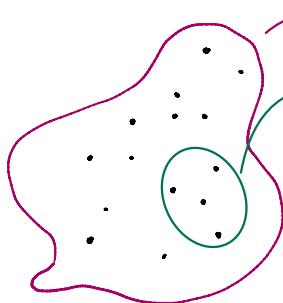


# Mathematical Statistics

## Notion of observations and populations



population (totality of samples) ... all students, cards in a deck  
sample (observation) ... several students, one card

Each observation is a value of a random variable  $X$  having some probability distribution  $f(x)$

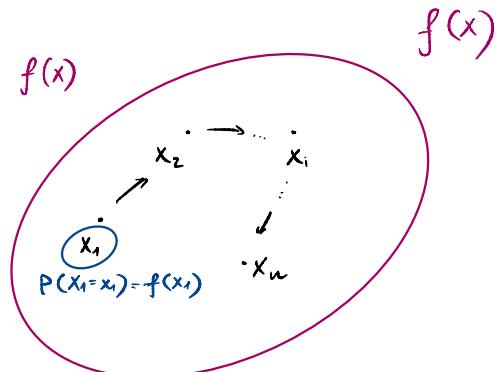
A **sample** is a subset of population

Since manually chosen samples might be **biased**, it is worth using **Random sampling**

In selecting random sample from population  $f(x)$  let us define **random variables**  $X_1, X_2, \dots, X_i, \dots, X_n$  that constitute this sample with values  $x_1, x_2, \dots, x_i, \dots, x_n$ .

Because of identical conditions, variables  $X_1, X_2, \dots, X_n$  are **independent** and **identically distributed**.

$$\text{So, } f(x_1, x_2, \dots, x_n) = f(x_1) f(x_2) \dots f(x_n)$$



Our main purpose in selecting random samples is to elicit information about **unknown population parameters**

Any function of the random variables constituting a random sample is called **statistic**.

Some important statistics:

- Sample Mean:  $\frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$
- Sample Median:  $\tilde{x} = \begin{cases} x_{(n+1)/2}, & n \text{ is odd} \\ \frac{1}{2}(x_m + x_{m+1}), & n \text{ is even} \end{cases}$
- Sample Mode: value that occurs more often  $\underset{n}{\text{occurs}}$
- Sample Variance:  $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = S^2$

If  $S^2$  is sample variance of a sample of size  $n$ , then we may write:

$$S^2 = \frac{1}{n(n-1)} \left( n \sum_{i=1}^n x_i^2 - \left( \sum_{i=1}^n x_i \right)^2 \right)$$

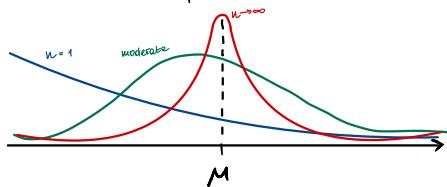
## Inference about populations

- The probability distribution of a statistic is called a **sampling distribution**. The sampling distribution depends on the distribution of the population, the size of samples, and the method of choosing samples.

## Sampling Distributions of Means and CLT

**Central Limit Theorem.** If  $\bar{X}$  is a mean of a sample taken from a population with mean  $\mu$  and finite variance  $\sigma^2$ , then the limiting form of the distribution of:

