

lec #6

$$\hat{y} = g(x) = \underbrace{\bar{y}_{\text{red}}}_{b_0} + \underbrace{(\bar{y}_{\text{green}} - \bar{y}_{\text{red}})}_{b_1} x, \quad \text{let } n_g = \sum x_i, \quad P_g = \bar{x} = \frac{n_g}{n}$$

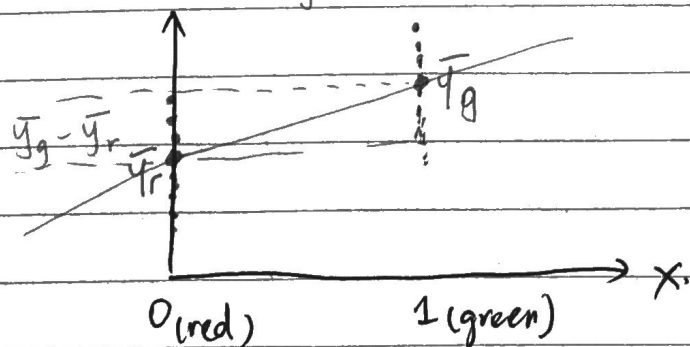
$$\bar{y} = \frac{1}{n} (\sum y_i) = \frac{1}{n} (\sum_{i=\text{green}} y_i + \sum_{i=\text{red}} y_i) = \frac{\sum_{i=\text{green}} y_i}{n} \cdot \frac{n_g}{n_g} + \frac{\sum_{i=\text{red}} y_i}{n} \cdot \frac{n_r}{n_r}$$

$$= P_g \frac{\sum y_i}{n_g} + (1 - P_g) \frac{\sum y_i}{n_r} = P_g \bar{y}_g + (1 - P_g) \bar{y}_r$$

$$b_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} = \frac{n_g \bar{y}_g - n P_g \bar{y}}{n_g - n P_g^2} \cdot \frac{\frac{1}{n}}{\frac{1}{n}} = \frac{P_g \bar{y}_g - P_g \bar{y}}{P_g - P_g^2} = \frac{\bar{y}_g - \bar{y}}{1 - P_g}$$

$$= \frac{\bar{y}_g - P_g \bar{y}_g - (1 - P_g) \bar{y}_r}{1 - P_g} = \frac{(1 - P_g) \bar{y}_g - (1 - P_g) \bar{y}_r}{1 - P_g} = \bar{y}_g - \bar{y}_r$$

$$b_0 = \bar{y} - b_1 \bar{x} = P_g \bar{y}_g + (1 - P_g) \bar{y}_r - (\bar{y}_g - \bar{y}_r) P_g = \bar{y}_r$$



What if $x \in \{\text{red, green, blue}\}$? This is then $p=2$ and we need an OLS solution for $p>1$. But intuitively...

$$g(x) = \begin{cases} \bar{y}_{\text{red}} & \text{if } x=\text{red} \\ \bar{y}_{\text{green}} & \text{if } x=\text{green} \\ \bar{y}_{\text{blue}} & \text{if } x=\text{blue} \end{cases} = \underbrace{\bar{y}_{\text{red}}}_{b_0} + \underbrace{(\bar{y}_{\text{green}} - \bar{y}_{\text{red}})}_{b_1} \underbrace{\mathbb{1}_{x=\text{green}}}_{x_1} + \underbrace{(\bar{y}_{\text{blue}} - \bar{y}_{\text{red}})}_{b_2} \underbrace{\mathbb{1}_{x=\text{blue}}}_{x_2}$$

How well does g predict? We need a "model performance metric". In the SVM this was accuracy or misclassification error. Here, it will can also be what we use internally in the algorithm:

$$\text{SSE} := \sum_{i=1}^n e_i^2 = \sum (y_i - g(x_i))^2$$

Is SSE interpretable? No. Let's take the mean at least, call that mean squared error (MSE)

$$\text{MSE} = \frac{1}{n-2} \text{SSE}$$

But this is still in the squared unit of the phenomenon so it's still uninterpretable. We can take the square root of MSE called root mean squared error (RMSE)

$$\text{RMSE} = \sqrt{\frac{1}{n-2} \text{SSE}} = \sqrt{\frac{1}{n-2} \sum e_i^2} = \sqrt{\text{MSE}}$$

RMSE is in the same unit as y (it is akin the s.d of the residuals S_e)

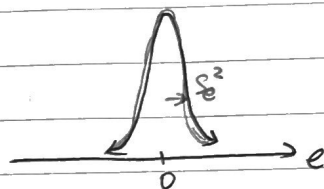
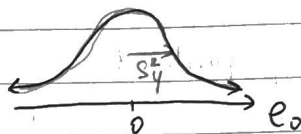
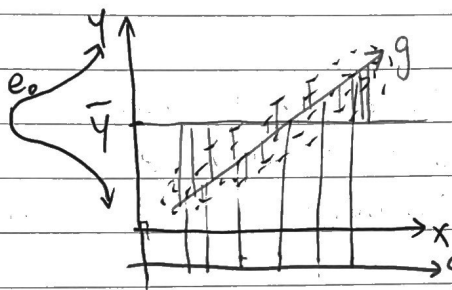
Also, from the CLT,

$[g(x) \pm 1.96 \cdot \text{RMSE}]$ is approx a 95% confidence interval for the true y at the x . RMSE is a very important metric in regression models.

