



Step 1. Fixing  $a > 0$ , find the best choice for  $b$ .

Let  $Z = Y - aX$ . We have  $E[(Y - aX - b)^2] = E[(Z - b)^2]$ . Notice that the best constant estimator of  $Z$  is  $\delta^* = E[Z]$ . Hence, the best choice of  $b \in \mathbb{R}$ , is  $b^* = E[Z] = \mu_Y - a\mu_X$

where  $\mu_Y = E[Y]$  and  $\mu_X = E[X]$ .

Step 2. Find  $a \in \mathbb{R}$  that minimizes  $E[(Y - \mu_X - a(X - \mu_X))^2]$ .

Notice that  $E[(Y - \mu_X - a(X - \mu_X))^2] = \text{Var}(Y - aX) = \text{Var}(Y) + a^2 \text{Var}(X) - 2a \text{Cov}(X, Y)$ .

By the first order condition,  $a \in \mathbb{R}$  that minimizes above should satisfy,

$$2a \text{Var}(X) - 2 \text{Cov}(X, Y) = 0 \Rightarrow a^* = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} = \frac{\sigma_Y}{\sigma_X} \rho_{X, Y}$$

where  $\rho_{X, Y}$  is correlation coefficient.

$\Rightarrow$  The best linear estimator of  $Y$  given  $X$  is  $L^*(X) = a^*X + b^* = \mu_Y + \sigma_Y \rho_{X, Y} \left( \frac{X - \mu_X}{\sigma_X} \right)$

Then mean square error of  $L^*(X)$  is

$$\begin{aligned} E[(Y - L^*(X))^2] &= \text{Var}(Y - a^*X) = \sigma_Y^2 + (a^*)^2 \sigma_X^2 - 2a^* \text{Cov}(X, Y) \\ &= \sigma_Y^2 + \sigma_Y^2 \rho_{X, Y}^2 - 2 \sigma_Y^2 \rho_{X, Y}^2 = \sigma_Y^2 (1 - \rho_{X, Y}^2) \end{aligned}$$

Important Remarks.

① We also use  $\hat{E}[Y|X]$  to denote best linear estimator of  $Y$  given  $X$ :

$$\hat{E}[Y|X] = L^*(X) = \mu_Y + \sigma_Y \rho_{X, Y} \left( \frac{X - \mu_X}{\sigma_X} \right)$$

$\hat{E}[Y|X]$  is called the wide sense conditional expectation of  $Y$  given  $X$ .

① The correlation coefficient is an integral part of the best linear estimator.

Notice that if instead of estimating  $Y$  using a linear function of  $X$ , we estimate standardized version of  $Y$  ( $\tilde{Y} = \frac{Y - \mu_Y}{\sigma_Y}$ ) using a linear function of standardized version of  $X$  ( $\tilde{X} = \frac{X - \mu_X}{\sigma_X}$ ), the best linear estimator is

$$\hat{E}[\tilde{Y}|\tilde{X}] = \rho_{\tilde{X}, \tilde{Y}} \tilde{X} = \rho_{X, Y} \tilde{X} = \rho_{X, Y} \left( \frac{X - \mu_X}{\sigma_X} \right)$$

② The class of constant estimators is a subset of class of linear estimators, which itself is a subset of class of unconstrained estimators. Hence, we have the following relation between the best estimators we have studied so far.

$$\text{MSE of } \delta^* = E(Y) \geq \text{MSE of } L^*(X) = \mu_Y + \sigma_Y \rho_{X, Y} \left( \frac{X - \mu_X}{\sigma_X} \right)$$

$$\geq \text{MSE of } g^*(X) = E[Y|X]$$

Equivalently, we have

$$\sigma_Y^2 \geq \sigma_Y^2 (1 - \rho_{XY}^2) \geq \sigma_Y^2 - E[(E[Y|X])^2]$$

③ If  $X$  &  $Y$  are independent, we have

$$E[Y] = \delta^* = L^*(X) = g^*(X)$$

as we expected (knowing  $X$  do not reduces our uncertainty about  $Y$ ).

