

روش حداقل مربعات با تنظیم خودکار

مهدی اکرمی و مرضیه عبدالحمیدی

تابستان ۱۴۰۱

فهرست موضوعی

- مسئله حداقل مربعات با تنظیم خودکار
- مسئله حداقل مربعات
- الگوریتم‌های برپایه گرادیان (Proximal Gradient)
- محاسبه گرادیان
- مسئله حداقل مربعات با قید خطی
- روش حداقل مربعات برای برازش داده‌ها
- تابع هدف حقیقی
- انواع توابع جریمه
- مجموعه داده‌های MNIST
- رگرسیون ریدج (Ridge Regression)
- پیاده‌سازی الگوریتم: مسئله رگرسیون ریدج
- نتایج
- منابع

مسئله حداقل مربعات با تنظیم خودکار

- **تفاوت هایپرپارامتر و پارامتر:** در یادگیری ماشین، هایپرپارامتر به پارامتری گفته می‌شود که مقدار آن، قبل از آغاز فرآیند یادگیری، تعیین می‌شود. برخلاف دیگر پارامترهای مدل، که حین فرآیند یادگیری تعیین می‌شوند. در یادگیری عمیق، هایپرپارامترها شامل متغیرهایی هستند که برای تنظیم شبکه عصبی استفاده می‌شوند، مانند رگولاریزاسیون و نرخ یادگیری.
- **هایپرپارامتر پیوسته:** هایپرپارامتری که مقادیر آن بطور پیوسته (مثلا در یک بازه) میتواند تغییر کند.
- **هایپرپارامتر در مسئله حداقل مربعات:** مسئله حداقل مربعات میتواند به هایپرپارامتر از طریق رگولاریزاسیون یا درجه چندجمله ای و ... مجهز شود.
مثال: رگرسیون ریدج
- **مزیت‌ها:**
 - ✓ انتخاب مدل بهتر
 - ✓ گذر به فراتر از تابع زیان مربعی

مسئله حداقل مربعات

- صورت کلی مسئله حداقل مربعات:

$$\underset{x}{\text{minimize}} \|Ax - b\|_2^2$$

که در آن A ماتریسی $m \times n$ و لاغر ($m > n$) است.

- پاسخ مسئله حداقل مربعات:

با فرض مستقل بودن ستون‌های ماتریس A ، جواب یکتا برای مسئله فوق بصورت زیر است:

$$\hat{x} = (A^T A)^{-1} A^T b = A^\dagger b$$

تابع هدف اصلی و عملگر proximal

- تابع هدف: یافتن ابرپارامتری که مسئله زیر را بهینه کند

$$\text{minimize } F(\omega) = \psi(\theta^{ls}(\omega)) + r(\omega)$$

- ایده حل: استفاده از روشی ابتکاری مانند proximal gradient که روشی تکراری بصورت زیر است

$$\omega^{k+1} = \text{prox}_{t^k r}(\omega^k - t^k \nabla_{\omega} \psi(\theta^{ls}(\omega^k)))$$

- تعریف عملگر proximal:

$$\text{prox}_{tr}(v) = \underset{\omega}{\operatorname{argmin}} (r(\omega) + \frac{1}{2t} \|\omega - v\|_2^2)$$

الگوریتم‌های بر پایه گرادیان (Proximal Gradient)

Given initial hyper-parameter vector $\omega^1 \in \Omega$, initial step size t^1 , number of iterations n_{iter} , tolerance ε .

for $k = 1, \dots, n_{iter}$

1. Solve the least squares problem. $\theta^{ls}(\omega^k) = \left(A^T(\omega^k) A(\omega^k) \right)^{-1} A^T(\omega^k) B(\omega^k)$

2. Compute the gradient. $g^k = \nabla_{\omega} \psi(\theta^{ls}(\omega^k))$

3. Compute the gradient step. $\omega^{k+\frac{1}{2}} = \omega^k - t^k g^k$

4. Compute the proximal operator. $\omega^{tent} = \text{prox}_{t^k r} \left(\omega^{k+\frac{1}{2}} \right)$

5. **if** $F(\omega^{tent}) \leq F(\omega^k)$

increase step size and accept update. $\omega^{k+1} = \omega^{tent}$, $t^{k+1} = 1.2t^k$

stopping criterion. **quit** if $\left\| \frac{\omega^k - \omega^{k+1}}{t^k} + (g^{k+1} - g^k) \right\|_2 \leq \varepsilon$

