

2/22/21

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

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

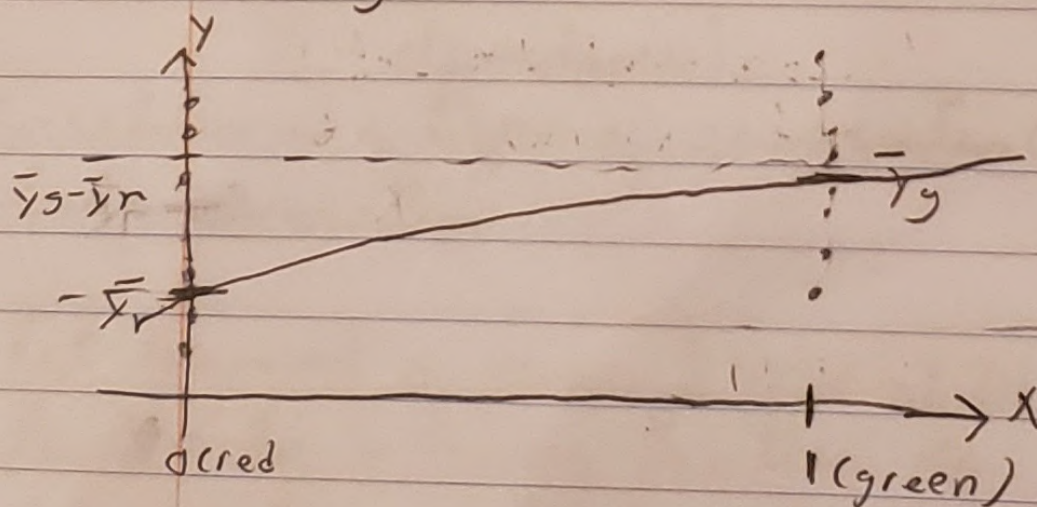
$$= p_g \frac{\sum y_i}{n_g} + (1 - p_g) \frac{\sum_{i \in red} y_i}{n_r} = p_g \bar{y}_g + (1 - p_g) \bar{y}_r$$

$$n_g \bar{y}_g = \sum y_i \quad p_g$$

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

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

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



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

$$\frac{2/22/21}{x = \frac{n_2}{n}}$$

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

$x_1 = \text{green}$ $x_2 = \text{blue}$

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:

$$\bar{y} = \frac{\bar{y}_g - \bar{y}}{1-p}$$

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

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

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

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);

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

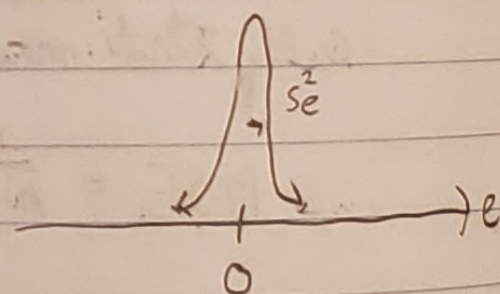
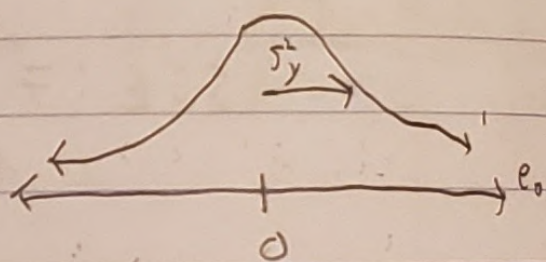
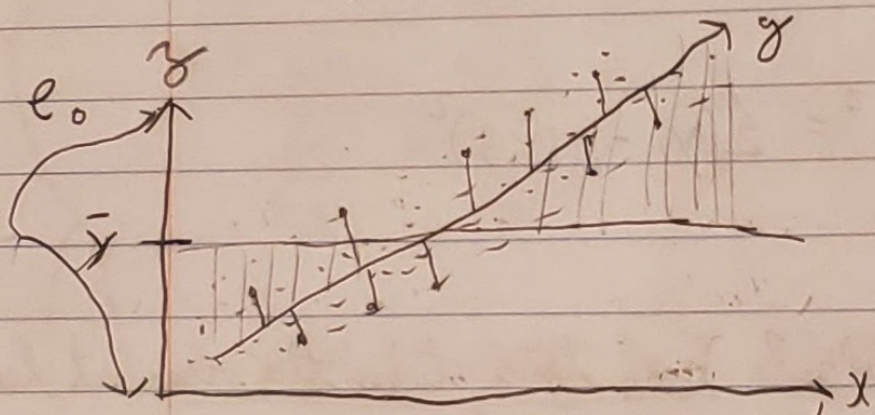
RMSE is in the same unit as y (just like the standard deviation is in the same unit as the random variable). Also, from the CLT, $[g(x) \pm 1.96 \cdot RMSE]$ (it is a bit the the SD of the residuals s_e)

is approx a 95% confidence interval for the true y at that 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", we will now explain this definition

