

مسئله رگرسیون خطی منظم سازی شد با نرم L_2 :

در این مسئله در واقع برای حل بیشترین میزان معیار شده است به این شکل که یک نرم جبریه به تابع هزینه مدل رگرسیون استاندارد اضافه می کنیم. نرم جبریه وزن های بزرگ w در مدل را جبریه می کند. هدف این است که مدلی پیدا کنیم که نه تنها به خوبی داده های آموزش را برازش کند بلکه وزن های آن نیز تا حد امکان کوچک باشند.

$$E(w) = J(w) = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$$

تابع هزینه در این مدل به صورت \leftarrow

$$E(w)_{Regu} = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2 + \alpha \cdot \text{Regularization Term}(w)$$

اعاد حالت منظم سازی شده داریم \leftarrow نوع نرم L_2 \leftarrow L_1

پارامتر تنظیم برای نرم \leftarrow مقدار هدف \leftarrow مقدار وزن های مدل \leftarrow مقدار وزن های نمونه

منظم سازی $L_2 \leftarrow$ مربع نرم L_2 وزن ها، نرم منظم سازی است:

$$RT(w) = \|w\| = \left(\sum_{j=1}^P w_j^2 \right)^{1/2} \quad L_2 = \sqrt{\sum_i x_i^2}$$

تعداد ویژگی ها است P

تابع هزینه به صورت زیر باز نویسی می شود:

$$E(w)_{Regu} = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2 + \frac{\alpha}{2} \sum_{j=1}^P w_j^2$$

الف) با کمک تکنیک کاهشی گرادینت تعادلی یک رابطه برای یادگیری یک رابطه برای یادگیری به خطا این مدل ارائه دهید. (sgd)

در این روش در هر گام وزن ها را در جهت منفی گرادینت تابع هزینه نسبت به وزن ها به روز رسانی می شوند. گرادینت جهت بیشترین افزایش تابع هزینه را نشان می دهد. بنابراین حرکت در جهت منفی آن باعث کاهش تابع هزینه می شود.

تکنیک کاهش گرادینت تعادلی یک نسخه از GD است که به جای استفاده از تمام مجموعه داده برای محاسبه گرادینت در هر تاپل تنها از یک نمونه داده تعادلی (یا یک دسته کوچکی از داده‌ها) استفاده می‌کند.

تا به خطی که در مرحله قبل موفق نگردیم را در نظر می‌گیریم، سپس مشتق آن را نسبت به بردار وزن w محاسبه می‌کنیم:

$$E(w) = \frac{1}{P} \sum_{i=1}^N (w^T x_i - y_i)^2 + \|w\|^2 \frac{\lambda}{2}$$

حالا برای λ کم + λ بزرگ

$$E_t(w) = \frac{1}{P} (w^T x_t - y_t)^2 + \|w\|^2 \frac{\lambda}{2}$$

$$\left. \begin{aligned} \frac{\partial}{\partial w} \left(\frac{1}{P} (w^T x_t - y_t)^2 \right) &= -(y_t - x_t (w^T x_t)) \\ \frac{\partial}{\partial w} \|w\|^2 \frac{\lambda}{2} &= \lambda w \end{aligned} \right\} \Rightarrow \nabla_w E_t(w) = x_t (w^T x_t - y_t) + \lambda w$$

حالا راجع به وزنهای وزن ها در SGD: در این الگوریتم وزن ها در جهت منفی گرادینت به وزنهای می‌روند. در SGD این به وزنهای به اساس گرادینت محاسبه شده برای نمونه تعادلی + انجام می‌شود:

$$w^{(t+1)} = w^{(t)} - \eta \nabla_w E_t(w^{(t)})$$

η (تا) نرخ یادگیری است یا طول گام است که اندازه

گام در جهت منفی گرادینت را کنترل می‌کند. انتخاب مناسب نرخ یادگیری برای همگرايي موفق الگوریتم بسیار مهم است.

حالا با جایگزینی عبارت گرادینت که بدست آوردیم، راجع به وزنهای برای رگرسیون خطی متغیر شده با λ با کمک SGD به صورت زیر است:

$$w^{(t+1)} = w^{(t)} + \eta (y_t - w^{(t)T} x_t) x_t - \eta \lambda w^{(t)}$$

با در تکلیف کاهش گرادیان، η ضریب یادگیری نامیده می‌شود و در محاسبه بهینه سازی حواله گام است. در رابطه ای که برای قسمت الف ارائه کرده‌اید، حواله گام بهینه را با حل معادلات ریاضی بدست آورید. هدف از حواله گام بهینه، مقداری برای ضریب یادگیری η است که به ازای آن، مقدار تابع هزینه در راستای بردار گرادیان به کمترین مقدار ممکن برسد. برای این کار ابتدا پارامتر وزن جدید را که به حساب η تعریف شده است محاسبه و وارد معادلات مدل، و تابع زیان کنید تا تابع زیان، تابعی به حساب η دلیله شود. حال نسبت به η مشتق بگیرید و برابر صفر قرار دهید و η را حل کنید.

