



دانشگاه صنعتی شریف
دانشکده ریاضی و علوم کامپیوتر

گزارش پروژه درس بهینه‌سازی در علوم داده
روش حداقل مربعات با تنظیم خودکار

استاد
دکتر مجتبی تفاق

اعضای تیم
مرضیه عبدالحمیدی (۹۹۲۰۴۸۰۲)
مهدی اکرمی (۹۹۲۰۴۷۳۸)

تابستان ۱۴۰۱

۳	مقدمه
۴	مسئله حداقل مربعات تنظیم شونده‌ی خودکار
۵	حل مسئله حداقل مربعات
۶	حل مسئله تنظیم حداقل مربعات
۸	محاسبه گرادیان
۱۱	تعمیم مسئله حالت مقید با قید تساوی
۱۴	روش حداقل مربعات برای برازش داده ها
۱۵	تابع هدف حقیقی
۱۶	تابع جریمه رگرسیون
۱۷	تابع جریمه طبقه بندی
۱۸	وزن دهی داده ها
۱۹	منظم ساز
۲۰	مهندسی کردن ویژگی ها
۲۱	توابع مهندسی سازی ویژگی‌های اسکالر
۲۲	توابع چند بعدی برای مهندسی سازی ویژگی ها
۲۳	مجموعه تست و توقف زود هنگام
۲۴	داده های MNIST
۲۷	نتایج پیاده سازی
۲۹	منابع

روش حداقل مربعات بیش از ۲۰۰ سال پیش توسط لژاندر و گاوس معرفی شد. این روش به عنوان یکی از پرکاربردترین روش‌های محاسباتی در زمینه‌های مختلف شامل یادگیری ماشین و آمار، پردازش سیگنال، کنترل، رباتیک و مالی می‌باشد. کاربرد این روش به این دلیل است که جواب تحلیلی ساده‌ای دارد و به راحتی قابل فهم است و همچنین الگوریتم‌های پایدار و کارآمدی برای محاسبه جواب‌های این مسئله توسعه یافته است.

برای در نظر گرفتن اختلاف بین هدف مسئله کمترین مربعات و هدف اصلی مسئله ما، یک روش رایج این است که مسئله کمترین مربعات را طوری تنظیم کنیم که جواب مسئله کمترین مربعات منجر به جواب خوبی برحسب هدف اصلی ما شود. ایده‌ها در این جا شامل اصلاح کردن داده‌ها، اضافه کردن جملات به تابع هزینه (منظم سازی) یا تغییر ابرپارامترها یا وزن‌ها در مسئله حداقل مربعات است.

هدف اصلی این است که فرایند تنظیم حداقل مربعات برای گستره‌ای از کاربردهای برازش داده‌ها بصورت خودکار باشد. مسئله حداقل مربعات توسط ابرپارامترها، پارامترلیزه می‌شود و سپس بصورت خودکار با استفاده از بهینه سازی بر پایه‌ی روش‌های گرادیان‌گونه تنظیم می‌شود تا بهترین عملکرد (یا حداقل عملکرد بهتری) را بدست آوریم. این فرایند به ما اجازه می‌دهد که بصورت خودکار، فضای ابرپارامترها را جستجو کنیم که این مزیت را دارد که پارامترهای بدست آمده، بهتر از پارامترهایی که بصورت دستی تنظیم شده است، باشد و همچنین کمک می‌کند که پارامترهای مناسب را خیلی سریع‌تر از روش‌های تنظیم دستی بدست آوریم. این روش به روش حداقل مربعات با تنظیم خودکار شناخته شده است.

به عنوان یک تست کیس خاص، الگوریتم مطرح شده را بر روی داده‌های MNIST با استفاده از روش ساختن classifier ها به همراه هایپر-پارامتر λ در مساله رگرسیون ریج پیاده‌سازی می‌کنیم.

مسئله حداقل مربعات تنظیم شوندهی خودکار

مسئلهی حداقل مربعات ماتریسی که وابسته به بردار هایپرپارامترهای $\omega \in \Omega \subseteq \mathcal{R}^p$ می باشد، بصورت زیر است:

$$\text{minimize } \|A(\omega)\theta - B(\omega)\|_F^2,$$

که در اینجا، $\theta \in \mathcal{R}^{n \times m}$ متغیر بهینه سازی حداقل مربعات یا ماتریس پارامتر می باشد و $A: \Omega \rightarrow \mathcal{R}^{k \times n}$ و $B: \Omega \rightarrow \mathcal{R}^{k \times m}$ هستند. تعبیر ساده ای از رابطه فوق، این است که ماتریس های A, B بردار هایپرپارامتر ω را به داده های مسئله حداقل مربعات می نگارند. $\|\cdot\|_F$ در اینجا همان نرم فروبنیوس یعنی جذر جمع مربعات درایه های یک ماتریس است. فرض بر این است که ستون های $A(\omega)$ مستقل خطی هستند و علاوه بر این، $A(\omega)$ یک ماتریس بلند (tall) می باشد ($k \geq n$). تحت این مفروضات، جواب حداقل مربعات، یکتا و بصورت زیر تعیین می شود:

$$\theta^{ls}(\omega) = A(\omega)^\dagger B(\omega) = (A(\omega)^T A(\omega))^{-1} A(\omega)^T B(\omega)$$

که $A(\omega)^\dagger$ شبه وارون (Moore-Penrose) می باشد. $\theta^{ls}(\omega)$ همان نگاشتی است که بردار هایپرپارامترها را به $\theta^{ls}(\omega) \in \mathcal{R}^{n \times m}$ می نگارد.

