# Math 4501 - Probability and Statistics II

6.5 - Regression $\longleftarrow$ we will employ MLE technique to determine the regression parameters

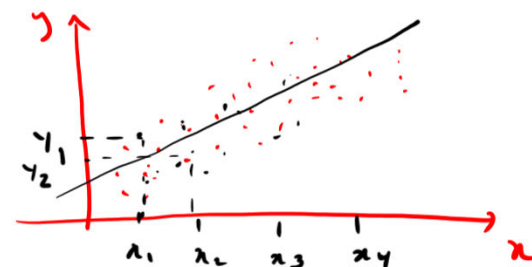# Simplest regression problem

$Y_i - \alpha_1 + \beta x_i = \epsilon_i$

Given the data points

$\rightarrow$ $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$

estimate the parameters $\alpha$ and $\beta$ of the linear model

$$E[Y|x] = \alpha_1 + \beta x ,$$

$\mu(x)$

that is, fit a straight line to the given set of data.

$Y - \underbrace{E[Y|x]}_{\alpha_1 + \beta Y} = \epsilon \sim N(0, \sigma^2)$

$\uparrow$ unknown

$E[Y|x] = \widetilde{\mu(x)} = \alpha_1 + \beta x$

**Assumptions**:

- for each particular value of $x$ the value of $Y$ differs from its mean by a random amount $\varepsilon$.

- the distribution of $\varepsilon$ is $N(0, \sigma^2)$.

$\sigma^2$ is another parameter to estimate

$\left.\begin{array}{l} \end{array}\right\} Y = \widetilde{\alpha_1 + \beta x} + \boxed{\varepsilon}$

**Consequence**: For the linear model described above, we have

$$\begin{array}{c} \mu(x_i) \\ Y_i = \alpha_1 + \beta x_i + \varepsilon_i , \end{array}$$

$\varepsilon_i \sim N(0, \sigma^2)$

where $\varepsilon_i$, $i = 1, 2, \ldots, n$, are independent $N(0, \sigma^2)$ random variables.

**GOAL**: Estimate

$\alpha_1, \beta, \sigma^2$

$\Rightarrow \left\{ Y_i \sim N(\alpha_1 + \beta x_i, \sigma^2) \right\}$

$\mu(x) = \alpha_1 + \beta x$

$y = mx + b \leftarrow$

- For convenience, we set

$\rightarrow \boxed{\alpha_1 = \alpha - \beta \bar{x}},$

$y = y_0 + m(x - x_0)$

where $\bar{x} = \dfrac{1}{n}\sum\limits_{i=1}^{n} x_i$ is the sample mean of the observations $x_1, \ldots, x_n$.

- For each $i = 1, 2, \ldots, n$, we have that

$\alpha_1 = \alpha - \beta\bar{x}$

$Y_i = \alpha_1 + \beta x + \varepsilon_i \longrightarrow \boxed{Y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i}$

is equal to a nonrandom quantity $\alpha + \beta(x_i - \bar{x})$ plus a mean-zero normal random variable $\varepsilon_i$.

- The random variables $Y_1, Y_2, \ldots, Y_n$ are mutually independent normal variables with respective means

estimate $\alpha, \beta, \sigma^2$

$$\alpha + \beta(x_i - \bar{x}), \quad i = 1, 2, \ldots, n$$

$Y_i \sim N\left(\alpha + \beta(x_i - \bar{x}), \sigma^2\right)$

and unknown variance $\sigma^2$.