6. **else** decrease step size and reject update. $\omega^{k+1} = \omega^k$, $t^{k+1} = \frac{t^k}{2}$

end for

محاسبه گرادیان

- هدف: محاسبه گرادیان $g = \nabla_{\omega} \psi(\theta^{ls}(\omega))$

- فرض: مقدار $\theta = \theta^{ls}(\omega)$ محاسبه شده

حال باید عبارت‌های زیر محاسبه شوند:

$$\nabla_{\omega} \psi(\theta)$$

$$C = (A^T A)^{-1} \nabla_{\omega} \psi(\theta) \in R^{n \times m}$$

پس از انجام محاسبات و استفاده از قاعده مشتق زنجیره‌ای داریم:

$$\nabla_A \psi = (B - A\theta)C^T - AC\theta^T$$

$$\nabla_B \psi = AC$$

و نهایتاً گرادیان بصورت زیر بدست می‌آید:

$$g = \sum_{i,j} (\nabla_A \psi)_{i,j} (\nabla_{\omega} A)_{i,j} + \sum_{i,j} (\nabla_B \psi)_{i,j} (\nabla_{\omega} B)_{i,j}$$

مسئله حداقل مربعات با قید خطی

- تعمیم مسئله حداقل مربعات به حالت مقید:

$$\begin{aligned} \min \quad & \frac{\|A(\omega)\theta - B(\omega)\|_F^2}{2} \\ \text{subject to} \quad & C(\omega)\theta = D(\omega) \end{aligned}$$

در مسئله فوق داریم: $\theta \in R^{n \times m}, C: \Omega \rightarrow R^{d \times n}, D: \Omega \rightarrow R^{d \times m}$

دوتایی متغیرهای اولیه و دوگان $(\theta, \nu) \in R^{n \times m} \times R^{d \times m}$ بهینه هستند اگر و تنها اگر در شرایط KKT بصورت زیر صدق کنند:

$$\underbrace{\begin{bmatrix} 0 & A(\omega)^T & C(\omega)^T \\ A(\omega) & -I & 0 \\ C(\omega) & 0 & 0 \end{bmatrix}}_{M(\omega)} \begin{bmatrix} \theta \\ q \\ \nu \end{bmatrix} = \begin{bmatrix} 0 \\ B(\omega) \\ D(\omega) \end{bmatrix}$$

مسئله حداقل مربعات با قید خطی

تابع هدف اصلی، تابعی از متغیرهای θ, ν خواهد شد. فرض میکنیم $\nabla_{\theta}\psi, \nabla_{\nu}\psi$ قابل محاسبه باشند. ابتدا محاسبه میکنیم:

$$\begin{bmatrix} g_1 \\ g_2 \\ g_3 \end{bmatrix} = M(\omega)^{-1} \begin{bmatrix} \nabla_{\theta}\psi \\ 0 \\ \nabla_{\nu}\psi \end{bmatrix}$$

حال گرادیانی نسبت به A, B بصورت زیر قابل محاسبه است:

$$\nabla_A \psi = -(r g_1^T + g_2 \theta^T), \nabla_B \psi = g_2$$

و مشابهها برای C, D داریم:

$$\nabla_C \psi = -(\nu g_1^T + g_3 \theta^T), \nabla_D \psi = g_3$$

روش حداقل مربعات برای برازش داده‌ها

✓ در مسئله برازش داده‌ها ← در دسترس بودن داده‌های ورودی به همراه مقدار مربوط به هر ورودی.

✓ هدف ← برازش پارامترهای تابع پیش بینی کننده زیر:

$$\hat{y} = \phi(u, \omega^{feat})^T \theta$$

✓ انتخاب مدل ← باید مسئله حداقل مربعات زیر حل شود:

$$A(\omega) = \begin{bmatrix} e^{\omega_1^{data}} \phi(u, \omega^{feat})^T \\ \vdots \\ e^{\omega_N^{data}} \phi(u, \omega^{feat})^T \\ e^{\omega_1^{reg}} R_1 \\ \vdots \\ e^{\omega_d^{reg}} R_1 \end{bmatrix}, B(\omega) = \begin{bmatrix} e^{\omega_1^{data}} y_1 \\ \vdots \\ e^{\omega_N^{data}} y_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