در خیلی از کاربردها، بیشتر از یک تابع هدف موجود است که این چنین مسئله ها بصورت زیر بازنویسی می شود:

$$\text{minimize } \lambda_1 \|A_1(\omega)\theta - B_1(\omega)\|_F^2 + \dots + \lambda_r \|A_r(\omega)\theta - B_r(\omega)\|_F^2$$

که در عبارت فوق، λ_i ها، وزن های نسبت بوده و مسئله به فرم استاندارد زیر قابل بیان است:

$$A(\omega) = \begin{bmatrix} \sqrt{\lambda_1} A_1 \\ \vdots \\ \sqrt{\lambda_r} A_r \end{bmatrix}, \quad B(\omega) = \begin{bmatrix} \sqrt{\lambda_1} B_1 \\ \vdots \\ \sqrt{\lambda_r} B_r \end{bmatrix},$$

$$\text{minimize } \|A(\omega)\theta - B(\omega)\|_F^2$$

حل مسئله حداقل مربعات

برای ω داده شده، راه‌های زیادی برای حل مسئله موجود است که از جمله می‌توان به تجزیه QR ماتریس چگال یا تنک (sparse) و روش‌های تکراری CG (گرادیان مزدوج) و LSQR اشاره کرد که کتابخانه‌های زیادی برای حل این مسئله با استفاده از چند CPU یا GPU توسعه یافته است؛ به این صورت که مسئله را با توجه به ستون‌ها (θ) جدا کرده و با ماتریس‌های A, B بصورت موازی حل میکنند.

فرض کنید A, B بصورت ماتریس‌های چگال ذخیره شده باشند. یک راه مناسب برای پیاده سازی توسط GPU استفاده از فرم گرام $G = A^T A$ در کنار $H = A^T B$ است که در مجموع knm و kn^2 فلاپ نیاز دارد. توجه کنید که این ضرب‌ها از عملگرهای سطح ۳ BLAS هستند که بصورت کارائی محاسبه میشود. برای محاسبه θ^{ls} می‌توان از طریق تجزیه چولسکی $G = LL^T$ با هزینه n^3 فلاپ و حل معادلات مثلثی $LY = H$ و $L^T \theta^{ls} = Y$ با هزینه n^3 و $n^2 m$ فلاپ مسئله را حل کرد که با توجه به $k > n$ هزینه محاسباتی کل از مرتبه $kn(n + m)$ میشود.

اگر $A(\omega)$ و $B(\omega)$ ماتریس‌های تنک و بزرگ باشند از نمایش عملگری استفاده میشود به این صورت که $A(\omega)u$ برای هر $u \in R^n$ و $v \in R^k$ برای $A^T(\omega)v$ قابل محاسبه است (نمایش غیرماتریسی). روش‌های CG و LSQR برای حل مسئله برای هر ستون θ بصورت موازی بکار برده میشود. پیچیدگی مسئله، متغیر نسبت به داده‌ها و ابعاد و میزان تنک بودن و فرض‌های مسئله و میزان دقت مورد نیاز است.

حال به مسئله اصلی باز میگردیم. هدف یافتن ابرپارامتری است که مسئله زیر را بهینه می‌کند:

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

که $\omega \in \Omega$ و $F: \Omega \rightarrow R \cup \{+\infty\}$ است و $\psi: R^{n \times m} \rightarrow R$ تابع هدف اصلی و $r: \Omega \rightarrow R \cup \{+\infty\}$ تابع منظم ساز بر حسب ω است. وقتی $\omega \notin \Omega$ فرض میکنیم $r(\omega) = \infty$ است.

مسئله توسط A, B, ψ و r مشخص میشود. توجه شود که $r(\omega)$ ممکن است تعدادی پارامتر داشته باشد که از طریق اثرگذاری بر روی هاپرپارامترهای مسئله بر روی پارامتر θ تاثیر بگذارند که به این پارامترها، هاپر-هاپر پارامتر گویند.

این مسئله را میتوان بصورت مسئله قیددار بصورت زیر بیان کرد که در آن θ و ω که مستقل هستند توسط قید زیر به هم مرتبط میشوند:

$$\begin{aligned} &\text{minimize} \quad \psi(\theta) + r(\omega) \\ &\text{s. t.} \quad A^T(\omega)A(\omega)\theta = A^T(\omega)B(\omega) \end{aligned}$$

حل مسئله تنظیم حداقل مربعات

این مسئله در حالت کلی غیرمحدب بوده و حل دقیق آن سخت یا حتی غیرممکن میشود (یک استثنا مهم حالتی است که ω اسکالر و Ω یک بازه باشد که در آن میتوان $F(\omega)$ روی نقاط مشبک Ω را راحت محاسبه کرد).

با توجه به آنچه گفته شد ناچاریم به روش ابتکاری (heuristic) یا بهینه سازی موضعی متوسل شویم. فرض را بر این میگذاریم که A و B نسبت به ω مشتق پذیر و در نتیجه $\theta^{ls}(\omega)$ نیز مشتق پذیر است. فرض کنید که ψ مشتق پذیر بوده یعنی جمله اول تابع F مشتق پذیر است. در مورد تابع $r(\omega)$ چنین فرضی نمیگذاریم. روشی که برای حل مسئله و بدست آوردن ω استفاده میکنیم روش proximal gradient است که روشی تکراری بصورت زیر است:

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

که در آن k شماره تکرار و عملگر proximal نیز بصورت زیر است:

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

در این روش فرض بر این است که argmin وجود دارد. وقتی که یکتا نباشد هر مینیمم کننده ای را میتوان انتخاب کرد فقط نیازمند به این است که tr به راحتی قابل محاسبه باشد.

وقتی که $r = 0$ و $\Omega = R^p$ باشد روش proximal gradient همان روش گرادیان معمولی است و حالت دیگر وقتی که $\Omega \subset R^p$ و $r(\omega) = 0 \forall \omega \in \Omega$ آنگاه عملگر proximal برای tr همان تصویر روی Ω و روش proximal gradient همان روش تصویر گرادیان میشود. برای تضمین طول قدم راههای متفاوتی وجود دارد اما روشی که توسط Lall and Boyd (2017) ارائه شده، در اینجا بکار میرود.

