

# Applied Statistics and Data Analysis

## Part I

# Statistical Process and Quality Control

## Introduction

### Quality and Process Control

#### Definition of Quality

- Quality means fitness for use. and
- Quality is inversely proportional to variability.
- My definition: *Is* and *should* are the same.

#### Statistical Process Control (SPC)

- Statistical process control is, first and foremost, a way of thinking which happens to have some tools attached.

#### The Magnificent Seven

1. histogram
2. check sheet
3. Pareto chart
4. defect concentration diagram
5. cause-and-effect diagram
6. control chart
7. scatter diagram

## Control Charts

The basis of a control chart is a statistical hypothesis test.

### Hypothesis test

#### Question:

is  $|\bar{x} - \mu_0|$  significant?

- $\mu_0$  target value
- $\bar{x}$  Arithmetic mean of the measurements

#### Hypothesis:

Two-sided statistical test to check the two alternative hypotheses:

$$\begin{aligned} H_0 : \mu_0 &= \bar{x} \text{ i.e. process is not disturbed} \\ H_1 : \mu_0 &\neq \bar{x} \text{ i.e. process is disturbed} \end{aligned} \quad (1)$$

#### Statistical test:

**Test statistic:** (z-test, since  $\sigma$  is known)

$$z = \frac{\bar{x} - \mu_0}{\sigma} \sqrt{n} \quad (2)$$

**Critical value:** i.e. q-quantile of the normal distribution

$$P(|z| \leq z_q) = q \quad (3)$$

with  $q = 1 - \frac{\alpha}{2}$ , where  $\alpha = 0.0027 \rightarrow z_q \approx 3$

#### Statistical conclusion:

- If  $\leq |z| z_q \rightarrow$  accept null hypothesis, i.e. process is not disturbed.
- If  $< |z| z_q \rightarrow$  reject null hypothesis, i.e. process is disturbed.

Reverse it and determine acceptance and rejection limits of the test!

#### Control limits:

$$UCL = \mu_0 + z_q \frac{\sigma}{\sqrt{n}} \quad LCL = \mu_0 - z_q \frac{\sigma}{\sqrt{n}} \quad (4)$$

#### Statistical conclusion:

- If  $LSL \leq \bar{x} \leq UCL \rightarrow$  process is not disturbed.
- If  $\bar{x} < LCL$  or  $UCL < \bar{x} \rightarrow$  process is disturbed.

Problem: In general, process standard deviation is unknown.

- Monitoring the mean and the variation of a process.
- First, monitoring the variation, then (if variation under control) monitoring the mean.

Solution: Control charts! :)

## The Control Chart

No.	sample values						mean	sd	range
1	$x_{11}$	$x_{12}$	$\cdots$	$x_{1j}$	$\cdots$	$x_{1n}$	$\bar{x}_1$	$s_1$	$R_1$
2	$x_{21}$	$x_{22}$	$\cdots$	$x_{2j}$	$\cdots$	$x_{2n}$	$\bar{x}_2$	$s_2$	$R_2$
$\vdots$	$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$
$i$	$x_{i1}$	$x_{i2}$	$\cdots$	$x_{ij}$	$\cdots$	$x_{in}$	$\bar{x}_i$	$s_i$	$R_i$
$\vdots$	$\vdots$	$\vdots$		$\vdots$		$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$x_{k1}$	$x_{k2}$	$\cdots$	$x_{kj}$	$\cdots$	$x_{kn}$	$\bar{x}_k$	$s_k$	$R_k$

Figure 1: Data Set with Mean, Standard Deviation and Range

#### Mean values

$$\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij} \quad (5)$$

#### Standard deviations

$$s_i = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2} \quad (6)$$

#### Ranges

$$R_i = \max \{x_{ij} | j \in \{1, \dots, n\}\} - \min \{x_{ij} | j \in \{1, \dots, n\}\} \quad (7)$$

for all  $i \in \{1, \dots, k\}$

### Control Chart for $\bar{x}$ and R

#### R Chart Centerline

$$\bar{R} = \frac{1}{k} \sum_{i=1}^k R_i \quad (8)$$

#### Control limits

$$UCL = D_4 \bar{R}; \quad LCL = D_3 \bar{R} \quad (9)$$

#### $\bar{x}$ Chart based on R chart

##### Control limits

$$UCL = \mu + 3 \frac{\sigma}{\sqrt{n}}; \quad LCL = \mu - 3 \frac{\sigma}{\sqrt{n}} \quad (10)$$

Problem:  $\mu$  and  $\sigma$  are in general unknown and must be estimated from the process data.

Two-stage process:

- Make sure that the process standard deviation (R chart) is under statistical control. That is, if some samples are out of bounds, it is recommended to omit these measurements and recalculate the limits.
- Use  $\bar{R}$  to estimate the process standard deviation.

#### Centerline

We know that for an independent sample  $x_1, \dots, x_n$  from a normal distribution with parameters  $\mu$  and  $\sigma$  the mean

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_j \quad (11)$$

satisfies

$$E(\bar{x}) = \mu \quad \text{and} \quad Var(\bar{x}) = \frac{\sigma^2}{n} \quad (12)$$

The mean is an unbiased estimator with the standard error

$$SE(\bar{x}) = \frac{\sigma}{\sqrt{n}} \quad (13)$$

Assumption: R chart is under statistical control.

