

## Useful Notes on Statistics

**Moment generating functions:** The moment generating function (mgf) is defined as:

$$M(t) = E[e^{tX}] = 1 + tX + \frac{1}{2!}t^2X^2 + \frac{1}{3!}t^3X^3 + \dots \quad (1)$$

For two independent random variables  $X$  and  $Y$ , with mgfs  $M_1(t)$  and  $M_2(t)$ , mgf of their sum is given by

$$E[X + Y] = E[e^{t(X+Y)}] = E[e^{tX}] E[e^{tY}] = M_1(t) M_2(t) \quad (2)$$

The mgf of the linear function  $a + bX$  is given by

$$E[a + bX] = E[e^{a+btX}] = e^{at} M(bt) \quad (3)$$

**Binomial Distribution:** The probability density function is given by

$$P(x) = \binom{n}{x} P^x Q^{1-x} \quad (4)$$

where  $P$  is the success probability,  $Q = 1 - P$  and  $x$  is the number of observations. The mgf of the binomial distribution is given by

$$M(t) = E[e^{tX}] = \sum_x \frac{n!}{x!(n-x)!} e^{tx} P^x Q^{1-x} = (e^t P + Q)^n \quad (5)$$

The moment generating function is useful for computing mean and variance:

$$M'(0) = E[X] = nP, \quad M''(0) = E[X^2] = nP + n(n-1)P^2 \quad (6)$$

thus,  $\mu = E[X] = nP$  and  $\sigma^2 = \text{Var}[X] = E[X^2] - E[X]^2 = nP(1-P) = nPQ$ .

**Normal Distribution:** The probability density function of the normal distribution is given by

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right] \quad (7)$$

Then, the mgf can be computed as

$$\begin{aligned} M(t) &= E[e^{tX}] = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} dx \exp\left[-\frac{1}{2\sigma^2} ((x-\mu)^2 - 2\sigma^2 tx)\right] \\ &= \exp\left[\mu t + \frac{1}{2}\sigma^2 t^2\right] \end{aligned} \quad (8)$$

Using the mgf of the normal distribution, we can show that the variable  $Z = (X - \mu)/\sigma$  is a normal variate, i.e normal distribution with mean 0 and standard deviation 1 ( $N(0, 1)$ ). The mgf for  $Z$  is given by

$$M_Z(t) = e^{-\mu t/\sigma} M(t/\sigma) = e^{\frac{1}{2} t^2} \quad (9)$$

which is the mgf of the normal variate.

**Central Limit Theorem:** The sum of a large number of independent random variables will be approximately normally distributed. Let us prove: Consider the sum of  $n$  independent random variables  $Y = X_1 + X_2 + \dots + X_n$ . Let  $\mu$  and  $\sigma^2$  be the mean and variance of  $Y$  and let  $M_i(t)$  be the mgf of  $X_i - \mu_i$ . Then, mgf of  $\sum_i (X_i - \mu_i)$  is

$$E[\exp(t(X_1 - \mu_1) + t(X_2 - \mu_2) \dots)] = \prod_i M_i(t) \quad (10)$$

The mgf of the variate  $(Y - \mu)/\sigma$  is then given by

$$M^*(t) = \prod_i M_i(t/\sigma) = \prod_i \left( 1 + \frac{\sigma_i^2}{2} \frac{t^2}{\sigma^2} + \frac{\mu_{3i}}{3!} \frac{t^3}{\sigma^3} + \dots \right) \quad (11)$$

since  $\mu = \mu_1 + \mu_2 + \dots + \mu_n$ . Taking the log of both sides

$$\begin{aligned} \log M^*(t) &= \sum_i^n \log M_i(t/\sigma) \simeq \sum_i^n \log \left( 1 + \frac{1}{2} \frac{\sigma_i^2 t^2}{\sigma^2} \right) \\ &\simeq \sum_i^n \frac{1}{2} \frac{\sigma_i^2 t^2}{\sigma^2} = \frac{1}{2} t^2 \end{aligned} \quad (12)$$

since for larger  $n$ ,  $\sigma^2$  will be large as well ( $\sigma^2 = \sum_i \sigma_i^2$ ), so the series expansion in  $\sigma_i/\sigma$  is convergent. Thus, the mgf of the variate  $(Y - \mu)/\sigma$  is  $e^{\frac{1}{2} t^2}$  which is the mgf of the standard normal variate.