این روش با قدم اولیه t^1 شروع شده و اگر تابع هدف اصلی کم شده یا تغییر نکند ($F(\omega^{k+1}) \leq F(\omega^k)$) در اینصورت، مقدار طول قدم را زیادهتر و ω جدید را میپذیریم. اگر $F(\omega^{k+1}) > F(\omega^k)$ طول قدم را کاهش داده و $\omega^{k+1} = \omega^k$ معمولاً برای افزایش طول قدم، $t^{k+1} = 1.2t^k$ و برای کاهش $t^{k+1} = \frac{1}{2}t^k$ استفاده میکنند. الگوریتم بصورت پیش فرض برای تعداد n_{iter} بار اجرا می شود.

اگر $F(\omega^{k+1}) \leq F(\omega^k)$ باشد، شرط توقف در مرحله $k+1$ بصورت

$$\left\| \frac{\omega^k - \omega^{k+1}}{t^k} + (g^{k+1} - g^k) \right\|_2 \leq \varepsilon$$

می شود. که $g^k = \nabla_{\omega^k} \psi(\theta^{ls}(\omega^k))$ و $\varepsilon > 0$ عدد کوچک دلخواهی است.

(اگر $r = 0$ و $\Omega = R^p$ روش پراکسیمال گرادیان همان روش گرادیان معمولی می شود).

الگوریتم کامل به صورت زیر است. توجه کنید که ω_1 ، t^1 ، n_{iter} مفروضات اولیه الگوریتم است و داده شده است:

برای n_{iter} ، \dots ، $k=1$ مراحل زیر را انجام بده.

$$1. \quad \theta^{ls}(\omega^k) = \left(A^T(\omega^k) A(\omega^k) \right)^{-1} A^T(\omega^k) B(\omega^k) \quad \text{را محاسبه کن.}$$

$$2. \quad g^k = \nabla_{\omega} \psi(\theta^{ls}(\omega^k)) \quad \text{را محاسبه کن.}$$

$$3. \quad \omega^{k+\frac{1}{2}} = \omega^k - t^k g^k \quad \text{قرار بده.}$$

$$4. \quad \omega^{tent} = \text{prox}_{t^k r} \left(\omega^{k+\frac{1}{2}} \right) \quad \text{محاسبه کن.}$$

$$5. \quad \text{اگر } F(\omega^{tent}) \leq F(\omega^k) \quad \text{اگر}$$

$$\left\| \frac{\omega^k - \omega^{k+1}}{t^k} + (g^{k+1} - g^k) \right\|_2 \leq \varepsilon \quad \text{قرار بده. اگر}$$

$$6. \quad \text{در غیر این صورت: } t^{k+1} = \frac{t^k}{2}, \quad \omega^{k+1} = \omega^k$$

اگر r و $\psi(\theta^{ls})$ محدب باشند روش پراکسیمال گرادیان به مینیمم سراسری همگرا خواهد شد و اگر چنین شرطهایی برقرار

نباشد، تحت فرضهای خاصی روش به یک نقطه پایا همگرا می شود.

محاسبه گرادیان

برای محاسبه گرادیان $g = \nabla_{\omega} \psi(\theta^{ls}(\omega))$ می توان از قاعده مشتق زنجیری استفاده کرد. فرض کنید که $\theta = \theta^{ls}(\omega)$ محاسبه شده است، در مرحله بعد باید $\nabla_{\omega} \psi(\theta)$ محاسبه و سپس $C = (A^T A)^{-1} \nabla_{\omega} \psi(\theta) \in R^{n \times m}$ تشکیل گردد.

اگر ماتریس A بصورت چگال ذخیره شده باشد، با استفاده از تجزیه چولسکی ماتریس گرام $G = A^T A = L^T L$ می توان ماتریس C را محاسبه نمود. اما اگر A بصورت عملگر ذخیره شده باشد برای هرستون C بصورت موازی می توان و با استفاده از روش های تکراری می توان C را محاسبه نمود.

مشتق نسبت به A به صورت زیر محاسبه می شود. قرار دهید $\theta = \phi(A, B) = (A^T A)^{-1} A^T B$.

$$D_A \phi(A, B) \Delta A = \phi(A + \Delta A, B) - \phi(A, B)$$

$$\begin{aligned} \phi(A + \Delta A, B) &= ((A + \Delta A)(A + \Delta A))^{-1} (A^T B + \Delta A^T B) \\ &\approx (A^T A)^{-1} (I - \Delta A^T A (A^T A)^{-1} - A^T \Delta A (A^T A)^{-1}) (A^T B + \Delta A^T B) \\ &\approx \phi(A, B) + (A^T A)^{-1} \Delta A^T B - (A^T A)^{-1} (\Delta A^T A + A^T \Delta A) \theta \end{aligned}$$

برای ساده سازی از رابطه تخمینی

$$(X + Y)^{-1} \approx (1 + X^{-1} Y)^{-1} X^{-1} \approx (1 - X^{-1} Y) X^{-1} \approx X^{-1} - X^{-1} Y X^{-1}$$

و رابطه $\theta = (A^T A)^{-1} A^T B$ استفاده شده است، بنابراین داریم.

$$D_A \phi(A, B) \Delta A = (A^T A)^{-1} \Delta A^T B - (A^T A)^{-1} (\Delta A^T A + A^T \Delta A) \theta$$

اکنون قرار دهید $C = (A^T A)^{-1} \nabla_{\omega} \psi(\theta)$ ، $f = \psi \circ \phi$ ، با توجه به تعریف گرادیان برای فضای ماتریس ها داریم

$$D_A f(A, B)(\Delta A) = \langle \nabla_A f, \Delta A \rangle = \text{tr}(\nabla_A f \Delta A)$$

با استفاده از خاصیت زنجیری مشتق و خاصیت چرخشی رد حاصل ضرب ماتریس ها داریم

$$\begin{aligned} D_A f(A, B)(\Delta A) &= D_{\phi^{ls}} \psi \circ D_A (\phi(A, B) \Delta A) \\ &= \text{tr} \left(\nabla_{\omega} \psi^T \left((A^T A)^{-1} \Delta A^T B - (A^T A)^{-1} (\Delta A^T A + A^T \Delta A) \theta \right) \right) \\ &= \text{tr} \left((BC^T - A\theta C^T - AC\theta^T)^T \Delta A \right) \end{aligned}$$