## Proposition

Under the conditions described above, the maximum likelihood estimators of $\alpha$, $\beta$ and $\sigma^2$ are given by:

$$\widehat{\alpha} = \bar{Y}$$

$$\widehat{\beta} = \frac{\sum_{i=1}^{n}\left(Y_i - \bar{Y}\right)\left(x_i - \bar{x}\right)}{\sum_{i=1}^{n}\left(x_i - \bar{x}\right)^2}$$

$$\widehat{\sigma^2} = \frac{1}{n}\sum_{i=1}^{n}\left[Y_i - \widehat{\alpha} - \widehat{\beta}\left(x_i - \bar{x}\right)\right]^2$$

Interpretation for:

estimates for $Y_i$

$$\widehat{Y}_i = \widehat{\alpha} - \widehat{\beta}\left(x_i - \bar{x}\right)$$

$$\widehat{\beta} = \frac{\frac{1}{n}\sum_{i=1}^{n}\left(y_i - \bar{y}\right)\left(x_i - \bar{x}\right)}{\frac{1}{n}\sum_{r=1}^{n}\left(x_i - \bar{x}\right)^2}$$

$= $ ratio of sample covariance to "sample" variance

$$\widehat{\sigma}^2 = \frac{1}{n}\sum_{i=1}^{n}\left[Y_i - \widehat{Y}_i\right]^2$$

average of the squares of the deviations between actual values and estimated values!

# Distributions of $\widehat{\alpha}$ and $\widehat{\beta}$

- As in the preceding discussion $x_1, x_2, \ldots, x_n$ are treated as nonrandom constants.
- Since the $x$-values are given, when determining the distributions of $\widehat{\alpha}$ and $\widehat{\beta}$, the only random variables are $Y_1, Y_2, \ldots, Y_n$.

## Proposition

*Under the conditions described earlier, we have that:*

*1)* $\widehat{\alpha}$ *is normally distributed with mean* $\alpha$ *and variance* $\dfrac{\sigma^2}{n}$, that is, $\boxed{\widehat{\alpha} \sim N\left(\alpha, \frac{\sigma^2}{n}\right)}$

*2)* $\widehat{\beta}$ *is normally distributed with mean* $\beta$ *and variance*, that is, $\widehat{\beta} \sim N\left(\beta, \dfrac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}\right)$

$$\text{Var}(\widehat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^{n}(x_i - \bar{x})^2}.$$

$\widehat{\alpha} = \bar{Y} = \frac{1}{n}\Sigma Y_i$

$\widehat{\beta} = \dfrac{\sum_{i=1}^{n}(Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^{n}(x_i - \bar{x})}$

CONSEQUENCE: $\widehat{\alpha}$ and $\widehat{\beta}$ are unbiased estimators for $\alpha$ and $\beta$, respectively, because $E[\widehat{\alpha}] = \alpha$ and $E[\widehat{\beta}] = \beta$

**Proof**: By assumption $Y_i = \alpha + \beta(x_i - \bar{x}) + \epsilon_i$, with $\epsilon_i \sim N(0, \sigma^2)$, that is, $Y_i \sim N(\overbrace{\alpha + \beta(x_i - \bar{x})}^{E[Y_i]}, \underbrace{\sigma^2}_{Var(Y_i)})$

1) Recall that $\hat{\alpha} = \bar{Y} = \frac{1}{m} \sum_{i=1}^{m} Y_i$

Since $\epsilon_1, \epsilon_2, \ldots, \epsilon_m$ are independent then $Y_1, Y_2, \ldots, Y_m$ are independent.

Since $\hat{\alpha}$ is a linear combination of independent normally distributed random variables, then $\hat{\alpha}$ is also normally distributed, that is $\hat{\alpha} \sim N(\mu_{\hat{\alpha}}, \sigma^2_{\hat{\alpha}})$

we need to find these two values

$$\mu_{\hat{\alpha}} = E[\hat{\alpha}] = E[\bar{Y}] = E\left[\frac{1}{m} \sum_{i=1}^{m} Y_i\right] \underset{linearity}{=} \frac{1}{m} \sum_{i=1}^{m} \underbrace{E[Y_i]}_{\alpha + \beta(x_i - \bar{x})} = \frac{1}{m} \sum_{i=1}^{m} \alpha + \beta(x_i - \bar{x})$$

