



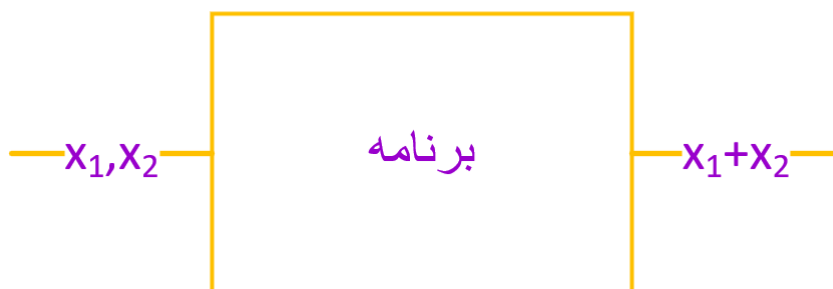
یادگیری ماشین

منابع درس

- James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). An introduction to statistical learning (Vol. 112, p. 18). New York: springer.
- Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). The elements of statistical learning: data mining, inference, and prediction (Vol. 2, pp. 1-758). New York: springer.

تفاوت یادگیری ماشین با منطق

تمایز یادگیری ماشین با منطق در این است که در یادگیری ماشین داده‌ها می‌توانند برنامه را عوض نمایند. مثلاً فرض نمایید که یک برنامه کامپیوتری نوشته شده است تا دو عدد را با یکدیگر جمع نماید



برنامه فوق یک برنامه مشخصی است که ورودیها را به خروجی انتقال می‌دهد. حال یک ماشین را در نظر بگیرید که یک برنامه دارد و در ضمن می‌تواند برنامه خود را نیز عوض کند.



برنامه متغیر معمولاً یک سری عدد می‌باشند که برنامه ثابت از آنها استفاده می‌نماید. این اعداد را معمولاً با β نمایش می‌دهیم.

مثال ۱. فرض کنید یک برنامه ثابت داریم که رابطه بین ورودی و خروجی آن به شکل زیر است:

$$y = \beta x$$

که در آن y ، x ، β همگی متعلق به \mathbb{R} می‌باشند. در یک مدل یادگیری مقدار β ثابت نمی‌باشد و با توجه به ورودیها و خروجیها که به آن داده گفته می‌شود می‌توانند تغییر نمایند. ممکن است پس از مدتی β ثابت گردد و دیگر آنرا تغییر ندهیم. در این صورت دیگر برنامه ثابت شده و یادگیری صورت نمی‌پذیرد. بنابراین هدف یادگیری استفاده از داده‌ها و رسیدن به روشی است که بتوان β را تنظیم نمود.

حفظ کردن یا تعمیم دادن

فرض نمایید داده‌ها به صورت زیر داده شده است:

$$X = \begin{matrix} & \text{سن} \\ \begin{matrix} \text{نمونه ۱} \\ \text{نمونه ۲} \\ \vdots \\ \text{نمونه } n \end{matrix} & \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \end{matrix} \quad Y = \begin{matrix} & \text{درآمد} \\ \begin{matrix} \text{نمونه ۱} \\ \text{نمونه ۲} \\ \vdots \\ \text{نمونه } n \end{matrix} & \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \end{matrix}$$

می‌خواهیم این داده‌ها را ذخیره نماییم و سپس هرگاه ورودی x_i را مشاهده کردیم y_i را در خروجی قرار دهیم. این عمل به سادگی قابل انجام است تنها یک مشکلی دارد و آن اینکه ممکن است برای یک سن خاص چندین درآمد وجود داشته باشد که می‌توان همه آنها را در خروجی ایجاد نمود.

اگر ورودی در داده‌ها وجود نداشت چه می‌توان کرد؟

به این مساله تعمیم می‌گوییم. چگونه می‌توانیم از مابقی داده‌ها به خروجی معقول برسیم و چگونه می‌توانیم نشان دهیم آنچه در خروجی ایجاد می‌شود مناسب است.

