# Lecture 18

## Probability Plots.

– Imagine we have some sample $\{x_1, ..., x_n\}$. (from some larger population)

– we may want to assume the data/population is distributed in a particular way.

<span style="color:red">↳ How can we know if this is a reasonable assumption? eg: given some data, can we assume a normal distribution?</span>

· Assumption more formally: we assume that the population consists of values of a random var $X$ with cdf $F(x)$.

– we use a probability plot to test this assumption.

– Minitab ( Graphs > probability plot ...)

## Step 1: Order the sample in increasing order and rename:

$$\{X_1, ..., X_n\} \longrightarrow X_{(1)} \leq X_{(2)} \leq ... \leq X_{(n)}$$

(so $X_{(1)}$ is the smallest, $X_{(n)}$ is largest.

$(--\quad\cdot(1)$                       $\cdot\quad\cdot(n)$            $\gamma\cdot\cdot$

__Now/ Idea:__   Want  $x_{(i)} \sim 100\left(\frac{i}{n}\right)^{th}$ percentile.

$\quad\quad \hookrightarrow$ __eg:__  If $\{x_{(1)} \leq \cdots \leq x_{(11)}\}$ is our sample, then

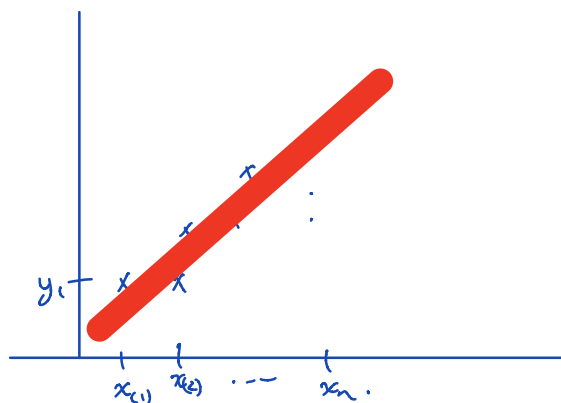$\quad\quad\quad$ we want  $x_{(3)} \sim$ median.

__Step 2:__  Find some points  $y_i$  such that:

$$P(X \leq y_i) = F(y_i) = \frac{i-0.5}{n}$$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ <span style="color:red">correction factor</span>

__Step 3:__  Plot  $(x_{(i)}, y_i)$.



$\quad\quad\quad$ straight.

__Step 4:__  Draw a line of best fit.

__Conclusions:__ Are all the points on or very near the

$\quad\quad$ line? $\longrightarrow$ __yes:__ your assumption is good.

$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ie. the distribution of $X$ is a good

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ reasonable fit to the data.

$\quad\quad\quad\quad\quad\quad\quad\quad$ __no:__ not so good.

most of the time we will be concerned with

the __normal__ distribution.

↳ here, $F(x) = \overline{\Phi}(x)$.
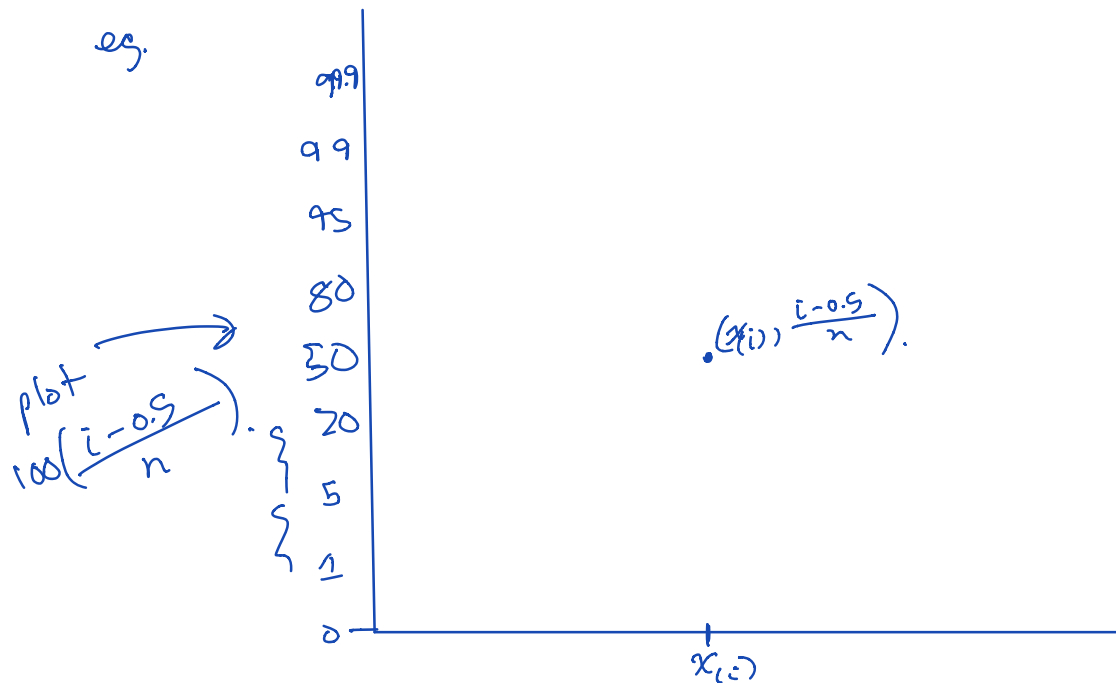
    ↳ to find $y_i$'s, use Z-score tables.