A corollary of the central limit theorem is that the distributions of the sample means is approximately normally distributed: Let  $X_i$  be a sample from a population. The sample mean of  $n$  samples is given by

$$\bar{X}_n = \frac{1}{n} \sum_i^n X_i \quad (13)$$

The expected value of  $\bar{X}_n$  is computed as

$$E[\bar{X}_n] = \frac{1}{n} \sum_i^n E[X_i] = \mu \quad (14)$$

where  $\mu$  is the population mean. Here, we assume that all of the  $X_i$ s are identical and independently distributed (iid). The variance of  $\bar{X}_n$  is computed as

$$\begin{aligned}
E[\bar{X}_n^2] &= E \left[ \left( \frac{1}{n} \sum_{i=1}^n X_i \right)^2 \right] \\
&= \frac{1}{n^2} E \left[ \sum_{i=1}^n X_i^2 + 2 \sum_{i=1}^n \sum_{j>i}^n X_i X_j \right] \\
&= \frac{1}{n} E[X_i^2] + \frac{2}{n} \frac{n(n-1)}{2} \mu^2 \\
&= \frac{1}{n} (\sigma^2 + \mu^2) + \frac{n-1}{n} \mu^2
\end{aligned} \tag{15}$$

where we have used  $E[X_i X_j] = E[X_i] E[X_j]$  when  $j \neq i$ . Then,  $\text{Var}[\bar{X}_n] = E[\bar{X}_n^2] - E[\bar{X}_n]^2$ , so

$$\text{Var}[\bar{X}_n] = \frac{\sigma^2}{n} \tag{16}$$

Thus the central limit theorem on the variate  $(\bar{X}_n - \mu)/(\sigma/\sqrt{n})$  gives us:

$$\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \sim N(0, 1) \tag{17}$$

**Poisson distribution:** The Poisson distribution is a limit of binomial distribution when the probability of success  $P = \mu/n$  is low but the number of trials  $n$  is very large. To arrive at Poisson distribution, we rewrite the probability density function of the binomial distribution as follows:

$$\begin{aligned}
P(x) &= \frac{n!}{x! (n-x)!} \left( \frac{\mu}{n} \right)^x \left( 1 - \frac{\mu}{n} \right)^{n-x} \\
&= \left[ \frac{n}{x} \cdot \frac{n-1}{x} \cdots \frac{n-x+1}{x} \right] \frac{\mu^x}{x!} \left( 1 - \frac{\mu}{n} \right)^{n-x}
\end{aligned} \tag{18}$$

The term in the brackets tend to 1 as  $n \rightarrow \infty$ , while  $(1 - \mu/n)^n \rightarrow e^{-\mu}$  and  $(1 - \mu/n)^x \rightarrow 1$ . Then, the above term reduces to

$$P(x) = \frac{e^{-\mu} \mu^x}{x!} \tag{19}$$

It is straightforward to show that the mgf of the Poisson distribution is given by  $e^{-\mu} e^{\mu t}$ .

**$\chi^2$  distribution:** The  $\chi^2$  distribution with  $f$  degrees of freedom (dom) is defined as a sum of  $f$  iid normal variates  $Z_i^2$ :

$$\Upsilon = Z_1^2 + Z_2^2 + \cdots + Z_f^2 \tag{20}$$

For 1 dof,  $E[\Upsilon] = E[Z^2] = 1$  and  $E[\Upsilon^2] = E[Z^4] = 3$ . The  $\chi_{[f]}^2$  distribution then satisfies:

$$E[\Upsilon] = f, \quad \text{Var}[\Upsilon] = 2f \quad (21)$$

Let's compute the mgf of  $\chi_{[f]}^2$ . First, for the  $\chi_{[1]}^2$  variate

$$M(t) = E[e^{tZ^2}] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} dz e^{-\frac{1}{2}z^2(1-2t)} = (1-2t)^{-1/2} \quad (22)$$

For  $\chi_{[f]}^2$ , the mgf will be the product of  $f$  mgfs which is

$$M_{\Upsilon}(t) = (1-2t)^{-\frac{1}{2}f} \quad (23)$$

Now, we can show that this mgf can be obtained from the density function  $f(y) = \frac{1}{A(f)} y^{\frac{1}{2}f-1} e^{-\frac{1}{2}y}$  where  $0 \leq y < \infty$  and  $A(f) = 2^{\frac{1}{2}f} \Gamma(\frac{1}{2}f)$ .

$$M_f(t) = E[e^{tY}] = \int_0^{\infty} e^{ty} f(y) dy = \frac{1}{A(f)} \int_0^{\infty} e^{ty-\frac{1}{2}y} y^{\frac{1}{2}f-1} dy \quad (24)$$

Using the substitution  $w = (1-2t)y$  we get

$$M_f(t) = (1-2t)^{-\frac{1}{2}f} \frac{1}{A(f)} \int_0^{\infty} dw w^{\frac{1}{2}f-1} e^{-\frac{1}{2}w} = (1-2t)^{-\frac{1}{2}f} \quad (25)$$

Using the properties of the  $\Gamma$ -function, one can show that (using partial integration)

$$A(f) = \begin{cases} 1 \cdot 3 \cdot 5 \cdots (f-2) \sqrt{2\pi}, & \text{odd } f \\ 2 \cdot 4 \cdot 6 \cdots (f-2) \cdot 2, & \text{even } f \end{cases} \quad (26)$$

