

همانطور که پیش تر گفته شد برای min یا max تعیین کردن یک تابع نیاز است که مشتق مرتبه دوم آن نیز بررسی شود. در حالت ماتریس مشتق مرتبه دوم با نام ماتریس Hessian شناخته می شود. به این ترتیب نیاز است مثبت یا منفی بودن یک ماتریس تعریف شود. اگر A یک ماتریس با ابعاد  $P \times P$  و بردار  $x$  با بعد  $P$  داشته باشیم آنگاه :

$$A \in \mathbb{R}^{P \times P}, x \in \mathbb{R}^P$$

*تبدیل است*

A is **positive** definite (PD) if  $x^T A x > 0$  for all nonzero  $x$

A is **semi positive** definite if  $x^T A x \geq 0$  for all nonzero  $x$

A is **negative** definite if  $x^T A x < 0$  for all nonzero  $x$

A is **semi negative** definite if  $x^T A x \leq 0$  for all nonzero  $x$

A is **indefinite** if  $\exists x_1, x_2 \Rightarrow \begin{cases} x_1^T A x_1 > 0 \\ x_2^T A x_2 < 0 \end{cases}$

\* می توان نشان داد که اگر یک ماتریس قطری تمام عناصرش مثبت باشد می توان گفت که مثبت معین است.

به عنوان چند تعریف از جبر خطی داریم (البته اگر A و B مربعی باشند)

$$(AB)^T = B^T A^T$$

$$(AB)^{-1} = B^{-1} A^{-1}$$

$$(A^{-1})^T = (A^T)^{-1}$$

یک مجموعه بردار وابسته خطی هستند اگر :

a set of vectors  $\{x_1, \dots, x_n\} \in \mathbb{R}^P$  is said to be linearly dependent

$$\text{if } \exists \alpha = (\alpha_1, \dots, \alpha_n) \neq 0 \Rightarrow \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n = 0$$

در جبر خطی linear combination به فرم زیر تعریف می شود :

$$\text{if } X = \begin{bmatrix} | & | & \dots & | \\ x^0 & x^1 & \dots & x^p \\ | & | & \dots & | \end{bmatrix} \text{ and } B = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$XB = \beta_0 x^0 + \beta_1 x^1 + \dots + \beta_p x^p = \sum \beta_i x^i = \text{linear combination}$$

سرایم

به عنوان یک مثال از تعریف بالا می توان به regression اشاره کرد. به این ترتیب ستون ها در واقع ویژگی ها هستند و سطرها مشاهدات هستند. به این ترتیب ترکیب خطی یک مشاهده ، حاصل ضرب آن سطر در بردار بتاها خواهد بود.

Rank یک ماتریس برابر است با تعداد سطرها یا ستون های آن ماتریس است که مستقل خطی هستند.

### Multi Linear Regression

در بحث رگرسیون خطی دیدم که در حالتی که نمونه ها multi variable می توان محاسبات را به فرم ماتریسی نشان داد. در این حالت نیز برای کمینه کردن تابع loss نیاز است ابتدا نقاط بحرانی را با محاسبه gradient محاسبه کرد و برای اطمینان از نقطه کمینه ماتریس Hessian باید در این نقطه positive definite باشد. به این ترتیب می توان محاسبات را به فرم زیر ادامه داد :

$$L(B) = \frac{1}{2} (Y - XB)^T (Y - XB)$$

$$\nabla_B L(B) = \frac{\partial L(B)}{\partial B} = 0 \quad \& \quad \nabla_B^2 L(B) > 0$$

$$L(B) = \frac{1}{2} (Y^T Y - \underbrace{Y^T X B}_{\text{یا } Y^T X B \text{ است پس برابر است با } Y^T X B} - \underbrace{(XB)^T Y}_{B^T X^T Y} + \underbrace{(XB)^T XB}_{B^T X^T X B})$$