η : ابتدا پارامتر وزن جدید را $(w^{(1)})$ به حساب η در تابع هزینه جایگزینی می‌کنیم و حذف می‌کنیم

$$w \leftarrow \text{وزن فعلی و } w' \text{ وزن جدید باشد: } w' = w + \eta (y - x_t^T w) - \eta \lambda w$$

$$E(w') = \frac{1}{P} (y - (w')^T x)^2 + \frac{\lambda}{P} \|w'\|^2 \quad (*) \quad \leftarrow \text{تابع هزینه}$$

$$E(\eta) = \frac{1}{P} (y - (w + \eta (y - x_t^T w) - \eta \lambda w)^T x)^2 + \frac{\lambda}{P} \|(w + \eta (y - x_t^T w) - \eta \lambda w)\|^2 \quad (*) \leftarrow \text{جایگزینی } w' \text{ در } (*)$$

حال برای η بهینه از رابطه بالا مشتق گرفته و برابر صفر قرار می‌دهیم:

$$\frac{\partial E(\eta)}{\partial \eta} = 0$$

$$\frac{\partial E(w')}{\partial w} = \nabla_w E_t(w') = - (y_t - w^T x_t) x_t + \lambda w$$

$$w' = w + - \eta \nabla_w E_t(w')$$

$$E(\eta) = \frac{1}{P} (y - (w - \eta \nabla_w E_t(w'))^T x)^2 + \frac{\lambda}{P} \|w - \eta \nabla_w E_t(w')\|^2$$

$$\frac{\partial E(\eta)}{\partial \eta} = (y - (w - \eta \nabla_w E_t(w'))^T x) (\nabla_w E_t(w'))^T x + \lambda (w - \eta \nabla_w E_t(w'))^T (-\nabla_w E_t(w')) = 0$$

$$y - \omega^T x + \eta (\nabla_{\omega} E_t(\omega'))^T x - \lambda \omega^T \nabla_{\omega} E_t(\omega') + \lambda \eta (\nabla_{\omega} E_t(\omega'))^T (\nabla_{\omega} E_t(\omega')) = 0$$

$$\Rightarrow (y - \omega^T x) ((\nabla_{\omega} E_t(\omega'))^T x) + \eta ((\nabla_{\omega} E_t(\omega'))^T x)^2 - \lambda (\omega^T (\nabla_{\omega} E_t(\omega'))^T) + \lambda \eta \|\nabla_{\omega} E_t(\omega')\|^2 = 0$$

$$\Rightarrow \eta ((\nabla_{\omega} E_t(\omega'))^T x)^2 + \lambda \|\nabla_{\omega} E_t(\omega')\|^2$$

$$= - (y - \omega^T x) ((\nabla_{\omega} E_t(\omega'))^T x) + \lambda \omega^T \nabla_{\omega} E_t(\omega')$$

$$\Rightarrow \eta^x = \frac{-(y - \omega^T x) ((\nabla_{\omega} E_t(\omega'))^T x) + \lambda \omega^T \nabla_{\omega} E_t(\omega')}{((-\lambda (y - \omega^T x) x + \lambda \omega)^T x)^2 + \lambda \|\lambda (y - \omega^T x) x + \lambda \omega\|^2}$$

ج) به طور مشابه به سوالات الف و ب، در صورتی که مسئله را می‌توان به صورت خطی بازنویس کرد، متغیرهای ساده‌شده، پاسخ‌دهی.

نرم ۱: $\omega = [x_1, x_2, \dots, x_n]$

$$\|\omega\| = L_1 = \sum_i |x_i|$$

$$E(\omega) = \frac{1}{2} \sum_{i=1}^N (y_i - \omega^T x_i)^2 + \lambda \|\omega\|_1$$

بردار علامت ω

$$\text{sgd} \Rightarrow \nabla_{\omega} E_t(\omega) = -(y_t - \omega^T x_t) x_t + \lambda \cdot \text{sgn}(\omega)$$

$$\text{sgn}(\omega_j) = \begin{cases} 1 & \omega_j > 0 \\ -1 & \omega_j < 0 \\ 0 & \text{other} \end{cases}$$

$$\omega^{(t+1)} = \omega^{(t)} - \eta \nabla_{\omega} E_t(\omega^{(t)})$$

$$\omega^{(t+1)} = \omega^{(t)} + \eta (y_t - (\omega^{(t)})^T x_t) x_t - \eta \lambda \cdot \text{sgn}(\omega^{(t)})$$

حالا برای پیدا کردن طول گام بهینه برای متغیرهای λ باید مشتق بگیریم. تابع هزینه ما به دلیل وجود قدر مطلق بهینه‌پذیر است و مشتق آن ممکن است ~~یافت~~ پیوسته نباشد. یک روش معمول در این حالت استفاده از روش `line search` است. اگر بخواهیم یک تقریب خوب برای طول گام بهینه در نظر بگیریم، می‌توانیم از تقریب مرتبه دوم تابع هزینه (روش نیوتن) کمک بگیریم.