راه حل اول: درون یابی

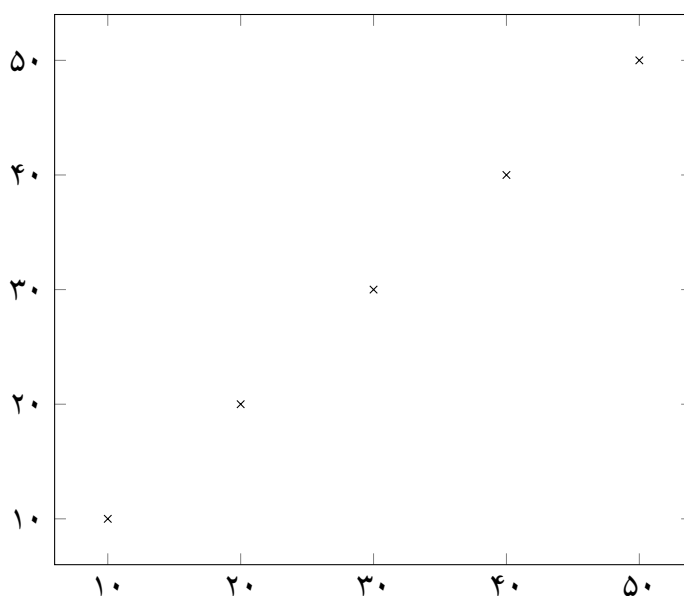
رابطه بین x و y را با یک تابع با پارامتر β نمایش می دهیم :

$$y = f_{\beta}(x)$$

حال می گوییم β را طوری بدست می آوریم که

$$y_i = f_{\beta}(x_i) \quad \forall i \in \{1, \dots, n\}$$

در این حالت از روی n معادله بایستی β را به دست آوریم . حال ممکن است β وجود نداشته باشد ، برای آنکه جواب β را امکان پذیر نماییم بایستی همزمان به داده ها و تابع توجه کنیم .



در شرایط بالا هرچند به یک دستگاه با چهار معادله روبه رو می باشیم اما یک تابع خطی با یک پارامتر می تواند از تمامی نقاط عبور کند . بنابراین کافیت در نظر بگیریم :

$$y = \beta x$$

و β را از روی یک نقطه بدست آوریم مثلاً

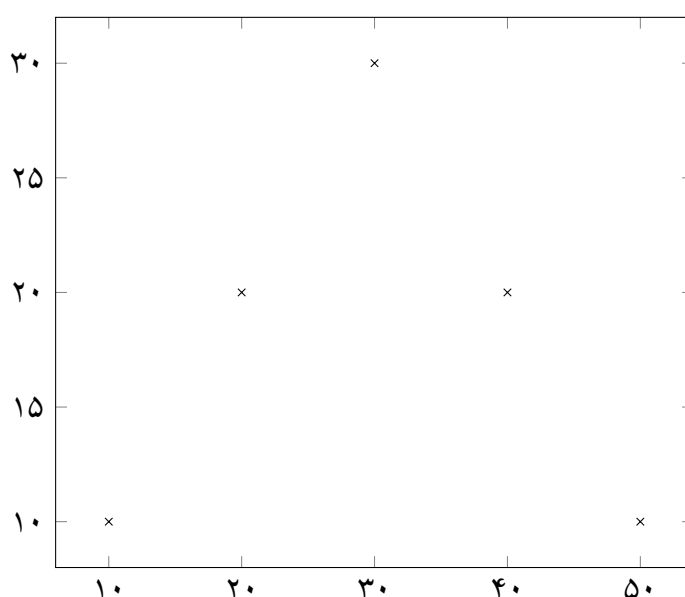
$$\beta = \frac{y_1}{x_1}$$

حال می توانیم برای یک x_{in} خروجی را به صورت

$$y_{out} = \beta x_{in}$$

در نظر بگیریم.

