

LSC 15

## Dealing with

### Scaling Data

$$\text{scaled value} = \frac{\text{Original}}{10^i}$$

i - smallest power (scaler factor) < 1  
 $10^i \rightarrow$  scaling factor.

- 1 - Find max absolute value
- a - find  $\varphi$  such that  $10^i$  becomes minimum absolute value below 1
- b - Scale data.

Value: 1200, 5600, 8900, 2100, 10,000

1 - max value = 10,000

$$\varphi = \frac{10,000}{10^4} = 1 \times 10^0$$

so  $\varphi = 10^0$

less scale

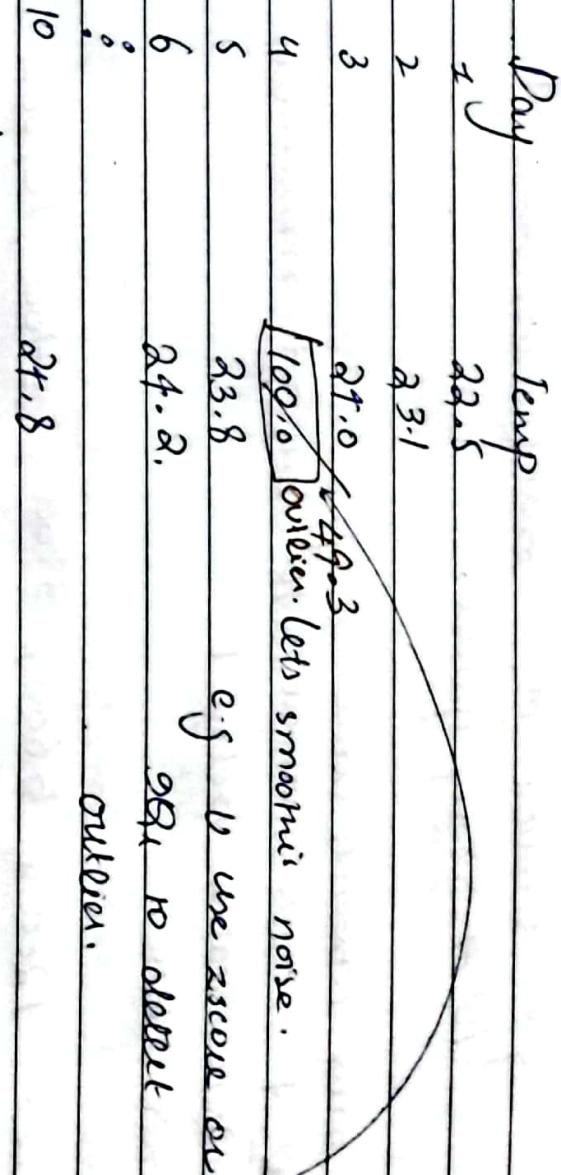
1200	0.012
5600	0.056
8900	0.089
2100	0.034
10000	0.1

Wish data set contain  
larger values  $\rightarrow$  Use min

Sensitive to outliers, so if  
may become very large

## Smoothing

Replacing noisy data.



Let's do Binning. Bins = 3 3 days.

Size,

$$10/3 = 3 \text{ or } 4.$$

Bin	Days	Temp	Mean
1	1, 2, 3	22.5, 23.1, 24.0	23.2
2	4, 5, 6	23.8, 24.2	24.0
3	7, 8, 9	23.5, 24.5, 25.0	24.5
4	10	24.8	24.8
	.		

## Normalization.

95 - 80  
95 - 80

- 1- Min max.  $\rightarrow$  when have diff  
want to bring to common

$$x_2 = 2, 5, 8, 12, \dots$$

$$x'_2 = 0, 0.06, 0.13, 0.21$$

Normalized data is easier to use  
2 - Z score.

$$\mu = \frac{2+5+8+12+\dots}{10}$$

~~sk~~

$$\sigma = \sqrt{\frac{(x-\mu)^2}{n}}$$

Z score.  $\rightarrow$

$$2 \rightarrow -1.3$$

$$5 \rightarrow -1.2$$

dis 0.7

$$8 \rightarrow -0.9$$

$$12 \rightarrow -0.6$$

$$18 \rightarrow -0.3$$

$\vdots$

## Normalization.

50      30      30 =

95 - 80

95 - 80

1 - Min max.  $\rightarrow$  When have diff scales in features - you

want to bring to common range.  $X' = \frac{X - \bar{X}}{S}$

$$X = 2, 2, 15, 8, 12, \dots \quad X' = \frac{X - \bar{X}}{S}$$

$$X' = 0, 0.08, 0.13, 0.21$$

$$S = 50^{-2}$$

Normalized data is easier to visualize.