- The value  $\bar{R}$  is a reliable estimate for the mean range.
- The value  $\bar{R}$  is a reliable estimate for the process standard deviation

$$\hat{\sigma} = \frac{\bar{R}}{d_2} \quad (14)$$

Any samples excluded for construction of the R chart should also be disregarded for construction of the  $\bar{x}$  chart. This results in a sample of  $k^*$  valid samples, (where  $k^*$  denotes the reduced number of samples). Mean values of  $\bar{x}_1, \dots, \bar{x}_{k^*}$  provide an estimate of  $\mu$ , i.e

$$\bar{\bar{x}} = \frac{1}{k^*} \sum_{i=1}^{k^*} \bar{x}_i \quad (15)$$

#### Control limits

$$UCL = \bar{\bar{x}} + 3 \frac{\bar{R}}{d_2} \frac{1}{\sqrt{n}} \approx \bar{\bar{x}} + A_2 \bar{R} \quad (16)$$

$$LCL = \bar{\bar{x}} - 3 \frac{\bar{R}}{d_2} \frac{1}{\sqrt{n}} \approx \bar{\bar{x}} - A_2 \bar{R}$$

### Control Chart with $\bar{x}$ and s

## s Chart

**Centerline** The centreline of the s chart is denoted by  $\bar{s}$  and is calculated from the arithmetic mean of the standard deviations

$$\bar{s} = \frac{1}{k} \sum_{i=1}^k s_i \quad (17)$$

### Control limits

$$UCL = B_4 \bar{s}; \quad LCL = B_3 \bar{s} \quad (18)$$

## $\bar{x}$ Chart based on s chart

Using an s chart of a process that is under control, the process standard deviation can be estimated by

$$\hat{\sigma} = \frac{\bar{s}}{c_4} \quad (19)$$

Any samples excluded for construction of the s chart should also be disregarded for construction of the  $\bar{x}$  chart. This results in a sample of  $k^*$  valid samples, (where  $k^*$  denotes the reduced number of samples). Mean values of  $\bar{x}_1, \dots, \bar{x}_{k^*}$  provide an estimate of  $\mu$ , i.e

$$\hat{\mu} = \bar{\bar{x}} = \frac{1}{k^*} \sum_{i=1}^{k^*} \bar{x}_i \quad (20)$$

### Control limits

$$UCL = \bar{\bar{x}} + 3 \frac{\bar{s}}{c_4} \frac{1}{\sqrt{n}} \approx \bar{\bar{x}} + A_3 \bar{s} \quad (21)$$
$$LCL = \bar{\bar{x}} - 3 \frac{\bar{s}}{c_4} \frac{1}{\sqrt{n}} \approx \bar{\bar{x}} - A_3 \bar{s}$$

## Individual Control Charts

Individual control charts have exactly one measurement per sample. Problem: You cannot estimate variability from a single measurement. Idea: Use variation of two adjacent measurements.

### Moving ranges

$$MR_i = |x_{i+1} - x_i| \quad (22)$$

for all  $i \in \{1, \dots, n-1\}$ .

### Arithmetic mean of the moving ranges

$$\overline{MR} = \frac{1}{n-1} \sum_{i=1}^{n-1} MR_i \quad (23)$$

### Estimated process standard deviation

$$\hat{\sigma} = \frac{\overline{MR}}{d_2} = \frac{\overline{MR}}{1.128} \quad (24)$$

Since two neighboring measurements were used to calculate the moving ranges we have  $d_2 = 1.128$ .

### Centerline

The centerline for the individuals control chart is the arithmetic mean of the measured values.

$$\bar{\bar{x}} = \frac{1}{k} \sum_{i=1}^k x_i \quad (25)$$

### Control limits

$$UCL = \bar{\bar{x}} + 3 \frac{\overline{MR}}{1.128}; \quad LCL = \bar{\bar{x}} - 3 \frac{\overline{MR}}{1.128} \quad (26)$$

## Control Charts for Attributes Data – p Chart

Number of defectives under number tested is a discrete random variable.

Given: Random sample of size n, of which D parts are defective We know: The number of defective D under n examined parts follows a binomial distribution with the unknown probability p of success.

### Estimated probability

$$\hat{p} = \frac{D}{n} \quad (27)$$

### Variance

$$Var(\hat{p}) = \frac{p(1-p)}{n} \quad (28)$$

### Given:

- k random samples with  $n_1, \dots, n_k$  values.
- Each of these samples contains  $d_1, \dots, d_k$  defective products.

### k relative frequencies

$$p_1 = \frac{d_1}{n_1}, \dots, p_k = \frac{d_k}{n_k} \quad (29)$$

### Centerline

The centreline and the control limits of a p chart are again determined from a stable trial run with  $k^*$  valid samples.

Again  $k^* \leq k$  is the reduced number of samples.

Distinguish 2 cases:

1. The sample sizes  $n_1, \dots, n_k$  are all equal to n.
2. The sample sizes are not all equal.

### Case 1

#### Centerline

$$\bar{p} = \frac{1}{k^*} \sum_{i=1}^{k^*} p_i \quad (30)$$

#### Control limits

$$UCL = \bar{p} + 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n}}; \quad LSL = \bar{p} - 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n}} \quad (31)$$

### Case 2

#### Centerline

$$\bar{p} = \frac{d_1 + \dots + d_{k^*}}{n_1 + \dots + n_{k^*}} \quad (32)$$

#### Control limits

$$UCL_i = \bar{p} + 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n_i}}; \quad LSL_i = \bar{p} - 3\sqrt{\frac{\bar{p}(1-\bar{p})}{n_i}} \quad (33)$$

The control limits now depend on the index i.

## Part II

# Multiple Regression

## Part III

# Design of Experiment