

Identifying and fitting distributions

Sunday, 09 October 2022 22:49

<ul style="list-style-type: none"><li>What's the approach we take for identifying and fitting distributions?</li></ul>	<p>Task: Identifying Distributions and fitting distributions</p> <p>Approach:</p> <ol style="list-style-type: none"><li>Look at visualizations</li><li>Use descriptive statistics to further identify distributions</li><li>Go for tests to confirm judgement</li></ol>
	Detailed Approach
What are the steps in the analysis?	<ol style="list-style-type: none"><li>Plot Histogram<ul style="list-style-type: none"><li>Whenever you want to find a distribution, the first thing you do is plot Histogram.</li><li>When you plot a histogram, you can look at it and tell what family of distributions a dataset may belong to.</li></ul></li><li>Do some descriptive Analysis</li><li>Plot Probability Plots to check goodness of fit<ol style="list-style-type: none"><li>P-P Plot</li><li>Q-Q Plot</li></ol></li><li>Apply statistical tests to check goodness of fit<ol style="list-style-type: none"><li>Chi square test<ul style="list-style-type: none"><li>For this test, we usually need two kind of frequencies:<ul style="list-style-type: none"><li>Observed frequency: sample frequency</li><li>Expected frequency: population frequency, the model distribution where we assume the data is coming from.</li></ul></li><li>Output: the chi-square statistic and the <math>p - value</math></li><li>Compare the <math>p - value</math> and <math>\alpha</math> (significance level) and, conclude about the null hypothesis.</li></ul></li><li>Calculated and tabular chi-square statistic comparison<ul style="list-style-type: none"><li>Input: confidence level, degrees of freedom (<math>df</math>)</li><li>Output: Tabular chi-square statistic</li><li>Compare tabular and calculated chi-square values and, conclude about the null hypothesis.</li></ul></li></ol></li></ol>
Defining Hypotheses	<p>Null Hypothesis: Sample distribution follows model distribution.</p> <p>Alternate Hypothesis: Sample distribution does not follow model distribution.</p>
Chi Square Test <ul style="list-style-type: none"><li>How to calculate expected frequency?</li></ul> <pre>stats.chisquare(obs_freq, expec_freq)</pre> <ul style="list-style-type: none"><li>What does this command return?</li><li>Chi-square formula?</li><li>What is p-value?</li><li>What is <math>\alpha</math>?</li><li>Comparing both these values, when do we accept the null hypothesis?</li></ul>	<pre>stats.chisquare(obs_freq, expec_freq)</pre> <p>This command returns two values:</p> <ol style="list-style-type: none"><li>Chi-Square Statistic</li><li>P-value</li></ol> <ul style="list-style-type: none"><li>P-value is the probability of observing this particular sample when the null hypothesis is assumed to be true.</li><li><math>\alpha</math> :: Significance level, that can take a value of 0.05, or 0.10, or 0.15.</li><li>If <math>P\text{-value} &gt; \alpha</math> :: We accept the null hypothesis.</li></ul>
<ul style="list-style-type: none"><li>Tabulated Chi-Square Statistic</li></ul> <pre>scipy.stats.chi2.ppf(0.95, df=9)</pre> <ul style="list-style-type: none"><li>What arguments does this command take?</li><li>How to calculate <math>df</math>?<ul style="list-style-type: none"><li>How do we calculate <math>df</math> in a contingency table?</li></ul></li></ul>	<pre>scipy.stats.chi2.ppf(0.95, df=9)</pre> <p>Above command gives us the tabulated chi-square statistic, that takes two arguments:</p> <ol style="list-style-type: none"><li>First argument is confidence level = 95% we have taken</li><li>Second argument is degrees of freedom, <math>df = k - p - 1</math> where,<ol style="list-style-type: none"><li><math>k</math> = number of classes/intervals/buckets</li><li><math>p</math> = number of parameters we estimate from the sample</li></ol></li></ol>
<ul style="list-style-type: none"><li>Comparing calculated and tabulated chi-square statistic, when do we accept the null hypothesis?</li></ul>	<p>If <math>\text{Tabulated value} \geq \text{Calculated value}</math></p> <p>» We accept the null hypothesis</p>
Some other observations:	<ul style="list-style-type: none"><li>In Python, whenever we run a code for any of the plots (Q-Q or P-P), by default, it's always comparing it with the <b>standard normal distribution</b>.</li><li>It converts the data that is fed to it to a standard normal, basically by doing <math>\frac{x - \mu}{\sigma}</math>.</li><li>It will scale it and then it will compare it with the standard normal(default), and see how the quantiles (or distributions) are fitting.</li><li>So, it normalizes the data, and then compare it with the standard normal.</li></ul>
	<ul style="list-style-type: none"><li>Poisson distribution has mean = variance (= <math>\lambda</math>)</li></ul>
	<ul style="list-style-type: none"><li>Generally for normal distributions the Kurtosis would be around 3.</li></ul>
	<ul style="list-style-type: none"><li>Degree of freedom = <math>k - p - 1</math></li><li><math>p = 0</math> for uniform distribution as there is no parameter.</li><li>This point is doubtful. Needs verification</li></ul>
	<input type="checkbox"/> Question: What range decides if the skewness is very small or very large to consider?