$$Z = \frac{1}{\sqrt{n}} \frac{\bar{X} - \mu}{\sigma}$$

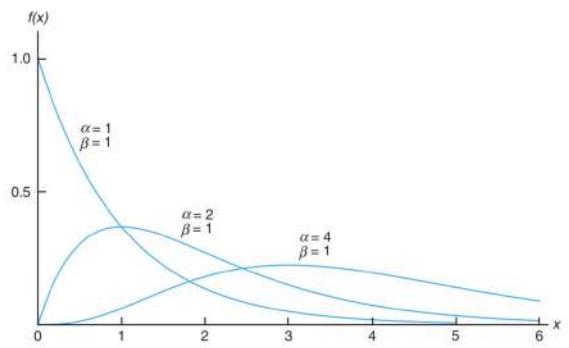


as  $n \rightarrow \infty$  is the standard normal distribution  $n(0,1)$

## Exercises

## Gamma Distribution

- Gamma Function:  $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$
- Properties:  $\Gamma(\alpha) = (\alpha-1)\Gamma(\alpha-1)$   
 $\Gamma(1) = 1$   
 $\Gamma(0.5) = \sqrt{\pi}$
- Gamma Distribution:  $f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-\frac{x}{\beta}} \cdot I_{x>0}$
- Exponential Distribution:  $f(x; \lambda) = \lambda e^{-\lambda x} \cdot I_{x>0}$



## Chi-squared Distribution

- A variable  $X$  is said to have chi-squared distribution with  $\gamma$  degrees of freedom if its p.d.f is defined as:
- The mean of such distribution is  $\gamma$
- The variance of such distribution is  $2\gamma$

$$f(x; \gamma) = \frac{1}{2^{\gamma/2} \Gamma(\gamma/2)} x^{\frac{\gamma}{2}-1} e^{-\frac{x}{2}}$$

## Sampling Distribution of $S^2$

If  $S^2$  is the variance of a random sample of size  $n$  taken from a normal population having the variance  $\sigma^2$ , then the statistic

$$\chi^2 = \frac{(n-1)S^2}{\sigma^2} = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma^2}$$

has chi-squared distribution with  $\gamma = n-1$  degrees of freedom

## t-Distribution

Let  $Z$  be a standard normal random variable and  $V$  a chi-squared random variable with  $\gamma$  degrees of freedom. If  $Z$  and  $V$  are independent, then the distribution of the random variable  $T$ , where

$$T = \frac{Z}{\sqrt{V/\gamma}}$$

is given by the density function  $h(t) = \frac{\Gamma(\frac{\gamma+1}{2})}{\Gamma(\frac{\gamma}{2}) \sqrt{\pi\gamma}} \left(1 + \frac{t^2}{\gamma}\right)^{-\frac{1}{2}(\gamma+1)}$

Let  $X_1, X_2, \dots, X_n$  be independent normally distributed random variables with mean  $\mu$  and deviation  $\sigma$ . Let  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  and  $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  then the random variable

$$T = \sqrt{n} \frac{\bar{X} - \mu}{S}$$

has a **t-distribution** with  $\gamma = n-1$  degrees of freedom.

The  $t$ -distribution is used extensively in problems with inferences about population mean or in problems that involve comparative samples

### F-Distribution

Let  $U$  and  $V$  be two independent random variables having chi-squared distributions with  $v_1$  and  $v_2$  degrees of freedom. Then the distribution of the random variable

$$F = \frac{v_2}{v_1} \frac{U}{V}$$

is given by the density function

$$h(f) = \begin{cases} \frac{\Gamma(\frac{v_1+v_2}{2}) (\frac{v_1}{v_2})^{\frac{v_1}{2}}}{\Gamma(\frac{v_1}{2}) \Gamma(\frac{v_2}{2})} \frac{f^{\frac{v_1}{2}-1}}{\left[1 + \frac{v_1}{v_2} f\right]^{\frac{v_1+v_2}{2}}}, & f > 0 \\ 0, & f \leq 0 \end{cases}$$

This is known as the **F-distribution** with  $v_1$  and  $v_2$  degrees of freedom

Writing  $f_{\alpha}(v_1, v_2)$  for  $f_\alpha$  with  $v_1$  and  $v_2$  degrees of freedom, we obtain:

$$f_{1-\alpha}(v_1, v_2) = \frac{1}{f_\alpha(v_1, v_2)}$$

### Quantile

A **quantile** of a sample,  $q(f)$  is a value for which a specified fraction  $f$  of the data values is less than or equal to  $q(f)$

**Example.** If  $q_X(0.8) = 4$  then  $P(X \leq 4) = 0.8$

## Estimation of Population Parameters

A **point estimate** of some population parameter  $\theta$  is a single value  $\hat{\theta}$  of a statistic  $\hat{\theta}$ .