↓  $(Y^T X B)^T = Y^T X B$

$$\mathcal{L}(\beta) = \frac{1}{2} (Y^T Y - 2\beta^T X^T Y + \beta^T \underbrace{X^T X}_{\text{symmetric}} \beta) \Rightarrow$$

$$(X^T X)^T = X^T (X^T)^T = X^T X \Rightarrow \text{symmetric}$$

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta} = \frac{1}{2} (-2 X^T Y + 2 X^T X \beta) = 0 \Rightarrow X^T X \beta = X^T Y$$

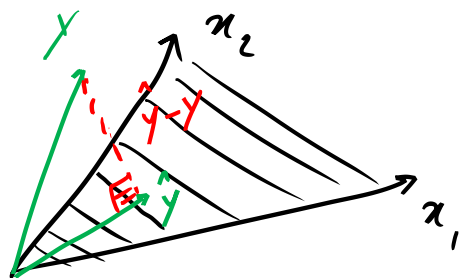
if  $(X^T X) \equiv \text{non singular}$   $\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T Y$

معکوس پذیر باشد

به این ترتیب این معادلات زمانی برقرار هستند که  $X^T X$  وارون پذیر باشد و ماتریس معکوس آن را بتوان حساب کرد. این در حالی است که می دانیم همیشه این وضعیت ممکن نیست.

معادلات مطرح شده را می توان از نگاه هندسی هم بیان کرد. در این دیدگاه می دانیم که  $\hat{Y}$  باید در فضای محاسباتی و پوشاننده  $X$  یا همان  $\text{span}(X)$  یا  $\text{column space}(X)$  قرار گیرد. این درحالی است که  $Y$  می تواند خارج از این فضا باشد. به این ترتیب می توان نشان داده که  $\hat{Y}$  در واقع تصویر  $Y$  در  $\text{span}(X)$  است. این موضوع به خوبی از  $\text{normal equation}$  که به فرم زیر تعریف می شود قابل برداشت است:

$$\nabla_{\beta} \mathcal{L}(\beta) = 0 \Rightarrow X^T Y - X^T X \beta = 0 \Rightarrow \underbrace{X^T (Y - X\beta)}_{\text{normal equation}} = 0$$



در مرحله بعد لازم است که ماتریس Hessian را حساب کنیم ببینیم که positive definite هست یا نه. تا متوجه شویم که نقطه کمینه را بدست آورده ایم یا نه. بنابراین داریم :

$$\nabla_{\beta}^2 \mathcal{L}(\beta) = X^T X$$

با یک مثال ساده می توان متوجه شد که این ماتریس همیشه positive definite است! فرض کنیم که داده ها تنها دو ویژگی دارند. به این ترتیب داریم :

$$X^T X = \begin{bmatrix} x_0[1] & \dots & x_0[n] \\ x_1[1] & \dots & x_1[n] \end{bmatrix} \begin{bmatrix} x_0[1] & x_1[1] \\ \vdots & \vdots \\ x_0[n] & x_1[n] \end{bmatrix} =$$

$$\begin{bmatrix} \sum_{i=1}^n x_0[i]^2 & \sum_{i=1}^n x_0[i] x_1[i] \\ \sum_{i=1}^n x_1[i] x_0[i] & \sum_{i=1}^n x_1[i]^2 \end{bmatrix}$$

if  $A \equiv PD$  then  $\forall Z \Rightarrow Z^T A Z > 0$  so,

$$Z^T = (Z_0, Z_1) \Rightarrow$$

$$Z^T (X^T X) Z = \begin{bmatrix} \sum_{i=1}^n z_0^2 x_0^2[i] + z_1^2 x_0[i] x_1[i] & \sum_{i=1}^n z_0 x_0[i] x_1[i] + z_1^2 x_1^2[i] \end{bmatrix} \begin{bmatrix} z_0 \\ z_1 \end{bmatrix}$$

$$\Rightarrow Z^T (X^T X) Z = \sum_{i=1}^n \left( z_0^2 x_0^2[i] + z_1^2 x_1^2[i] + 2 z_0 z_1 x_0[i] x_1[i] \right) \Rightarrow$$

$$\underbrace{\hspace{10em}}_{(a+b)^2 \text{ اینجا}}$$

$$\text{if } \begin{cases} \omega_0[i] = z_0 x_0[i] \\ \omega_1[i] = z_1 x_1[i] \end{cases} \Rightarrow Z^T (X^T X) Z = \sum_{i=1}^n (\omega_0[i] + \omega_1[i])^2$$