2 - Z score.

$$Z = \frac{X - \bar{X}}{S}$$

$$11 = 2 + 5 + 8 + 12 + \dots$$

-  $B_3$  measure

$$B_1 = M = \overline{a+s+8} = \overline{[5, 5, 15]} \leftarrow S = \overline{[5, 5, 15]} , B_2 = , [a1.25, a1.25, a1.25, a1.25]$$

3

DR

change,

65-80 Hertz

$$B_1 = 18 + \frac{1}{20 - 20} = 18 + 0.05 = 18.05$$

3

-  $B_3$  measure

$$B_3 = (44.33, 44.33, 44.33)$$

2- Bin boundaries

$$[88, 88, 88]$$

$$[13, 13, 30, 30]$$

$$B_1 = [a, s, 8] \rightarrow [a, a, 8]$$

## Bin Range Data

8 - 18	2,15, 8, 12, 18
19 - 34	25, 30
35 - 50	38, 45, 50

- Equi-depth (On <sup>order of</sup> ~~selected data~~)  
Total value = 10

$$\text{Bin size} = 10/3$$

Bin	Values
1st	2, 5, 8
2nd	12, 18, 25, 30
3rd	38, 45, 50

↳ Sloping customers purchase more bins where each bin contains several smooth no. of customers

1 - Bin means

$$B.M = M =$$

center over

$$\text{Range} = 50 - 2 = 48$$

$$\text{Bin width} = \frac{48}{3} = 16$$

Bin range      Values

2 - 18      2, 15, 8, 12, 1

19 - 34      25, 30

35 - 50      38, 45, 5

- Equi depth      Bin width  
fixed

Total value  $\Rightarrow$

Bin size = 1

Bin      Values

$$\textcircled{2} \approx \sqrt{10} = 3.16 \approx \textcircled{3}$$

$$k = \sqrt{n}$$

Solve more to understand.

$$100 \text{ or } b^2 n = 3$$

1 - Even number (on natural numbers)

2, 5, 8, 12, 18, 25, 30, 38, 45, 50

lets take alternate

$7.396 < 7.81$  we fail to reject  $H_0$  at  
 $\alpha = 0.05$ .

We can conclude what actual affects energy usage.

stats based on

for a whole dataset

using a test?

Water: Eneg

$$\frac{(5-7.25)^2}{7.25} = \frac{5.06}{7.25} = 0.698$$

Water: Troch

$$\frac{(5-2.75)^2}{2.75} = 10.824$$

$$\frac{\text{Coffee} - E}{7.25} = \frac{(8-7.25)^2}{7.25} = 0.072$$

$$C - T = 0.204$$

$$\text{Tea } E = 0.215, T - T = 0.567$$

$$E-D-E = 1045, E-D-T = 2.75$$

$$6.1 \quad \chi^2 = 0.698 + 1.84 + 0.072 + 0.204 + 0.215 + 0.567 +$$

$$\boxed{\chi^2 = 7.396}$$

To Compare with Chi Table.

df	P=0.10	P=0.05 (sig)	P=0.01 (v. sig.)
3	6.25	7.81	11.34

$$df = (4-1) + (3-1) = 3$$

lets calculate 3 expected value manually.

$$E = \frac{R.T \times C.T}{C.p.T} =$$

3 - Energic.

$$1) - E = \frac{10 \times 29}{40} = \frac{290}{40} = 7.25$$

$$3) \quad E = \frac{10 \times 29}{40} = 7.25$$

7.25

No rest will be determined automatically.

4 - Expected.

