



یادگیری ماشین

نیم سال دوم ۱۴۰۱-۱۴۰۲

مدرس: دکتر سید ابوالفضل مطهری

تمرین سوم

۱. در عمل، $P_{\mathcal{X} \times \mathcal{Y}}$ عموماً ناشناخته است و برای مدل کردن آن از کمینه کننده خطای تجربی ^۱ (ERM) استفاده می کنیم. ما مسئله را به صورت یک مسئله رگرسیون خطی d بعدی بازنویسی می کنیم. در ابتدا در نظر داشته باشید که توابع درون \mathcal{H}_d ^۲ توسط بردار $\mathbf{b} = [b_0, b_1, \dots, b_d]^T$ نمایش داده می شوند، به همین دلیل برای نمایش آن ها از نماد $f_{\mathbf{b}}$ استفاده می کنیم. به طرز مشابهی بردار $a \in \mathbb{R}^3$ برداری است که تابع $g(x) = f_a(x)$ را نمایش می دهد. در نهایت داده های آموزشی را فرم ماتریسی به صورت مقابل نمایش می دهیم:

$$X = \begin{bmatrix} 1 & x_1 & \dots & x_1^d \\ 1 & x_2 & \dots & x_2^d \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_N & \dots & x_N^d \end{bmatrix}, \quad \mathbf{y} = [y_0, y_1, \dots, y_N]^T$$

استفاده از این نمادها اجازه استفاده از قضایا و تکنیک های جبرخطی را به ما می دهند. ماتریس X را ماتریس طراحی ^۳ می نامیم.

- نشان دهید که کمینه کننده خطای تجربی (ERM) $\hat{\mathbf{b}}$ توسط کمینه سازی رابطه مقابل بدست می آید.

$$\hat{\mathbf{b}} = \arg \max_b \|X\mathbf{b} - \mathbf{y}\|_2^2$$

- اگر $N > d$ باشد و ماتریس X دارای رتبه ^۴ کامل باشد، نشان دهید $\hat{\mathbf{b}} = (X^T X)^{-1} X^T \mathbf{y}$ (راهنمایی: باید گرادیان تابع loss بالا را نسبت به \mathbf{b} بگیرید). چرا شروط $N > d$ و کامل بودن رتبه ماتریس X لازم هستند؟

۲. n داده آموزش با m ویژگی ^۵ فرض کنید. فرض کنید برچسب داده ها بردار $[y^{(1)}, \dots, y^{(n)}]$ و $X = [x^{(1)}, \dots, x^{(n)}]$ داده های آموزش باشند. ($X \in \mathbb{R}^{n \times m}$).

- نشان دهید اگر یک مدل خطی تنها بر روی یکی از ویژگی ها آموزش دهیم، آنگاه $w_j = \frac{x_j^T y}{x_j^T x_j}$
- فرض کنید ستون های ماتریس X بر هم عمود ^۶ باشند. نشان دهید پیدا کردن وزن های بهینه در این حالت با پیدا کردن هر وزن به صورت مستقل تفاوتی ندارد.

¹Empirical Risk Minimizer

²Hypothesis Space (d dimension linear regression functions)

³Design Matrix

⁴Rank

⁵Feature

⁶Orthogonal

۳. فرض کنید \hat{g}_1 و \hat{g}_2 به صورت زیر تعریف شده اند.

$$\hat{g}_1 = \underset{g}{\operatorname{argmin}} = \left(\sum (y_i - g(x_i))^2 + \lambda \int [g^{(3)}(x)]^2 dx \right)$$

$$\hat{g}_2 = \underset{g}{\operatorname{argmin}} = \left(\sum (y_i - g(x_i))^2 + \lambda \int [g^{(4)}(x)]^2 dx \right)$$

توجه کنید که $g^{(m)}$ مشتق m ام تابع g است.

- اگر $\lambda \rightarrow \infty$ کدام یک از توابع RSS کمتری بر روی داده های آموزش دارند؟
- اگر $\lambda \rightarrow \infty$ کدام یک از توابع RSS کمتری بر روی داده های تست دارند؟

Spline

به موارد زیر پاسخ دهید:

۱. مزیت Spline نسبت به رگرسیون چندجمله ای با درجه آزادی برابر چیست؟
۲. تفاوت دو مورد Natural Spline و Restricted Cubic Spline را توضیح دهید.
۳. مفهوم Curse of Dimensionality یا نحسی ابعاد بالا را توضیح دهید و توضیح دهید این چه اثری بر روی استفاده از Spline و GAM ها در ابعاد بالا خواهد گذاشت؟

GAM

به موارد زیر پاسخ دهید:

۱. در استفاده از Spline ها تعداد knot ها چگونه انتخاب می شوند؟ یا به طور مشابه در GAM ها تعداد توابع پایه را چگونه انتخاب می کنید؟ مزایا و معایب تعداد بیشتر یا کمتر در هریک چیست؟ توضیح دهید.
۲. مزایا و معایب استفاده از GAM ها در مقایسه با بقیه تکنیک های رگرسیون غیر خطی چیست؟ توضیح دهید.