به این ترتیب جواب نهایی همواره بزرگتر یا مساوی صفر است. می توان نشان داد که اگر  $X$  یک ماتریس full rank باشد آن گاه جواب معادله هیچگاه صفر نخواهد بود و به این ترتیب بردار بتا تخمین زده شده نقطه کمینه در loss function را نمایش می دهد.

$$\text{if } \omega_0[i] \neq -\omega_1[i] \Rightarrow Z^T (X^T X) Z = \sum_{i=1}^n (\omega_0[i] + \omega_1[i])^2 > 0 \Rightarrow$$

$$\text{if } \omega_0[i] = -\omega_1[i] \Rightarrow Z^T (X^T X) Z = 0$$

$$\omega_0[i] = -\omega_1[i] \Rightarrow z_0 x_0[i] = -z_1 x_1[i] \Rightarrow x_0[i] = -x_1[i]$$

$$\Rightarrow X \neq \text{full Rank matrix}$$

این اثبات را در حالت کلی هم می توان نشان داد :

$$Z^T (X^T X) Z = \underbrace{(XZ)^T}_{a} XZ = a^T a = \sum_{i=1}^P a_i^2 \geq 0 \Rightarrow$$

$$\text{for } a^T a = 0 \rightarrow XZ = 0 \rightarrow X \neq \text{full Rank} \Rightarrow$$

$$\text{if } X \equiv \text{full Rank} \Rightarrow X^T X > 0$$

پس در کل برای آنکه بردار بتا تخمین زده شده نقطه کمینه در  $\text{loss function}$  را نمایش دهد ، باید دو شرط معکوس پذیر بودن  $x^T x$  و  $\text{full rank}$  بودن ماتریس  $x$  برقرار باشد.

نزدیک ترین نقطه به  $y$  را برمی گرداند

پس در کل برای آنکه بردار بتا تخمین زده شده نقطه کمینه در loss function را نمایش دهد، باید دو شرط معکوس پذیر بودن  $x^T x$  و full rank بودن ماتریس  $X$  برقرار باشد.

به این ترتیب نشان دادیم که می توان بردار بتا را با استفاده از فرمول زیر تخمین زد :

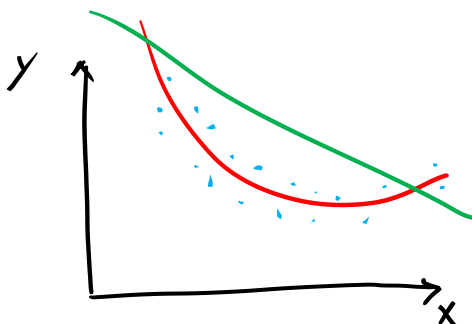
$$\hat{\beta} = (x^T x)^{-1} x^T y \Rightarrow \hat{y} = X \hat{\beta} = \underbrace{x (x^T x)^{-1} x^T}_H y = H y$$

با استفاده از این تعریف می توان ماتریس  $H$  را که با نام های Hat Matrix و Projection Matrix شناخته می شود معرفی کرد. علاوه بر این چون این ماتریس نزدیک ترین نقطه به  $y$  را در  $\text{span}(x)$  برمی گرداند (یعنی عمود است) به آن Orthogonal projection matrix هم گفته می شود. به این ترتیب می توان خواص زیر را برای آن برشمرد:

$$\hat{y} = H y \Rightarrow \begin{cases} H^2 = H \\ H^T = H \end{cases}$$

## Non-linear Regression

تمامی محاسباتی که تاکنون بررسی شد مربوط به رگرسیون خطی بود. حالت هایی وجود دارد که در آن ها یک رگرسیون غیر خطی می تواند تخمین بهتری و خطای کمتری داشته باشد. در شکل زیر مشخص است که رگرسیون غیر خطی که با رنگ قرمز نشان داده شده است، از رگرسیون خطی که با رنگ سبز مشخص شده است، تخمین بهتری را برمی گرداند. به صورت شهودی می توان نتیجه گرفت که اگر در معادلات رگرسیون خطی جملات با درجه بالاتر به کار برده شود می توان مدل های منعطف تر و پیچیده تری تولید کرد. برای مثال در این شکل با اضافه کردن جمله درجه دو می توان رگرسیون قرمز رنگ را تولید کرد. این درحالی است که در فضای درجه دو هنوز رابطه میان  $y$  و ترم درجه دو خطی است.



$$y \approx \beta_0 + \beta_1 x + \beta_2 x^2$$

$$\Rightarrow X = \begin{bmatrix} 1 & x[1] & x^2[1] \\ \vdots & \vdots & \vdots \\ 1 & x[n] & x^2[n] \end{bmatrix}$$

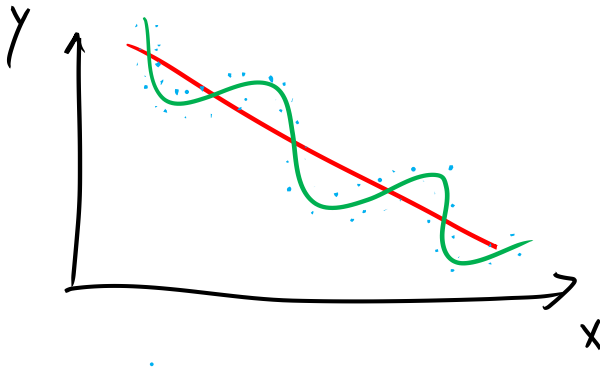
به این ترتیب می توانیم در حالت کلی یک polynomial regression را به فرم زیر تعریف کرد:

$$Y \simeq \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_p x^p$$

مشخص است که اگر این مدل را  $M_2$  بنامیم می توان به خطای کمتری از حالت  $M_1$  بدست آورد.

$$L(M_2) < L(M_1)$$

نکته ای که در اینجا وجود دارد این است که اگر  $P$  را افزایش بدهیم ، انعطاف پذیری مدل در تطبیق ما داده های train بالا می رود و training error کاهش می یابد؛ اما ممکن است این موضوع برای prediction اصلا مناسب نباشد.



به عبارت دیگر ممکن است با پیچیده شدن مدل مشکل overfitting اتفاق بیافتد.

آنچه که بیان شد را می توان به صورت کلی تر نیز نمایش داد. به عبارت دیگر از basis function ها برای ترم های دیگر استفاده کنیم می توانیم یک non-linear regression را در حالت کلی تر داشته باشیم که ترم های آن در واقع توابعی از بردارها هستند. بنابراین می توان رگرسیون را به شکل زیر تعریف کرد:

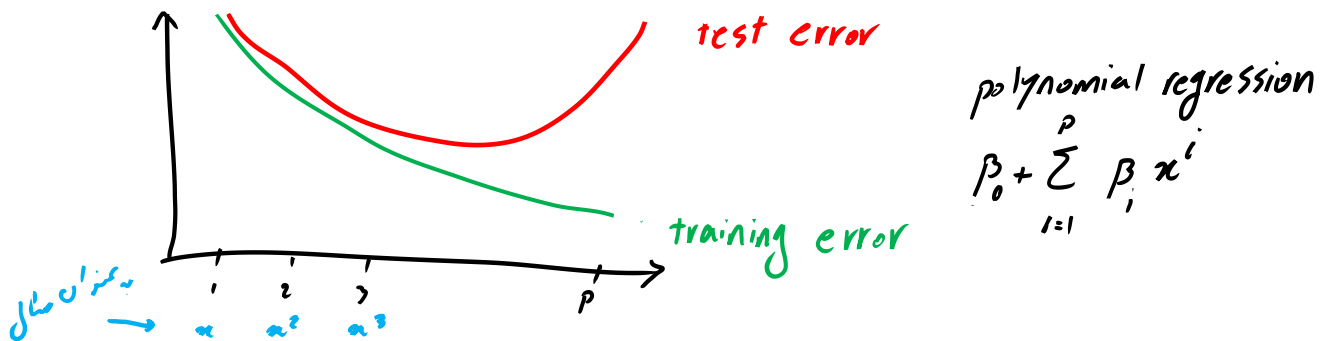
$$Y \simeq \beta_0 + \beta_1 h_1(x) + \beta_2 h_2(x) + \dots + \beta_m h_m(x)$$