$$= \frac{1}{m} \underbrace{\sum_{i=1}^{m} \alpha}_{m\alpha} + \frac{\beta}{m} \underbrace{\sum_{i=1}^{m} (x_i - \bar{x})}_{m\bar{x} - m\bar{x} = 0} = \frac{1}{m} \cdot m\alpha + \frac{\beta}{m} \cdot 0 = \alpha$$

$$\sigma_{\hat{\alpha}}^2 = Var(\hat{\alpha}) = Var\left(\frac{1}{m}\sum_{i=1}^{m}Y_i\right) = \frac{1}{m^2}Var\left(\sum_{i=1}^{m}Y_i\right) = \frac{1}{m^2}\sum_{i=1}^{m}\underbrace{Var(Y_i)}_{\sigma^2}$$

$$\uparrow$$
$$Y_1, Y_2, ..., Y_m$$
$$\text{independence}$$

$$= \frac{1}{m^2}\underbrace{\sum_{i=1}^{m}\sigma^2}_{m\sigma^2} = \frac{1}{m^2} \cdot m\sigma^2 = \frac{\sigma^2}{m}$$

CONCLUSION:

$$\boxed{\hat{\alpha} \sim N\left(\alpha, \frac{\sigma^2}{m}\right)}$$

$E[Y|x]$

DETOUR:

$$\underline{\underline{x}}$$

$\boxed{Y(x) = ???}$ ← not easy

$\boxed{E[Y|x] = \alpha + \beta(x - \bar{x})}$

$$Y_i - \boxed{E[Y|x_i]} = \overbrace{Y_i - \alpha - \beta(x_i - \bar{x})} = \boxed{\varepsilon_i} \leftarrow \text{error or deviation (random)}$$

$$\boxed{Y_i} = \alpha + \beta(x_i - \bar{x}) + \boxed{\varepsilon_i}$$

Proof of item 2:

Recall that $\boxed{Y_i} = \overbrace{\alpha + \beta(x_i - \bar{x})} + \varepsilon_i \sim N(\overbrace{\alpha + \beta(x_i - \bar{x})}^{E[Y_i]}, \overbrace{\sigma^2}^{Var(Y_i)})$, $i = 1, \dots, m$, independent

and that

$$\hat{\beta} = \boxed{\frac{\sum_{i=1}^{m} (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_{i=1}^{m} (x_i - \bar{x})^2}} = \frac{\sum_{i=1}^{m} Y_i (x_i - \bar{x})}{\sum_{i=1}^{m} (x_i - \bar{x})^2}$$

last time

because $\sum_{i=1}^{m} \bar{Y}(x_i - \bar{x}) = \bar{Y} \underbrace{\sum_{i=1}^{m}(x_i - \bar{x})}_{=0} = \bar{Y} \cdot 0 = 0$

Since $\hat{\beta}$ is a linear combination of independent normally distributed r.vs

$$\hat{\beta} = Y_1 \cdot \overbrace{\frac{(x_1 - \bar{x})}{\sum (x_i - \bar{x})^2}}^{a_1} + Y_2 \cdot \overbrace{\frac{(x_2 - \bar{x})}{\sum_{i=1}^{m} (x_i - \bar{x})^2}}^{a_2} + \dots + Y_m \overbrace{\frac{(x_m - \bar{x})}{\sum_{i=1}^{m} (x_i - \bar{x})^2}}^{a_m}$$

these are the $Y_i$s

then $\hat{\beta}$ is itself normally distributed, that is

$$\hat{\beta} \sim N\left(\mu_{\hat{\beta}}, \sigma^2_{\hat{\beta}}\right)$$

to be determined.

