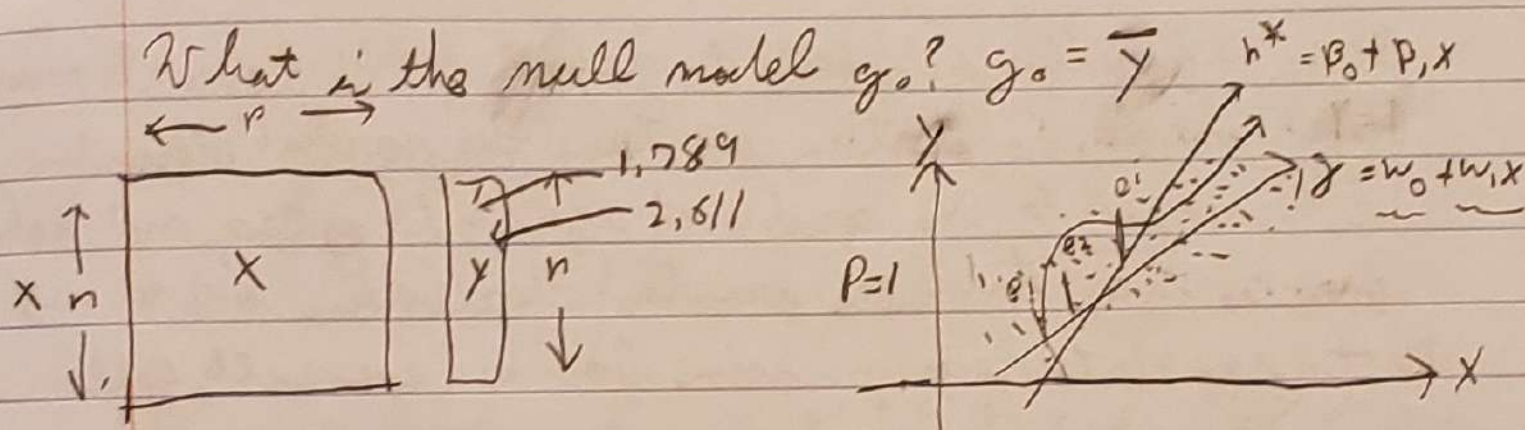


2/17/21

So far, the response space was  $\{0, 1\}$  and the models were "binary classification" models. What if  $y = \mathbb{R}$  or  $y \subset \mathbb{R}$ ? This means the response is continuous and our prediction will be continuous. These models are called "regression" models. The word "regression" is used because of historical circumstances only (see lab).



$$\mathcal{H} = \left\{ \vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{p+1} \right\}$$

$w_0 + w_1 x_1 + \dots + w_p x_p$

Like before, this candidate set, requires, a "1" appended to each of the original  $p$ -length  $x$ -vectors

$$h^*(\vec{x}) = w_0^* + w_1^* x_1 + \dots + w_p^* x_p = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \epsilon$$

Standard notation for the best / "true" value of the linear coefficients.

We have training data and the candidate set of linear models. We need an algorithm that will compute  $w_0$  and  $w_1$  for us. We first need an "objective function" or "error function" or "loss function" which gauges the degree of our model mistakes



Let  $e_i := y_i - \hat{y}_i$ . Consider the loss function:

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

(sum of squared error)

Our algorithm will seek to  $\arg\min \{SSE\}$  over all possible  $w_0, w_1$  values. To do this, we take the partial derivative with respect to  $w_0$  and set equal to zero and solve for  $b_0$ , then take the partial derivative w.r.t  $w_1$  and set equal to zero and solve for  $b_1$ . we will call  $g(x) = b_0 + b_1 x$  the "least squares" regression model or "ordinary least squares" (OLS)

$$\begin{aligned} & \sum y_i^2 + w_0^2 + w_1^2 \sum x_i^2 - 2 \sum y_i w_0 - 2 \sum y_i w_1 x_i + 2 w_0 w_1 \sum x_i \\ &= \sum y_i^2 + n w_0^2 + w_1^2 \sum x_i^2 - 2 w_0 n \bar{y} - 2 w_1 \sum x_i y_i + 2 w_0 w_1 n \bar{x} \end{aligned}$$

$$\frac{\partial}{\partial w_0} (SSE) = \frac{\partial}{\partial w_0} [\quad] = 2 n w_0 - 2 n \bar{y} + 2 w_1 n \bar{x} \stackrel{\text{set}}{=} 0 \Rightarrow b_0 =$$

$$b_0 = \frac{n \bar{y} - w_1 n \bar{x}}{n} = \bar{y} - b_1 \bar{x}$$

$$\frac{\partial}{\partial w_1} (SSE) = \frac{\partial}{\partial w_1} [\quad] = 2 w_1 \sum x_i^2 - 2 \sum x_i y_i + 2 w_0 n \bar{x} \stackrel{\text{set}}{=} 0$$

$$\Rightarrow b_1 \sum x_i^2 = \sum x_i y_i - b_0 n \bar{x} = \sum x_i y_i - (\bar{y} - b_1 \bar{x}) n \bar{x}$$

$$\Rightarrow b_1 \sum x_i^2 = \sum x_i y_i - n \bar{y} \bar{x} + n \bar{x}^2 b_1 \Rightarrow b_1 \sum x_i^2 - b_1 n \bar{x}^2 = \sum x_i y_i - n \bar{x} \bar{y}$$