A statistic  $\hat{\theta}$  is said to be an **unbiased estimator** for a parameter  $\theta$  if

$$\mu_{\hat{\theta}} = E(\hat{\theta}) = \theta$$

### Variance of Estimators

If we consider all unbiased estimators of some population parameter  $\theta$ , the one with the smallest variance is called the **most efficient estimator** of  $\theta$ .

### Interval Estimation

There are many situations in which it is preferable to determine an interval within we would expect to find the value of the parameter. Such an interval is called **interval estimate**. Such estimate can be expressed as two numbers  $\theta_L$  and  $\theta_U$  such that  $\theta_L \leq \theta \leq \theta_U$

### Estimating the mean

Recalling **Central Limit Theorem**, we can obtain a confidence interval for a sample mean estimation.

$$P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$$

where

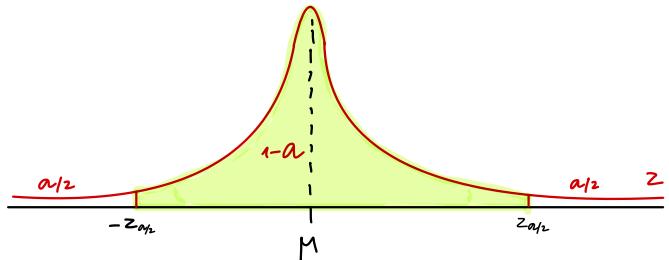
$$Z = \sqrt{n} \frac{\bar{X} - \mu}{\sigma}$$

$$\Phi_0(x) = \int_0^x \frac{e^{-t^2/2}}{\sqrt{2\pi}} dt = \frac{\alpha}{2}$$

$$x = \Phi_0^{-1}\left(\frac{\alpha}{2}\right)$$

We obtain:

$$P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$



If  $\bar{X}$  is the mean of a random sample of size  $n$  from a population with known variance  $\sigma^2$ , a  $100(1-\alpha)/$  confidence interval for  $\mu$  is given by:

$$\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

### Example Exercise I.

The average zinc concentration recovered from a sample of measurements taken in 36 different locations in a river is found to be 2.6 g per mm. Find the 95% and 99% confidence intervals for the mean zinc concentration. Assume that population std. deviation is 0.3 g. per mm.

a)  $2.6 - (1.96)\left(\frac{0.3}{\sqrt{36}}\right) < \mu < 2.6 + (1.96)\left(\frac{0.3}{\sqrt{36}}\right)$

$$2.5 < \mu < 2.7$$

b)  $2.6 - (2.575)\left(\frac{0.3}{\sqrt{36}}\right) < \mu < 2.6 + (2.575)\left(\frac{0.3}{\sqrt{36}}\right)$

$$2.47 < \mu < 2.73$$

Gorodetskiy's Notes

### Estimator Biasness

$$\text{Bias } \hat{\theta} = E\hat{\theta} - \theta$$

### Mean Square Error

$$E(\hat{\theta}(x) - \theta)^2 = \text{Var}\hat{\theta}(x) + \text{bias}^2\hat{\theta}(x)$$

### Consistency

Estimator  $\hat{\theta}$  is said to be consistent if  $\forall \epsilon > 0 \lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| \geq \epsilon) = 0$

$$P(|\hat{\theta} - \theta| \geq \epsilon) \stackrel{\text{E}\hat{\theta}(\text{unbiased})}{\leq} \frac{\text{Var}\hat{\theta}}{\epsilon^2} \longrightarrow 0$$

If  $\bar{X}$  is used as an estimate of  $\mu$ , we can be  $100(1-\alpha)\%$  confident that the error will not exceed a specified amount  $e$  when the sample size is

$$n = \left( \frac{z_{\alpha/2} \sigma}{e} \right)^2$$

### Example Exercise II:

How large a sample is required if we want to be  $95\%$  confident that our estimate of  $\mu$  in previous exercise is off by less than  $0.05$ .