هایپرپارامترهای وزنی $\omega^{data} \in \Omega^{data} \subseteq R^N \rightarrow$

ماتریس‌های منظم‌ساز با سایز مناسب $R_1, \dots, R_d \rightarrow$

هایپرپارامترهای منظم‌ساز $\omega^{reg} \in \Omega^{reg} \subseteq R^d \rightarrow$

مجموعه کل هایپرپارامترها $\omega = (\omega^{feat}, \omega^{data}, \omega^{reg}) \in \Omega^{feat} \times \Omega^{data} \times \Omega^{reg} \rightarrow$

تابع هدف حقیقی

فرض: داده‌های اعتبارسنجی از ورودی‌های $u_1^{val}, \dots, u_{N_{val}}^{val} \in U$ و $y_1^{val}, \dots, y_{N_{val}}^{val} \in R^m$ تشکیل شده‌اند.

پیش‌بینی: ابتدا پیش‌بینی زیر را انجام می‌دهیم:

$$\hat{y}_i^{val} = \phi(u_i^{val})\theta^{ls}(\omega)$$

در اینجا تابع $\phi(u, \omega^{feat})$ مشخص و ثابت در نظر گرفته می‌شود.

تابع هدف حقیقی: تابع هدف حقیقی ψ در مسئله برازش داده‌ها با حداقل مربعات را میانگین زیان برای پیش‌بینی داده‌های اعتبارسنجی بصورت زیر در نظر می‌گیریم:

$$\psi(\theta) = \frac{1}{N_{val}} \sum_{i=1}^{N_{val}} l\left(\hat{y}_i^{val}, y_i^{val}\right)$$

در عبارت فوق، $l: R^m \times R^m \rightarrow R$ تابع جریمه است که فرض می‌شود نسبت به متغیر اول خود مشتق پذیر است.

تابع جریمه رگرسیون

صورت کلی:

$$l\left(\hat{y}, y\right)=\pi(r)$$

که $r=\hat{y}-y$ مانده و $R \rightarrow \pi: R^m$ تابع جریمه اعمال شده روی مانده است.

رایج ترین توابع جریمه

$$\pi(r)=\begin{cases} \frac{M^2}{6}\left(1-\left(1-\frac{\|r\|_2^2}{M^2}\right)^3\right) & \|r\|_2 \leq M \\ \frac{M^2}{6} & \|r\|_2 > M \end{cases}$$

تابع دو مربعی

$$\pi(r)=\begin{cases} \|r\|_2^2 & \|r\|_2 \leq M \\ M\left(2\|r\|_2^2-M\right) & \|r\|_2 > M \end{cases}$$

تابع هوبر

$$\pi(r)=\|r\|_2^2$$

تابع مربعی

تابع جریمه طبقه‌بندی

برای طبقه‌بندی، پیش‌بینی $\hat{y} \in R^m$ متناسب با توزیع احتمال روی m برچسب بصورت زیر داده می‌شود:

$$\Pr(y = e_i) = \frac{e^{\hat{y}_i}}{\sum_{j=1}^m e^{\hat{y}_j}}, \quad i = 1, 2, \dots, m$$

پیش‌بینی \hat{y} به عنوان یک توزیع احتمال شرطی بر روی برچسب‌های \mathcal{Y} برای x داده شده، تعبیر می‌شود.

✓ تابع جریمه پیشنهادی در طبقه‌بندی:

$$l(\hat{y}, y) = -\hat{y}_i + \log\left(\sum_{j=1}^m e^{\hat{y}_j}\right) \quad i = 1, 2, \dots, m$$