Let us consider an important application of the  $\chi^2$  distribution. Let  $Z_1, Z_2, \dots, Z_n$  be iid normal variates and  $\Upsilon_1, \Upsilon_2, \dots, \Upsilon_n$  be linear functions of them:

$$\begin{aligned} \Upsilon_1 &= a_1 Z_1 + a_2 Z_2 + \cdots + a_n Z_n \\ \Upsilon_2 &= b_1 Z_1 + b_2 Z_2 + \cdots + b_n Z_n \\ &\dots\dots\dots \\ \Upsilon_n &= u_1 Z_1 + u_2 Z_2 + \cdots + u_n Z_n \end{aligned} \quad (27)$$

with  $a_i, b_i, \dots, u_i$  being a set of **orthonormal vectors**. Then, it is clear that

$$\sum_i \Upsilon_i^2 = \sum_i Z_i^2 \quad (28)$$

The orthonormality also implies that all the  $\Upsilon_i$ 's are independently distributed; for example

$$\begin{aligned}\text{Cov}[\Upsilon_1, \Upsilon_2] &= E[(a_1 Z_1 + \cdots + a_n Z_n) \cdot (b_1 Z_1 + \cdots + b_n Z_n)] \\ &= \sum_{i,j} (a_i b_j) E[Z_i Z_j] \\ &= \sum_i a_i b_i E[Z_i^2] = \sum_i a_i b_i = 0\end{aligned}\tag{29}$$

since for  $i \neq j$ ,  $E[Z_i Z_j] = E[Z_i] E[Z_j] = 0$ . Using the identity, it is possible to work out the sampling distribution of the variance.

**Sampling distribution of variance:** Consider the sum of squared deviations from the mean

$$S^2 = \sum_i (x_i - \bar{x})^2\tag{30}$$

where  $x_i$ 's are sample distributions and  $\bar{x}$  is the mean of the sample distribution, from a population with mean  $\mu$  and variance  $\sigma^2$ . Then, the following identity holds:

$$\begin{aligned}\sum_i (x_i - \bar{x})^2 &= \sum_i ((x_i - \mu)^2 - (\bar{x} - \mu)^2) \\ &= \sum_i (x_i - \mu)^2 - n(\bar{x} - \mu)^2\end{aligned}\tag{31}$$

Then, the expectation value of  $S^2$  is given by

$$\begin{aligned}E[S^2] &= \sum_i E[(x_i - \mu)^2] - n E[(\bar{x} - \mu)^2] \\ &= \sum_i \text{Var}[x_i] - n \text{Var}[\bar{x}] \\ &= n \sigma^2 - n \frac{\sigma^2}{n} = (n-1) \sigma^2\end{aligned}\tag{32}$$

Thus, the **unbiased** estimator for the variance is given by

$$s^2 = \frac{1}{n-1} S^2\tag{33}$$

Now, we have from the above equations

$$\frac{\sum_i (x_i - \mu)^2}{\sigma^2} = \frac{S^2}{\sigma^2} + \frac{n(\bar{x} - \mu)^2}{\sigma^2}\tag{34}$$

Assuming that  $x_i$  being **normally** distributed, notice that  $\left(\frac{\bar{x} - \mu}{\sigma/n}\right)^2$  becomes a  $\chi^2_{[1]}$  variate. At the same time,  $\sum_i \left(\frac{x_i - \mu}{\sigma}\right)^2$  becomes a  $\chi^2_{[n]}$  variate. Thus, the above equation implies that  $S^2/\sigma^2$  to be a  $\chi^2_{[n-1]}$  variate, as long as  $S^2$  and  $\bar{x}$  are independently distributed.

We establish this independence by constructing an orthonormal transformation from a set of normal variates  $Z_i$  to their linear combinations  $\Upsilon_i$ :

$$\begin{aligned} Z_i &= \frac{x_i - \mu}{\sigma} \quad , \quad i = 1, 2, \dots, n \\ \Upsilon_1 &= \frac{1}{\sqrt{n}} Z_1 + \dots + \frac{1}{\sqrt{n}} Z_n \end{aligned} \quad (35)$$

Notice that

$$\begin{aligned} \Upsilon_1 &= \frac{1}{\sqrt{n}} \sum_i \left( \frac{x_i - \mu}{\sigma} \right) = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \\ \rightarrow \Upsilon_1^2 &= \frac{n(\bar{x} - \mu)^2}{\sigma^2} \sim \chi_{[1]}^2 \\ \sum_i Z_i^2 &= \frac{\sum_i (x_i - \mu)^2}{\sigma^2} \end{aligned} \quad (36)$$