$$n = \left[ \frac{1.96 \cdot 0.3}{0.05} \right]^2 = 138.3 \rightarrow 139$$

### One-sided intervals

$$P(\mu > \bar{X} - z_{\alpha} \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

$$P(\mu < \bar{X} + z_{\alpha} \frac{\sigma}{\sqrt{n}}) = 1 - \alpha$$

### The case of $\sigma$ unknown

If the particular value of  $\sigma$  is unknown, it can be replaced with statistic  $S$ . Then, recalling that statistic (considering normal sample)

$$T = \sqrt{n} \frac{\bar{X} - \mu}{S}$$

has a **t-distribution** with  $n-1$  degrees of freedom. Hence, we can use that fact in the same manner as we used CLT.

$$P(-t_{\alpha/2} < T < t_{\alpha/2}) = 1 - \alpha$$

$$P\left(-t_{\alpha/2} < \sqrt{n} \frac{\bar{X} - \mu}{S} < t_{\alpha/2}\right) = 1 - \alpha$$

$$P\left(\bar{X} - t_{\alpha/2} \frac{S}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{S}{\sqrt{n}}\right) = 1 - \alpha$$

### Large-Sample Confidence Interval

Often even when **normality** cannot be assumed, if  $n \geq 30$ ,  $s$  can replace  $\sigma$  and the confidence interval

$$\bar{X} = \pm z_{\alpha/2} \frac{s}{\sqrt{n}}$$

may be used

### Standard Error of Point Estimate

A standard error of a point estimate of  $\bar{X}$  is given as  $\frac{\sigma}{\sqrt{n}}$

### Prediction Intervals

For a **normal** distribution of measurements with unknown mean and known variance  $\sigma^2$ , a  $100(1-\alpha)\%$  prediction interval of a future observation  $x_0$  is:

$$\bar{X} - z_{\alpha/2} \sigma \sqrt{1 + \frac{1}{n}} < x_0 < \bar{X} + z_{\alpha/2} \sigma \sqrt{1 + \frac{1}{n}}$$

where  $z_{\alpha/2}$  is the **z-value** leaving an area of  $\alpha/2$  to the right.  
or

$$\bar{X} - t_{\alpha/2} s \sqrt{1 + \frac{1}{n}} < x_0 < \bar{X} + t_{\alpha/2} s \sqrt{1 + \frac{1}{n}}$$

where  $t_{\alpha/2}$  is the **t-value** leaving an area of  $\alpha/2$  to the right

## Tolerance Intervals

For a normal distribution of measurements with unknown mean  $\mu$  and unknown std. deviation  $\sigma$ , tolerance limits are given by  $\bar{x} \pm ks$ , where  $k$  is determined such that one can assert with  $100(1-\alpha)\%$  confidence that the given limits contain at least the proportion  $1-\alpha$  of the measurements. (Table values)

## Estimating the differences between two means

If we have two populations with means  $\mu_1$  and  $\mu_2$  and variances  $\sigma_1^2$  and  $\sigma_2^2$  respectively, a point estimator of the difference between  $\mu_1$  and  $\mu_2$  is given by the statistic  $\bar{x}_1 - \bar{x}_2$  that can be approximated using the normal distribution

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

## Estimating the Variance

An interval estimate of  $\sigma^2$  can be established by using the statistic:

$$\chi^2 = \frac{(n-1)s^2}{\sigma^2}$$

If samples are taken from normal population, then  $\chi^2 \sim \chi_{n-1}^2$ .  
So,

$$P\left(\frac{(n-1)s^2}{\chi_{\alpha/2}^2} < \sigma^2 < \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2}\right) = 1-\alpha$$

## Sufficient Statistic

The statistic  $T(\vec{x})$  is said to be sufficient for parameter  $\theta$  if  $\forall \theta \in \Theta \quad P(\vec{x} \in D | T(\vec{x}))$  is not dependent on  $\theta$ .