پس داریم

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

مشابها داریم

$$\begin{aligned}\phi(A, B + \Delta B) &= (A^T A)^{-1} A^T (B + \Delta B) \\ &= (A^T A)^{-1} A^T B + (A^T A)^{-1} A^T \Delta B\end{aligned}$$

درنتیجه

$$D_B \phi(A, B) \Delta B = (A^T A)^{-1} A^T \Delta B$$

با استفاده از خاصیت زنجیری مشتق و خاصیت چرخشی رد حاصلضرب ماتریس‌ها و تعریف گرادیان داریم:

$$\begin{aligned}D_B \phi(A, B) \Delta B &= D_{\theta^T} \psi \circ D_B (\phi(A, B) \Delta B) \\ &= \text{tr} \left(\nabla_{\theta} \psi^T (A^T A)^{-1} A^T \Delta B \right) \\ &= \text{tr} \left((AC)^T \Delta B \right)\end{aligned}$$

و نتیجه مطلوب عبارت است از:

$$\nabla_B f = AC$$

اگر ماتریس A چگال باشد $\nabla_A \psi$ و $\nabla_B \psi$ هم ابعاد ماتریس‌های A و B می‌شوند و می‌توان به صورت صریح محاسبه کرد. پیچیدگی کل محاسبات گرادیان‌های فوق در حالت چگال $kn(n+m)$ می‌شود. وقتی که A ماتریس تنک باشد محاسبه گرادیان $\nabla_B \psi$ امکان پذیر ولی محاسبه و ذخیره گرادیان $\nabla_A \psi$ به صورت صریح امکان پذیر نیست. درحالتی که متغیر ω فقط زیر مجموعه از درایه‌های A مانند Γ را تحت تاثیر قرار می‌دهد. یعنی

$$A_{ij}(\omega) = 0, \quad i, j \notin \Gamma$$

گرادیان $\nabla_A \psi$ همان الگوی تنک بودن ماتریس A را خواهد داشت و فقط نیاز است عبارت زیر حساب شود.

$$\nabla_A \psi = \begin{cases} (b_i - \theta a_i) c_j^T - a_i^T (C\theta^T)_j & i, j \in \Gamma \\ 0 & o.w \end{cases}$$

که در اینجا b_i و a_i به ترتیب i امین سطر ماتریس B و A و c_j و $(C\theta^T)_j$ به ترتیب j امی ستون ماتریس‌های C و $C\theta^T$ هستند.

با دانستن $\nabla_{\omega} A_{ij}, \nabla_{\omega} B_{ij} \in R^p$ می‌توان $g = \nabla_{\omega} \psi(\theta^L(\omega))$ را به صورت زیر با استفاده از مشتق زنجیری بدست آورد.

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

اگر که ماتریس ها تنک باشند جمع فوق روی درایه های غیر صفر بسته می شود.

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

تعمیم مسئله حالت مقید با قید تساوی

در این قسمت مساله را به حالت مقید توسعه می دهیم.

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

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

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

که $q = A(\omega)\theta - B(\omega)$ بردار مانده است. قرار دهید

$$M(\omega) = \begin{bmatrix} 0 & A(\omega)^T & C(\omega)^T \\ A(\omega) & -I & 0 \\ C(\omega) & 0 & 0 \end{bmatrix}$$

اگر ماتریس های A, C چگال باشند با استفاده از تجزیه LDL^T می توان مساله فوق را حل کرد. در صورت تنک بودن ماتریس های نامبرده با استفاده از روش های تکرار یا سالورهای مخصوص ماتریس اسپارس LDL^T مساله قابل حل است.

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

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

را محاسبه کنید. گرادیان ψ نسبت به A, B به صورت زیر قابل محاسبه است:

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

و مشابهها برای C, D داریم.

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

محاسبه گرادیان نگاشت جواب، نیازمند به حل یک معادله خطی است، بنابراین پیچیدگی آن به صورت تقریباً به اندازه پیچیدگی محاسبه خود جواب است. در صورتی که فاکتورگیری ها ذخیر شده باشد پیچیدگی خیلی کمتر خواهد شد. اگر ماتریس های A, C

تنک باشند فقط نیازمند این هستیم که روی اعضای غیرصفر گرادیان را حساب کنیم.

در زیر به محاسبه دقیق فرمول‌های ذکر شده می‌پردازیم:

$$M = \begin{bmatrix} 0 & A^T(\omega) & C^T(\omega) \\ A(\omega) & -I & 0 \\ C(\omega) & 0 & 0 \end{bmatrix}$$

$$\phi(A) = M^{-1}S = Z, \quad S = \begin{bmatrix} 0 \\ B(\omega) \\ D(\omega) \end{bmatrix}$$

$$D_A \phi(A) \Delta A = \phi(A + \Delta A) - \phi(A)$$

$$\phi(A + \Delta A) = (M + \Delta M)^{-1}S = (I + M^{-1}\Delta M)M^{-1}S$$

$$\simeq (I - M^{-1}(\Delta M)M^{-1}S) = \phi(A) - M^{-1}(\Delta M)Z$$

در نتیجه

$$D_A \phi(A) \Delta A = M^{-1}(\Delta M)Z$$

قرار دهید $C = M^{-1}\nabla_{\theta}\psi, F = \psi \circ \phi$

$$D_A f(A) \Delta A = D_{\theta^T \psi} D_A \phi(A) \Delta A = \text{tr}(\nabla_{\theta} \psi^T M^{-1}(\Delta M)Z),$$

$$\begin{cases} \Delta M^T = \Delta M \\ (M^{-1})^T = M^{-1} \end{cases}$$

داریم

$$\Delta M C = \begin{bmatrix} 0 & \Delta A^T & 0 \\ \Delta A & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} g_1 \\ g_2 \\ g_3 \end{bmatrix} = \begin{bmatrix} \Delta A^T g_2 \\ \Delta A g_1 \\ 0 \end{bmatrix}, \quad C = \begin{bmatrix} g_1 \\ g_2 \\ g_3 \end{bmatrix}$$

