



یادگیری ماشین

نیم سال اول ۱۴۰۳-۱۴۰۲

مدرس: دکتر سید ابوالفضل مطهری

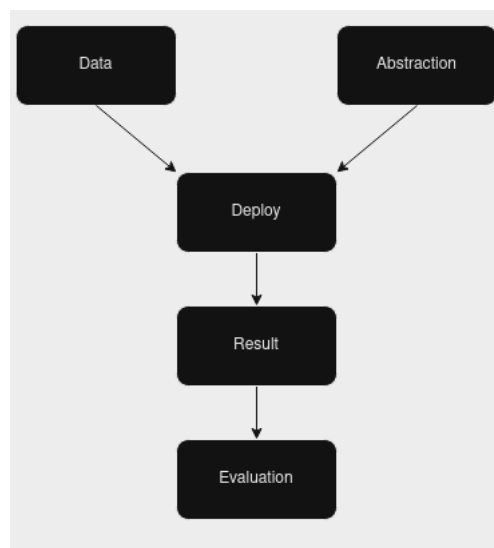
درسنامه یکم

فرآیند یادگیری ماشین

یادگیری ماشین با داده شروع می شود. سپس آن داده با استفاده از Abstraction که نشان دهنده طرز تفکر و چگونگی نگاه به داده است، به پیاده سازی یا Deployment می رسد. از این پیاده سازی نتایجی گرفته می شود و در نهایت این نتایج ارزیابی می شوند.

داده

تعدادی عدد است، مانند یک ماتریس. اگر عدد نباشد هم آن را به عدد تبدیل می کنیم.



شکل ۱: فرآیند یادگیری

اهداف یادگیری

در یادگیری ماشین ما سعی می کنیم یک الگو یا pattern را در داده شناسایی کنیم. همچنین داده ها را بر اساس یک دسته از ویژگی ها دسته بندی کنیم. گاهی هم نیاز داریم که بفهمیم چگونه می توانیم یک داده جدید تولید کنیم؟ بنابراین با توجه به نوع Abstraction ما، مدل ها به دو دسته Discriminative و Generative تقسیم بندی می شوند. در مدل های Generative برای تولید خروجی های مختلف نیاز به Random seed داریم.

Karyotyping

برای مثال، ما نمونه‌های تصاویر کروموزوم‌ها را از آزمایشگاه گرفته و به همراه label خود، در یک ماتریس قرار می‌دهیم. پس از دسته‌بندی ویژگی‌های داده‌ها یا همان pre-processing، اگر نمونه جدیدی به ما داده شود، می‌توانیم آن را دسته‌بندی کنیم و مدل ما Discriminative است. به عبارت دیگر، مدل Generative توزیع $\mathbb{P}(x)$ را مدل می‌کند و مدل Discriminative توزیع $\mathbb{P}(y|x)$.

حفظ کردن یا تعمیم دادن

فرض کنید ما داده‌های سن و درآمد افراد را به صورت زیر داریم:

سن	درآمد
20	0
30	10
⋮	⋮

30 → box → 10

32 → → ?

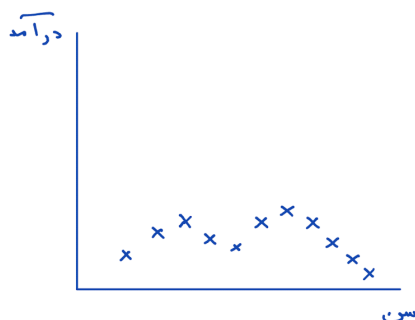
تعمیم

هدف ساختن جعبه‌ای است که با دریافت یک x_i از ورودی y_i صحیح را در خروجی بدهد. یک روش حفظ کردن داده‌ها مثل یک جدول است ولی با استفاده از آن با دو مشکل رو به رو هستیم، اول اینکه ممکن است یک داده ورودی چند خروجی در جدول ما داشته باشد و دوم اینکه اگر داده‌ی جدیدی به عنوان ورودی بیاید چه اتفاقی می‌افتد؟

بنابراین سعی می‌کنیم از ترکیب تعمیم دادن Generalization و حفظ کردن Memorization استفاده کنیم. راه حل ما برای ترکیب این دو، درونیابی است.

