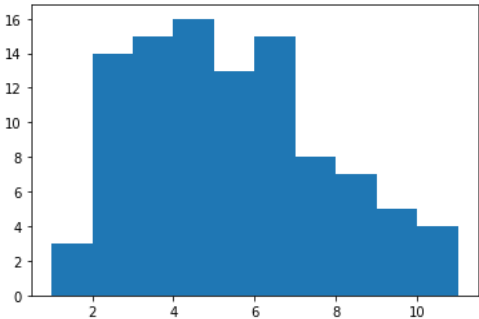
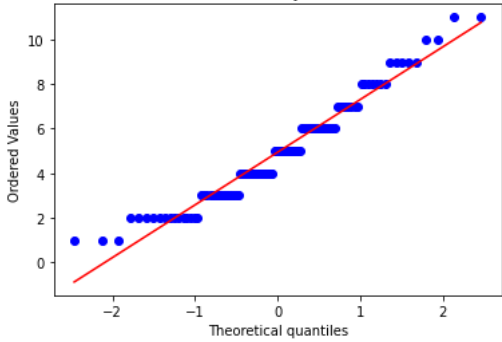
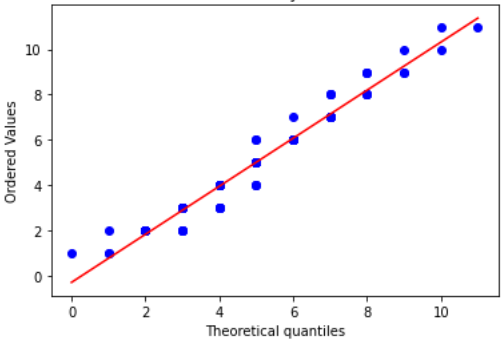


2.4 Guessing the Distribution of Dataset 1

Sunday, 09 October 2022 14:09

	Poisson Distribution
Histogram	<div></div> <div><ul style="list-style-type: none"><li>• Since there are two peaks, it could be a non-Gaussian distribution.</li><li>• But, we cannot say whether it's a symmetric or a non-symmetric distribution.</li></ul></div>
Descriptive Statistics	<div><div><div>count</div><div>100.000000</div></div><div><div>mean</div><div>4.940000</div></div><div><div>std</div><div>2.381834</div></div><div><div>min</div><div>1.000000</div></div><div><div>25%</div><div>3.000000</div></div><div><div>50%</div><div>5.000000</div></div><div><div>75%</div><div>6.000000</div></div><div><div>max</div><div>11.000000</div></div></div> <div><div>{('Variance Observed', 5.67), ('Mean Observed', 4.94)},</div><div>('Skew Observed', 0.51), ('Kurt Observed', -0.38)</div></div> <div><ul style="list-style-type: none"><li>• Variance and mean are close to each other. Think about what distributions have mean and variance close to each other. Ans: Poisson Distribution</li><li>• Observing the skewness and the kurtosis, we can conclude that we can eliminate certain symmetric distributions such as normal distribution, uniform distribution.</li><li>• Mean = Median = 5. So, 50% points are lying below 5 and 50% lying above 5. So it's still a symmetric distribution. But, Poisson is not a symmetric distribution.</li></ul></div>
P-P Plot	<div><div>Probability Plot</div><div></div></div> <div><ul style="list-style-type: none"><li>• Points are not lying on the red line. It means the data is not coming from the normal.</li><li>• Since P-P plot deals with the CDF, and we have data similar to a step function, which also hints that this <b>distribution could be discrete in nature.</b></li></ul></div>
	<div>P-P w.r.t Poisson distribution with lambda = 4.94</div> <div><div>Probability Plot</div><div></div></div> <div><ul style="list-style-type: none"><li>• Now the points are almost around the red line. So, with a certain degree of certainty, we can say that the cumulative distribution from which this data is coming is a Poisson distribution which has a lambda of 4.94.</li></ul></div>
Statistical tests	<div><div><div>frequency</div><div>OBS_PROBA</div><div>POISSON_PMF</div><div>POISSON_FREQ</div></div><div>obs</div><div><div><div>1</div><div>3</div><div>0.03</div><div>0.035344</div><div>3.53</div></div><div><div>2</div><div>14</div><div>0.14</div><div>0.087299</div><div>8.73</div></div><div><div>3</div><div>15</div><div>0.15</div><div>0.143752</div><div>14.38</div></div><div><div>4</div><div>16</div><div>0.16</div><div>0.177534</div><div>17.75</div></div><div><div>5</div><div>13</div><div>0.13</div><div>0.175404</div><div>17.54</div></div><div><div>6</div><div>15</div><div>0.15</div><div>0.144416</div><div>14.44</div></div><div><div>7</div><div>8</div><div>0.08</div><div>0.101916</div><div>10.19</div></div><div><div>8</div><div>7</div><div>0.07</div><div>0.062933</div><div>6.29</div></div><div><div>9</div><div>5</div><div>0.05</div><div>0.034543</div><div>3.45</div></div><div><div>10</div><div>2</div><div>0.02</div><div>0.017064</div><div>1.71</div></div><div><div>11</div><div>2</div><div>0.02</div><div>0.007663</div><div>0.77</div></div></div></div> <div><ul style="list-style-type: none"><li>• What we first did here, we made a frequency chart (shown in frequency column).</li></ul></div>