Now, complete the transformation to from  $Z_i$  to  $\Upsilon_i$  by adding  $\Upsilon_2, \dots, \Upsilon_n$  on  $\Upsilon_1$ , such that  $\text{Cov}[\Upsilon_i, \Upsilon_j] = \delta_{ij}$ . Otrhonormality requires

$$\begin{aligned} \sum_{i=1}^n Z_i^2 &= \Upsilon_1^2 + \sum_{i=2}^n \Upsilon_i^2 \\ \rightarrow \sum_{i=2}^n \Upsilon_i^2 &= \frac{\sum_i (x_i - \mu)^2}{\sigma^2} - \frac{n(\bar{x} - \mu)^2}{\sigma^2} = \frac{S^2}{\sigma^2} \chi_{[n-1]}^2 \end{aligned} \quad (37)$$

where each  $\Upsilon_i$  is iid and  $\chi_{[1]}^2$  by orthonormality. Thus,  $S^2$  is  $\chi_{[n-1]}^2$  variate.

**t-Distribution:** The Student t-distribution is defined through the random variable T as follows:

$$T = \frac{Z}{\sqrt{\Upsilon/f}} \quad (38)$$

where  $\Upsilon$  is a  $\chi^2$  variate with  $f$  dof and  $Z$  is a normal variate. This distribution is useful when the population variance  $\sigma^2$  is unknown, but estimated by  $s^2$ . Since  $\frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$  is a normal variate (assuming  $x_i$ 's are iid and normally distributed) and  $S^2/\sigma^2$  is a  $\chi_{[n-1]}^2$ , then

$$\frac{(\bar{x} - \mu)}{\sigma/\sqrt{n}} \bigg/ \sqrt{\frac{S^2}{(n-1)\sigma^2}} = \frac{\bar{x} - \mu}{s/\sqrt{n}} \quad (39)$$

will follow the t-distribution with  $n - 1$  dof.

Let us know calculate the probability density function of the t-distribution. First consider the distribution  $V = \sqrt{\Upsilon/f}$  so that  $T = Z/V$ . The cumulative density function for  $V$  is given by

$$F_V(v) = P[V \leq v] = P[\sqrt{\Upsilon/f} \leq v] = P[\Upsilon \leq f v^2] = F_\Upsilon(f v^2) \quad (40)$$

Then, the probability density function is given by

$$\begin{aligned} f_V(v) &= \frac{dF_V(v)}{dv} = \frac{dF_{\Upsilon}(f v^2)}{dv} = 2f v f_{\Upsilon}(f v^2) \\ &= \frac{1}{A(f)} 2f v (f v^2)^{\frac{1}{2}f-1} e^{-\frac{1}{2}f v^2} \end{aligned} \quad (41)$$

Now, for the  $T$  distribution, we have

$$\begin{aligned} F_T(t) &= P[Z/V \leq t] = \sum_v (P[Z \leq vt] \cup P[v \leq V \leq v + dv]) \\ &= \int F_Z(vt) f_V(v) dv \\ f_T(t) &= \int f_Z(vt) v f_v(v) dv \\ &= \frac{2 f^{f/2}}{\sqrt{2\pi} A(f)} \int_0^\infty e^{-\frac{1}{2}v^2 t^2 - \frac{1}{2}f v^2} v^f dv \end{aligned} \quad (42)$$

The integral can be evaluated by the substitution  $\xi = v^2 (f + t^2)$ , and after some algebra

$$f_T(t) = \frac{A(f+1)}{\sqrt{2\pi} f A(f)} \left(1 + \frac{t^2}{f}\right)^{-\frac{1}{2}(f+1)} \quad (43)$$

t-distribution can be applied to test whether the means of two distributions are the same. Suppose that we have  $m$  observations on a random variable  $X$  ( $x_1, \dots, x_m$ ) and  $n$  observations on another random variable  $Y$  ( $y_1, \dots, y_n$ ). Assuming that  $X$  and  $Y$  are normally distributed with the same variance  $\sigma^2$ , but different means  $\mu_1$  and  $\mu_2$ , we test the hypotheses  $\mu_1 = \mu_2$ . We define

$$\begin{aligned} S_1^2 &= \sum_{i=1}^m (x_i - \bar{x})^2 & S_2^2 &= \sum_{i=1}^n (y_i - \bar{y})^2 \\ S^2 &= S_1^2 + S_2^2 & s^2 &= \frac{S^2}{m+n-2} \end{aligned} \quad (44)$$

Then,  $\bar{x} - \bar{y}$  will be normally distributed with mean  $\mu_1 - \mu_2$  and variance  $(1/n + 1/m) \sigma^2$ , and  $S^2$  will follow a  $\chi^2$  distribution with  $m+n-2$  dof. Then,

$$t = \frac{\bar{x} - \bar{y}}{s \sqrt{\frac{1}{n} + \frac{1}{m}}} \quad (45)$$

will follow a t-distribution with  $m+n-2$  dof.

## Linear Regression

We first consider the case of two variables, multivariable extension is straightforward. The observations  $Y_i$  are modeled by the following linear function:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad (46)$$

where  $\epsilon_i \sim N(0, \sigma^2)$  are assumed to be iid normal variates. The precision is denoted by hatted symbols, i.e.

