

Regression & Correlation

§ 14.1

Let $f(x, y)$ be joint density of r.v.s X and Y .① determine conditional density of Y given $X=x$.

② Conditional mean $\mu_{Y|X} = E(Y|X) = \int_{-\infty}^{\infty} y w(y|x) dy$.

regression equation of Y on X

Remarks: 1 can define regression of Y on X_1 and X_2 (or X_1, \dots, X_n).

$$\mu_{Y|X_1, X_2} = E(Y|X_1, X_2) = \int_{-\infty}^{\infty} y w(y|X_1, X_2) dy$$

2: in the discrete case, we can replace pdf w/ pmf. $\int \rightarrow \text{sum}$.3: regression eqn of Y on X , $\mu_{Y|X}$, is often used for prediction of Y when we observe only X .Ex: Let $(X_1, X_2, X_3) \sim \text{pdf } f(x_1, x_2, x_3) = \begin{cases} (x_1 + x_2) e^{-x_3} & \text{for } x_1, x_2 \in (0, 1), x_3 > 0 \\ 0 & \text{otherwise.} \end{cases}$ find $\mu_{X_2|X_1, X_3}$ sol: marginal density of X_1, X_3 .

$$m(x_1, x_3) = \begin{cases} \int_0^1 (x_1 + x_2) e^{-x_3} dx_2 & x_1 \in (0, 1), x_3 > 0 \\ 0 & \text{elsewhere.} \end{cases}$$

$$= \begin{cases} (x_1 + \frac{1}{2}) e^{-x_3} & x_1 \in (0, 1), x_3 > 0 \\ 0 & \text{elsewhere.} \end{cases}$$

$$\mu_{X_2|X_1, X_3} = E(X_2|X_1, X_3) = \int x_2 \cdot \frac{f(x_1, x_2, x_3)}{m(x_1, x_3)} dx_2$$

$$\begin{aligned}
 \mu_{y_2|x_1, x_3} &= E(x_2 | x_1, x_3) = \int_0^1 x_2 \cdot \frac{f(x_1, x_2, x_3)}{m(x_1, x_3)} dx_2 \\
 &= \int_0^1 x_2 \cdot \frac{x_1 + x_2}{x_1 + \frac{1}{2}} dx_2 \\
 &= \frac{1}{x_1 + \frac{1}{2}} \int_0^1 (x_1 x_2 + x_2^2) dx_2 \\
 &= \frac{\frac{x_1}{2} + \frac{1}{3}}{x_1 + \frac{1}{2}} = \frac{x_1 + \frac{2}{3}}{2x_1 + 1}
 \end{aligned}$$

§14.2 Linear Regression.

$\mu_{y|x}$ is a function of x . In particular, when $\mu_{y|x} = \alpha + \beta x$ is linear, we call it "linear regression". This is good.

Theorem 14.1: If $\mu_{y|x}$ is linear, then:

$$\mu_{y|x} = \mu_2 + \rho \frac{\sigma_2}{\sigma_1} (x - \mu_1)$$

$$\text{where } \mu_2 = EY \quad \mu_1 = EX \quad \sigma_2^2 = \text{Var}(Y) \quad \sigma_1^2 = \text{Var}(X)$$

$$\rho = \frac{\sigma_2 \rightarrow \text{cov}(x, y)}{\sigma_1 \sigma_2}$$

correlation coefficient

Proof: $\mu_{y|x} = \alpha + \beta x$

①

$$\Rightarrow \int y w(y|x) dy = \alpha + \beta x$$

$$\begin{aligned}
 &\nearrow \cdot g(x) \\
 \Rightarrow &\int y \underbrace{w(y|x) g(x)}_{f(x, y)} dy = (\alpha + \beta x) g(x)
 \end{aligned}$$

marginal density on x

$$\Rightarrow \iint y f(x, y) dy dx = \int (\alpha + \beta x) g(x) dx$$

$$\begin{aligned}\Rightarrow EY &= \alpha \int g(x) dx + \beta \int x g(x) dx \\ &= \alpha + \beta EX\end{aligned}$$

$$\Rightarrow \mu_2 = \alpha + \beta \mu_1 \quad (*)$$

$$(2) \quad \stackrel{\cdot x g(x)}{\Rightarrow} \int xy f(x,y) dy = (\alpha x + \beta x^2) g(x)$$

$$\Rightarrow \int xy f(x,y) dy dx = \int (\alpha x + \beta x^2) g(x) dx$$

$$\begin{aligned}\Rightarrow E(XY) &= \alpha E(X) + \beta \int x^2 g(x) dx \\ &= \alpha E(X) + \beta (\text{Var}(X) + (EX)^2)\end{aligned}$$

$$\Rightarrow \sigma_{12} + \mu_1 \mu_2 = \alpha \mu_1 + \beta (\sigma_1^2 + \mu_1^2) \quad (**)$$

Now solve the system $(*, **)$

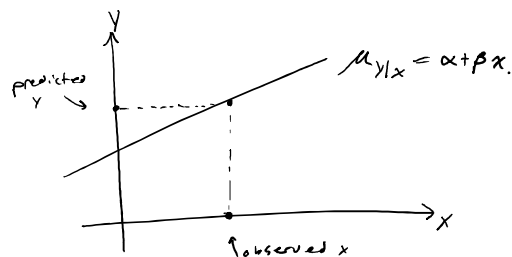
$$\begin{aligned}\Rightarrow \alpha &= \mu_1 - \rho \frac{\sigma_2}{\sigma_1} \mu_1 \\ \beta &= \rho \frac{\sigma_2}{\sigma_1}\end{aligned}$$

HW: if Regression of Y on X is linear, is the regression of X on Y also linear?

Remarks: ① $\rho = 0 \Leftrightarrow \mu_{Y|X}$ is constant (does not depend on x).

when $\rho = 0$, X, Y are uncorrelated, but not necessarily independent. Ex 14.9

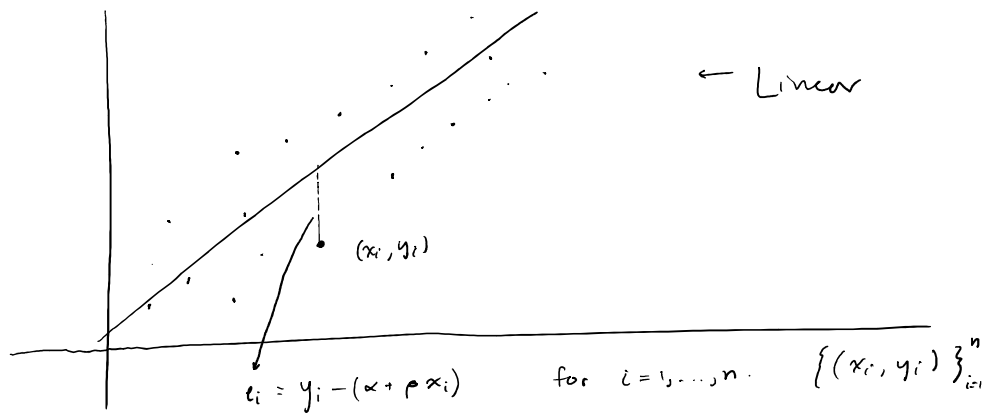
② $-1 \leq \rho \leq 1$ and is slope of regression line



§14.3 Method of Least Squares.

When only paired data is given. "curve fitting."





not linear

is a set of paired data.

The least-squares estimate $\hat{\alpha}^*$ and $\hat{\beta}^*$ minimize $\sum_{i=1}^n e_i^2$.