تابع آنتروپی متقابل

مجموعه تست و توقف زودهنگام

افزایش تعداد هایپرپارامترها \Leftarrow افزایش ریسک برازش بیش از حد

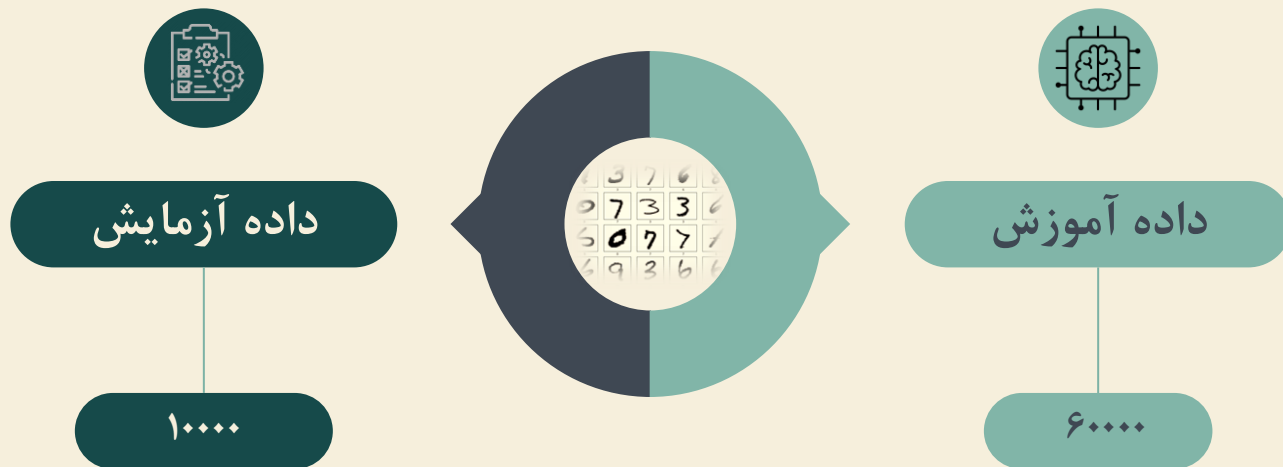
- چگونگی تشخیص برازش بیش از حد: در نظر گرفتن یک مجموعه سوم از داده‌ها و محاسبه روی آن \Leftarrow مجموعه تست

توجه!

مقدار زیان مجموعه اعتبارسنجی لزومی ندارد که تخمین خوبی از اندازه دقیق عملکرد مدل بر روی داده‌های دیده نشده باشد.

- روشی متفاوت برای رویایی با مسئله برازش بیش از حد: محاسبه مقدار زیان بر روی داده تست در هر تکرار و توقف الگوریتم در صورت افزایش مقدار زیان \Leftarrow توقف زودهنگام
 - پیش‌نیاز این روش، مجموعه داده چهارم است که پس از توقف الگوریتم، روی این داده محاسبات صورت گرفته و بعنوان عملکرد نهایی الگوریتم، گزارش می‌شود \Leftarrow مجموعه داده تست نهایی

MNIST داده‌های



رگرسیون ریج (Ridge Regression)

- مسئله رگرسیون ریج با یک هایپرپارامتر:

$$\|A\theta - B\|_F^2 + \exp(2\lambda) \|\theta\|_F^2$$

- حل مسئله با استفاده از: $\left. \begin{array}{l} \checkmark \text{ تابع آنتروپی متقابل} \\ \checkmark \text{ تابع هدف مسئله حداقل مربعات} \end{array} \right\}$

- جواب مسئله :

$$X = \begin{bmatrix} A \\ \exp(\lambda)I \end{bmatrix} \quad C = \begin{bmatrix} B \\ 0 \end{bmatrix} \quad X^T X \theta^{ls} = X^T C$$

- محاسبه مشتق θ^{ls} بر حسب λ

$$\frac{d\theta^{ls}}{d\lambda} = D^{-1} \left(\frac{dX^T}{d\lambda} C - \frac{dD}{d\lambda} (D^{-1} X^T C) \right)$$

رگرسیون ریدج (Ridge Regression)

- در نظر گرفتن تابع حقیقی آنتروپی:

$$l\left(\hat{y}, y = e_i\right) = -\hat{y}_i + \log\left(\sum_{j=1}^m e^{\hat{y}_j}\right)$$

- اعمال تابع آنتروپی متقابل روی داده‌های اعتبارسنجی و میانگین گیری:

$$\hat{y}_{val} = A_{val} \theta^{ls}, \quad \psi(\theta) = \frac{1}{N_{val}} \sum_{i=1}^{N_{val}} l\left(y_i^{val}, \hat{y}_i^{val}\right)$$

- به منظور محاسبه $g = \nabla_{\lambda} \psi$ کافی است عبارت زیر محاسبه شود:

$$\frac{dl}{d\lambda}\left(\hat{y}_j, y = e_i\right) = -\left(A_{val} \frac{d\theta^{ls}}{d\lambda}\right)_{ji} + \frac{\sum_{l=1}^m e^{\hat{y}_{li}} (A_{val} \frac{d\theta^{ls}}{d\lambda})_{li}}{\sum_{l=1}^m e^{\hat{y}_{li}}}$$