Rem:

$$\mu_{\hat{\beta}} = E\left[\hat{\beta}\right] = E\left[\frac{\sum_{i=1}^{n} Y_i (x_i - \bar{x})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}\right] = \frac{1}{\sum_{r=1}^{n} (x_i - \bar{x})^2} E\left[\sum_{i=1}^{n} Y_i (x_i - \bar{x})\right]$$

linearity

$$\underset{\text{linearity}}{\downarrow} = \frac{\sum_{i=1}^{n} (x_i - \bar{x}) \overbrace{E[Y_i]}^{\alpha + \beta(x_i - \bar{x})}}{\sum_{i=1}^{n} (x_i - \bar{x})^2} = \frac{\sum_{i=1}^{n} (x_i - \bar{x})\left(\alpha + \beta(x_i - \bar{x})\right)}{\sum_{i=1}^{n} (x_i - \bar{x})^2}$$

$$= \frac{\sum_{i=1}^{m} \left( \alpha \left( x_i - \bar{x} \right) + \beta \left( x_i - \bar{x} \right)^2 \right)}{\sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2} = \frac{\alpha \overbrace{\sum_{i=1}^{m} \left( x_i - \bar{x} \right)}^{= 0} + \beta \sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2}{\sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2}$$

$$= \frac{\beta \cdot \sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2}{\sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2} = \beta \qquad , \text{ that is } \quad E\left[ \hat{\beta} \right] = \beta$$

Finally, let us compute $\nabla^2_{\hat{\beta}}$ :

$$\nabla^2_{\hat{\beta}} = \text{Var} \left( \hat{\beta} \right) = \text{Var} \left( \frac{\sum_{i=1}^{m} Y_i \left( x_i - \bar{x} \right)}{\sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2} \right) = \frac{1}{\left( \sum_{i=1}^{m} \left( x_i - \bar{x} \right)^2 \right)^2} \text{Var} \left( \sum_{i=1}^{m} Y_i \left( x_i - \bar{x} \right) \right)$$

$Y_1, \ldots, Y_n$ independent

$$\overset{\downarrow}{=} \frac{1}{\left(\sum_{i=1}^{m} (x_i - \bar{x})^2\right)^2} \sum_{i=1}^{m} \text{Var}\left(Y_i (x_i - \bar{x})\right) = \frac{\sum_{i=1}^{m} (x_i - \bar{x})^2 \overbrace{\text{Var}(Y_i)}^{\sigma^2}}{\left(\sum_{i=1}^{m} (x_i - \bar{x})^2\right)^2}$$

$$= \frac{\sum_{i=1}^{m} (x_i - \bar{x})^2 \sigma^2}{\left(\sum_{i=1}^{m} (x_i - \bar{x})^2\right)^2} = \frac{\sigma^2 \sum_{i=1}^{m} (x_i - \bar{x})^2}{\left(\sum_{i=1}^{m} (x_i - \bar{x})^2\right)^2} = \frac{\sigma^2}{\sum_{i=1}^{m} (x_i - \bar{x})^2}$$

CONCLUSION: $\hat{\beta} \sim N\left(\beta, \boxed{\dfrac{\sigma^2}{\sum_{i=1}^{m} (x_i - \bar{x})^2}}\right)$