رابطه فوق را ممکن است به شکل زیر باشد :



ولیکن همچنان بتوانیم با تعداد پارامترهای کم تابعی بدست آوریم که از تمام نقاط بگذرد. مثلاً در مثال فوق می‌توانیم یک درجه دوم در نظر بگیریم که از تمام داده‌ها عبور می‌کند.

$$y = \beta_1 x + \beta_2 x^2$$

اگر دو نقطه از داده را انتخاب نماییم، می‌توانیم β_1 و β_2 را از یک دستگاه معادله بدست آوریم :

$$y_1 = \beta_1 x_1 + \beta_2 x_1^2$$

$$y_2 = \beta_1 x_2 + \beta_2 x_2^2$$

$$\Rightarrow \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} x_1 & x_1^2 \\ x_2 & x_2^2 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

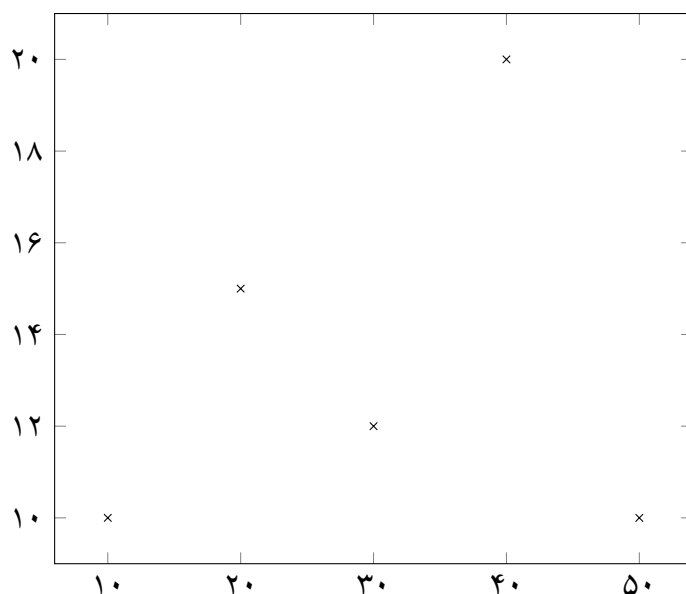
بنابراین خواهیم داشت:

$$\Rightarrow \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} x_1 & x_1^2 \\ x_2 & x_2^2 \end{pmatrix}^{-1} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

مثال‌های فوق نشان می‌دهند که در حالتیکه شرایط مناسب باشد با تعداد کمی نمونه می‌توانیم رابطه بین ورودی و خروجی را بدست آوریم.

در بسیاری از مواقع به این میزان خوش شانس نمی‌باشیم و بایستی از تمام نمونه‌ها برای بدست آوردن پارامترها

استفاده نمود. یعنی تعداد پارامترها با تعداد نمونه ها بایستی برابر باشد.



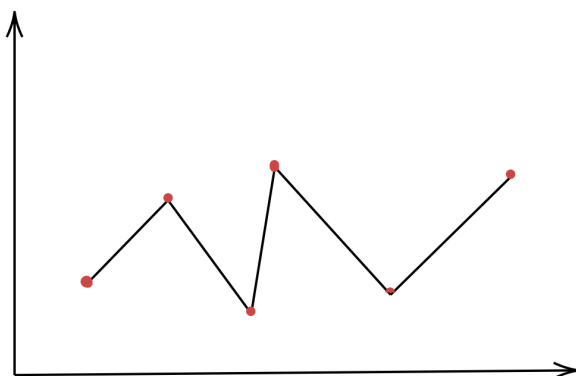
برای حل مسئله فوق دو راه پیش رو داریم:

۱. تابعی را برای کل R در نظر بگیریم و سعی می‌کنیم نقاط را fit نماییم، مثلاً اگر یک تابع درجه $n - 1$ در نظر بگیریم که چند جمله ای به صورت زیر باشد:

$$f(x) = \sum_{i=0}^{n-1} \beta_i x^i$$

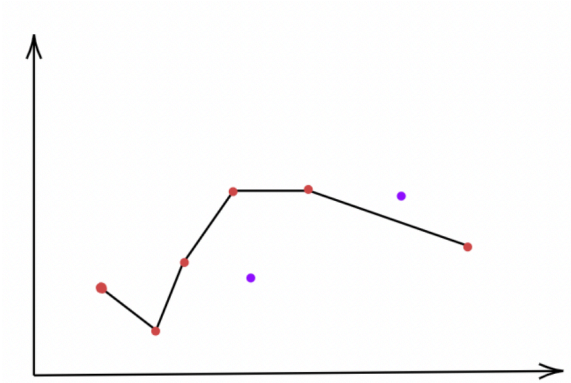
آنگاه می‌توانیم با داشتن n معادله و n مجهول سعی نماییم β_i ها را بدست آوریم. (دیدگاه جهانی)