$$\text{tr}(AB) = \text{tr}(BA), \quad \text{tr}(X^T) = \text{tr}(X)$$

در نتیجه داریم

$$\text{tr}(Z^T \Delta M C)^T = \text{tr}((\theta^T, q^T, v^T) \cdot (\Delta A^T g_2, \Delta A g_1, 0))$$

$$= \text{tr}((\theta^T \Delta A^T g_2 + q^T \Delta A g_1)^T)$$

$$= \text{tr}((g_2 \theta^T + q^T g_1) \Delta A), \quad q^T = r_{\text{مانده}}$$

$$\nabla_A \psi = \nabla_A f = g_2 \theta^T + q^T g_1 = g_2 \theta^T + r g_1$$

دوباره قرار دهید

$$\phi(B) = M^{-1}S$$

$$\begin{aligned} D_B \phi(B) \Delta B &= \phi(B + \Delta B) - \phi(B) \\ &= M^{-1}(S + \Delta S) - M^{-1}(S) = M^{-1} \Delta S \end{aligned}$$

$$f = \psi \circ \phi \quad \text{داریم}$$

$$D_B f(B) \Delta B = D_{\theta^T S} \psi D_B \phi(B) \Delta B = \text{tr}(\nabla_{\theta} \psi^T M^{-1} \Delta S) = \text{tr}(\Delta S^T C)$$

$$\Delta S^T = (0, \Delta B^T, 0)$$

$$\Delta S^T C = (0, \Delta B^T, 0) \cdot (g_1, g_2, g_3) = \Delta B^T g_2$$

$$\text{tr}(\Delta S^T C) = \text{tr}(g_2^T \Delta B)$$

$$\nabla_B f = \nabla_B \psi = g_2$$

مشابهها برای C, D نیز می‌توان محاسبات فوق را انجام داد.

روش حداقل مربعات برای برازش داده ها

در مساله برازش داده ها، باید داده های آموزشی که شامل داده های ورودی به همراه مقدار مربوط به ورودی در دسترس باشد.

هدف برازش کردن پارامترهای، تابع پیش-بینی کننده زیر است.

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

که $\theta \in R^{n \times m}$ متغیرهای مدل و $\omega^{feat} \in \Omega^{feat} \subseteq R^{p^{feat}}$ هایپر-پارامترهای مهندسی کردن پارامترها است. توجه کنید که Ω^{feat} دامنه هایپر-پارامترهای مهندسی کردن ویژگی ها است و $\phi: U \times \Omega^{feat} \rightarrow R^n$ تابع ویژگی سازاست، که مشتق پذیر نسبت به هایپر-پارامترها در نظر گرفته می شود. تابع پیش-بینی کننده نسبت به بردار ویژگی ها خطی است.

برای انتخاب یک مدل باید مساله حداقل مربعات با داده های زیر حل شود.

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

که $\omega^{data} \in \Omega^{data} \subseteq R^N$ هایپر-پارامترهای وزنی هستند. R_1, \dots, R_d ماتریس های منظم ساز با سایز مناسب و $\omega^{reg} \in \Omega^{reg} \subseteq R^d$ هایپر-پارامترهای منظم ساز هستند. مجموعه کل هایپر-پارامترهای به صورت زیر نمایش داده می شود.

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

فرض براین است که تابع منظم ساز جداپذیر است یعنی

$$r(\omega) = r^{feat}(\omega^{feat}) + r^{data}(\omega^{data}) + r^{reg}(\omega^{reg})$$

این بدین معناست که عملگر پراکسیمال جداپذیر است. توابع فوق به ترتیب توابع منظم ساز برای هایپر-پارامترهای ویژگی ها و وزن داده ها و خود تابع منظم ساز است.

تابع هدف حقیقی

فرض کنید که داده های اعتبار سنجی از ورودی های $u_1^{val}, \dots, u_{N_{val}}^{val} \in U$ و $y_1^{val}, \dots, y_{N_{val}}^{val} \in R^m$ تشکیل شده باشند. ابتدا پیش-بینی را انجام می دهیم.

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

که در اینجا ویژگی ها یا تابع $\phi(u, \omega^{feat})$ مشخص و ثابت در نظر گرفته می شود. تابع هدف حقیقی ψ در مساله برازش داده ها با حداقل مربعات را میانگین زیان برای پیشبینی داده های اعتبار سنجی بصورت زیر در نظر میگیریم.

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

که $l: R^m \times R^m \rightarrow R$ تابع جریمه است، و فرض می شود که نسبت به متغیر اول خود مشتق پذیر است. توابع حقیقی پیچیده مانند تابع زیان اعتبار سنجی متقاطع با K لایه نیز استفاده کرد.

روش فوق برای خیلی از مسائل برازش داده از جمله رگرسیون و طبقه بندی کاربرد دارد. در رگرسیون خروجی عددی حقیقی و در رگرسیون چند-وظیفه ای خروجی یک بردار با درایه های حقیقی، و در طبقه بندی خروجی متغیر بولی و در طبقه بندی چند-وظیفه ای خروجی مشخص کننده کلاس برای داده ورودی است.

تابع جریمه رگرسیون

تابع جریمه در رگرسیون اکثر اوقات به صورت زیر است

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

که $r = \hat{y} - y$ مانده و $\pi: R^m \rightarrow R$ تابع جریمه اعمال شده بر روی مانده است. برخی از توابع رایج برای تابع جریمه بصورت زیر هستند.

- تابع مربعی: تابع جریمه مربعی، $\pi(r) = \|r\|_2^2$
- تابع هوبر: تابع جریمه هوبر، جریمه قدرتمندی به شکل تابع جریمه مربعی برای مانده های کوچک وبصورت ۲-نرمی برای مانده های بزرگ می شود.

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

که متغیر M خود یک هایپر-هایپر-پارامتر است.

- دو مربعی: جریمه دو مربعی، جریمه قدرتمند با مقدار ثابت برای مانده های بزرگ است. متغیر M خود یک هایپر-هایپر-پارامتر است

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

تابع جریمه طبقه بندی

برای طبقه بندی، پیش بینی $\hat{y} \in R^m$ متناسب با توزیع احتمال روی m برچسب به صورت زیر داده می شود.