Interpolation یا درونیابی

فرض کنید داده‌ها را به صورت نقاطی در شکل زیر نشان داده‌ایم:



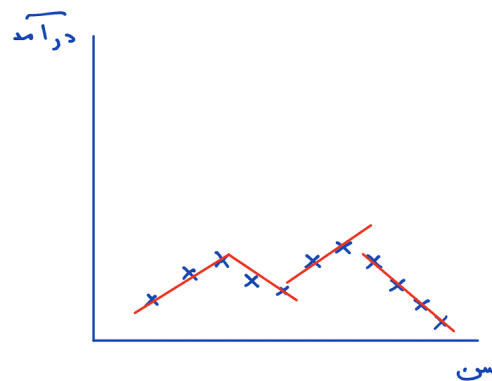
اگر داشته باشیم

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

سعی می‌کنیم پارامترهای β_i و تابع f را به گونه‌ای پیدا کنیم تا از تمام نقاط بگذرد. یعنی

$$f_{\beta}(x) = \sum_{i=0}^5 (\beta_i x^i)$$

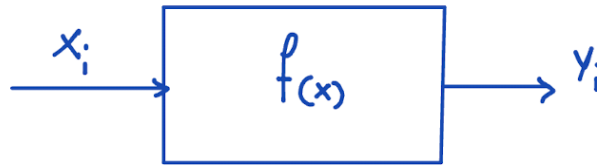
به این دیدگاه که پارامترها را برای کل دامنه پیدا کرده و ثابت در نظر بگیریم، درونیابی Global می‌گوییم. در مقابل دیدگاه گلوبال دیدگاه Local را داریم که مانند شکل زیر دامنه را بازه‌بندی کرده و سپس برای هر بازه پارامترهای جداگانه در نظر می‌گیریم.



تفاوت Regression و Classification و توصیف مدل و ارزیابی آن

اگر بتوانیم برای y_i ها فاصله تعرف کنیم و از آن استفاده کنیم می‌توانیم از Regression استفاده کنیم و در غیر این صورت از classification باید استفاده کنیم. X می‌تواند هر چیزی باشد و در نهایت می‌توانیم آن را به عدد تبدیل کنیم.

x	y
x_1	y_1
x_2	y_2
x_3	y_3
x_4	y_4
\vdots	\vdots



اگر $f(x)$ را $f_\beta(x)$ در نظر بگیریم و داده ها را بصورت $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ هدف ما این است که از روی نمونه های داده شده که D هستند β را بیابیم.

برای این کار سه رویکرد وجود دارد که شامل memorization و prediction و inference می شود که در memorization تعمیم یا generalization وجود ندارد در صورتی که در prediction وجود دارد و تفاوت prediction و inference این است که prediction نگاه black box به موضوع دارد و به مکانیزم کاری ندارد و مهم این است که برای این ورودی این خروجی را می گیریم ولی در inference نگاه white box وجود دارد و فرایند و رابطه ی میان داده های موجود را مورد بررسی قرار میدهد.

حال فرض کنید D که نمونه های ماست از یک توزیع احتمالاتی مانند $p(x, y)$ که در جمعیت وجود داشته است آمده باشند. یک مدل مانند $y = f_\beta(x) + \epsilon$ برای predict کردن در نظر می گیریم و باید یک تابع loss برای ارزیابی مدل ارائه شده تعریف کنیم و نمونه ها را به دو قسمت train و test تقسیم می کنیم و برای ارزیابی، امید ریاضی تابع loss را بر روی نمونه های train و test محاسبه می کنیم اگر این امید ریاضی برای نمونه های train و test نزدیک به هم باشد، یعنی مدل ما توانایی تعمیم خوبی دارد.

به عنوان مثال اگر مدل ما $y = \omega x$ باشد، $l(x, y) = (y - \omega x)^2$ تابع loss می باشد و برای ارزیابی مدل برای تعمیم یافته بودن از $\mathbb{E}((y - \omega x)^2)$ استفاده می کنیم.

همچنین طبق قانون اعداد بزرگ برای محاسبه امید ریاضی بر روی داده ها می توانیم از $\mathbb{E}(l) = \frac{1}{m} \sum (l_i)$ استفاده کنیم.

Linear Regression

رگرسیون خطی

در جمعیت:

$$Y = \beta_1^* X + \beta_0^* + \epsilon$$

ϵ از توزیع نرمال $\epsilon \sim N(0, \sigma^2)$ پیروی می کند و X از توزیع $P[X]$ می آید.