۲. تعداد نقاط کمی در کنارهم را در نظر گرفته و یک تابع ساده را روی آنها برازش نماییم و سپس آنها را به یکدیگر وصل کنیم. (دیدگاه محلی)



اگر دقت نماییم درونیابی خاصیت ذخیره سازی را نیز دارد، یعنی به ازای X_i همان Y_i را در خروجی قرار می‌دهد. سوال اصلی را البته بایستی جواب دهیم، اینکه برای یک نقطه جدید X_{in} خروجی Y_{out} چقدر معتبر می‌باشد؟ جواب

این مسئله را می‌توان اینگونه داد که تعدادی از نقاط را در مسئله درونیابی مورد استفاده قرار نمی‌دهیم و پس از مشخص کردن پارامترهای مسئله، روی نقاط کنار گذاشته شده چک می‌کنیم که خروجی مدل به ازای این نقاط چه میزان به آنچه نیاز داریم نزدیک است. در مثال زیر نقاط استفاده شده در فرایند درونیابی با رنگ قرمز و نقاط آزمون با رنگ آبی مشخص شده‌اند که مشاهده می‌شود که نقاط آزمون فاصله زیادی تا خروجی بدست آمده دارند.



راه حل دوم: یادگیری آماری

در یادگیری آماری رابطه بین ورودی و خروجی را یک توزیع مشترک در نظر می‌گیریم. به عنوان مثال در مسئله سن و درآمد، یک توزیع مشترک می‌تواند به شکل زیر باشد:

$$f_{XY}(x, y) = \frac{1}{\sqrt{\det(2\pi Q)}} e^{-\frac{1}{2} \begin{pmatrix} x & y \end{pmatrix} Q^{-1} \begin{pmatrix} x \\ y \end{pmatrix}}$$

که یک توزیع گوسی می‌باشد و می‌توان آن را به صورت $(X, Y) \sim N(\mu, Q)$ نشان داد. از آنجاییکه هدف، پیدا کردن خروجی Y بر پایه ورودی X می‌باشد، می‌توانیم از روی توزیع مشترک مقدار Y را بدست آوریم. می‌توان نشان داد در شرایط خاص بهترین خروجی به صورت زیر است:

$$E[Y|X] = f(X)$$

که در صورت وجود توزیع مشترک می‌توان آن را محاسبه نمود. به عنوان مثال در رابطه فوق داریم:

$$f(X) = \beta X$$

که در آن

$$\beta = \frac{E[XY]}{E[Y^2]} = \frac{Q_{12}}{Q_{22}}$$

که ماتریس Q را به صورت زیر فرض کرده‌ایم:

$$Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix} = \begin{pmatrix} E[X^2] & E[XY] \\ E[XY] & E[Y^2] \end{pmatrix}$$

داشتن توزیع مشترک عملاً هم آنچه بدان نیاز داریم را در اختیار قرار می‌دهد و هم نیازی به داده برای تعمیم وجود ندارد. در بسیاری از مواقع توزیع مشترک وجود ندارد و تنها داده‌های ورودی و خروجی در دسترس است. دو کار می‌توان انجام داد:

۱. توزیع مشترک را بیابیم.
۲. مستقیماً $E[Y|X]$ را پیدا کنیم.