Another important error / performance metric is "R-squared" which is the "proportion of variance explained".

Consider the null model $g_0 = \bar{y}$. What is the SSE of this model? Let's call it SSE_0 .

$$SSE_0 = \sum_{i=1}^n e_{0i}^2 = \sum_{i=1}^n (y_i - \bar{y})^2 = \underbrace{SST}_{\text{sum of squared total}} = (n-1)S_y^2$$

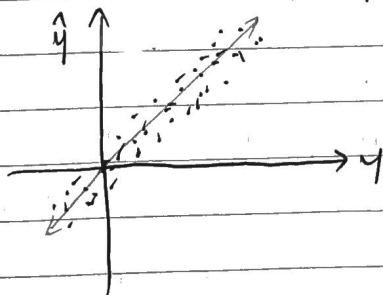


$$\frac{SSE}{SST} = \frac{(n-1)S_e^2}{(n-1)S_y^2} = \frac{S_e^2}{S_y^2}$$

$$R^2 = \frac{SST - SSE}{SST} = \frac{(n-1)S_y^2 - (n-1)S_e^2}{(n-1)S_y^2} = \frac{\overbrace{S_y^2 - S_e^2}^{\Delta S^2}}{S_y^2} = \frac{\Delta S^2}{S_y^2}$$

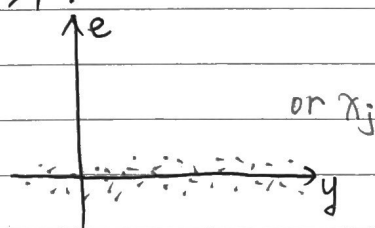
R^2 can never be more than 100%. But R^2 can be negative. This occurs when $S_e^2 > S_y^2$ meaning the model is predicting worse than $g_0 = \bar{y}$.

Here is some other useful plot especially when $p > 1$:



$$R^2 = 1 \Leftrightarrow RMSE = 0$$

$$R^2 \uparrow \Leftrightarrow RMSE \downarrow$$



or x_j

Q: If R^2 is 99%, does this mean the model is for sure "good"?

No. Because if the initial variance was so large, even a 99% reduction wouldn't result a small residual variance i.e. RMSE still could be high after 99% variance reduction

We now would like to generalize the least squares estimation algorithm to cases where $p > 1$. Let's begin with $p=2$.

$$\mathcal{H} = \{w_0 + w_1 x_1 + w_2 x_2 : w_0, w_1, w_2 \in \mathbb{R}\}$$

$$SSE = \sum_{i=1}^n e_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - w_0 - w_1 x_{1,i} - w_2 x_{2,i})^2$$

$$b_0 = \underset{w_0 \in \mathbb{R}}{\operatorname{argmin}} \{SSE\}, \quad b_1 = \underset{w_1 \in \mathbb{R}}{\operatorname{argmin}} \{SSE\}, \quad b_2 = \underset{w_2 \in \mathbb{R}}{\operatorname{argmin}} \{SSE\}$$

This problem can be solved more simply with matrix algebra and a matrix equation:

$$\mathcal{D} = \langle X, \vec{y} \rangle, \text{ Let } X = [\vec{1}_n \quad \vec{x}_{.1} \quad \vec{x}_{.2}] = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ 1 & x_{31} & x_{32} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} \end{bmatrix}$$

$$\text{e.g. } \hat{y}_i = \vec{x}_{i.} \cdot \vec{w}$$

$$\vec{\hat{y}} = X \cdot \vec{w} = \begin{bmatrix} w_0 + w_1 x_{11} + w_2 x_{12} \\ w_0 + w_1 x_{21} + w_2 x_{22} \\ \vdots \\ w_0 + w_1 x_{n1} + w_2 x_{n2} \end{bmatrix}$$

define: $\vec{e} := \vec{y} - \vec{\hat{y}}$

$$SSE = \sum_{i=1}^n e_i^2 = \vec{e}^T \vec{e} = (\vec{y} - \vec{\hat{y}})^T (\vec{y} - \vec{\hat{y}}) = (\vec{y}^T - \vec{\hat{y}}^T) (\vec{y} - \vec{\hat{y}})$$

$$= \vec{y}^T \vec{y} - \vec{y}^T \vec{\hat{y}} - \vec{\hat{y}}^T \vec{y} + \vec{\hat{y}}^T \vec{\hat{y}} = \vec{y}^T \vec{y} - 2 \vec{\hat{y}}^T \vec{y} + \vec{\hat{y}}^T \vec{\hat{y}}$$

$$= \vec{y}^T \vec{y} - 2 (\overset{1 \times n}{X} \overset{n \times 1}{\vec{w}})^T \vec{y} + (\overset{1 \times n}{X} \vec{w})^T \overset{(p+1) \times n}{X} \vec{w} = \vec{y}^T \vec{y} - 2 \vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w}$$

$$\frac{\partial SSE}{\partial \vec{w}} := \begin{bmatrix} \frac{\partial SSE}{\partial w_0} \\ \frac{\partial SSE}{\partial w_1} \\ \vdots \\ \frac{\partial SSE}{\partial w_p} \end{bmatrix}$$

set $\vec{0}_{p+1}$ and solve for b_0, b_1, \dots, b_p .