
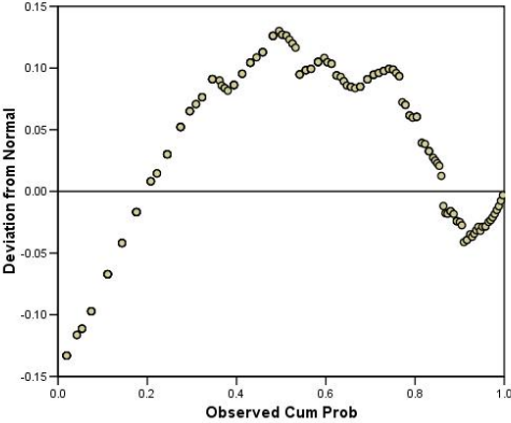
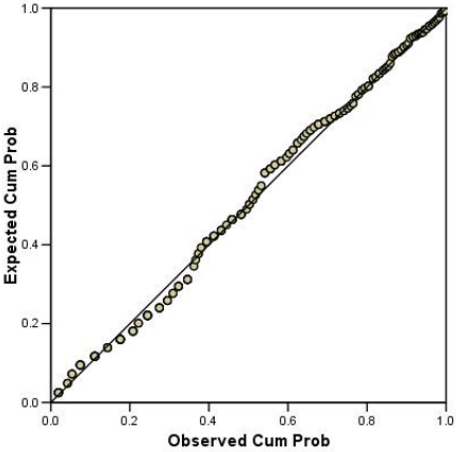
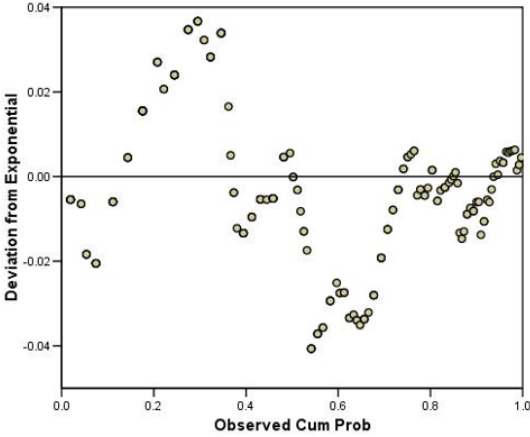
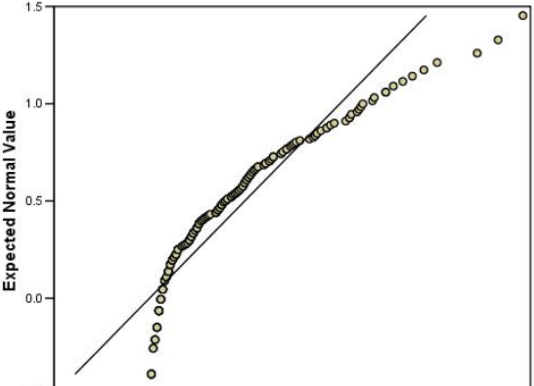