ϵ و X مستقل از همند.

$$P[X, \epsilon] = P[X].P[\epsilon]$$

$$P[X, Y] = P[Y|X].P[X] \rightarrow P[(\epsilon = Y - \beta_1^* X + \beta_0^*) | X] P[X]$$

در حد جمعیت می دانیم که چه اتفاقی افتاده است ولی آن را در نظر نمی گیریم و فکر می کنیم که فقط دیتا از جمعیت داریم.

$$L = E[|y - \hat{y}|^2]$$

آیا می توان loss function را صفر کرد؟

اگر دیتای کمی در نظر بگیریم مثلا $p+1$ تا feature یعنی $p+1$ ستون داده داشته باشیم، می توانیم β را طوری بیابیم که تابع loss دقیقا صفر شود اما این خوب نیست چون قابلیت تعمیم پذیری ندارد. یکی از مسایل مهم شناختن ابعاد است و در این شناخت اشتباهات زیادی رخ می دهد. باید در فرایند training حواسمان به دیتا و پارامترها باشد. این که کدام در سطح population است و کدام در سطح sample باید دقت کرد. دیتا را به n, m تقسیم می کنیم و یکی بر روی population است و دیگری بر روی sample.

$$E[|y - \hat{y}|^2] \Rightarrow \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|^2$$

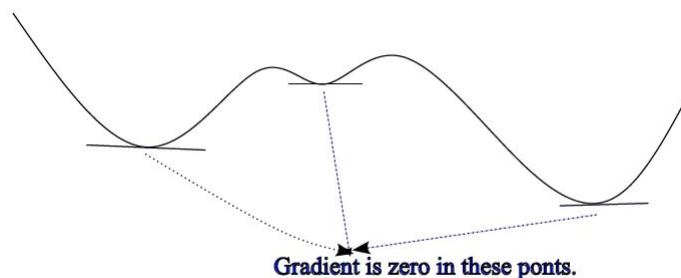
با تخمین از عبارت بالایی β را پیدا می کنیم و از روی آن مدل را ارزیابی می کنیم.

$$L_{test} = \frac{1}{m} \sum_{i=1}^m |y - \hat{y}|^2$$

که m تعداد جمعیت اولیه است.

$$\arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (y_j - X_j^T \beta)^2 = \arg \min g(\beta) \Rightarrow \nabla g(\beta) = 0 \Rightarrow \hat{\beta}$$

با β عبارت بالایی را مینیمم می کنیم. آیا $\hat{\beta}$ که بدست آورده ایم، یک Global Minimum است؟ اگر تابع ما convex باشد، یک مینیمم جهانی دارد. اما اگر non-convex باشد، β لزوما مینیمم جهانی نیست. مساله رگرسیون خطی، قطعا convex بوده و β قطعا جهانی است.



عوض کردن پی در پی مدل برای به جواب رسیدن اشتباه است. تفکر درست این است که ابتدا فضای مدل را مشخص کنیم. مثلا مشخص کنیم که می خواهیم از مدل Neural Net استفاده کنیم یا Random Forest، بعد β ای را پیدا کنیم که loss کمینه شود. بحث fine tune کردن هم مطرح می شود که باید ببینیم برای مساله چه

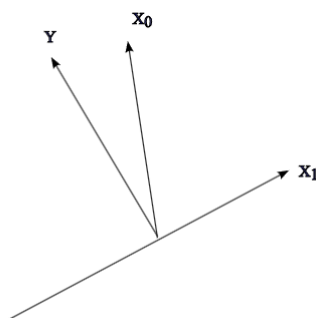
مناسب است.

$$\begin{aligned}\langle \nabla g \cdot \mathbf{d} \rangle &= \lim_{\varepsilon \rightarrow 0} \frac{g(\beta + \varepsilon \mathbf{d}) - g(\beta)}{\varepsilon} \Rightarrow \\ \sum_j \left[(y_j - X_j^T \beta - X_j^T \varepsilon \mathbf{d})^2 - (y_j - X_j^T \beta)^2 \right] &\Rightarrow \\ -2 \left[\sum_j (y_j - X_j^T \beta) X_j^T \right] \mathbf{d} &\Rightarrow\end{aligned}$$

$$\nabla g(\beta) = \frac{-2}{n} \left[\sum_j (y_j - X_j^T \beta) X_j^T \right] = \frac{-2}{n} \left[\sum_j X_j (y_j - \beta^T X_j) \right] \Rightarrow \sum_j X_j Y_j = \sum_j X_j X_j^T \beta$$

$X_j X_j^T$ ، ضرب داخلی دو بردار است و بیانگر شباهت بین دو بردار است. و نیز رابطه هر feature با feature های دیگر را نشان می دهد.

