

Today, we will:

- Do some more example problems – regression analysis
- Review the pdf module: **Outlier Points** and do some example problems

Example: Regression Analysis

Given: Twenty data points of (x,y) pairs with lots of scatter (see Excel spreadsheet on website for the raw data).

Data point	x	y
1	0	1.924
2	0.1	2.377
3	0.2	2.088
4	0.3	3.245
5	0.4	3.031
6	0.5	2.779
7	0.6	3.186
8	0.7	4.102
9	0.8	4.278
10	0.9	3.701
11	1	3.654
12	1.1	3.991
13	1.2	3.891
14	1.3	3.790
15	1.4	3.437
16	1.5	4.080
17	1.6	3.204
18	1.7	3.130
19	1.8	2.614
20	1.9	2.066

See EXCEL & MATLAB
files posted on the course
website →

Regression_nonlinear_fits.xls
" " .m

To do: Perform regression analysis – linear, quadratic, and cubic – and compare how the fitted curves fit to the data points.

Solution:

See Excel spreadsheet – I will show in class how to do the regression analysis in Excel.

Conclusions

- As polynomial order $m \uparrow$, multiple R also \uparrow
(correlation coefficient)
(i.e. the fit improves with increasing m)
- As polynomial order $m \uparrow$, standard error S \downarrow
(better overall agreement with the data points)

But - bigger m
is not always
better!

In this example, the cubic fit ($m=3$) is the best. If m gets too large, the curve fit will have too many "wiggles"



Example: Outliers – single set of data

Given: Four data points are measured.

Data point	x
1	38
2	42
3	44
4	53

$$\text{sample mean} = 44.25 = \bar{x}$$

$$\text{sample standard deviation} = 6.3443 = S$$

[We suspect either the min or the max to possibly be an outlier]

To do: Eliminate any “official” outliers, one at a time.

Solution:

- @ $n=4$, look up Thompson Tau $\rightarrow \bar{\tau} = 1.4250$ (either from the table or the eqn.)
- @ min x , $f = |38 - 44.25| = 6.25$ This one is biggest so this is the one we suspect may be an outlier
- @ max x , $f = |53 - 44.25| = 8.75$
- Compare $f_{\max} \div \bar{\tau}S \rightarrow f_{\max} = 8.75$ } Since $f_{\max} < \bar{\tau}S$, There are $(8.75 < 9.0406)$ no outliers
 $\bar{\tau} \cdot S = (1.4250)(6.3443) = 9.0406$ }

Example: Outliers – single set of measurements

Given: Janet takes 12 temperature measurements ranging from 23.0°C (lowest reading) to 25.7°C (highest reading).

- The sample mean of all 12 readings is 24.88°C .
- The sample standard deviation is 1.04°C .

To do: Is there an “official” outlier?

Solution:

- @ $n=12$, look up or calculate $\bar{\tau} \rightarrow \bar{\tau} = 1.8290$
- $\bar{\tau}S = (1.8290)(1.04^{\circ}\text{C}) = 1.902^{\circ}\text{C}$
- @ minimum reading, $f = |23.0 - 24.88| = 1.88^{\circ}\text{C}$ $f_{\max} \rightarrow$ The minimum reading is the one we suspect may be an outlier
- @ maximum reading, $f = |25.7 - 24.88| = 0.82^{\circ}\text{C}$
- Compare: $f_{\max} = 1.88^{\circ}\text{C} < \bar{\tau}S = 1.902^{\circ}\text{C}$

Thus, there are no outliers



Example: Outliers – single set of measurements

Given: Eleven measurements of pump efficiency are taken (listed in increasing order below):

Point # 1 2 3 4 5 6 7 8 9 10 11
58, 72, 74, 74, 74, 76, 77, 80, 80, 82, 85, and 92%

(a) **To do:** Are there any “official” outliers? If so, remove them. How many “good” measurements are left?

(b) **To do:** Based on the measurements that are left, estimate the population mean and its confidence interval to 95% confidence.

(c) **To do:** Based on the measurements that are left, estimate the population standard deviation and its confidence interval to 95% confidence.

Solution:

(a). First calculate sample mean \bar{x} : Sample standard deviation: $\bar{x} = 77.273$
 $s = 8.580$

• At x_{\min} , $s_1 = |58 - 77.273| = 19.27$ → Point X, is the suspected outlier since $s_{\max} = s_1 = 19.27$
Point #1
• At x_{\max} , $s_{11} = |92 - 77.273| = 14.73$
Point #11

• TABLE → @ $n=11$, $\bar{C} = 1.8153 \rightarrow \bar{CS} = (1.8153)(8.580) = 15.575 = \bar{CS}$

• Compare: $s_{\max} = s_1 = 19.27 > \bar{CS} = 15.575 \rightarrow$ So, Point 1 is an outlier

NOTE: Other data points may or may not also have $s > \bar{CS}$, but it does not matter. WE CAN REMOVE ONLY ONE OUTLIER AT A TIME

• Remove point 1: start all over, now with 10 points instead of 11 pts.

★ NOTE: Since we removed a data point, we have to re-calculate \bar{x} & s ★

ROUND 2: . New $\bar{x} = 79.20$ (sample mean of points x_2, x_3, \dots, x_{11})
 $s = 6.0332$ (sample standard deviation of points x_2, x_3, \dots, x_{11})

• \bar{C} is also different, since now $n=10$ instead of 11

TABLE @ $n=10$, $\bar{C} = 1.7984 \rightarrow \bar{CS} = 10.850$

• The maximum deviation turns out to be point 11:

$$s_{\max} = s_{11} = |92 - 79.20| = 12.80 = \underline{s_{\max}}$$

• Compare: $s_{\max} = s_{11} = 12.80 > \bar{CS} = 10.850 \rightarrow$ Point 11 is also an outlier!



ROUND 3 → Remove Point 11 & we are left with only 9 data points: (points 2 through 10)

$x \rightarrow$	72	74	74	76	77	80	80	82	85	11
Point #	1	2	3	4	5	6	7	8	9	10

• New $\bar{x}, s, \hat{C} \rightarrow \bar{x} = 77.778, s = 4.2655, \hat{C} = 1.7770$
 $\rightarrow \hat{CS} = 7.5798$

• Must calculate the deviations again, based on the revised value of \bar{x}

$$\min x \rightarrow \delta_2 = |72 - 77.778| = 5.778$$

$$\max x \rightarrow \delta_{10} = |85 - 77.778| = 7.222 \rightarrow \delta_{\max} = \delta_{10} \text{ this time}$$

• Compare: $\delta_{\max} = \delta_{10} = 7.222 < \hat{CS} = 7.5798 \rightarrow$ No MORE OUTLIERS

Bottom line: We removed two official outliers (points 1 & 11), and are left with the remaining 9 (points 2 through 10)

(b) Estimate μ based on the remaining data points, to 95% confidence

Solution: This is review $\rightarrow \bar{x} = 77.778 \quad df = n-1 = 9-1 = 8$
 $s = 4.2655 \quad t_{\alpha/2} = 2.3060 \text{ (TABLE)}$

$$\text{So, we predict } \mu = \bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}} = 77.778 \pm \frac{(2.306)(4.2655)}{\sqrt{9}}$$

$\mu = 77.78 \pm 3.28 \text{ % to 95% confidence}$

(c) Similarly, for review, estimate σ to 95% confidence

* DO THIS ON YOUR OWN FOR PRACTICE

ANSWER:

$2.88 \leq \sigma \leq 8.17 \text{ to 95% confidence}$