② Law of large numbers

Notice that given  $n$  samples  $X_1, \dots, X_n$  that are independent & identically distributed, the sample average  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$  is unbiased estimator of  $\mu$ . The question is: how accurate is our estimator. Intuitively, more samples should result in a more accurate estimator.

From basic calculus, when we say a sequence of numbers  $y_1, y_2, \dots$  converges to  $y$ , we mean for any  $\delta > 0$ , there is  $N > 0$  s.t.  $|y - y_n| < \delta$  for any  $n > N$ . However, in the case of sample average  $\hat{\mu}_n$ , we are dealing with random number.

Consider the set of bad estimates:  $\left\{ \left| \frac{X_1 + \dots + X_n}{n} - \mu \right| > \delta \right\} = \left\{ \omega \in \Omega : \left| \frac{X_1(\omega) + X_2(\omega) + \dots + X_n(\omega)}{n} - \mu \right| > \delta \right\}$

One way to characterize the accuracy of our estimate is to say the probability of observing bad estimates is small.

This is exactly what the following Law of large numbers say.

**Proposition 4.10.1.** (Law of Large Numbers) Suppose that  $X_1, X_2, \dots, X_n$  are uncorrelated random variables, they have the same mean  $\mu = E[X_1] = \dots = E[X_n]$ , and their variance is less than  $C$ , i.e.,  $\text{Var}(X_i) < C$  for all  $i \in \{1, 2, \dots, n\}$ . For any

$\delta > 0$ , we have: 
$$P\left(\left|\frac{S_n}{n} - \mu\right| > \delta\right) \leq \frac{C}{n\delta^2} \xrightarrow[n \rightarrow \infty]{} 0$$

where  $S_n = X_1 + \dots + X_n$ .

**Proof:** Notice that by Chebyshev's inequality:

$$P\left(\left|\frac{S_n}{n} - \mu\right| > \delta\right) \leq \frac{\text{Var}\left(\frac{S_n}{n}\right)}{\delta^2} = \frac{\text{Var}(S_n)}{n^2\delta^2}$$

Since  $X_i$ 's are uncorrelated, we have

$$\text{Var}(S_n) = \text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n) \leq nC.$$

Hence,

$$P\left(\left|\frac{S_n}{n} - \mu\right| > \delta\right) \leq \frac{nC}{n^2\delta^2} = \frac{C}{n\delta^2}$$

③ Central limit theorem

Suppose that we have  $n$  random variables  $X_1, X_2, \dots, X_n$  which are independent & identically distributed. Law of large numbers

### Central Limit Theorem

Suppose that we have  $n$  random variables  $X_1, X_2, \dots, X_n$  which are independent & identically distributed. Law of large numbers characterizes the accuracy of sample mean estimator. What if we want more, i.e., we want to estimate the distribution of  $S_n$ ?

Notice that  $E[S_n] = n\mu$  and  $\text{Var}(S_n) = n\sigma^2$  where  $\mu = E[X_1]$  and  $\sigma^2 = \text{Var}(X_1)$ . Hence,  $S_n$  itself is not well-behaved, i.e., its mean & variance blows up. Hence, it makes more sense to study the standardized version of  $S_n$  instead, i.e.,  $\tilde{S}_n = \frac{S_n - n\mu}{\sqrt{n\sigma^2}}$ . The central limit theorem characterizes the asymptotic distribution of  $\tilde{S}_n$ .

**Theorem.** (Central Limit Theorem) Suppose that  $X_1, X_2, \dots, X_n$  are independent & identically distributed random variables. Let  $\mu = E[X_1]$ ,  $\sigma^2 = \text{Var}(X_1)$ . Set  $S_n = X_1 + \dots + X_n$ . We have

$$\lim_{n \rightarrow \infty} P(\tilde{S}_n \leq c) = \lim_{n \rightarrow \infty} P\left(\frac{S_n - n\mu}{\sqrt{n\sigma^2}} \leq c\right) = \Phi(c)$$

### Remarks.

- ① We can use Gaussian approximation for sum of independent & identically distributed random variables.
  - ② CLT is called a universal law since it has nothing to do with distribution of  $X_1$ , other than existence of  $\mu$  &  $\sigma^2$ .
-