27  $b_1 = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2}$  this is the answer and now we simplify it... using Maths 241-like notation

$$S_y^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 = \frac{1}{n-1} \left( \sum x_i^2 - 2 \bar{x} \sum x_i + n \bar{x}^2 \right) = \frac{1}{n-1} (\sum x_i^2 - 2n \bar{x}^2 + n \bar{x}^2)$$

$$= \frac{1}{n-1} (\sum x_i^2 - n \bar{x}^2)$$

all possible  
derivatives  
and solve  
w. and  
call  
del or

$$e' = \text{Corr}[X, Y] := \frac{\text{Cov}[X, Y]}{\text{SE}[X] \text{SE}[Y]} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sqrt{\text{Var}[Y] \text{Var}[X]}}$$

Covariance is estimated with

$$S_{xy} = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n-1} \left( \sum x_i y_i - \bar{y} \sum x_i - \bar{x} \sum y_i + n \bar{x} \bar{y} \right)$$

$$= \frac{1}{n-1} (\sum x_i y_i - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y}) = \frac{1}{n-1} (\sum x_i y_i - n \bar{x} \bar{y})$$

$$b_1 = \frac{(n-1) S_{xy}}{(n-1) S_x^2} = \frac{S_{xy}}{S_x^2} \quad r = \frac{S_{xy}}{S_x S_y} \Rightarrow S_{xy} = r S_x S_y$$

$$= \frac{r S_x S_y}{S_x^2} = r \frac{S_y}{S_x} \Rightarrow b_1 = \bar{y} - r \frac{S_y}{S_x} \bar{x}$$

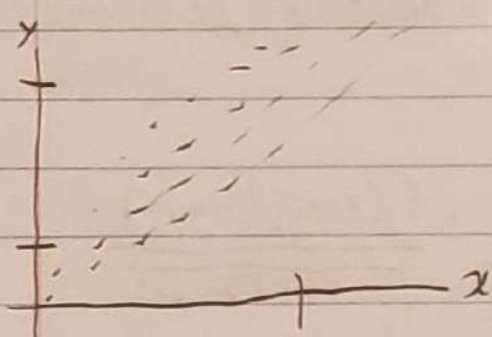
$$\sum x_i y_i - n \bar{x} \bar{y}$$



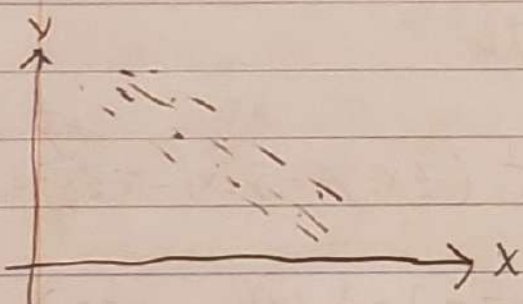
$$\text{Cor}[X, Y] = 0$$



Covariance measures change in expected value of the second rv if the first rv changes



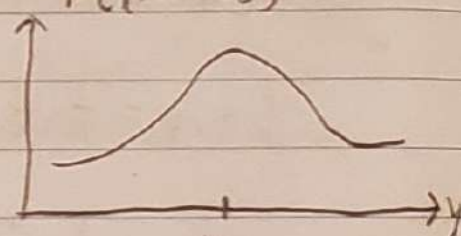
$$\text{Cor}[X, Y] > 0$$



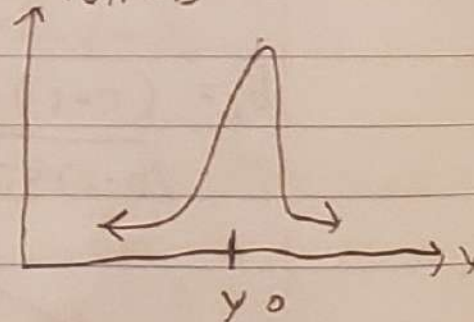
$$\text{Cor}[X, Y] < 0$$

Are  $X, Y$  independent?

$$P(Y|X=x_0)$$



$$P(Y=x_1)$$



The word "association" just means "dependence". Correlation means linear dependence (and covariance mean linear dependence).

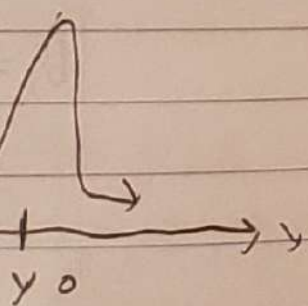
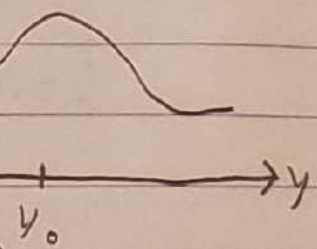
Correlation is a type of association (it is linear association).

Let examine a special case of OLS where  $p=1$ . Let the only feature be a binary feature eg.  $x_i$  is either "red" or "green". Let's create a new  $x$  which is a dummy / binary variable which is 0 if red and 1 if green. What is a good model for prediction



assumes  
predicted value  
is the  
mean

independent?  
( $x_0$ )



the only  
on "green".  
ble which  
on prediction

$$\left. \begin{aligned} g(\text{red}) &= \bar{y}_{\text{red}} \\ g(\text{green}) &= \bar{y}_{\text{green}} \end{aligned} \right\} \text{OLS model}$$

x	y
0	3.71
1	8.43
0	6.72
1	1.07
1	7.11
1	7.10

