# 2.6 Guessing the Distribution of Dataset 3

| | Normal Distribution |
|---|---|
| |  <br><br> • Somewhat Gaussian in nature. <br> • Side bit is sitting between 90 to 100, so it's inclining towards something. It's telling us that there are significant number of observations that are towards one side. Meaning the distribution has a tail which is on the right side. So it could be a Gaussian with the right tail. |
| | count    48.000000 <br> mean    56.958333 <br> std    14.871018 <br> min    27.000000 <br> 25%    48.000000 <br> 50%    56.000000 <br> 75%    68.000000 <br> max    101.000000 <br><br> {('Mean Observed', 56.96), ('Variance Observed', 221.15), ('Skew Observed', 0.37), ('Kurt Observed', 0.6)} <br><br> • The difference between the minimum and the 25th quantile around 20. <br> • Between 25th and 50th is around 10. <br> • 50th to 75th is again around 10. <br> • 75th to the maximum is more than 30, that means it has a tail on the right side. <br><br> • Also, mean and the median are almost the same. Meaning, it's symmetric. So, huge symmetry with a slight tail. <br><br> • Positive skew, means that the distribution is right tailed. |
| | **Q-Q Plot** <br>  <br><br> **P-P Plot** <br>  <br><br> • Both plots suggest that it's a normal distribution. |
| | NULL HYPOTHESIS: The given data follows Normal distribution. <br><br> ALTERNATE HYPOTHESIS: The given data does not follow Normal distribution |
| | We split the data into 6 equal intervals (you can split it in as many intervals you want). <br><br> As the area under the curve is one, the area under each bucket will be $\frac{1}{6} = 0.167$. <br><br> This the command to do that: <br><br> ```python<br>n = 1/6<br>for i in range(1,6):<br>    prob_intervals = [scipy.stats.norm.ppf(i*n,df['obs'].mean(), df['obs'].std())]<br>``` |

| print(prob_intervals) | So the 6 buckets we have now are: |
|---|---|
| [42.571790240767264] <br> [50.552980104100826] <br> [56.958333333333336] <br> [63.363686562565846] <br> [71.3448764258994] | 1st bucket  : 42.57 and below <br> 2nd bucket : 42.57 to 50.55 <br> 3rd bucket : 50.55 to 56.95 <br> 4th bucket : 56.95 to 63.36 <br> 5th bucket : 63.36 to 71.34 <br> 6th bucket : 71.34 and above |

| | |
|---|---|
| | Since, we're splitting into six buckets, and we had $48$ observations, so the expected frequency under each bucket $\frac{48}{6} = 8$.<br><br>• Expected freq = [8, 8, 8, 8, 8, 8]<br>• Observed freq = [9, 7, 9, 7, 9, 7]<br><br>• Observed frequencies can be obtained by first sorting the dataset in excel in ascending order, and then count the number of values fall in each buckets.<br>   ○ e.g., here, there are 9 values in the 1st bucket, i.e., 9 values in the dataset below 42.57 |
| | Calculated chi square statistic = 0.75<br>p-value = 0.98<br><br>• p-value > α<br>  We accept the null. |
| | $df = k - p - 1 = 6 - 2 - 1 = 3$<br><br>$k = 6$ : number of buckets<br><br>$p = 2$ : for normal distribution |
| | Tabulated Chi Square value $= 7.81$<br><br>    Tabulated value $>$ Calculated value<br>      We accept the null. |
| | |
| Business Cases | • Most of our life events can be assumed to be normal.<br>  • e.g., Our test scores over the years could be a normal distribution.<br>  • Or test scores of an entire class will be normally distributed.<br><br>• The data that we were working on comes from automotive sales. There were around 48 salesman who have sold tractors on a particular month. So, ideally you will see a couple of people who are performing extremely well, and there will be some salesman who may not be performing well, and in between you will have the majority. |
| | |
| | |