اگر حالتی پیش بیاید که ببینیم یک feature خیلی با y شباهت داشته باشد و بقیه ستون ها این شباهت را نداشته باشند، حذف کردن بقیه ستون ها کار اشتباهی است چون ممکن است ستون ها با y correlation نداشته باشند ولی شاید با همدیگر correlation داشته باشند و این اطلاعات مهمی باشد. مثلاً در شکل زیر X_0 با y ، correlation دارد اما X_1 ندارد. اگر X_1 را دور بیندازیم کار اشتباهی است چون X_1 در مورد X_0 به ما اطلاعات میدهد.



مساله Linear Regression بردار ضرایب β با ترکیب خطی ستون های ماتریس X ، y را می سازد (تخمین می زند).

$$\begin{aligned}\left[X, \quad X_1 \quad \dots \quad X_p \right], \quad \beta &= \left[\beta_0 \quad \beta_1 \quad \dots \quad \beta_p \right] \Rightarrow \hat{y} = X \beta \\ \hat{\beta} &= \arg \min_{\beta} \|y - \hat{y}\|^2 = \|y - X \beta\|^2 \xrightarrow{\frac{\partial}{\partial \beta}} 2 X^T (y - X \hat{\beta}) = 0 \longrightarrow \hat{\beta} = (X^T X)^{-1} X^T y\end{aligned}$$

در نتیجه $\hat{\beta}$ یک تابعی از داده است.

حال فرض کنید یک نمونه جدید به صورت
$$X^{(\cdot)} = \begin{bmatrix} X^{(\cdot)} \\ \vdots \\ X_p^{(\cdot)} \end{bmatrix}$$
 داریم. در این صورت:

$$\hat{y}^{(\cdot)} = X^{(\cdot)T} \beta \longrightarrow L = \mathbb{E} [(y - \hat{y})^2] = \mathbb{E} \left[\left(y^{(\cdot)} - X^{(\cdot)T} \beta \right)^2 \right]$$

در prediction مقادیر β اهمیتی ندارد و مینیمم بودن امید ریاضی برای ما مهم است.

$$y = \langle \beta^*, X \rangle + \epsilon : \quad \text{Var}(\epsilon) = \sigma^2$$

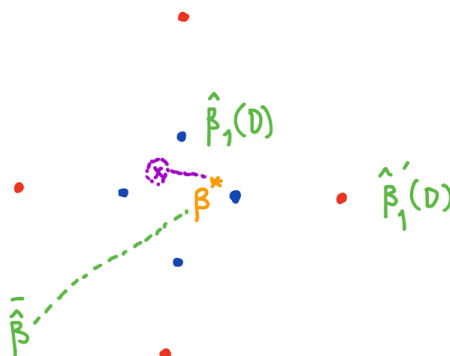
$$\mathbf{L} = \mathbb{E} \left[\left(X^T \beta^* + \epsilon - X^T \hat{\beta}(D) \right)^2 \right] = \mathbb{E} \left[\left\{ \left(X^T (\beta^* - \mathbb{E}(\hat{\beta})) \right) - X^T (\hat{\beta}(D) - \mathbb{E}(\hat{\beta})) \right\} + \epsilon \right\}^2 \right]$$

در عبارت بالا با توجه به اینکه X , ϵ از یکدیگر مستقل هستند می‌توان مقدار Loss را به صورت زیر نوشته

$$\begin{aligned} \mathbf{L} &= \mathbb{E} \left[\left\{ X^T (\beta^* - \mathbb{E}[\hat{\beta}]) \right\}^2 \right] + \mathbb{E} \left[\left\{ X^T (\hat{\beta} - \mathbb{E}[\hat{\beta}]) \right\}^2 \right] + \mathbb{E}[\epsilon^2] \\ &= \underbrace{\mathbb{E}_X \left[X^T (\beta^* - \mathbb{E}[\hat{\beta}]) \right]^2}_{\text{bias}^2} + \underbrace{\mathbb{E} \left[\left\{ X^T (\hat{\beta} - \mathbb{E}[\hat{\beta}]) \right\}^2 \right]}_{\text{Variance}} + \sigma^2 \end{aligned}$$