	Obsv.	Expected.	Total.
N	7.25	$10 - 7.25 = 2.75$	10
C	7.25	2.75	10
T	7.25	2.75	10
C.D	7.25	2.75	10
Total.	29.	11	40

5 - Chi-Square.

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

Cell by cell calculation.

### 3 d) Example.

A researcher wants to test whether the type of drink affects energy level. They conduct experiment of 40 people. They divide them in 4 groups based on drink they consumed.

Water, coffee, Tea, Energy drink

$H_0 \rightarrow$  type of drink has no effect - all due to random chance.

$H_1 \rightarrow$  type of drink consumed effect on person energy.

d) Actual table.

	Energy drink	Total
Water	5	10
Coffee	8	10
Tea	6	10
E.D	10	10
Total	29	40

"Yes, Charlie! This means we do not have enough evidence to say carrots actually help. The difference could have happened just by chance."

The Grand Professor nodded.

"Yes, Charlie! This means we do not have enough evidence to say carrots actually help. The difference could have happened just by chance."

Charlie now understood:

- df controls which Chi-Square Table we use.
- p-value tells us if the result is significant.
- If  $p < 0.05$ , we reject  $H_0$  (the null hypothesis) and say it's significant.
- If  $p > 0.05$ , we fail to reject  $H_0$ , meaning there's not enough proof.

Calculate Carter's V strength of association b/w variables

$$V = \sqrt{\frac{X^2}{\max(k-1)}} = \sqrt{\frac{3.84}{20 \times (2-1)}} = 0.48$$



Charlie ran back to Mathville and told everyone:

"Statistics is like detective work! You can't just guess—you have to prove if something is real or just random!"

And from that day on, Mathville became the most logical place in the world.

The End.

## Key Takeaways:

Term	Meaning
p-Value	Probability that your result happened just by chance
Significant Value (Critical Value)	The cutoff number that decides if your result is "real"
$p < 0.05$	We call the result statistically significant <input checked="" type="checkbox"/>
$p > 0.05$	The result is not significant—it could be just luck <input checked="" type="checkbox"/>

$$df = (\text{Rows} - 1) \times (\text{Columns} - 1)$$

Here, we have 2 rows and 2 columns, so:

$$df = (2 - 1) \times (2 - 1) = 1$$

Charlie checked  $df = 1$  in the Chi-Square Table:

df	p = 0.10	p = 0.05 (Significant)	p = 0.01 (Very Significant)
1	2.71	3.84	6.63

His chi-square value was 3.34, which was less than 3.84.

p

- The p-value is greater than 0.05!

## Final Level: What Does This Mean?

Charlie had solved the case! 🎉

"Because our chi-square number (3.34) is less than 3.84, the p-value is greater than 0.05, meaning we CANNOT say for sure that carrots help running speed!"

The Grand Professor nodded.

"Yes, Charlie! This means we do not have enough evidence to say carrots actually help. The difference could have happened just by chance."

Charlie now understood:

- df controls which Chi-Square Table we use.
- p-value tells us if the result is significant.
- If  $p < 0.05$ , we reject  $H_0$  (the null hypothesis) and say it's significant.
- If  $p > 0.05$ , we fail to reject  $H_0$ , meaning there's not enough proof.

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

Where O is observed and E is expected."

Charlie calculated:

$$\begin{aligned}\chi^2 &= \frac{(8 - 6)^2}{6} + \frac{(2 - 4)^2}{4} + \frac{(4 - 6)^2}{6} + \frac{(6 - 4)^2}{4} \\ &= \frac{4}{6} + \frac{4}{4} + \frac{4}{6} + \frac{4}{4} = 0.67 + 1 + 0.67 + 1 = 3.34\end{aligned}$$

"Great job, Charlie! Now, to get the p-value, we need the Chi-Square Table!"

## Level 3: The p-Value and the Significance Test ⚖

Charlie ran to the Chi-Square Table and found his degrees of freedom (df).