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

پیش بینی \hat{y} بعنوان یک توزیع احتمال شرطی بر روی برچسب های y برای x داده شده تعبیر می شود. تابع زیان آنتروپی متقابل نیز یک تابع پیشنهادی بعنوان تابع جریمه در طبقه بندی محسوب می شود و بصورت زیر است:

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

وزن دهی داده ها

در مسئله برازش داده با حداقل مربعات هایپر پارامتر w_i^{data} به خطای مربعی داده (u_i, y_i) وزن $e^{2w_i^{data}}$ را نسبت می دهد. اگر w_i^{data} کوچک باشد داده متناظر نقش کمی در محاسبه پارامترهای مدل دارد و بالعکس. وزن دهی هر داده در مسئله با تابع زیان مربعی همانند این است که تابع زیان غیرمربعی در مسئله استفاده شده است به همین دلیل تنظیم خودکار هایپر-پارامترهای وزن دار کار انتخاب تابع زیان را ساده می کند. اما دامنه هایپر-پارامتر وزن دهی به صورت زیر در نظر گرفته می شود.

$$\Omega^{data} = \{x | 1^T x = 0\}$$

و این بدین معناست که میانگین هندسی $\exp(w^{data})$ مساوی یک است. هایپر-پارامتر را می توان به سمت وزن های یکسان با استفاده از منظم ساز $r^{data}(w) = \frac{\lambda}{2} \|w\|_2^2$ یا $r^{data}(w) = \lambda \|w\|_1$ سوق داد، که مجدداً یک هایپر-هایپر-پارامتر مثبت است. عملگر پراکسیمال برای $r^{data}(w) = \frac{\lambda}{2} \|w\|_2^2$ با $\Omega = \Omega^{data}$ در v با طول گام t متناظر با حل مساله بهینه سازی زیر می شود:

$$\begin{aligned} \min \quad & t \frac{\lambda}{2} \|w\|_2^2 + \frac{1}{2} \|w - v\|_2^2 \\ \text{subject to} \quad & 1^T w = 0 \end{aligned}$$

که جواب تحلیلی آن به صورت زیر خواهد شد.

$$\begin{aligned} L &= t \frac{\lambda}{2} \|w\|_2^2 + \frac{1}{2} \|w - v\|_2^2 + k 1^T w \\ \nabla_w L &= t \lambda w^* + (w^* - v) + k 1 = 0 \\ w^* &= \frac{1}{1 + t \lambda} (v - k 1) \\ 1^T w^* &= \frac{1}{1 + t \lambda} (1^T v - k 1^T 1) = 0 \\ k &= \frac{1^T v}{p} \end{aligned}$$

پس داریم

$$w^* = \frac{1}{1 + t \lambda} \left(v - \frac{1^T v}{p} 1 \right)$$

هایپر-پارامتر ω^{reg} از طریق $\sum_{i=1}^d \exp(2\omega_i^{reg}) \|R_i \theta\|_F^2$ روی محاسبه پارامترها تاثیر می گذارد. هر جمله $\|R_i \theta\|_F^2$

یک معیار برای میزان پیچیدگی θ در نظر گرفته می شود. ساده ترین شکل ماتریس $R_i = \text{diag}(e_i)$ یا همان ماتریس همانی است که منجر به رگرسیون ریج می شود. انتخاب دیگر برای R_i می تواند ماتریس وقوع گراف برای اعضای هر ستون θ باشد که هایپر-پارامتر منظم ساز در اینجا میزان اهمیت منظم سازی گراف را نشان می دهد.

مهندسی کردن ویژگی ها

هایپر-پارامتر ω^{feat} ، ویژگی ساز ها را پارامتریزه می کند. هدف از انتخاب ω^{feat} اینست که خروجی y را بصورت تقریباً خطی از $\phi(u, \omega^{feat})$ می کند. که دامنه ورودی اول فرض میشود که فضای برداری U است.

معمولاً تابع ϕ بصورت زنجیره ای از تولید ویژگی ها ساخته می شود. یعنی به صورت ترکیبی از توابع مهندسی سازی ویژگی ها به صورت ϕ_1, \dots, ϕ_l است پس داریم.

$$\phi = \phi_l \circ \dots \circ \phi_1$$

توجه کنید معمولاً آخرین تابع مهندسی سازی یک ثابت به ترکیب توابع قبلی اضافه میکند که نتیجه ترکیب به تابع آفین می شود.

توابع مهندسی سازی ویژگی‌های اسکالر

در این بخش به توابع مهندسی سازی ویژگی‌ها $R \rightarrow R: \phi$ با این فرض که می‌توان آنها را به مولفه‌های یک بردار با هایپر-پارامترهای متفاوت اعمال کرد، خواهیم پرداخت.

- مقیاس بندی: یکی از ساده‌ترین توابع مهندسی سازی ویژگی‌ها مقیاس بندی خطی است که به صورت زیر داده می‌شود.

$$\phi(x, (a, b)) = ax + b$$

یک رویه رایج در برازش داده این است که داده‌ها را با استفاده از مقیاس بندی به پارامترهای $a = \frac{1}{\sigma}, b = -\frac{\mu}{\sigma}$ استاندارد می‌کنند. در اینجا μ میانگین و σ انحراف معیار استاندارد متغیر تصادفی x هستند. در مساله حداقل مربع با تنظیم خودکار، b و a بر اساس داده‌ها تعیین می‌شوند.

- تبدیل توانی: تبدیل توانی به صورت $\phi(x, (c, \gamma)) = \text{sgn}(x - c) |x - c|^\gamma$ تعریف می‌شود. تابع sgn همان تابع علامت و مقیاس‌بند $\gamma \in R$ و پارامتر $c \in R$ هایپر-پارامترهای مساله هستند. برای مقادیر متفاوت γ و c تابع تبدیل متفاوتی بدست می‌آید. مثلاً اگر $\gamma = 1, c = 0$ تابع تبدیل همانی و اگر $\gamma = 0$ تابع تبدیل تعیین‌کننده این است که x سمت راست یا چپ مقدار c قرار دارد. و اگر $\gamma = \frac{1}{2}, c = 0$ این تبدیل همان ریشه قدر مطلق می‌شود. توجه کنید که تبدیل توانی همه جا به غیر از $\gamma = 0$ مشتق پذیر است.