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i \quad (47)$$

where  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are obtained through least-squares, by minimizing the sum of squared errors:

$$S^2 = \sum_i (Y_i - \beta_0 - \beta_1 X_i)^2 \quad (48)$$

Minimizing with respect to  $\beta_0$  and  $\beta_1$ , we obtain

$$\hat{\beta}_1 = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} = \beta_1 + \frac{\sum_i \epsilon_i (X_i - \bar{X})}{\sum_i (X_i - \bar{X})^2} \quad , \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \quad (49)$$

Equivalently, it is easy to show that

$$\begin{aligned} \hat{\beta}_1 &= \text{Cor}(Y, X) \frac{S_Y}{S_X} \quad , \quad \text{Cor}(Y, X) = \frac{\text{Cov}(Y, X)}{S_X S_Y} \\ \text{Cov}(X, Y) &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \\ S_X^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \end{aligned} \quad (50)$$

**Prediction and confidence intervals:** For simplifying the algebra, let us subtract  $\bar{X}$  from  $X_i$  to normalize:

$$Y_i = \beta_0 + \beta_1 (X_i - \bar{X}) + \epsilon_i \quad (51)$$

this is equivalent to redefining  $\beta_0 \rightarrow \beta_0 - \beta_1 \bar{X}$ , which results in  $\hat{\beta}_0 = \bar{Y}$ . Now it is easy to see the following apply:

$$\begin{aligned} \bar{Y} &= \beta_0 + \frac{1}{n} \sum_i \epsilon_i \quad , \quad \bar{Y} = \hat{\beta}_0 \sim N(\beta_0, \frac{\sigma^2}{n}) \\ E[Y_i] &= \beta_0 + \beta_1 (X_i - \bar{X}) \quad , \quad \text{Var}[Y_i] = \text{Var}[\epsilon_i] = \sigma^2 \quad , \quad Y_i \sim N(\beta_0 + \beta_1 (X_i - \bar{X}), \sigma^2) \end{aligned} \quad (52)$$



now,

$$\begin{aligned}\text{Var}[\hat{\beta}_1] &= \frac{1}{[\sum_i (X_i - \bar{X})^2]^2} \sum_i (X_i - \bar{X})^2 \text{Var}[\epsilon_i] \\ &= \frac{\sigma^2}{\sum_i (X_i - \bar{X})^2}\end{aligned}\tag{53}$$

thus, in summary:

$$\boxed{\hat{\beta}_0 \sim N\left(\beta_0, \frac{\sigma^2}{n}\right) \quad , \quad \hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum_i (X_i - \bar{X})^2}\right)}\tag{54}$$

If we were to add the  $\bar{X}$  term back,  $\beta_0 \rightarrow \beta_0 + \beta_1 \bar{X}$  so  $\text{Var}[\hat{\beta}_0] \rightarrow \bar{X}^2 \text{Var}[\hat{\beta}_1] + \text{Var}[\hat{\beta}_0]$ , which is more widely used. Now, consider a **new prediction**  $\hat{Y}_{n+1}$  from a **new predictor**  $X_{n+1}$ , i.e.

$$Y_{n+1} = \beta_0 + \beta_1 (X_{n+1} - \bar{X}) + \epsilon_{n+1} \quad , \quad \hat{Y}_{n+1} = \hat{\beta}_0 + \hat{\beta}_1 (X_{n+1} - \bar{X})\tag{55}$$

consider the quantity  $W = Y_{n+1} - \hat{Y}_{n+1}$ . It is easy to see that  $E[W] = 0$ . Now, we compute the variance of  $W$ , which will be related to the **prediction interval**. To proceed, let us first note that  $\text{Cov}[Y_{n+1}, \hat{\beta}_0 + \hat{\beta}_1 (X_{n+1} - \bar{X})] = 0$ , since  $Y_{n+1}$  is a new data uncorrelated to  $Y_i$  from which  $\hat{\beta}_{0,1}$  are determined. In addition,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are uncorrelated as well:

$$\text{Cov}[\hat{\beta}_0, \hat{\beta}_1] = \text{Cov}\left[\hat{\beta}_0, \beta_1 + \frac{\sum_i \epsilon_i (X_i - \bar{X})}{\sum_i (X_i - \bar{X})^2}\right]\tag{56}$$

since  $E[\hat{\beta}_0] = E[\bar{Y}] = \beta_0$  and  $E[\hat{\beta}_1] = \beta_1$ , we get

$$\begin{aligned}\text{Cov}[\hat{\beta}_0, \hat{\beta}_1] &= E\left[\frac{1}{n} \sum_i \epsilon_i \times \frac{\sum_i \epsilon_i (X_i - \bar{X})}{\sum_i (X_i - \bar{X})^2}\right] \\ &= \frac{1}{n \sum_i (X_i - \bar{X})^2} \sum_{i,j} (X_i - \bar{X}) E[\epsilon_i \epsilon_j] \\ &= \frac{1}{n \sum_i (X_i - \bar{X})^2} \sum_i (X_i - \bar{X}) \sigma^2 = 0\end{aligned}\tag{57}$$

