

2/28/18

Lecture 9

$$\mathcal{H} = \{w_0 + w_1 \tilde{x} : w_0 \in \mathbb{R}, w_1 \in \mathbb{R}\} = \{w_0 + w_1 \mathbb{1}_{x=\text{green}} : w_0 \in \mathbb{R}, w_1 \in \mathbb{R}\}$$

$$\Rightarrow g(x) = b_0 + b_1 \mathbb{1}_{x=\text{green}}$$

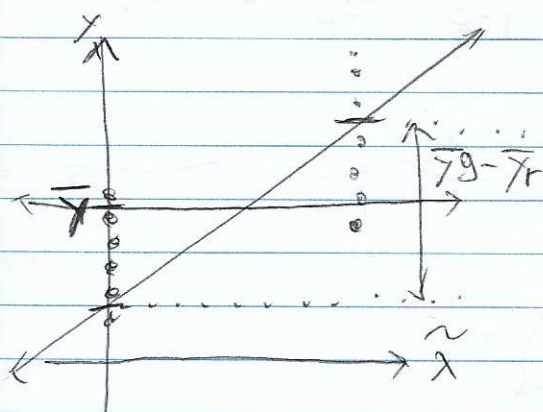
$$g_{\text{null}} = \bar{y}$$

$$b_0 = \bar{y}_r$$

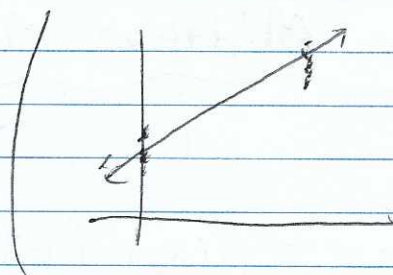
$$b_1 = \bar{y}_g - \bar{y}_r$$

$$g(x) = \begin{cases} \bar{y}_r & \text{if } x = \text{Red} \\ \bar{y}_g & \text{if } x = \text{Green} \end{cases}$$

$$= \bar{y}_r + (\bar{y}_g - \bar{y}_r) \mathbb{1}_{x=g}$$



Here the value of R^2 at \bar{y} is 0



Better R^2 if points close to \bar{y}_g and \bar{y}_r

$$b_0 = \bar{y} - b_1 \bar{x} \quad \text{Let's assume } b_1 = \bar{y}_g - \bar{y}_r \text{ and compare to } b_0.$$

$$b_1 = \frac{r S_y}{S_x} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$$

$$\bar{y} = \frac{y_1 + y_2 + \dots + y_n}{n} = \frac{y_{g_1} + \dots + y_{g_{n_g}} + y_{r_1} + \dots + y_{r_{n_r}}}{n}$$

$$\text{let } p = \frac{n_g}{n}$$

$$= \frac{\sum y_{g_i}}{n} \frac{n_g}{n_g} + \frac{\sum y_{r_i}}{n} \frac{n_r}{n_r} = \bar{y}_g \frac{n_g}{n} + \bar{y}_r \frac{n_r}{n} = p \bar{y}_g + (1-p) \bar{y}_r$$

$$b_0 = (p \bar{y}_g + (1-p) \bar{y}_r) - (\bar{y}_g - \bar{y}_r) = \cancel{p \bar{y}_g + (1-p) \bar{y}_r} - p \bar{y}_g + p \bar{y}_r = \bar{y}_r (1-p) + p \bar{y}_r = \bar{y}_r \checkmark$$

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} = \frac{X_{g1} + \dots + X_{gn} + X_{r1} + \dots + X_{rn}}{n} = \frac{n_g}{n} = p$$

$$\sum x_i y_i = \sum y_i g_i = n_g \bar{y}_g \quad n \bar{x} \bar{y} = n p \bar{y} \quad \sum x_i^2 = n_g \quad n \bar{x}^2 = n p^2$$

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

$$= \frac{\bar{y}_g - (p \bar{y}_g + (1-p) \bar{y}_r)}{1-p} = \frac{\bar{y}_g}{1-p} - \frac{p \bar{y}_g}{1-p} - \bar{y}_r = \bar{y}_g - \bar{y}_r \quad \checkmark$$

Midterm Material End

$$\hat{y} = g(x)$$

$$y = g(x) + e = g(x) + \underbrace{(h^*(x) - g(x))}_{\text{estimation error}} + \varepsilon$$

$$= g(x) + \underbrace{(h^*(x) - g(x))}_{\text{misspecification error}} + (f(x) - h^*(x)) + (f(x) - f(x))$$

$h^* \in \mathcal{H}$ where $h^* :=$ "closest" element $\in \mathcal{H}$ to f
 $g \in \mathcal{H}$

$$h^* = \beta_0 + \beta_1 x \quad g = b_0 + b_1 x$$

To minimize estimation, we need a larger n (law of large numbers)
 (also better \mathcal{H})

To minimize misspecification, better \mathcal{H} (then also \mathcal{H})

To minimize ignorance, choose better X_s that reflect true causal inputs

SVM model (Support vector machine)

1D if linearly separable...

$$\min \|\vec{w}\| \text{ subj to } \forall i \quad (y_i - \frac{1}{2})(\vec{w} \cdot \vec{x}_i + b) \geq \frac{1}{2}$$

~~1D~~ If 1D is not linearly separable...

$$\text{minimize: } \frac{1}{n} \sum_{i=1}^n \max \left\{ 0, \underbrace{-\frac{1}{2} - (y_i - \frac{1}{2})(\vec{w} \cdot \vec{x}_i + b)}_{d_i} \right\} + \lambda \|\vec{w}\|^2$$

d_i : distance away from line on the wrong side of the line
→ Parking Tickets
worse mistakes are more costly