$T(\vec{x})$  is sufficient for  $\theta$  if

$$f_{\vec{x}}(\vec{x}, \theta) = g(\vec{x}) h(T(\vec{x}), \theta) \quad (\text{factorization criterion})$$

## Example

1)  $\bar{x}$  is sufficient for  $\theta$  if  $X_1, X_2, \dots, X_n \sim \text{Pois}(\theta)$

$$P(X=x) = \frac{e^{-\theta} \theta^x}{x!} \quad f_{\vec{x}}(\vec{x}, \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!} = \underbrace{\frac{1}{\prod_{i=1}^n x_i!}}_{g(\vec{x})} \underbrace{e^{-n\theta} \theta^{n\bar{x}}}_{h(\bar{x}, \theta)}$$

## Method of maximum likelihood

Likelihood function:  $L(x_1, \dots, x_n, \theta) = f(x_1, \dots, x_n, \theta) = \prod_{i=1}^n f_i(x_i, \theta)$

The likelihood of the sample is the following joint probability:

$$P(X_1=x_1, \dots, X_n=x_n | \theta)$$

or

$$\underbrace{P(X_1=x_1 | \theta) \dots P(X_n=x_n | \theta)}$$

This function is ought to be maximized

### Examples

- I) Suppose  $X_1, X_2, \dots, X_n$  are i.i.d random variables distributed as  $N(\mu, \sigma^2)$ . Find good estimators for parameters  $\mu$  and  $\sigma^2$

$$L(\vec{x}; \mu, \sigma^2) = \prod_{x_i \in \vec{x}} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} = \frac{1}{(\sqrt{2\pi}\sigma)^n} \exp\left(-\frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2\right)$$

Likelihood function

Since  $\ln L$  is a monotone increasing function it is possible to maximize  $\ln L$  to acquire maximum for  $L$

$$\ln L = -\frac{n}{2} \ln(2\pi) - n \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \mu)^2$$

$$\begin{cases} \frac{\partial \ln L}{\partial \mu} = \sigma^{-2} \sum_{j=1}^n (x_j - \mu) = 0 \\ \frac{\partial \ln L}{\partial \sigma^2} = -\frac{n}{\sigma^2} + \frac{1}{\sigma^4} \sum_{j=1}^n (x_j - \mu)^2 = 0 \end{cases}$$

$$\begin{cases} \frac{1}{n} \sum_{j=1}^n x_j = \bar{x} \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 \end{cases} \quad \left| \begin{array}{l} \hat{\mu} = \bar{x} \\ \hat{\sigma}^2 = S^2 \end{array} \right\} \text{good estimators for } N(\mu, \sigma^2)$$

## Tests of Hypotheses

A statistical hypothesis is an assertion or conjecture concerning one or more populations.

The structure of hypothesis testing will be formulated with the use of the term **null hypothesis**, which refers to **any hypothesis we wish to test** and is denoted as  $H_0$ . The rejection of  $H_0$  leads to acceptance of  $H_1$  (**alternative hypothesis**)

### Errors

**Type I Error.** Rejection of the null hypothesis when it is true.

**Type II Error.** Non-rejection of the null hypothesis when it is false

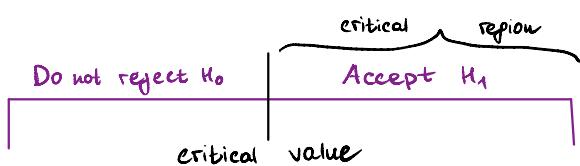
The probability of committing a type I error is called the **level of significance**, and is denoted via Greek letter  $\alpha$

The probability of committing a type II error is denoted via letter  $\beta$

- The type I error and type II error are related. A decrease in the probability of one generally results in an increase in the other.
- The size of the critical region and therefore the probability of committing a type I error, can always be reduced by adjusting the critical values.
- An increase in the sample size  $n$  will reduce both  $\alpha$  and  $\beta$ .
- The greater the distance between the true value and hypothesized value, the smaller  $\beta$  will be.