now we have shown that  $Y_{n+1}, \hat{\beta}_0, \hat{\beta}_1$  are uncorrelated, it is straightforward to compute  $\text{Var}[W]$ :

$$\begin{aligned}\text{Var}[W] &= \text{Var}[Y_{n+1}] + \text{Var}[\hat{\beta}_0] + (X_{n+1} - \bar{X})^2 \text{Var}[\hat{\beta}_1] \\ &= \sigma^2 + \frac{\sigma^2}{n} + (X_{n+1} - \bar{X})^2 \frac{\sigma^2}{\sum_i (X_i - \bar{X})^2}\end{aligned}\tag{58}$$

$$\boxed{\text{Var}[W] = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum_i (X_i - \bar{X})^2} \right]} \quad (59)$$

the quantity  $\text{Var}[W]$  is a measure of the error in prediction, so it determines the **prediction interval** via:

$$\hat{Y}_{n+1} \pm t_{\alpha/2, n-2} \hat{\sigma} \left[ 1 + \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum_i (X_i - \bar{X})^2} \right]^{1/2} \quad (60)$$

where  $t_{\alpha/2, n-2}$  is the  $\alpha$ th  $t$ -quantile with  $n - 2$  dof, and  $\hat{\sigma}$  is the unbiased estimate of  $\sigma$ :

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (61)$$

instead, the intrinsic error we make in  $\hat{Y}_{n+1}$  is simply  $\text{Var}[\hat{Y}_{n+1}]$ , which determines the **confidence interval**. the confidence interval is always smaller than the prediction interval, and is characterized by:

$$\boxed{\text{Var}[\hat{Y}_{n+1}] = \sigma^2 \left[ \frac{1}{n} + \frac{(X_{n+1} - \bar{X})^2}{\sum_i (X_i - \bar{X})^2} \right]} \quad (62)$$

**Some properties of the residuals:** The residuals are defined by  $e_i = Y_i - \hat{Y}_i$ , and are used in the definition of the unbiased estimate of  $\text{Var}[\epsilon_i] = \sigma^2$ . The residuals satisfy the following properties:

$$\boxed{\sum_i e_i = 0 \quad , \quad \sum_i e_i X_i = 0} \quad (63)$$

Here are the proofs:

$$\begin{aligned} \sum_i e_i &= \sum_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) \\ &= \sum_i (\beta_0 + \beta_1 X_i + \epsilon_i - \bar{Y} + \hat{\beta}_1 X_i - \hat{\beta}_1 X_i) \\ &= \sum_i \left( \beta_0 + \beta_1 X_i + \epsilon_i - \beta_0 - \beta_1 \bar{X} - \frac{1}{n} \sum_j \epsilon_j + \hat{\beta}_1 \bar{X} - \hat{\beta}_1 X_i \right) \\ &= (\beta_1 - \hat{\beta}_1) \sum_i (X_i - \bar{X}) = 0 \\ \sum_i e_i X_i &= \sum_i e_i (X_i - \bar{X}) = \sum_i (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i) (X_i - \bar{X}) \\ &= \sum_i \epsilon_i (X_i - \bar{X}) + (\beta_0 - \hat{\beta}_0) \sum_i (X_i - \bar{X}) + (\beta_1 - \hat{\beta}_1) \sum_i X_i (X_i - \bar{X}) \\ &= \sum_i \epsilon_i (X_i - \bar{X}) + (\beta_1 - \hat{\beta}_1) \sum_i (X_i - \bar{X})^2 \\ &= \sum_i \epsilon_i (X_i - \bar{X}) - \sum_i \epsilon_i (X_i - \bar{X}) = 0 \end{aligned} \quad (64)$$

where we used the least squares fit for  $\hat{\beta}_1$  in the last line.

Now, let's compute the **variance of the residuals** (using the normalized  $X_i \rightarrow X_i - \bar{X}$ ):

$$\begin{aligned}\text{Var}[e_i] &= \text{Var}[Y_i - \hat{\beta}_0 - \hat{\beta}_1 (X_i - \bar{X})] \\ &= \text{Var}[Y_i] + \text{Var}[\hat{\beta}_0] + (X_i - \bar{X})^2 \text{Var}[\hat{\beta}_1] - 2 \text{Cov} \left[ Y_i, \hat{\beta}_0 + \hat{\beta}_1 (X_i - \bar{X}) \right]\end{aligned}\quad (65)$$

where we have used the fact that  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are uncorrelated and  $\text{Var}[A \pm B] = \text{Var}[A] + \text{Var}[B] \pm 2 \text{Cov}[A, B]$ . The covariance term was not present when computing  $\text{Var}[W]$  since there,  $Y_{n+1}$  was new data and not correlated to  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . Now, notice that

$$\begin{aligned}Y_i - E[Y_i] &= \epsilon_i \\ \hat{\beta}_0 + \hat{\beta}_1 (X_i - \bar{X}) - E \left[ \hat{\beta}_0 + \hat{\beta}_1 (X_i - \bar{X}) \right] &= (\hat{\beta}_0 - \beta_0) + (\hat{\beta}_1 - \beta_1) (X_i - \bar{X})\end{aligned}\quad (66)$$