"Professor, how do I get df?" 🤔

The professor explained,

"You count how much your data can change before one value becomes fixed. Use this formula:

$$df = (\text{Rows} - 1) \times (\text{Columns} - 1)$$

Here, we have 2 rows and 2 columns, so:

$$df = (2 - 1) \times (2 - 1) = 1$$

Charlie checked  $df = 1$  in the Chi-Square Table:

## INFLUENCE OF POLYMER ON PROPERTIES OF RESINS (1)

W. J. KELLY AND R. L. HARRIS, JR., E. I. DU PONT DE NEMOURS & COMPANY, WILMINGTON, DELAWARE

Received June 1, 1966; revised August 1, 1967

Editor

Editorial

Editorial

Editorial

Editorial, Technical, and Editorial Assistant: G. S. SAWYER

Editorial

# The Game of Numbers: The Secret of Degrees of Freedom and p-Value

Once upon a time, in the land of Mathville there was a fun game played by scientists and detectives. This game helped them find out if something was real or just random luck.

Our hero Charlie Square  loved solving number mysteries. One day, the Grand Professor of Numbers gave him a new case:

✓ "Charlie, people are claiming that eating carrots helps kids run faster! But we don't know if this is actually true or just a random coincidence. Your mission: find the truth!"

Charlie was excited! But how could he prove whether carrots really make kids faster? 

He Null Eating carrots does not improve running speed.  
 $H_0 \rightarrow$  opposite

## Level 1: The Running Race

Charlie gathered 20 kids and made them run a race.

- 10 kids ate carrots before running
- 10 kids did not eat carrots before running

Charlie wrote down their times and wanted to see if carrot-eaters really ran faster. But how? 

The Grand Professor handed Charlie a magic tool called the Chi-Square Test! 

✓ "Charlie, you can't just look at the times and guess. You need to test if the difference is real or just random chance!"

Charlie was ready! 

## Level 2: The Secret of Degrees of Freedom

## Step deviation.

Step h.10

Since values have constant gap of 10.

$$\Delta x = \frac{X - 30}{h}$$

assumed

Uniform gap from  
more no were  
sample.

$$\Delta y = Y - 33$$

$\Delta x$	$\Delta y$	$\Delta x \cdot \Delta y$	$\Delta x^2$	$\Delta y^2$
-2	-2.1	4.2	4	4.41
-1	-0.8	0.8	1	0.64
0	0	0	0	0
1	1.2	1.2	1	1.44
2	2.2	4.4	4	4.84

$$h = \frac{6.6}{\sqrt{10 \times 11.43}}$$

$$= 0.997$$

Variance in grouped data or when nos are  
larger.

## Assumed Mean

- For larger dataset.
- Excluding mean calculation by using assumption of mean like by median or mode.
- If value is decimal and you want whole number.

$$A_x = 30$$

$$A_y = 33$$

$x - 30$	$d_y = y - 33$	$d_x \cdot d_y$	$d_x^2$	$d_y^2$
20	-2	40	400	4
10	-13	80	100	169
0	0	0	0	0
10	12	120	100	144
20	22	440	400	484

$$\lambda = \frac{1060}{\sqrt{1000 + 1133}} = 0.996$$

$$dx = X - 30 \Rightarrow -80, -10, 0, 10, 20$$

$$dy = Y - 34 \Rightarrow -22, -9, -1, 11, 21$$

$$dx \cdot dy \Rightarrow -140, 90, 0, 110, 120$$

$$\begin{aligned} dx^2 &\Rightarrow 400, 100, 0, 100, 400 \\ dy^2 &\Rightarrow 484, 81, 1, 121, 441 \end{aligned}$$

$$k = \frac{\sum dx \cdot dy}{\sqrt{\sum dx^2 \cdot \sum dy^2}}$$

$$= \frac{-1060}{\sqrt{1000 \times 1128}} \Rightarrow \frac{-1060}{1062.07}.$$

$$\boxed{k = 0.998}$$

## Correlation.

Actual mean method.

For small datasets due to complex calculations.

$x$	10	20	30	40	50
$y$	12	25	33	45	55

$$1 - \bar{x} = 30 \quad \text{it can be decimal.}$$

$$\bar{y} = 34 \quad \text{integer}$$

2 - Compute deviations:

$$\Delta x = x - \bar{x}$$