The **power** of a test is the probability of rejecting  $H_0$  given that alternative is true.

The power of a test may be computed as  $1 - \beta$



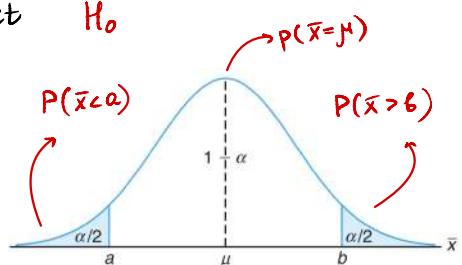
## Tests Concerning Sample Mean

Suppose  $X_1, X_2, \dots, X_n$  represent a random sample from a distribution with mean  $\mu$  and variance  $\sigma^2 > 0$ .

$$H_0: \mu = \mu_0 \\ H_1: \mu \neq \mu_0$$

Using the property of standard normal distribution we can derive a critical region for hypothesis  $H_0$  with probability  $\alpha$  to reject it when its true

If  $-Z_{\alpha/2} < Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} < Z_{\alpha/2}$  do not reject  $H_0$   
 or reject if  $\bar{X} < \mu_0 - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$  or  $\bar{X} > \mu_0 + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$



## Test on Two Means

The two-sided hypothesis on two means can be written generally as

$$H_0: \mu_1 - \mu_2 = d_0$$

The test-statistics is given as:

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - d_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Reject  $H_0$  if  $Z > Z_{\alpha/2}$  or  $Z < -Z_{\alpha/2}$

### Pooled t-test

$$H_0: \mu_1 = \mu_2 \quad (\text{assume } \sigma_1 = \sigma_2) \\ H_1: \mu_1 \neq \mu_2$$

$$\text{Test statistics: } t = \frac{(\bar{X}_1 - \bar{X}_2) - d_0}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \\ S_p = \sqrt{\frac{S_1^2(n_1-1) + S_2^2(n_2-1)}{n_1 + n_2 - 2}}$$

Reject  $H_0$  if  $|t| > t_{\alpha/2, n_1+n_2-2}$  or  $S_p^2 < t_{\alpha/2, n_1+n_2-2}$

## Unknown and Unequal Variances

The test statistic is given as

$$T' = \frac{(\bar{X}_1 - \bar{X}_2) - d_0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

$$\text{and has } \nu = \left( \frac{S_1^2}{n_1} \right)^2(n_1-1) + \left( \frac{S_2^2}{n_2} \right)^2(n_2-1) \quad \text{degrees of freedom}$$

As a result, we do not reject  $H_0$  when  $-t_{\alpha/2, \nu} < t' < t_{\alpha/2, \nu}$

Table 10.3: Tests Concerning Means

$H_0$	Value of Test Statistic	$H_1$	Critical Region
$\mu = \mu_0$	$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$ ; $\sigma$ known	$\mu < \mu_0$	$z < -z_\alpha$
		$\mu > \mu_0$	$z > z_\alpha$
		$\mu \neq \mu_0$	$z < -z_{\alpha/2}$ or $z > z_{\alpha/2}$
$\mu = \mu_0$	$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$ ; $v = n - 1$ , $\sigma$ unknown	$\mu < \mu_0$	$t < -t_\alpha$
		$\mu > \mu_0$	$t > t_\alpha$
		$\mu \neq \mu_0$	$t < -t_{\alpha/2}$ or $t > t_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$ ; $\sigma_1$ and $\sigma_2$ known	$\mu_1 - \mu_2 < d_0$	$z < -z_\alpha$
		$\mu_1 - \mu_2 > d_0$	$z > z_\alpha$
		$\mu_1 - \mu_2 \neq d_0$	$z < -z_{\alpha/2}$ or $z > z_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{1/n_1 + 1/n_2}}$ ; $v = n_1 + n_2 - 2$ , $\sigma_1 = \sigma_2$ but unknown, $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$	$\mu_1 - \mu_2 < d_0$	$t < -t_\alpha$
		$\mu_1 - \mu_2 > d_0$	$t > t_\alpha$
		$\mu_1 - \mu_2 \neq d_0$	$t < -t_{\alpha/2}$ or $t > t_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$t' = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$ ; $v = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$ , $\sigma_1 \neq \sigma_2$ and unknown	$\mu_1 - \mu_2 < d_0$	$t' < -t_\alpha$
		$\mu_1 - \mu_2 > d_0$	$t' > t_\alpha$
		$\mu_1 - \mu_2 \neq d_0$	$t' < -t_{\alpha/2}$ or $t' > t_{\alpha/2}$
paired observations	$t = \frac{\bar{d} - d_0}{s_d/\sqrt{n}}$ ; $v = n - 1$	$\mu_D < d_0$	$t < -t_\alpha$
		$\mu_D > d_0$	$t > t_\alpha$
		$\mu_D \neq d_0$	$t < -t_{\alpha/2}$ or $t > t_{\alpha/2}$