Thus, the covariance term becomes

$$\begin{aligned}\text{Cov} \left[ Y_i, \hat{\beta}_0 + \hat{\beta}_1 (X_i - \bar{X}) \right] &= E \left[ \epsilon_i \times \left( (\hat{\beta}_0 - \beta_0) + (\hat{\beta}_1 - \beta_1) (X_i - \bar{X}) \right) \right] \\ &= E \left[ (\hat{\beta}_0 - \beta_0) \epsilon_i \right] + (X_i - \bar{X}) E \left[ (\hat{\beta}_1 - \beta_1) \epsilon_i \right]\end{aligned}\quad (67)$$

Now,  $E[(\hat{\beta}_0 - \beta_0) \epsilon_i] = E[(\bar{Y} - \beta_0) \epsilon_i] = E[\frac{1}{n} \sum_j \epsilon_j \epsilon_i] = \sigma^2/n$ . Similarly,

$$E \left[ (\hat{\beta}_1 - \beta_1) \epsilon_i \right] = \frac{1}{\sum_j (X_j - \bar{X})^2} E \left[ \sum_j (X_j - \bar{X}) \epsilon_j \epsilon_i \right] = \frac{(X_i - \bar{X}) \sigma^2}{\sum_j (X_j - \bar{X})^2}\quad (68)$$

Thus,

$$\text{Cov} \left[ Y_i, \hat{\beta}_0 + \hat{\beta}_1 (X_i - \bar{X}) \right] = \sigma^2 \left[ 1 + \frac{(X_i - \bar{X})^2}{\sum_j (X_j - \bar{X})^2} \right]\quad (69)$$

The calculation of  $\text{Var}[Y_i] + \text{Var}[\hat{\beta}_0] + (X_i - \bar{X})^2 \text{Var}[\hat{\beta}_1]$  follows similarly to  $\text{Var}[W]$ , so finally combining with the covariance term we get

$$\boxed{\text{Var}[e_i] = \sigma^2 \left[ 1 - \frac{1}{n} - \frac{(X_i - \bar{X})^2}{\sum_j (X_j - \bar{X})^2} \right]}\quad (70)$$

**Multiple features:** Consider the case of two features (can easily be generalized to more). Assume that  $\beta_0 = 0$  (i.e.  $\rightarrow Y - \beta_0$  is being fitted) then, the least-squares minimization is performed on

$$S^2 = \sum_{i=1}^n (Y_i - \beta_1 X_{1i} - \beta_2 X_{2i})^2\quad (71)$$

Evaluating partial derivatives with respect to  $\beta_1$ ,  $\beta_2$  and setting them to zero, we obtain

$$\begin{bmatrix} \sum_i X_{1i}^2 & \sum_i X_{1i} X_{2i} \\ \sum_i X_{1i} X_{2i} & \sum_i X_{2i}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} \sum_i Y_i X_{1i} \\ \sum_i Y_i X_{2i} \end{bmatrix} \quad (72)$$

This equation is easily solved for  $\hat{\beta}_1$  and  $\hat{\beta}_2$ . The solution is given by

$$\begin{aligned} \Delta &= \left( \sum_i X_{1i}^2 \right) \left( \sum_i X_{2i}^2 \right) - \left( \sum_i X_{1i} X_{2i} \right)^2 \\ \hat{\beta}_1 &= \left[ \left( \sum_i X_{2i}^2 \right) \left( \sum_i Y_i X_{1i} \right) - \left( \sum_i X_{1i} X_{2i} \right) \left( \sum_i Y_i X_{2i} \right) \right] / \Delta \\ \hat{\beta}_2 &= \left[ \left( \sum_i X_{1i}^2 \right) \left( \sum_i Y_i X_{2i} \right) - \left( \sum_i X_{1i} X_{2i} \right) \left( \sum_i Y_i X_{1i} \right) \right] / \Delta \end{aligned} \quad (73)$$

However, we would like to express the solutions in terms of residuals. We define:

$$e_{i, X_1|X_2} = X_{1i} - \left( \frac{\sum_j X_{2j} X_{1j}}{\sum_j X_{2j}^2} \right) X_{2i} \quad (74)$$

which is **the residual having fit  $X_2$  on  $X_1$** . The term in the paranthesis is the regression coefficient (via least squares) if we were to fit  $X_{1i} = \beta X_{2i}$  and the residual is  $X_{1i} - \hat{\beta} X_{2i}$ . The residuals  $e_{i, Y|X_2}$ ,  $e_{i, Y|X_1}$  are similarly defined.