روش دوم را فعلاً انتخاب می‌کنیم و فرض می‌کنیم که $E[Y|X]$ از یک دسته تابع مشخص با پارامتر β درست شده است.

$$E[Y|X] = f_{\beta}(x)$$

در این شرایط رابطه بین X و Y را میتوان به صورت زیر نوشت:

$$y = f_{\beta}(x) + \epsilon$$

که در آن ϵ خود یک متغیر تصادفی است که دارای توزیع مشترک با X بوده و

$$E[\epsilon|X] = 0$$

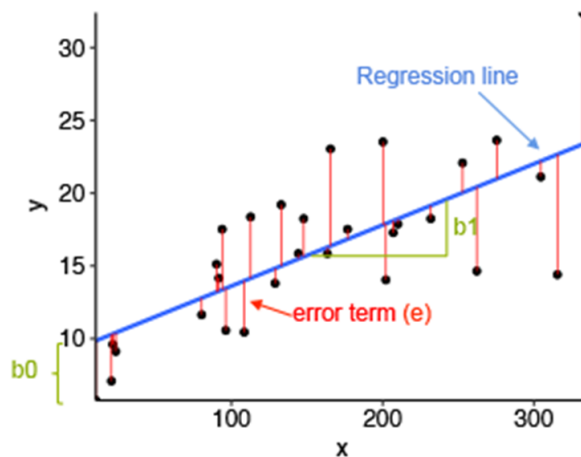
در اغلب موارد و برای سادگی فرض میشود که X و ϵ از هم مستقل بوده و در نتیجه:

$$E[\epsilon|X] = E[\epsilon] = 0$$

نکته ۱. اگر Y گسسته باشد، آنگاه فرض استقلال نمی‌تواند برقرار باشد. جز در شرایط خاص که مثلاً تابع f را خود گسسته بدانیم.

در دیدگاه آماری چون هدف دیگر برازش نیست آنگاه میتوانیم تابع f را با آزادی بیشتری انتخاب کنیم. به عنوان مثال در داده‌های زیر تابع خطی را در نظر میگیریم و خواهیم داشت:

$$E[Y|X] = \beta X$$



شکل ۱: $f_{\beta}(x) = \beta x$

حال باید از روی داده مقدار β را محاسبه کنیم:

$$y_1 = \beta x_1 + \epsilon_1$$

$$y_2 = \beta x_2 + \epsilon_2$$

...

$$y_n = \beta x_n + \epsilon_n$$

که به دلیل وجود $\epsilon_1, \dots, \epsilon_n$ عملاً نمیتوانیم بی نهایت جواب پیدا کنیم. برای رهایی از این مشکل باید فرضیاتی داشته باشیم که برای رسیدن به این فرضیات از نظریه احتمالات بهره میبریم

قضیه ۱ (قانون اعداد بزرگ). اگر متغیرهای تصادفی X_1, \dots, X_n متغیرهای تصادفی *i.i.d* باشند، داریم:

$$\frac{\sum_{i=1}^n X_i}{n} \rightarrow \mathbb{E}[X]$$

با بهره گیری از قانون اعداد بزرگ و فرض $\mathbb{E}[\epsilon] = 0$ خواهیم داشت:

$$\sum_{i=1}^n y_i = \beta \sum_{i=1}^n x_i + \sum_{i=1}^n \epsilon_i$$

بنابراین:

$$\hat{\beta} = \frac{\sum_{i=1}^n y_i}{\sum_{i=1}^n x_i}$$

میتواند تقریب خوبی از β اصلی باشد. حال اگر $\sum x_i = 0$ باشد، $\hat{\beta}$ قابل محاسبه نیست. فرض بعدی این است که واریانس متغیر تصادفی ϵ کمینه باشد یعنی β طوری انتخاب شود که σ_{ϵ}^2 کمینه شود.

به این ترتیب، طبق قانون اعداد بزرگ داریم:

$$\sigma_\epsilon^2 \approx \frac{\sum_{i=1}^n \epsilon_i^2}{n} = \frac{1}{n} \sum_{i=1}^n (y_i - \beta x_i)^2$$

در این حالت β که تقریب اعداد بزرگ مذکور را کمینه میکند به صورت زیر خواهد بود:

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

حال مشابه قبل این سوال مطرح است که مدل بدست آمده را چطور ارزیابی کنیم؟

راه حل پیشنهادی قبل را میتوانیم دوباره مدنظر قرار دهیم: تعدادی از داده ها را کنار گذاشته و در فرآیند یادگیری از آنها استفاده نکنیم. سپس در مرحله ازمون کیفیت مدل را بر حسب آنچه خروجی داده میشود ارزیابی کنیم.