## Goodness-of-fit test

This kind of the test is ought to determine the distribution of the given population

We suppose that given population has a certain distribution and therefore sampling vector is expected as  $\vec{o} = (o_1, o_2, \dots, o_n)^T$  where  $o_i$  is expected number of occurrences of  $i$ -th outcome.

Suppose  $\vec{e}$  is a vector of observed number of occurrences.

Then

$$\chi^2 = \sum_{i=1}^k \frac{(o_i - e_i)^2}{e_i}$$

approx. has chi-squared distribution with  $k-1$  degrees of freedom.

## Test for Independence

$\chi^2$  test can be used for testing hypothesis of independence

Let our observations consists of pairs of i.i.d.r.v  $(X_1; Y_1), (X_2; Y_2) \dots (X_k; Y_k)$

$H_0: X_i$  and  $Y_i$  are independent

$X \setminus Y$	$b_1, b_2, \dots, b_m$	$b_m$	$\sum$
$a_1$	$\tilde{v}_{11}, \tilde{v}_{12}$	$\tilde{v}_{1m}$	$\tilde{v}_{1+}$
$a_2$	$\tilde{v}_{21}$	$\tilde{v}_{2m}$	$\tilde{v}_{2+}$
$\vdots$			
$a_n$	$\tilde{v}_{n1}$	$\tilde{v}_{nm}$	$\tilde{v}_{n+}$
$\sum$	$\tilde{v}_{++}, \tilde{v}_{+2}$	$\tilde{v}_{+m}$	$\tilde{v}_{++}$
	$P_{1+} P_{+2}$	$P_{+m}$	

$$H_{ij}: 1 \leq i \leq n \quad 1 \leq j \leq m \quad P(X=a_i, Y=b_j) = P(X=a_i)P(Y=b_j)$$

independence condition

Let  $\tilde{v}_{ij}$  be a number of outcomes such that  $\begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} a_i \\ b_j \end{pmatrix}$

$H_0: P(X=a_i, Y=b_j) = P_{i+} \cdot P_{+j} \quad H_{ij}$

1) Using the method of maximum likelihood:  $L = \prod_{\substack{i=1 \\ 1 \leq i \leq n}}^{m} (P_{i+} P_{+j})^{\tilde{v}_{ij}} = \prod_{i=1}^n P_{i+}^{\tilde{v}_{i+}} \prod_{j=1}^m P_{+j}^{\tilde{v}_{+j}}$

$$\ln L = \sum_{i=1}^n \tilde{v}_{i+} \ln P_{i+} + \sum_{j=1}^m \tilde{v}_{+j} \ln P_{+j}$$