Bias-Variance trade off

هدف اصلی Linear Regression این است که بهترین $\hat{\beta}$ را پیدا کنیم. یک روش قدیمی به این صورت است که ابتدا bias و سپس Variance را کم کنیم. اما در روش‌های جدیدتر جمع این دو را کم می‌کنیم.



$$\mathbb{E} [\hat{\beta}] = \mathbb{E} \left[(X^T X)^{-1} X^T y \right] = \mathbb{E} \left[(X^T X)^{-1} X^T (X \beta^* + \epsilon) \right] = \mathbb{E} \left[\beta^* + (X^T X)^{-1} X^T \epsilon \right] = \beta^*$$

با توجه به عبارت بالا داریم که unbiased, Linear Regression است.

- **مدل پیچیده:** در این شرایط فضای جست و جو بزرگتری داریم، شانس رسیدن به β^* بیشتر می شود. در نتیجه bias کمتر و variance بیشتر می شود.
- **مدل ساده:** در این شرایط فضای جست و جو کوچکتری داریم، شانس رسیدن به β^* کمتر می شود. در نتیجه bias بیشتر و variance کمتر می شود.
- اگر دیتا به اندازه کافی زیاد نباشد، مدل پیچیدگی کافی را ندارد و شانس کمتری برای رسیدن به جواب وجود دارد و در نهایت مدل ما overfit می شود.

فرض کنید که داده های ما متعامد باشند

$$\hat{\beta} = (X^T X)^{-1} X^T y = \begin{bmatrix} \frac{1}{\|X_{\cdot}\|^2} & & \\ & \ddots & \\ & & \frac{1}{\|X_p\|^2} \end{bmatrix} \begin{bmatrix} \langle X_{\cdot}, y \rangle \\ \vdots \\ \langle X_p, y \rangle \end{bmatrix} \Rightarrow \begin{bmatrix} \beta_{\cdot} \\ \vdots \\ \beta_p \end{bmatrix} = \begin{bmatrix} \frac{1}{\|X_{\cdot}\|^2} \langle X_{\cdot}, y \rangle \\ \vdots \\ \frac{1}{\|X_p\|^2} \langle X_p, y \rangle \end{bmatrix}$$

در نتیجه β_p تنها مربوط به دیتای ستون p است. اگر دیتاها متعامد به یکدیگر نباشند نیز می توانیم با استفاده از الگوریتم Gram-Schmidt داده ها را متعامد کنیم و مقادیر β را بدست آوریم.

استنتاج

در بخش های قبلی دیدیم که هدف برای یک prediction خوب کم تر شدن امید ریاضی loss function است.

$$y = \langle \beta, x \rangle + \epsilon$$

$$l(y, \hat{y})$$

$$\beta(D)$$

$$\mathbb{E}(l(y, \hat{y})) \downarrow ?$$

اما در استنتاج، سوال هایی وجود دارند که قصد داریم به آن ها پاسخ بدهیم:

۱. آیا مقادیر زیر صفر هستند؟ (مثال: دومین سطر به ما نشان می دهد که آیا ستون اول داده با y رابطه ای داشته و informative بوده است یا خیر.)

$$\beta_{\cdot}^* = \bullet ?$$

$$\beta_1^* = \bullet ?$$

$$(\beta_{\cdot}^*, \beta_1^*, \beta_p^*) = \bullet ?$$

۲. آیا مدل خطی مناسب و درست می باشد؟

۳. آیا x_i ها از یک توزیع آمده اند؟

۴. آیا ϵ_i ها از هم مستقل هستند؟

۵. آیا واقعا x_i ها و ϵ_i ها با یکدیگر رابطه ای ندارند؟

در ادامه در ۵ بخش به سوالات مذکور پاسخ می دهیم.

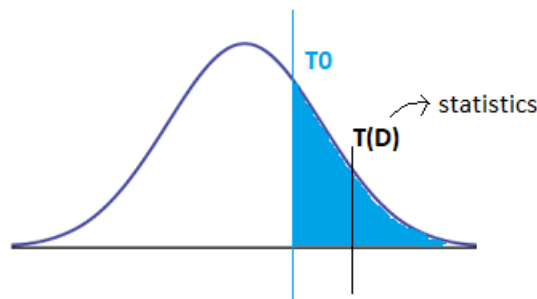
۱

می خواهیم پاسخ این سوال را پیدا کنیم که آیا θ^* برابر صفر هست یا خیر.
می دانیم که بین θ^* و داده رابطه وجود دارد:

$$\theta^* \leftrightarrow D$$

T تابعی از داده ها با توزیعی مستقل از θ^* است که با کمک آن قصد داریم به سوال مورد نظر پاسخ بدهیم. بر اساس این که $T(D)$ از یک threshold مشخص بیشتر یا کمتر باشد می توانیم تصمیم گیری خود را انجام دهیم. پس پاسخ سوال ما به نتیجه مقایسه زیر بررسی دارد:

$$T(D) \leq T.$$



شکل ۲: در مثال مورد نظر خروجی تابع در بازه مورد نظر قرار گرفته و بنابراین فرض اولیه ($\theta^* = 0$) برقرار است.

برای فهم بهتر به مثال زیر توجه کنید:

مثال

فرض کنید توزیع و داده های زیر را در دست داریم:

$$X \sim N(\mu, \sigma^2)$$

$$D(x_1, x_2, \dots, x_n)$$

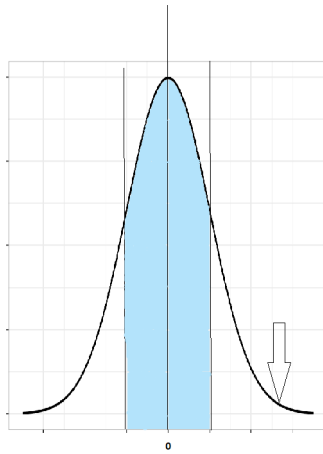
می خواهیم بررسی کنیم که اگر δ^2 برابر ۱ باشد، آیا می توان نتیجه گرفت μ برابر ۰ است:

$$\delta^2 = 2 \Rightarrow \mu = 0 ?$$

داریم:

$$T(D) = \frac{\sigma x_i}{n} \Rightarrow T(D) \sim N(\mu, \frac{1}{n}) \Rightarrow N(0, \frac{1}{n})$$

اگر تست را تعداد بسیار بالا تکرار کنیم مشاهده می کنیم که μ حول و اطراف ۰ قرار نمی گیرد. (فرض می کنیم بازه اطمینان ما ۹۹ درصد است.) بنابراین نتیجه گرفته می شود که فرض $\mu = 0$ نادرست بوده است.



شکل ۳: داده در بازه مد نظر قرار نمی گیرد و بنابراین فرض ما صحیح نیست.

اما فرض کنید که ما مقدار δ^2 را در اختیار نداشتیم. در این صورت آماره ما به δ^2 وابسته می شد:

$$T(D) = \frac{\sigma x_i}{n\delta(\hat{D})} \rightarrow \frac{\sigma(x_i - \bar{x})^2}{n-1}$$

در این صورت توزیع به دست آمده نه از توزیع نرمال، بلکه از توزیع T-Student پیروی می کند.

پس به طور خلاصه، اگر توزیعی داشته باشیم:

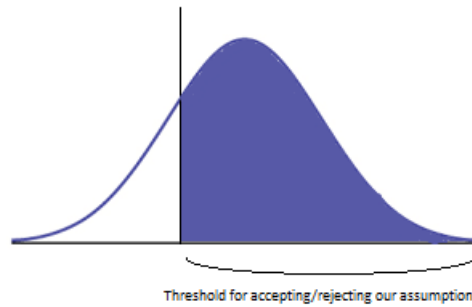
$$p_{\theta_1, \theta_2}(x)$$

$$D = (x_1, x_2, \dots, x_n)$$

برای به دست آوردن پاسخ پرسشی مشابه:

$$\theta_1 = 0 ?$$

آماره $T(D)$ را روی داده تعریف می کنیم و با استفاده از توزیع آن به پاسخ سوال می رسیم. با توجه به این نکته می توان دریافت که برای ما بسیار راحت تر خواهد بود اگر آماره دارای توزیع مشخصی باشد. توجه: $T(D)$ زمانی توزیع خوبی خواهد بود که از θ_2 مستقل باشد.



شکل ۴: توزیع $T(D)$

در ادامه بحث، فرض کنید قصد داریم به سوال اول پاسخ بدهیم و بررسی کنیم آیا مثلاً مقدار β^* برای y رابطه ای داشته یا خیر.