	<ul style="list-style-type: none"><li>• <math>OBS\_PROB = P(x_i = k) = \frac{frequency\ of\ k}{total\ frequency}</math> .<ul style="list-style-type: none"><li>◦ <i>e.g.</i>, Total frequency = 100 <math>P(x_i = 1) = \frac{3}{100} = 0.03</math> .</li></ul></li><li>• Then we calculate Poisson pmf<ul style="list-style-type: none"><li>◦ <math>P(x_i = k) = e^{-\lambda} \frac{\lambda^k}{k!}</math> .</li><li>◦ Here, <math>\lambda = 4.94</math> .</li></ul></li><li>• After this step, we calculate Poisson frequency,<ul style="list-style-type: none"><li>◦ <math>Poission\ freq = Poisson\ pmf \times total\ freq</math> <math>\implies = Poisson\ pmf \times 100</math> .</li></ul></li></ul>
Observed and Expected frequency	<ul style="list-style-type: none"><li>• Observed probability tells you the probability of an observation occurring from the sample.</li><li>• The Poisson pmf tells you, the probability of getting the same observation from the population (assuming that the population is Poisson) with the mean of 4.94.</li><li>• So now we have both, the observed frequency and the expected frequency (Poisson freq).</li></ul>
Hypothesis	NULL HYPOTHESIS: The given data follows Poisson distribution.  ALTERNATE HYPOTHESIS: The given data does not follow Poisson distribution
Chi-square Test	Calculated chi square statistic = 7.92 p-value = 0.64 <ul style="list-style-type: none"><li>• <math>p\text{-value} &gt; \alpha</math> This is said to be almost coming from the distribution. So we can accept the null.</li></ul>
Degrees of freedom	$df = k - p - 1 = 11 - 1 - 1 = 9$ .  $k = 11 \implies$ total number of classes.
Tabulated Chi-Square	= 16.92 <ul style="list-style-type: none"><li>• Tabulated value &gt; Calculated value  » We accept the null hypothesis</li></ul>
Business Cases	<ul style="list-style-type: none"><li>• suppose the given data is from a traffic signal in a city, where the number represents the number of times the signal was violated on a given day, i.e., on day #1 the signal was violated 5 times, on day #2 the signal was violated 4 times, and so on. This is a <b>count information</b>.</li><li>• Another example would be that let's say everyday you're manufacturing a thousand products. Out of those thousand products, how many are defective. So, 5 could be the number of defective items on day 1, 4 could be on day #2, and so on.</li><li>• <b>Wherever you have count information, those places could be ideal for Poisson distributions to happen.</b> So, a number of discrete events happening in continuous space is Poisson distribution.</li><li>• Other examples could include customer arrivals and so on.</li></ul>

•