$$\text{under condition } P_{1+} + P_{2+} + \dots + P_{n+} = 1$$

$$\begin{aligned} P_{i+} &= \frac{\tilde{v}_{i+}}{\tilde{v}_{++}} & \tilde{P}_{ij} &= \frac{\tilde{v}_{i+} \tilde{v}_{+j}}{\tilde{v}_{++}} \\ P_{+j} &= \frac{\tilde{v}_{+j}}{\tilde{v}_{++}} & \sum \frac{\left( \tilde{v}_{ij} - \frac{\tilde{v}_{i+} \tilde{v}_{+j}}{\tilde{v}_{++}} \right)^2}{\frac{\tilde{v}_{i+} \tilde{v}_{+j}}{\tilde{v}_{++}}} &= \tilde{v}_{++} \sum_{ij} \frac{\left( \tilde{v}_{ij} - \frac{\tilde{v}_{i+} \tilde{v}_{+j}}{\tilde{v}_{++}} \right)^2}{\tilde{v}_{i+} \tilde{v}_{+j}} \sim \chi^2_{(n-1)(m-1)} \end{aligned}$$

2) Calculate the sum and decide about the hypothesis.

## Linear Regression

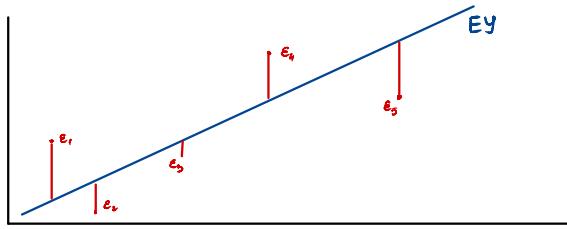
### Simple Linear Regression Model

The model must include the set  $\{(x_i, y_i); i=1, 2, \dots, n\}$  of data involving  $n$  pairs of  $(xy)$  values. The value  $y_i$  depends on  $x_i$  via a linear structure that also has the random component involved. The basis for the use of a statistical model relates to how the random variable  $y$  moves with  $x$  and the random variable.

$$y = \beta_0 + \beta_1 x + e \quad \text{Random variable (error)}$$

where  $\beta_0$  and  $\beta_1$  are unknown intercept and slope, while  $e$  is a random variable that is distributed with  $E(e) = 0$  and  $\text{Var}(e) = \sigma^2$

Since  $E(e)$  is assumed to be 0, then  $Ey = \beta_0 + \beta_1 x$  is a **regression line**



The purpose of the regression analysis is to estimate unknown coefficients  $\beta_0$  and  $\beta_1$ .

The estimated regression line is given as:  $\hat{y} = b_0 + b_1 x$

Residual: Error in Fit

Given a set of regression data  $\{(x_i, y_i); i=1, 2, \dots, n\}$  and a fitted model  $\hat{y}_i = b_0 + b_1 x_i$ , the  $i$ -th residual is given as  $e_i = y_i - \hat{y}_i, i=1, 2, \dots, n$

### The Method of Least Squares

$$\text{SSE} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - b_0 - b_1 x_i)^2$$

Differentiating SSE yields:

$$\text{SSE}_{b_0} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) \quad \text{SSE}_{b_1} = -2 \sum_{i=1}^n (y_i - b_0 - b_1 x_i) x_i$$

Setting both derivatives to 0, we obtain values for  $b_0$  and  $b_1$ :

$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \text{and} \quad b_0 = \bar{y} - b_1 \bar{x}$$

### Properties of Estimators

Denote  $R = \sum_{i=1}^n (y_i - \hat{a} - \hat{\beta}(x_i - \bar{x}))^2$  as a **residual sum of squares** and  $\epsilon_i \sim N(0, \sigma^2)$

The following statements hold:

- 1)  $\hat{a} \sim N(\alpha, \frac{\sigma^2}{n})$ ,  $\hat{\beta} \sim N(\beta, \frac{\sigma^2}{S_{xx}})$
- 2)  $\frac{R}{\sigma^2} \sim \chi^2_{n-2}$
- 3)  $\hat{a}, \hat{\beta}, R$  are independent
- 4)  $(n-2)^{-1}R$  is an unbiased estimator for  $\sigma^2$