- چند جمله‌ای‌های اسپلاین: چندجمله‌ای اسپلاین یک تابع قطعه چند جمله‌ای است. چند جمله‌ای اسپلاین متشکل از یک آرایه صعودی از اعداد $z \in R^{k+1}$ که گره‌های اسپلاین را مشخص می‌کند به همراه درجه d و ضرایب چندجمله‌ای $f_0, \dots, f_{k+1} \in R^{d+1}$ است.

$$\phi(x, (z, f_1, \dots, f_{k+1})) = \begin{cases} \sum_{i=0}^d (f_0)_i x^i & x \in (-\infty, z_1) \\ \sum_{i=0}^d (f_j)_i x^i & x \in (z_j, z_{j+1}) \quad j=1, \dots, k \\ \sum_{i=0}^d (f_{k+1})_i x^i & x \in (z_{k+1}, \infty) \end{cases}$$

توجه داشته باشید که r^{feat}, Ω^{feat} برای اعمال شرط پیوستگی و پیوستگی مشتق‌ها در نقاط گره‌ای به بکار برده می‌شود.

توابع چند بعدی برای مهندسی سازی ویژگی‌ها

- رتبه پایین: ویژگی ساز ϕ می تواند یک تبدیل رتبه پایین به صورت زیر باشد.

$$\phi(x, T) = Tx$$

که $T \in R^{r \times n}$, $r < n$. در عمل یک انتخاب رایج برای عملگر T انتخاب چند بردار اول در تجزیه SVD ماتریس داده است که اب استفاده از تنظیم خود کار انتخاب T به صورت خود کار صورت می گیرد.

- شبکه های عصبی: تابع ϕ می تواند یک شبکه عصبی باشد، در این صورت ω^{feat} همان پارامتر های شبکه عصبی می شوند.

- انتخاب ویژگی: کسر f از ویژگی ها را می توان با $\phi(x) = diag(a)x$ انتخاب کرد که

$$\Omega^{feat} = \left\{ a \mid a \in \{0,1\}^{n_1}, 1^T a = \lfloor fn_1 \rfloor \right\}$$

که $\lfloor \rfloor$ همان تابع جز صحیح کف و f هایپر-هایپر-پارامتر مساله است.

مجموعه تست و توقف زود هنگام

وقتی که تعداد هایپر-پارامترها زیاد می‌شود ریسک برازش بیش از حد وجود دارد زیرا تنظیم خودکار حداقل مربعات بصورت مستقیم تابع زیان روی مجموعه اعتبارسنجی را کمینه می‌کند. برای تشخیص این برازش بیش از حد، یک مجموعه سوم از داده‌ها به اسم مجموعه داده تست در نظر گرفته می‌شود که در آخر الگوریتم روی این مجموعه محاسبه می‌شود. مقدار زیان مجموعه اعتبارسنجی لزومی ندارد که تخمین خوبی از اندازه دقیق عملکرد مدل بر روی داده‌های دیده نشده باشد.

روشی کمی متفاوت برای رویارویی با مساله برازش بیش از حد این است که مقدار زیان روی مجموعه داده تست در هر تکرار محاسبه کرده و وقتی که مقدار زیان روی مجموعه داده تست شروع به افزایش کرد، الگوریتم را متوقف می‌کنند. به این روش توقف زود هنگام می‌گویند. در این روش نیاز به مجموعه داده چهارمی است که به عنوان مجموعه داده تست نهایی شناخته شده است، وقتی که الگوریتم متوقف شد، روی مجموعه داده نهایی الگوریتم حساب شده و بعنوان عملکرد الگوریتم در نظر گرفته می‌شود.

یک روش مناسب و قدرتمند استفاده از تابع زیان روش اعتبارسنجی متقابل است. به این ترتیب که گرادیان روش اعتبارسنجی متقابل همان میانگین گرادیان‌های حساب شده در هر مرحله اعتبارسنجی متقابل است.

داده های MNIST

مجموعه داده MNIST شامل ۶۰۰۰۰ نقطه آموزشی به همراه برچسب که هر کدام آن یک بردار از سایز ۷۸۴ (یک بعدی شده ماتریس های تصویری خاکستری ۲۸ در ۲۸) هستند. برچسبها شامل ده کلاس متناسب با اعداد ۰ تا ۹ می شوند. داده های تست شامل ۱۰۰۰۰ نقطه آموزشی به همراه برچسب هستند.

مساله ای که با استفاده از این روش حل می شود رگرسیون ریج با یک هایپرپارامتر است.

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

برای حل این مساله می توان از تابع آنتروپی متقابل و یا خود تابع هدف مساله حداقل مربعات بعنوان تابع حقیقی هدف استفاده کرد.

همانطور که قبلا دیدیم جواب مساله به شکل زیر است:

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

توجه کنید که ماتریس C همان مقادیر برچسب به ازای هر نمونه آموزشی است و هر سطر آن یک بردار ۱۰ تایی که دارای یک مولفه غیر صفر و با مقدار یک است و جایگاه آن مولفه نشان دهنده برچسب مورد نظر است. مثلا اگر مولفه دهم یک باشد آن برچسب مساوی عدد ۹ و اگر جایگاه اول مساوی یک باشد برچسب نشان دهنده برچسب صفر است.