$\hookrightarrow R^p(x_1, \dots, x_p)$

موضوع بعدی مربوط به overfitting و generalization است. در واقع آنچه که برای ما مهم است ، پیدا کردن مدلی است که بتواند prediction بهتری روی داد هایی که هنوز دیده نشده انجام دهد و نیز تعمیم بهتری پیدا کند به گونه ای که خطای تعمیم کمی داشته باشد. این موضوع وقتی به وقوع می پیوندد که مدل به اندازه کافی پیچیده باشد تا شرط لازم برای تعمیم یافتگی وجود



داشته باشد. اما در این حالت مسئله overfitting می تواند مشکل ساز باشد. پس کاری که در این شرایط باید انجام داد ، ایجاد ساز و کارهایی است که از overfitting جلوگیری کند. از آنجا که خطای تعمیم روی داده های هنوز دیده نشده اتفاق می افتد راه کار دقیقی برای بدست آوردن آن وجود ندارد و تنها می توان آن را تخمین زد. این کار می تواند با استفاده از validation data صورت گیرد؛ که یعنی بخشی از داده ها را برای تست نگه داریم. راه کار دیگر استفاده از روش cross validation است که ساختار پیچیده تری دارد. شکل زیر خطای تعمیم و مفهوم overfitting را نمایش می دهد.



یکی دیگر از راه کارهای جلوگیری از overfitting استفاده از مفهوم regularized least squares کردن است. فرض کنید فرم زیر معادله رگرسیون ما را تشکیل دهد :

$$Y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_p x^p$$

بدیهی است که در این حالت رگرسیون بسیار پیچیده خواهد بود و امکان overfitting بسیار محتمل می باشد. اما اگر ترم های بتا را صفر کنیم (به غیر از intercept) ، آن گاه مسئله دوباره به حالت بسیار ساده و حتی خطی کاهش پیدا می کند. به این ترتیب یک راه کار جالب می تواند کوچک نگه داشتن ترم های بتا (مثلا تا نزدیک صفر) باشد. در این شرایط می توان از پیچیده شدن رگرسیون جلوگیری کرد. برای این کار کافی است با استفاده از یک متغیر trade off یک loss function جدید تعریف کرد که میزان پیچیدگی رگرسیون به صورت یک پنالتی در loss function ظهور می کند. تنها نکته مهم این است که در ترم پنالتی نباید هیچ پنالتی متوجه intercept باشد. چرا که در غیر این صورت خط رگرسیون همیشه از مبدا خواهد گذشت و این چیزی نیست که ما به دنبال آن باشیم.

$$\sum_{i=1}^p \beta_i^2 = \beta^T \beta \equiv \text{complexity of model} \quad \beta_0 \text{ is not included}$$

$$\mathcal{L}(\beta) = \underbrace{(Y - X\beta)^T (Y - X\beta)}_{\text{fitting the data}} + \underbrace{\lambda \beta^T \beta}_{\text{penalize complexity}}$$

control trade off between fitness & complexity of model

Ridge Regression

به این ترتیب هدف ما عبارتند از :

$$\hat{\beta} = \arg \min_{\beta} \mathcal{L}'(\beta) \quad \beta^T = [\beta_1, \dots, \beta_p]$$