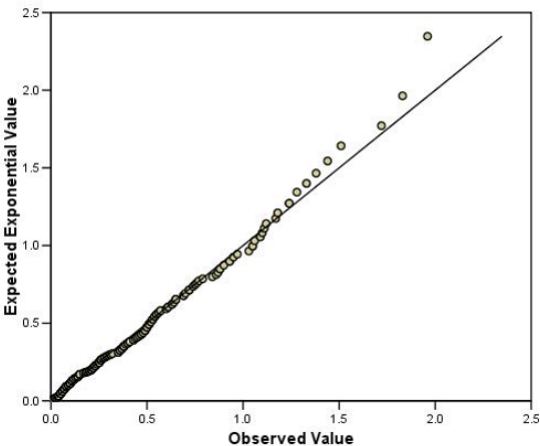


## 2.3 Guessing the Distribution

Sunday, 09 October 2022 12:01

Summary	<ul style="list-style-type: none"><li>• Q-Q plot</li><li>• P-P plot</li><li>• Chi-Square test</li></ul>
<ul style="list-style-type: none"><li>• What are the two probability plots we use to check goodness-of-fit?</li><li>• How do you interpret Q-Q plot?</li><li>• What points to we plot in Q-Q plot?</li></ul>	<h3>Probability plots</h3> <p><b>Q-Q plot: Quantile-quantile plot</b></p> <ul style="list-style-type: none"><li>• Graph of the <math>q_i</math>-quantile of a fitted (model) distribution versus the <math>q_i</math>-quantile of the sample distribution.</li></ul> $x_{q_i}^M = \hat{F}^{-1}(q_i)$ $x_{q_i}^S = \tilde{F}_n^{-1}(q_i) = X_{(i)}, i = 1, 2, \dots, n.$ <ul style="list-style-type: none"><li>• If <math>F^\wedge(x)</math> is the correct distribution that is fitted, for a large sample size, then <math>F^\wedge(x)</math> and <math>F_n(x)</math> will be close together and the Q-Q plot will be approximately linear with intercept 0 and slope 1.</li><li>• For small sample, even if <math>F^\wedge(x)</math> is the correct distribution, there will some departure from the straight line.</li></ul> <ul style="list-style-type: none"><li>• <math>F^\wedge</math> :: indicates the distribution that we're trying to fit</li><li>• <math>F\sim_n</math> :: indicates the distribution that come from the sample</li></ul> <ul style="list-style-type: none"><li>• If <math>x^M_{qi}</math> that comes from the fitted/model distribution matches with the <math>x^S_{qi}</math> that comes from the sample distribution, then you're going to get a line.</li></ul> <div></div> <ul style="list-style-type: none"><li>• You'll plot all the <math>x</math> that comes from <math>X^M</math>, that is model distribution on <math>x</math>-axis, and plot that against the <math>x</math> that comes from sample distribution, <math>X^S</math>.</li><li>• If <math>x_i</math> from model distribution matches with the sample distribution, you'll get a nice <math>45^\circ</math> line in Q-Q plot.</li><li>• This line will have an intercept of 0 and a slope of 1.</li></ul> <ul style="list-style-type: none"><li>• Practically we'll get a line that is around this <math>45^\circ</math> line. And how far away we are from this <math>45^\circ</math> line determines how good is the model distribution.</li><li>• If very far from this line, then the model distribution doesn't match the sample distribution.</li></ul>
<ul style="list-style-type: none"><li>• How do you interpret P-P plot?</li><li>• What points do we plot in P-P plot?</li></ul>	<h3>Probability plots</h3> <ul style="list-style-type: none"><li>• <b>P-P plot:</b> Probability-Probability plot.</li></ul> <p>A graph of the model probability <math>\hat{F}(X_{(i)})</math> against the sample probability <math>\tilde{F}_n(X_{(i)}) = q_i, i = 1, 2, \dots, n.</math></p> <ul style="list-style-type: none"><li>• It is valid for both continuous as well as discrete data sets.</li><li>• If <math>F^\wedge(x)</math> is the correct distribution that is fitted, for a large sample size, then <math>F^\wedge(x)</math> and <math>F_n(x)</math> will be close together and the P-P plot will be approximately linear with intercept 0 and slope 1.</li></ul> <ul style="list-style-type: none"><li>• Here, we compare probability distributions: <math>F^\wedge</math> and <math>F\sim</math> (CDF).</li></ul>
<ul style="list-style-type: none"><li>• Q-Q plot amplifies the differences between the _____.</li><li>• P-P plot amplifies the differences between the _____.</li></ul>	<h3>Probability plots</h3> <ul style="list-style-type: none"><li>• The <b>Q-Q</b> plot will amplify the <b>differences between the tails</b> of the model distribution and the sample distribution.</li><li>• Whereas, the <b>P-P</b> plot will amplify the <b>differences at the middle portion</b> of the model and sample distribution.</li></ul>
<ul style="list-style-type: none"><li>• In general, for both probability plots, what do we plot on the:<ul style="list-style-type: none"><li>1. x-axis?</li><li>2. y-axis?</li></ul></li></ul>	<h3>Probability plots: Dataset</h3> <div></div>

	<div><p>Observed Cum Prob</p></div> <div><ul style="list-style-type: none"><li>On the observed data Var1, we've tried to fit Normal distribution using P-P Plot.</li><li>X-axis :: Observed cumulative frequency</li><li>Y-axis :: Expected(model) cumulative frequency</li><li>So, both x- and y-axis will go from 0 to 1.</li><li>Here, observed points don't seem to be close to expected points.</li></ul></div>
	<div>Probability plots: Dataset</div> <div><p>Deviation from Normal</p><p>Observed Cum Prob</p></div> <div><ul style="list-style-type: none"><li>Here is shown the deviation of observed points from Normal.</li><li>The deviation seems to be high in P-P plot, of the order of 0.15.</li></ul></div>
	<div>Probability plots: Dataset</div> <div><p>Exponential P-P Plot of VAR00001</p><p>Expected Cum Prob</p><p>Observed Cum Prob</p></div> <div><ul style="list-style-type: none"><li>Here, we're trying to fit exponential distribution in the dataset using P-P plot.</li><li>The observed points seem to be very close to the 45° line.</li></ul></div>
	<div>Probability plots: Dataset</div> <div><p>Deviation from Exponential</p><p>Observed Cum Prob</p></div> <div><ul style="list-style-type: none"><li>Here, the deviation of the dataset points from the exponential distribution is shown.</li><li>Notice the scale of Y-axis, it's of the order of 0.04, which is small.</li></ul></div>
	<div><ul style="list-style-type: none"><li>Conclusion: In P-P plots, the exponential distribution seem to be a better fit than the normal distribution.</li></ul></div>
	<div>Probability plots: Dataset</div> <div><p>Normal Q-Q Plot of VAR00001</p><p>Expected Normal Value</p></div>

