Overfitting و تعادل بهتر بین خطای آموزشی و تست شده است.