به این ترتیب می توان برای تابع loss جدید محاسبات را تکرار کرد :

$$\nabla_{\beta} \mathcal{L}'(\beta) = \frac{\partial}{\partial \beta} \left( y^T y - 2 \beta^T x^T y + \beta^T x^T x \beta + \lambda \beta^T \beta \right) = 0$$

$$\Rightarrow -2 x^T y + 2 x^T x \beta + 2 \lambda \beta = 0 \Rightarrow$$

$$(x^T x + \lambda I) \beta = x^T y \Rightarrow \hat{\beta}_{ridge} = (x^T x + \lambda I)^{-1} x^T y$$

در این حالت دیگر لازم نیست که عبارت داخل پرانتز non-singular یا همان معکوس پذیر باشد. به عبارت دیگر می توان مقدار لاندا را به گونه تعیین کرد (به مقدار لازم بزرگ در نظر گرفت) که ماتریس حاصل همواره مشتق پذیر و معکوس پذیر باشد. به عبارت دیگر اگر تعداد مشاهدات از تعداد ویژگی ها خیلی کمتر باشد ، می توان نشان داد که ماتریس  $x^T x$  مشتق پذیر نیست. علاوه براین می توان نشان داد که اگر مقادیر قطر اصلی یک ماتریس را از حد مشخصی بیشتر کنیم آن ماتریس مشتق پذیر خواهد شد. بنابراین با انتخاب درست مقدار لاندا می توان ماتریس مذکور را مشتق پذیر کرد.

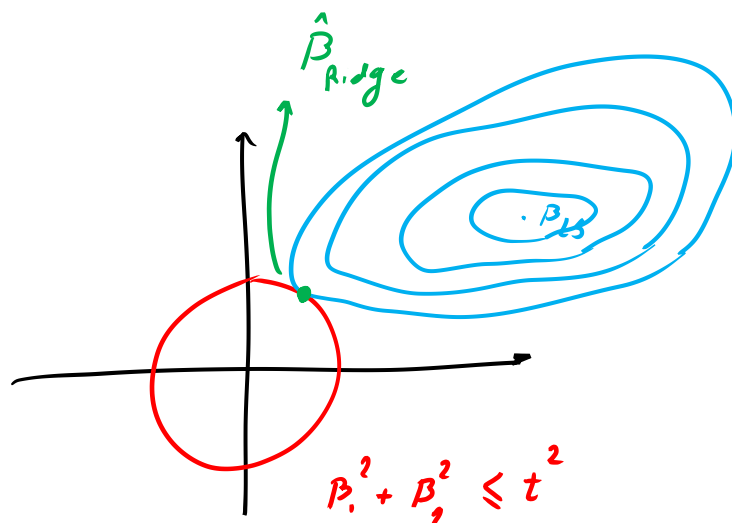
پس مرحله بعدی انتخاب مقدار لاندا است. مقدار لاندا عموماً از طریق cross validation مشخص می شود. اما به صورت شهودی می توان مشاهده کرد که اگر مقدار آن را صفر در نظر بگیریم ، مقدار پناستی را از بین برده و همان فرم linear regression را قبل را خواهیم داشت. در حالتی هم که مقدار آن را بسیار بزرگ در نظر بگیریم (بی نهایت) چون ترم پناستی باعث صفر شدن تمام ترم های بتا می شود پس تخمین ما از  $y$  همان میانگین  $y$  ها ( $\bar{y}$ ) می شود.