	<div>  </div> <ul style="list-style-type: none"> <li>In this Q-Q plot, we're trying to fit normal distribution to the dataset Var1.</li> <li>X-axis :: X-values from the sample Y-axis :: Y-values from the model</li> <li>The observed points seem to be deviating from the 45° line.</li> </ul>
	<div> <h3>Probability plots: Dataset</h3> <div>  </div> <ul style="list-style-type: none"> <li>Q-Q plot with the exponential distribution.</li> <li>Seems very close to the 45° line. There are some deviations in the upper portion, these deviations are for higher observed values.</li> </ul> </div>
	<ul style="list-style-type: none"> <li>Conclusion: In Q-Q plots as well, the exponential distribution seem to be a better fit than the normal distribution.</li> </ul>
	<div> <p>Q. Do we now conclude that the exponential distribution is a good fit for the data?</p> <p>A. No. We also need to look at statistical goodness-of-fit tests.</p> </div>
<ul style="list-style-type: none"> <li>Name the two famous statistical goodness-of-fit tests.</li> </ul>	<div> <h3>Goodness-of-fit tests</h3> <ul style="list-style-type: none"> <li>A goodness-of-fit test is a <b>statistical hypothesis test</b> that is used to assess formally whether the observations <math>X_1, X_2, X_3...X_n</math> are an independent sample from a particular distribution with function <math>F^\wedge</math>.</li> </ul> <p><math>H_0</math>: The <math>X_i</math>'s are IID random variables with distribution function <math>F^\wedge</math>.</p> <ul style="list-style-type: none"> <li>Two famous tests:</li> </ul> <ol style="list-style-type: none"> <li>Chi-square test</li> <li>Kolmogorov - Smirnov test</li> </ol> </div>
<ul style="list-style-type: none"> <li>How do you calculate chi-square test? <ul style="list-style-type: none"> <li>Look <a href="#">here</a>.</li> </ul> </li> </ul>	<div> <h3>Chi-square test</h3> <ul style="list-style-type: none"> <li><b>Applicable for both</b>, continuous as well as discrete, distributions.</li> <li>Method of calculating chi-square test statistic:</li> </ul> <ol style="list-style-type: none"> <li>Divide the entire range of fitted distribution into k adjacent intervals -- <math>[a_0, a_1), [a_1, a_2), ... [a_{k-1}, a_k)</math>, where it could that <math>a_0 = -\infty</math> in which case the first interval is <math>(-\infty, a_1)</math> and/or <math>a_k = \infty</math>.</li> </ol> <p><math>N_j</math> = # of <math>X_i</math>'s in the <math>j</math>th interval <math>[a_{j-1}, a_j), j= 1, 2...n</math>.</p> <ol style="list-style-type: none"> <li>Next, we compute the expected proportion of <math>X_i</math>'s that would fall in the <math>j</math>th interval if we were sampling from fitted distribution</li> </ol> </div>
	<div> <h3>Chi-square test</h3> <div> <p>For continuous distributi ons : <math>p_j = \int_{a_{j-1}}^{a_j} \hat{f}(x)dx</math></p> <p>For discrete distributi ons : <math>p_j = \sum_{a_{j-1} \leq x_j &lt; a_j} \hat{p}(x_j).</math></p> </div> <ul style="list-style-type: none"> <li>Finally the test statistic is calculated as:</li> </ul> <math display="block">\chi^2 = \sum_{j=1}^k \frac{(N_j - np_j)^2}{np_j}.</math> </div>

	<ul style="list-style-type: none"> <li>Better formula is given <a href="#">here</a>.</li> </ul>
<ul style="list-style-type: none"> <li>Based on chi-square test statistic when do we: <ul style="list-style-type: none"> <li>Accept <math>H_0</math>?</li> <li>Reject <math>H_0</math>?</li> </ul> </li> </ul>	<div>Chi-square test</div> <ul style="list-style-type: none"> <li>This calculated value of the test statistic is compared with the tabulated value of chi-square distribution with <math>k-1</math> df at <math>1-\alpha</math> level of significance.</li> </ul> <div> <math display="block">\text{If } \chi^2 &gt; \chi^2_{k-1, 1-\alpha} \text{ Reject } H_0</math> <math display="block">\text{If } \chi^2 \leq \chi^2_{k-1, 1-\alpha} \text{ Do not Reject } H_0</math> </div>
	<ul style="list-style-type: none"> <li>The data given to us was a time data. Time to get a service in a bank.</li> <li>The exponential distribution has a strong association with the <u>queuing theory</u>.</li> </ul>