رگرسیون ریج (Ridge Regression)

- در نظر گرفتن تابع هدف مسئله حداقل مربعات به عنوان تابع حقیقی:

$$\psi(\theta^{ls}) = \|X\theta^{ls} - b\|_F^2 + \exp(2\lambda) \|\theta^{ls}\|_F^2$$

- مشتق گیری:

$$g = \frac{d\psi}{d\lambda} = 2(X\theta^{ls} - b)^T \frac{d\theta^{ls}}{d\lambda} + 2\exp(2\lambda) \|\theta^{ls}\|_F^2 + 2\exp(2\lambda)(\theta^{ls})^T \frac{d\theta^{ls}}{d\lambda}$$

پیاده‌سازی الگوریتم: رگرسیون ریج

رویکرد اول
تابع هدف حداقل مربعات

ساخت classifier + الگوریتم Proximal Gradient


رویکرد اول
آنتروپی متقابل

ساخت classifier + الگوریتم Proximal Gradient

رویکرد دوم
آنتروپی متقابل

ساخت بردار باینری به طول ۱۰ از هر برچسب + الگوریتم Proximal Gradient

نتایج

- مقایسه زمان محاسبه θ با استفاده از دو روش Least Squares و Convex Optimization:
- Convex Optimization: 15.594337 seconds
- Least Squares: 7.582242 seconds 
- مهندسی ویژگی‌ها:
 - ✓ حذف ویژگی‌هایی که کمتر از ۶۰۰ درایه فعال دارند
 - ✓ افزودن ویژگی تصادفی

نتایج

- نتایج پیاده‌سازی با رویکرد اول + تابع هدف حداقل مربعات بعنوان تابع هدف حقیقی:

Random features	Iterations	Accuracy
0	10	0.86
1200	5	0.9555
1500	5	0.9587

- نتایج پیاده‌سازی با رویکرد اول + تابع حقیقی آنتروپی متقابل:

Train set	Random features	Validation set	Test set	Iteration	Accuracy
20000	0	1000	10000	2	0.85
40000	1000	4000	10000	2	0.94
40000	1000	4000	10000	6	0.95

نتایج

- نتایج پیاده‌سازی با رویکرد دوم + تابع حقیقی آنتروپی متقابل:

Train set	Random features	Validation set	Test set	Iteration	Accuracy
40000	1200	4000	10000	5	0.9573
50000	1200	5000	10000	10	0.9576
55000	1500	5000	10000	10	0.9581
59000	2000	1000	10000	5	0.9637
59000	5000	1000	10000	5	0.9726

- نتایج پیاده‌سازی با استفاده از شبکه‌های عصبی ساده در ۱۰ تکرار:

test_loss = 0.22957107, test_accuracy = 0.9352

- نتایج پیاده‌سازی با استفاده از شبکه عصبی پیچشی در ۱۰ تکرار:

Test: (loss = 0.0445f0, acc = 98.53)

همکاری در پروژه

نام	تئوری و مرور ادبیات	پیاده سازی	گزارش نویسی	ساخت اسلاید
مرضیه عبدالحمیدی	۵۰	۵۰	۳۰	۷۰
مهدی اکرمی	۵۰	۵۰	۷۰	۳۰

منابع

- 1) S. Barratt. On the differentiability of the solution to convex optimization problems. arXiv preprint arXiv:1804.05098, 2018.
- 2) Shane Barratt, Stephen Boyd: Least Squares Auto-Tuning. arXiv:1904.05460v1, 2019.
- 3) A. Baydin, B. Pearlmutter, A. Radul, and J. Siskind. Automatic differentiation in machine learning: a survey. Journal of Machine Learning Research, 18:1-43, 2018.
- 4) R. Eigenmann and J. Nossek. Gradient based adaptive regularization. In Proc. Neural Networks for Signal Processing, pages 87-94, 1999.
- 5) Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. 2009. The elements of statistical learning: data mining, inference, and prediction. Springer Science & Business Media.
- 6) Parikh, N., and S. Boyd. 2014. "Proximal algorithms." Foundations and Trends R in Optimization 1 (3): 127-239.

سپاس از توجه شما