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

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

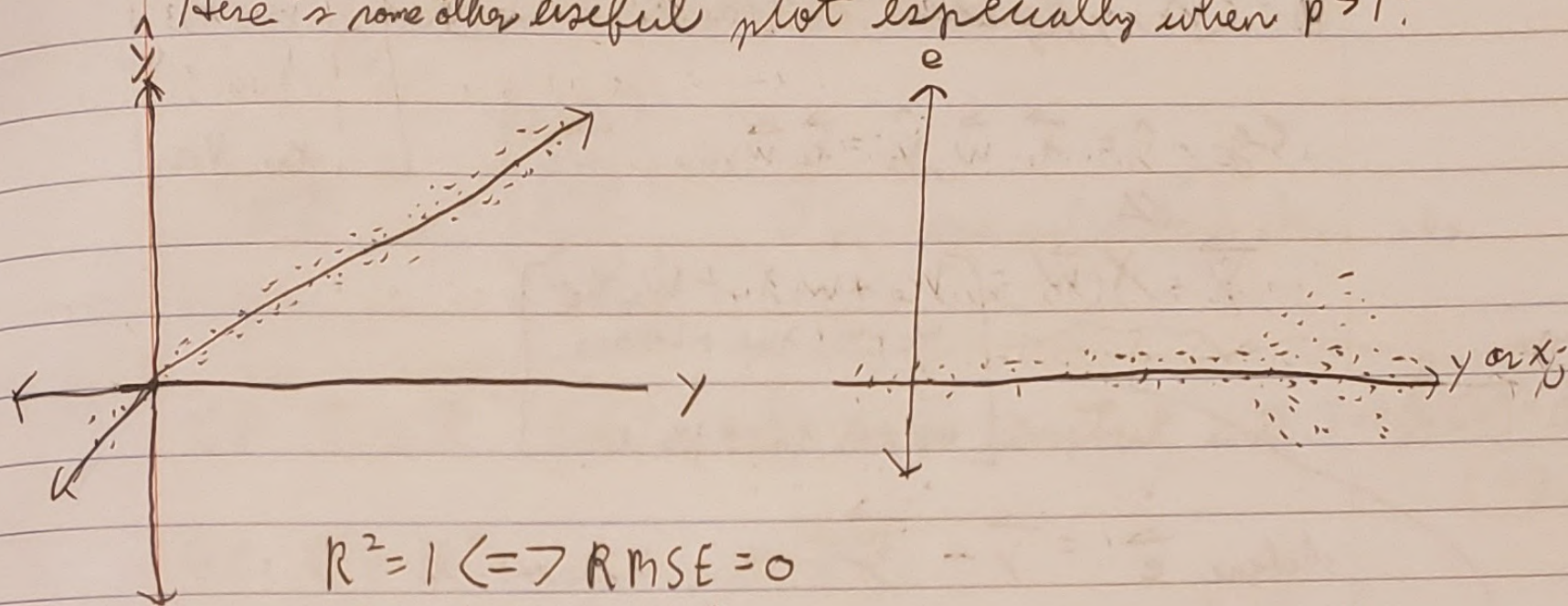


$$\frac{SSE}{SST} = \frac{(n-1)S_e^2}{(n-1)S_y^2} = \frac{S_e^2}{S_y^2}$$

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

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

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



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

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

If $R^2 = 99\%$, does this mean the model is for sure "good"? No, because if the initial variance was so very large, even a 99% reduction wouldn't result in 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 $p = 2$.

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

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

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

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

$$D = \langle X, \vec{y} \rangle, \text{ let } X = [\vec{1}_n \quad \vec{X}_{\cdot 1} \quad \vec{X}_{\cdot 2}] = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ 1 & x_{31} & x_{33} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} \end{bmatrix}$$

$$\text{eg. } \hat{y}_1 = \vec{x}_1 \vec{w}, \hat{y}_2 = \vec{x}_2 \vec{w}, \dots$$

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

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

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

$$\vec{\hat{y}}^T \vec{\hat{y}} = (\vec{y}^T - \vec{\hat{y}}^T) (\vec{y} - \vec{\hat{y}})$$

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

$$= \vec{y}^T \vec{y} - 2 (X \vec{w})^T \vec{y} + (X \vec{w})^T X \vec{w} = \vec{y}^T \vec{y} - 2 \vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w}$$