این حالت را می توان از نگاه هندسی هم مورد بحث قرار داد. اگر فرمول های بیان شده را به فرم زیر بنویسیم :

$$\hat{\beta}_{Ridge} = \arg \min_{\beta} (y - x\beta)^T (y - x\beta)$$

$$\text{subject to : } \sum_{j=1}^p \beta_j^2 \leq t$$

در این حالت می توان نشان داد که نقطه مورد نظر ما نقطه تلاقی contour plot و دایره constraint است. به عبارت دیگر ما به دنبال نقطه کمینه در تابع loss نیستیم. چرا که این نقطه ممکن است نقطه ای باشد که over fit در آن اتفاق افتاده باشد. به این ترتیب ما با اعمال یک محدودیت یا پناستی سعی می کنیم از آن نقطه فاصله بگیریم و شهودی هندسی آن را می توان به فرم زیر نشان داد.



تا کنون سعی کردیم راه هایی را معرفی کنیم که پیچیدگی مدل را کنترل کنیم. موضوع بعدی بررسی درجه آزادی موثر مدل است که در واقع همان پیچیدگی مدل را نشان می دهد. می توان نشان داد که درجه آزادی یک مدل برابر است با تریس Hat matrix به این ترتیب داریم :

$$\hat{y} = Hy \Rightarrow \text{Degree of freedom (Dof)} = \text{tr}(H)$$

linear regression

$$\text{Dof} = \text{tr}(H) = \text{tr}(X(X^T X)^{-1} X^T) \quad \left\{ \begin{array}{l} \Rightarrow \text{tr}(X^T X)^{-1} X^T X = \\ \text{if } AB \equiv \text{square matrix then } \text{tr}(AB) = \text{tr}(BA) \end{array} \right. = \text{tr}(I_{p \times p}) = p$$

در این تعریف intercept لحاظ نشده است که با احتساب آن درجه آزادی  $P+1$  می شود. مشابه همین محاسبات را می توان برای Ridge در نظر گرفت که به این ترتیب داریم :

$$Dof(\lambda) = \text{tr} \left( X (X^T X + \lambda I)^{-1} X^T \right)$$

برای حساب کردن تریس در این حالت نیاز است که مفهوم singular value decomposition را بررسی کنیم. این مفهوم بیان می کند که هر ماتریس را می توان توسط حاصلضرب سه ماتریس به گونه ای نشان داد که :

$$X_{n \times p} = U_{n \times n} D_{n \times p} V_{p \times p} \quad U^T U = I \quad V^T V = I \quad D = \begin{bmatrix} d_{11} & & \\ & d_{22} & \\ & & \ddots \\ & & & d_{pp} \end{bmatrix}$$

(البته فرمول های نوشته شده مال حالت خاصی است که در آن ها ماتریس U مربعی نیست ولی ماتریس های D و V مربعی هستند)

در این تعریف ماتریس D یک ماتریس قطری است که به مقادیر روی قطر آن singular value گفته می شود. می توان نشان داد که با استفاده از مفهوم SVD ،  $Dof(\lambda)$  برابر است با :

$$Dof(\lambda) = \sum_{i=1}^p \frac{d_{ii}^2}{d_{ii}^2 + \lambda}$$

این تعریف نشان می دهد که اگر مقدار لاندا را صفر در نظر بگیریم درجه آزادی همان P می شود و اگر مقدار آن را خیلی زیاد در نظر بگیریم درجه آزادی به سمت صفر میل می کند.

یکی دیگر از مدل هایی که به دلیل خیلی بزرگ بودن مقدار ویژگی ها از تعداد مشاهدات مورد بحث قرار می گیرد مدل lasso است. در این مدل نرم پناستی نرم اول بردار بتا است. چون نرم اول را می توان معادل قدر مطلق در نظر گرفت ، مشتق نمی توان برای آن محاسبه کرد. به این ترتیب در این مدل راه حل analytical وجود ندارد و از راه حل های convex optimization استفاده می کنیم.