After some algebra, we get

$$\begin{aligned} \hat{\beta}_1 &= \frac{\sum_i e_{i, Y|X_2} e_{i, X_1|X_2}}{\sum_i e_{i, X_1|X_2}^2} \\ \hat{\beta}_2 &= \frac{\sum_i e_{i, Y|X_1} e_{i, X_2|X_1}}{\sum_i e_{i, X_2|X_1}^2} \end{aligned} \quad (75)$$

In other words, the regression estimate for  $\beta_1$  is the regression through the origin estimate having regressed  $X_2$  out of both the response and the predictors (similar for  $\beta_2$ ). Moreover, if we were to fit the regression through the origin between the residuals  $e_{i, Y|X_2}$  and  $e_{i, X_1|X_2}$ :

$$e_{i, Y|X_2} = \gamma e_{i, X_1|X_2} \quad (76)$$

we would get

$$\hat{\gamma} = \frac{\sum_i e_{i, Y|X_2} e_{i, X_1|X_2}}{\sum_i e_{i, X_1|X_2}^2} \quad (77)$$

which is equivalent to  $\hat{\beta}_1$ . So the residuals having regressed out  $X_2$  contains the information on the  $X_1$  dependence. This is the reason why we see left over variability in residual plots if there is a feature which is not fitted.

### Bias-Variance trade-off

Suppose we know the true model  $f(X)$  that characterizes the response  $Y$ , i.e.

$$Y = f(X) + \epsilon \quad (78)$$

Here,  $\epsilon$  represents the irreducible error, i.e. the error that is still there even if we knew the exact model  $f(X)$ . Let's assume that we have a estimate of the true model by some function  $\hat{f}(X)$ . Then, we can compute the total error we make (residual sums squared) as

$$\begin{aligned} E \left[ (Y - \hat{f}(X))^2 \right] &= E \left[ \epsilon + (f(X) - \hat{f}(X))^2 \right] \\ &= \text{Var}(\epsilon) + E \left[ (f(X) - \hat{f}(X))^2 \right] \end{aligned} \quad (79)$$

where we have used the fact that  $E[\epsilon] = 0$  so  $\text{Var}(\epsilon) = E[\epsilon^2]$ , and that  $\epsilon$  and  $f(X) - \hat{f}(X)$  are uncorrelated. The first term above is the irreducible error. Let us look into the second term:

$$\begin{aligned} E \left[ (f(X) - \hat{f}(X))^2 \right] &= E[f(X)^2] - E[\hat{f}(X)]^2 + E[\hat{f}(X)]^2 - 2 f(X) E[\hat{f}(X)] + f(X)^2 \\ &= \text{Var}(\hat{f}(X)) + \left( E[\hat{f}(X)] - f(X) \right)^2 \\ &= \text{Var}(\hat{f}(X)) + \text{Bias}(\hat{f}(X))^2 \end{aligned} \quad (80)$$

where the bias is defined as

$$\text{Bias}(\hat{f}(X)) = E[\hat{f}(X)] - f(X) \quad (81)$$

The first term above is the variance of  $\hat{f}(X)$  and the second one is bias, which measures how much  $\hat{f}(X)$  deviates from the true  $f(X)$ . By including more flexibility in model  $\hat{f}(X)$  we can reduce the bias, but the variance increases (the bias-variance trade-off).

### Maximum Likelihoods:

The determination of linear regression parameters through minimization of residual sums squared can be thought of a maximizing the likelihood function of  $\mathbf{Y}$ . Given  $Y_i$ 's are random variables with identical variances  $\sigma^2$ , their conditional probability distribution (given data  $X_i$  and parameters  $\beta = (\beta_0, \beta_1)$ ) can be written as

$$P(\mathbf{Y}|\mathbf{X}; \beta) = \prod_{i=1}^n P(Y_i|X_i; \beta) \quad (82)$$

The optimal parameters  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$  can be obtained from maximizing this joint probability (i.e. likelihood), or equivalently its logarithm. The probability distribution function for each  $Y_i$  is Gaussian centered around  $\hat{Y}_i$  (see above) with variance  $\sigma^2$ , so the maximum likelihood problem reduces to

$$\begin{aligned}\hat{\beta} &= \operatorname{argmax}_{\beta} \left[ \sum_{i=1}^n \left( -\frac{1}{2} \log(2\pi) + \log \sigma - \frac{|\hat{Y}_i - Y_i|^2}{2\sigma^2} \right) \right] \\ &= \operatorname{argmin}_{\beta} \left[ \sum_{i=1}^n |\hat{Y}_i - Y_i|^2 \right]\end{aligned}\tag{83}$$