اکنون می خواهیم مشتق تابع θ^{ls} بر حسب λ را بدست بیاوریم. برای سادگی $D = X^T X$ نشان می دهیم پس رابطه جواب مساله به صورت زیر است:

$$D\theta^{ls} = X^T C$$

با استفاده از خاصیت ضرب داریم:

$$\frac{dD}{d\lambda} \theta^{ls} + D \frac{d\theta^{ls}}{d\lambda} = \frac{dX^T}{d\lambda} C$$

برای محاسبه $\frac{d\theta^{ls}}{d\lambda}$ احتیاج به محاسبه مشتق $\frac{dX^T}{d\lambda}$ و $\frac{dD}{d\lambda}$ داریم زیرا

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

حال ماتریس $X^T X$ را در نظر بگیرید.

$$D = X^T X = \begin{bmatrix} A^T & \exp(\lambda)I \end{bmatrix} \begin{bmatrix} A \\ \exp(\lambda)I \end{bmatrix} =$$

$$\begin{bmatrix} A^T A + \exp(2\lambda)I \end{bmatrix} \rightarrow \frac{dD}{d\lambda} = 2\exp(2\lambda)I$$

و مشابهها برای X^T داریم

$$X^T = \begin{bmatrix} A^T & \exp(\lambda)I \end{bmatrix} \rightarrow \frac{dX^T}{d\lambda} = \begin{bmatrix} 0 & \exp(\lambda)I \end{bmatrix}$$

پس با توجه به محاسبات بالا $\frac{d\theta^{ls}}{d\lambda}$ را می توان حساب کرد. حال تابع هدف حقیقی را تابع آنتروپی متقابل در نظر بگیرید:

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

با اعمال تابع آنتروپی متقابل روی داده های اعتبارسنجی و میانگین گرفتن داریم:

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

در الگوریتم نیاز به محاسبه $g = \nabla_{\lambda} \psi$ داریم. و چون $\psi(\theta)$ جمع بر روی تمام داده ها اعتبارسنجی است فقط کافی است که $\frac{dl}{d\lambda}$ محاسبه شود.

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

که در اینجا اندیس اول اندیس مربوط به نمونه و اندیس دوم مربوط به مولفه های نمونه مربوطه است. که با جمع بستن روی تمام نمونه های اعتبارسنجی و میانگین گرفتن تابع g بدست می آید.

اگر تابع هدف حقیقی همان تابع هدف مساله حداقل مربعات باشد داریم:

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

و با مشتق گیری داریم:

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

که در اینجا X ماتریس ساخته شده از نمونه ها و b بردار برچسب هاست.

نتایج پیاده سازی

برای پیاده سازی ۲ تابع هدف حقیقی آنتروپی متقابل و تابع هدف حداقل مربعات را در نظر گرفتیم. رویکرد اول، استفاده از روش‌های ساخت classifier ها و ترکیب آن با الگوریتم تنظیم خودکار است؛ در اینجا ابتدا تابع هدف حقیقی را همان تابع هدف مسئله حداقل مربعات قرار داده و در پیاده سازی دیگری، تابع هدف را تابع زیان بر روی داده های اعتبارسنجی قرارداده و با استفاده از روش proximal gradient، تابع هدف را مینیمم سازی می‌کنیم. در کنار آن به مهندسی ویژگی نیز پرداخته‌ایم و تعداد مختلفی ویژگی تصادفی اضافه کرده‌ایم.

اما رویکرد دیگر، ابتدا از هر برچسب، یک بردار سطری به طول ۱۰ ایجاد می‌کنیم و جایگاه متناظر با مقدار برچسب را مقدار ۱ می‌دهیم. برای این بخش از میانگین تابع آنتروپی متقابل بر روی داده‌های اعتبارسنجی، به عنوان تابع هدف اصلی استفاده کرده‌ایم که محاسبات مربوط به آن قبلاً در توضیحات آورده شده است.

قابل توجه است که در ادبیات، پیاده سازی های مرتبط با روش تنظیم خودکار، با استفاده از شبکه های عصبی صورت پذیرفته اما در این پروژه، تلاش بر این بوده است که با استفاده از روش‌های ساده بحث شده در کلاس به حل مسئله رگرسیون رنج پردازیم. در ادامه نتایج برای پیاده سازی مختلف آورده شده است.

در ابتدا به مقایسه محاسبه θ با استفاده از دو روش convex optimization و least squares پرداختیم :

Convex Optimization:

15.594337 seconds

Least Squares:

7.582242 seconds

در نتیجه از روش حداقل مربعات که زمان کمتری می‌گیرد، در الگوریتم خود استفاده می‌کنیم.

• نتایج پیاده‌سازی با رویکرد اول:

Code: first_app_1.ipynb/.jl

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

Code: first_app_2.ipynb/.jl

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 پیاده‌سازی با رویکرد دوم:

Code: second_app_1.ipynb/.jl

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 شبکه‌های عصبی ساده در ۱۰ تکرار (code: simple NN MNIST.ipynb/.jl):

test_loss = 0.22957107, test_accuracy = 0.9352

- نتایج پیاده‌سازی با استفاده از شبکه‌های عصبی پیچشی در ۱۰ تکرار (code: CNN on MNIST.ipynb):

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

همانگونه که از نتایج پیاده‌سازی مشخص است، تاثیر افزودن ویژگی‌های تصادفی بیشتر از منظم ساز در رگرسیون ریدج است. با پیچیده تر کردن تابع هدف مسئله حداقل مربعات، بطور مثال با افزودن منظم ساز متناسب با نرم فروبنیوسی ماتریس وقوع گراف پیکسل‌ها، می‌توان اطلاعات مربوط به همواری عکس‌ها را نیز به تابع هدف اضافه نمود. همچنین می‌توان تابع هدف حقیقی را اینگونه تغییر داد که به جای استفاده از تابع زیان بر روی داده‌های اعتبارسنجی، با روش cross validation، ابتدا داده‌های آموزشی را به k قسمت تقسیم کرده و $k - 1$ قسمت را بعنوان آموزش و باقی مانده را بعنوان داده‌ی اعتبارسنج استفاده کرد. تابع هدف حقیقی را می‌توان میانگین توابع زیان در هر مرحله از cross validation در نظر گرفت. مجدداً می‌توان تابع زیان برای یک داده از داده‌های اعتبارسنجی، را همان تابع آنتروپی متقابل قرار داد.

- [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.