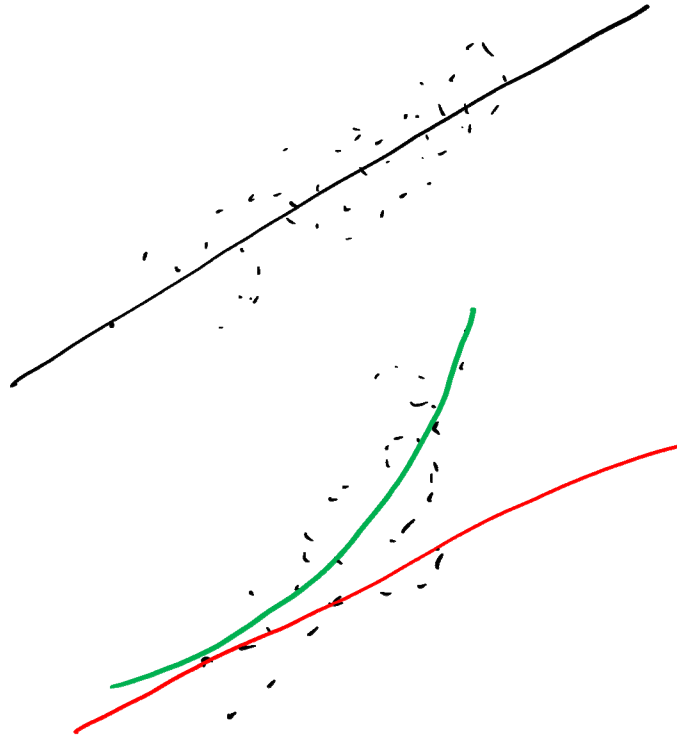$Y_1 = y_1$

$Y_2 = y_2$

$Y_3 = y_3$

$Y_n = y_n$

$x_1$   $x_2$   $x_3$   $- - - -$   $x_n$

$\alpha + \beta x$

$\alpha e^{\beta x}$

$$y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i$$

$$\alpha e^{\beta x}$$

$$\alpha + \beta x + \gamma x^2$$

$$y = \alpha e^{\beta x}$$

$$y = e^x$$

$$\rightarrow \quad z = \ln y = \ln(\alpha e^{\beta x}) = \ln \alpha + \ln e^{\beta x} = \overbrace{\ln \alpha}^{\tilde{\alpha}} + \beta x$$

$$z = \tilde{\alpha}_i + \beta x$$

In general, when we want to estimate an unknown parameter $\underline{\theta} \in \mathbb{R}$.

for a prob. distr. $f(x, \theta)$

We take a random sample $\{X_1, X_2, \ldots, X_m\}$ and determine

$$\underbrace{i.i.d. \ r.v.s}$$

an estimator $\hat{\theta} = \underline{u}(X_1, X_2, \ldots, X_m)$  $\leftarrow$ can are MLE or method of moments

When we actual collect data, we get $X_1 = x_1, \ X_2 = x_2, \ldots, \ X_m = x_m$

$\uparrow$

$x_1, x_2, \ldots, x_m$ are real numbers!!

we can use $x_1, x_2, \ldots, x_m$ to find a point estimate $\hat{\theta} = \underbrace{u(x_1, x_2, \ldots, x_m)}_{\text{a number}}$

With regression: pick values $(x_1), (x_2), \ldots, (x_m)$ non random.

For each value we propose to observe a r.v. $Y_i$, $i = 1, 2, \ldots m$

$Y_1, Y_2, \ldots, Y_m$ are mutually independent

We assume that $Y_i = \alpha + \beta(x_i - \bar{x}) + \varepsilon_i$, $\varepsilon_i \sim N(0, \sigma^2)$

unknown (pointing to $\alpha$, $\beta$)    unknown (pointing to $\sigma^2$)

Use MLE to find estimators

$$\hat{\alpha}(Y_1, Y_2, \ldots, Y_m, x_1, x_2, \ldots, x_m) = \frac{1}{m}\sum_{i=1}^{m} Y_i$$

$$\hat{\beta}(Y_1, Y_2, \ldots, Y_m, x_1, x_2, \ldots, x_m) = \frac{\sum_{i=1}^{m}(Y_i - \bar{Y})(x_i - \bar{x})}{\sum(x_i - \bar{x})^2}$$

$$\widehat{\sigma^2}(Y_1, Y_2, \ldots, Y_m, x_1, x_2, \ldots x_m) = \frac{1}{m}\sum_{i=1}^{m}\left[Y_i - \hat{\alpha} - \hat{\beta}(x_i - \bar{x})\right]^2$$

estimators for $\alpha, \beta,$ and $\sigma^2$ $\leftarrow$

functions of the sample as r.v.

Finally, after observing $Y_1 = y_1, Y_2 = y_2, \ldots, Y_m = y_m$ we can get point estimates

($\varepsilon_1$, $\varepsilon_2$, $\varepsilon_m$)

for $\hat{\alpha}, \hat{\beta},$ and $\widehat{\sigma^2}$ $\leftarrow$ numerical values