    ↳ want $y_i$ s.t.

$$\Phi(y_i) = \frac{i - 0.5}{n}$$

<u>Remark:</u> Sometimes prob plots are done on particular graphing paper

e.g.

plot $100\left(\dfrac{i-0.5}{n}\right)$.



On the vertical axis, values are marked: 99.9, 99, 95, 80, 50, 20, 5, 1, 0. The point plotted is $\left(x_{(i)}, \dfrac{i-0.5}{n}\right)$. The horizontal axis is labelled $x_{(i)}$.

# Chapter 7.    Points of Parameters.

__Recall:__ the general goal of statistics/statistical inference is to make predictions/draw conclusions about population, especially based on limited data.

- a major component of statistical inference is called __parameter estimation__.

  $\hookrightarrow$ eg: maybe you have some data set and you want to estimate the mean or variance.

- In general, a __parameter__ (usually denoted by a lowercase $\theta$) is any numerical property/feature of the data being studied.

- Recall that we assume a sample $\{x_1,...,x_n\}$ is a particular instance of independent and identically distributed random variables $\{X_1, ..., X_n\}$.

- A __statistic__ is any function of random variables.

  $\hookrightarrow$ eg. - $\overline{X} = \frac{1}{n}(X_1 + ... + X_n)$ $\mp$ sample mean
  - $S^2$
  - $S$

- Given a particular parameter, $\theta$, an __estimator__ for $\theta$ is a statistic $\wedge$

$$\widehat{(H)} = h(X_1, ..., X_n)$$

↳ capital theta w/ hat

used to estimate $\theta$.

↳ eg. $\overline{X}$ is a estimator for $\mu$.

↗ statistic.  ↗ parameter

- If $\widehat{(H)} = h(X_1 ..., X_n)$ is an estimator for $\theta$, then a particular value $h(x_1, ..., x_n) = \hat{\theta} \Leftarrow$ a <u>point estimate</u> for $\theta$.

## Descriptive Statistics: Sample

### Statistics

| Variable | N | N* | Mean | SE Mean | StDev | Minimum | Q1 | Median | Q3 | Maximum |
|---|---|---|---|---|---|---|---|---|---|---|
| Sample | 20 | 0 | 82.500 | 2.620 | 11.718 | 55.560 | 73.610 | 86.110 | 88.890 | 100.000 |

| | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 | C11 | C12 | C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Test 1 Scores | Sample | | | | | | | | | | | |
| 1 | 72.22 | 77.78 | | | | | | | | | | | |
| 2 | 61.11 | 55.56 | | | | | | | | | | | |
| 3 | 66.67 | 88.89 | | | | | | | | | | | |
| 4 | 83.33 | 83.33 | | | | | | | | | | | |
| 5 | 83.33 | 94.44 | | | | | | | | | | | |
| 6 | 94.44 | 66.67 | | | | | | | | | | | |
| 7 | 88.89 | 94.44 | | | | | | | | | | | |
| 8 | 88.89 | 88.89 | | | | | | | | | | | |
| 9 | 61.11 | 66.67 | | | | | | | | | | | |
| 10 | 66.67 | 77.78 | | | | | | | | | | | |
| 11 | 55.56 | 88.89 | | | | | | | | | | | |
| 12 | 83.33 | 100.00 | | | | | | | | | | | |
| 13 | 94.44 | 83.33 | | | | | | | | | | | |
| 14 | 72.22 | 83.33 | | | | | | | | | | | |
| 15 | 77.78 | 88.89 | | | | | | | | | | | |
| 16 | 72.22 | 94.44 | | | | | | | | | | | |
| 17 | 77.78 | 66.67 | | | | | | | | | | | |
| 18 | 38.89 | 88.89 | | | | | | | | | | | |
| 19 | 83.33 | 88.89 | | | | | | | | | | | |
| 20 | 61.11 | 72.22 | | | | | | | | | | | |
| 21 | 100.00 | | | | | | | | | | | | |
| 22 | 83.33 | | | | | | | | | | | | |

<u>Ex:</u> The mean for test 1 was $\mu = 79.15$.

Taking a sample of size 20.

$\overline{x} = 82.5 \sim 79.15$.

↳ point estimate for $\mu$.

Similarly: the sample variance

$\widehat{\sigma}^2 = s^2$ is an estimator for $\sigma^2$.

## <u>Sample Distributions & Central Limit Theorem.</u>

- Recall that a statistic is a function of random variables, $h(X_1, ..., X_n)$

— So a statistic is itself a random variable.
  ↳ therefore, each statistic has a prob. distrib.
    ↳ such distributions are called
         sampling distributions.
      ↳ eg: $-\overline{X} = \frac{1}{n}(X_1 + \cdots + X_n) \leftarrow$ sampling distribution
                                              of the mean $(\mu)$

      $- s^2$ is the sample distribution of $\sigma^2$.

Next: Central limit theorem.