جدول داده زیر را در نظر بگیرید. فرض کنید همه مقادیر y برابر میانگین مقدار آن یعنی \bar{y} بود. (تمامی مقدار ها ثابت). در آن صورت loss برابر ۰ می شد و دیگر برای محاسبه y نیازی به β ها و بررسی داده های x نبود. اما در شرایطی که چنین شرطی برقرار نباشد، ما نیازمند بررسی داده های x می باشیم. در این صورت دو حالت وجود دارد. اگر با استفاده از داده های x امکان توصیف یا describe کردن y به صورت کامل وجود داشته باشد، مجدداً مقدار loss برابر ۰ خواهد بود. اما اگر چنین نباشد، یعنی y به مقدار دیگری به جز x وابسته است که آن را noise می نامیم. اصل کار ما این است که بتوانیم variability در y را با داده describe کنیم. به مثال زیر توجه کنید.

قصد داریم با استفاده از ستون های مذکور رابطه بین x و y را به دست بیاوریم. اگر مقدار ستونی را برابر ۰ قرار دهیم و مشاهده کنیم که تغییری در تخمین ما از y توسط x ها به وجود نمی آید، متوجه می شویم که این ستون تاثیر چندانی ندارد. (informative نمی باشد). مدل زیر را در نظر بگیرید:

$\hat{\beta}_0$	$\hat{\beta}_1$		$\hat{\beta}_p$	y	\hat{y}	$y - \hat{y} = \hat{\epsilon}$
		...				

شکل ۵: داده ها

پس ابتدا با فرض این که هیچ کدام از β ها صفر نیستند محاسبات را انجام می دهیم و \hat{y} و $\hat{\epsilon}$ به دست می آید. سپس برای بررسی هر β مقدار آن را نگه داشته و باقی را صفر قرار می دهیم، و به این ترتیب مقادیر \hat{y} و $\hat{\epsilon}$ به دست می آیند.

$\hat{\beta}_0$	$\hat{\beta}_1 = 0$						$\hat{\beta}_p = 0$	y	$\hat{y}' = \bar{y}$	$y - \hat{y}' = \epsilon'$
		...								

شکل ۶:

داریم:

$$RSS = \sum_n^{i=1} (y - \hat{y})^2$$

$$RSS = \sum_n^{i=1} (y - \hat{y}')^2 = \sum_n^{i=1} (y - \bar{y})^2 = \sum_n^{i=1} \epsilon'^2$$

برای یکی کردن scale ها RSS را به DoF یا degree-of-freedom تقسیم میکنیم. اگر همه داده ها از هم مستقل بودند، DoF برابر n می شد، اما چون $\hat{\beta}$ ها را از روی تمام داده ها ساخته ایم، dependency داریم. پس درنهایت:

$$RSS = \sum_n^{i=1} \frac{(y - \hat{y})^2}{n - (p + 1)}$$

$$RSS = \sum_n^{i=1} \frac{(y - \hat{y}')^2}{n - 1}$$

با مقایسه نسبت RSS و RSS* میتوان به سوال های مورد نظر و به طور کلی سری اول سوال ها پاسخ داد.

۲

اکنون به پاسخ سوال دوم می پردازیم. ما فرض کرده ایم که مدل خطی باشد. یعنی ما مقدار \hat{y} را تنها با استفاده از x به دست آورده ایم. این گزاره زمانی درست است که مدل تنها به x وابستگی داشته باشد. اگر چنین نباشد مدل واقعی دیگر خطی نیست. فرض کنیم مدل واقعی به صورت زیر باشد:

$$y = f(x) + \epsilon$$

و ما تخمین خطی زیر را روی تابع f انجام داده باشیم:

$$f(x) = \langle x, \beta \rangle$$

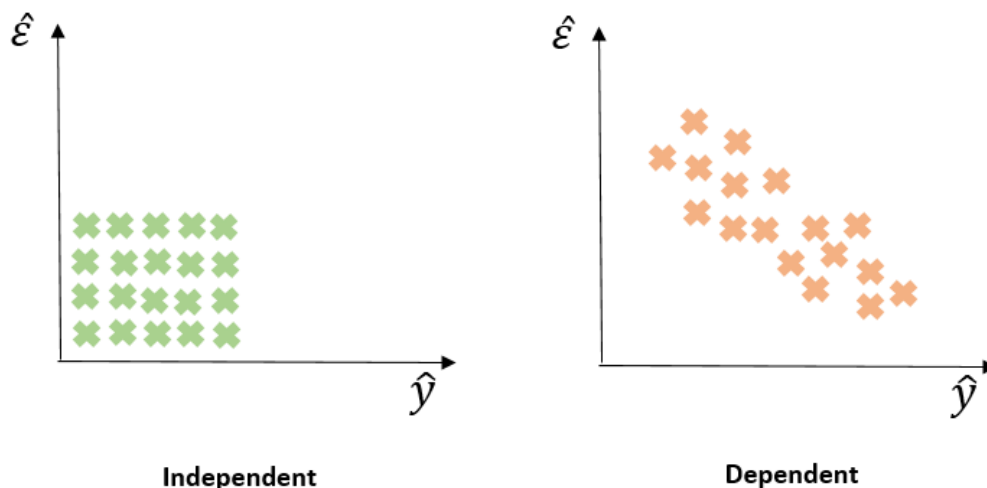
می توان نوشت:

$$y = f(x) - \langle x, \beta \rangle + \langle x, \beta \rangle + \epsilon$$

پس:

$$y - \langle x, \beta \rangle = g(x) + \epsilon \quad (g(x) = f(x) - \langle x, \beta \rangle)$$

$$\Rightarrow \hat{\epsilon} = g(x) + \epsilon$$

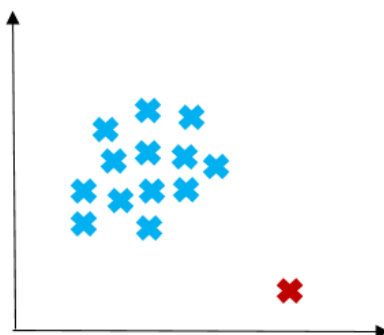


شکل ۷: رابطه $\hat{\epsilon}$ و \hat{y}

اگر مقدار $g(x)$ برابر صفر باشد یعنی مدل خطی بوده است. در غیر این صورت مدل خطی نبوده و \hat{y} و $\hat{\epsilon}$ از یکدیگر مستقل نبوده اند.

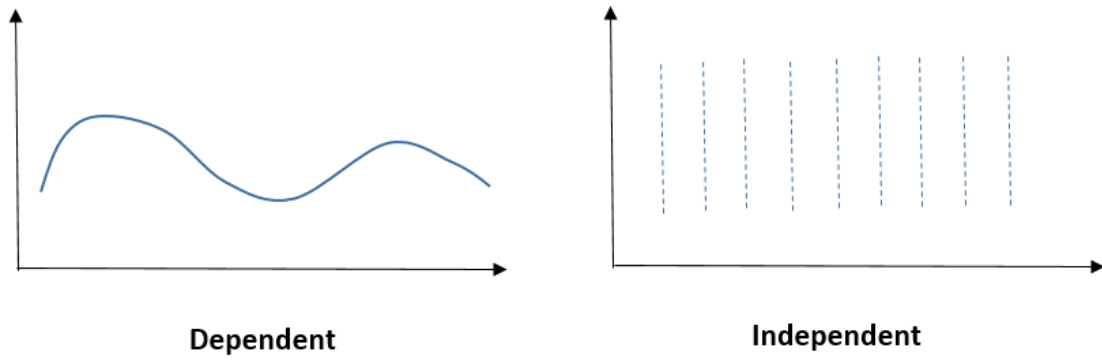
۳

سوال بعدی این است که آیا x_i ها از یک توزیع آمده اند یا خیر. مثلا فرض کنید تستی انجام شده است و یکی از داده ها بسیار پرت از سایر مقادیر است. یا این دیتای پرت valid هست و از یک توزیع دیگر آمده و یا valid نیست و مثلا اندازه گیری غلط بوده است. یک راه این است که برای سایر داده ها distribution ای فرض کنیم و سپس احتمال این که این داده عضوی از این توزیع باشد را محاسبه کرده و اگر این احتمال خیلی پایین بود از این داده صرف نظر کنیم. در چنین شرایطی median بسیار کاربردی است. median مقدار وسطی بعد از sort داده ها می باشد. (از مزایای median این است که robust می باشد).



شکل ۸: داده پرت

برای بستگی مستقل بودن یا نبودن ϵ_i ها از یکدیگر باید رابطه آن ها را با زمان بررسی کنیم. اگر پراکندگی آن ها در طول زمان تغییر نکند یعنی ϵ_i ها از یکدیگر مستقل هستند. (اگر مستقل نباشند یعنی مدل خوبی نداشته ایم).



شکل ۹: ارتباط ϵ_i ها