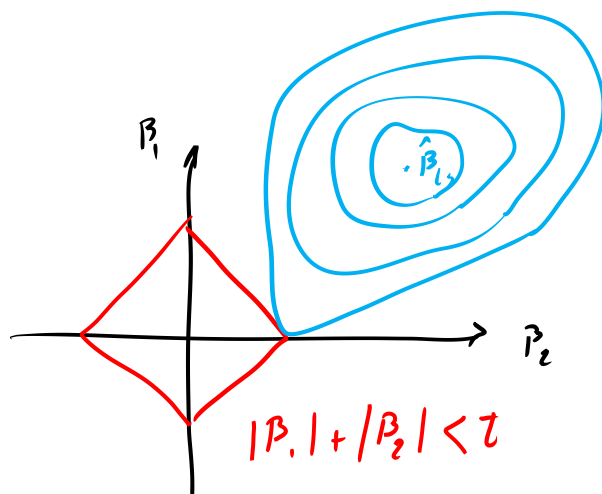
$$L(\beta) = (y - X\beta)^T (y - X\beta) + \lambda \|\beta\|_1 \quad \sum_{i=1}^p |\beta_i|$$

در این حالت هم می توان تفسیر هندسی برای این مدل در نظر گرفت که در شکل نشان داده شده است. نکته ای که در این مدل وجود دارد این است که می توان نشان داد نقطه مورد نظر در گوشه های لوزی بدست می آید. این یعنی اینکه یکی از ویژگی ها

صفر خواهد شد. از این نکته می توان در PCA استفاده کرد. در این مدل درجه آزادی برابر است با تعداد ترم های non zero بردار بتا.

$$\arg \min_{\beta} (y - X\beta)^T (y - X\beta)$$

$$\text{s.t. } \|\beta\|_1 < \tau \equiv |\beta_1| + |\beta_2| < \tau$$



بحث بعدی در ارتباط با خطای تعمیم است. آنچه که ما داریم یک مدل محاسباتی است که بر مبنای داده های train آموزش داده شده است. اما آنچه که برای ما اهمیت دارد این است که این داده ها چقدر روی داده های دیده نشده خوب می توانند عمل کنند. مشکل در اینجا است که خروجی ما روی داده های دیده نشده حکم یک پیش بینی را دارد. به این ترتیب نمی توانیم خطای داده های دیده نشده را حساب کنیم. کاری که می شود کرد این است که روی بخشی از داده های موجود مرحله یادگیری را انجام داد و روی بخش دیگری از داده ها که باقی مانده است، خطای مرحله پیش بینی را اندازه گرفت. بدیهی است که چون داده ها و نوع توزیع آن ها را نمی دانیم، باز هم خطای در نظر گرفته شده در این حالت نمی تواند معادل خطای واقعی باشد. به همین دلیل خطای حاصل یک میانگین یا امید ریاضی خطای اصلی است.

کاری که در این مرحله می توان انجام داد این است که داده ها را به چند دسته یا fold تقسیم کنیم. فرض کنید داده ها را به n دسته تقسیم کردیم. سپس در هر بار از اجرا یک دسته را برای تست و n-1 دسته باقی مانده را برای آموزش مورد استفاده قرار می گیریم. به این ترتیب ما n خطای متفاوت خواهیم داشت. اگر میانگین خطای بدست آمده را در نظر بگیریم، آنگاه cross validation انجام داده ایم.

یک نکته بسیار مهم در استفاده از CV نوع استفاده از داده ها است. به عبارت دیگر وقتی ما داده ها را به چند دسته تقسیم می کنیم و یک دسته را برای تست در نظر می گیریم، این دسته نباید در هیچ یک از مراحل آموزش و ساخت مدل مورد استفاده قرار بگیرند. به عبارت دیگر باید کاملاً برای ما ناشناخته باشند. مثلاً اگر بخواهیم از داده های خود ویژگی استخراج کنیم، ابتدا باید

دسته مربوط به تست را خارج کرده و از دسته های باقی مانده ویژگی استخراج کنیم. سپس میزان خطای مدل را از داده های دسته خارج شده حساب نماییم. به این ترتیب اگر ویژگی ها را با استفاده از تمام داده ها انجام دهیم و بعد دسته ها را تشکیل داده و مراحل تست و یادگیری را انجام دهیم ، خطای بدست آمده ، دارای دقت و اعتبار علمی نخواهد بود.