

Summary	<ul style="list-style-type: none">Interpreting Descriptive StatisticsGoodness-of-fit																																																
	<h3>Business example</h3> <ul style="list-style-type: none">Data points: 217.For these data points, we need to fit a probability distribution.																																																
<ul style="list-style-type: none">What is the relation between mean, mode and median for symmetric distributions?Try to interpret the given statistics.How do you interpret skewness?<ul style="list-style-type: none">If skewness is positive, the data is skewed to the _____?<ul style="list-style-type: none">How will its shape look like?	<h3>Summary statistics</h3> <table border="1"><thead><tr><th colspan="3">VAR00001</th></tr><tr><th>N</th><th>Valid</th><th>217</th></tr><tr><th></th><th>Missing</th><th>1</th></tr></thead><tbody><tr><td>Mean</td><td></td><td>.4012</td></tr><tr><td>Median</td><td></td><td>.2800</td></tr><tr><td>Mode</td><td></td><td>.05^a</td></tr><tr><td>Std. Deviation</td><td></td><td>.38093</td></tr><tr><td>Variance</td><td></td><td>.145</td></tr><tr><td>Skewness</td><td></td><td>1.466</td></tr><tr><td>Std. Error of Skewness</td><td></td><td>.165</td></tr><tr><td>Range</td><td></td><td>1.95</td></tr><tr><td>Minimum</td><td></td><td>.01</td></tr><tr><td>Maximum</td><td></td><td>1.96</td></tr><tr><td>Percentiles</td><td>25</td><td>.1000</td></tr><tr><td></td><td>50</td><td>.2800</td></tr><tr><td></td><td>75</td><td>.5500</td></tr></tbody></table> <p>a. Multiple modes exist. The smallest value is shown</p> <ul style="list-style-type: none">For symmetric distributions: Mean = Mode = Median (or very close to each other) (eg, normal distribution)But here, mean, mode and median are different. It means it's not a symmetric distribution. So it rules all the symmetric theoretical distributions for us.Looking at the min-max of the data points, it can be observed that none of the data point takes on negative value. Rules out all the distributions that go to the negative side of the line.Skewness tells us about the symmetry of the distribution. The value of skewness is positive, so it's a positive skew, i.e., the data is skewed to the right. Right Tail > Left Tail <div></div> <ul style="list-style-type: none">So consider all the positive skewed distributions as the potential distributions here.	VAR00001			N	Valid	217		Missing	1	Mean		.4012	Median		.2800	Mode		.05 ^a	Std. Deviation		.38093	Variance		.145	Skewness		1.466	Std. Error of Skewness		.165	Range		1.95	Minimum		.01	Maximum		1.96	Percentiles	25	.1000		50	.2800		75	.5500
VAR00001																																																	
N	Valid	217																																															
	Missing	1																																															
Mean		.4012																																															
Median		.2800																																															
Mode		.05 ^a																																															
Std. Deviation		.38093																																															
Variance		.145																																															
Skewness		1.466																																															
Std. Error of Skewness		.165																																															
Range		1.95																																															
Minimum		.01																																															
Maximum		1.96																																															
Percentiles	25	.1000																																															
	50	.2800																																															
	75	.5500																																															
	<h3>Box plot</h3>																																																
	<ul style="list-style-type: none">Below is given the same plot tilted 90 degrees, to get a better picture of the points on x-axis.																																																
	<h3>Box plot</h3>																																																

	<div data-bbox="642 71 1213 356"> </div> <ul style="list-style-type: none"> It's clear from the box plot that this distribution has a positive skew.
	<div data-bbox="514 463 728 513"> <h3>Histograms</h3> </div> <div data-bbox="663 557 1148 1041"> </div> <ul style="list-style-type: none"> Large number of values close to 0. Very few values more than 1.5.
	<div data-bbox="522 1166 730 1213"> <h3>Histograms</h3> </div> <div data-bbox="655 1246 1157 1748"> </div> <ul style="list-style-type: none"> We increased the width of the bars here. Frequency seems to be dropping as you go to the right.
	<div data-bbox="512 1863 724 1914"> <h3>Histograms</h3> </div> <div data-bbox="648 1941 1165 2457"> </div> <ul style="list-style-type: none"> Even thicker bars here.
<ul style="list-style-type: none"> What is Coefficient of variation, cv? <ul style="list-style-type: none"> This statistic is for continuous or discrete distributions? cv = 1 for? cv > 1 for? 	<div data-bbox="518 2543 1008 2588"> <h3>Clues from summary statistics</h3> </div> <ul style="list-style-type: none"> For the symmetric distributions mean and median should match. In the sample data, if these values are sufficiently close to each other, we can think of a symmetric distribution (e.g. normal). Coefficient of variation (cv): (ratio of std dev and the mean) for continuous distributions. The $cv = 1$ for exponential dist. If the histogram looks like a slightly right-skewed curve with $cv > 1$, then lognormal could be better approximation of the distribution.

	<p>Note: For many distributions cv may not even be properly defined. When?</p> <p>Examples?</p> <ul style="list-style-type: none">• <i>Coefficient of variation</i>, $(cv) = \frac{Standard\ Dev}{Mean} = \frac{\sigma}{\mu}$• $cv = 1$:: Exponential Distribution for slightly right-skewed curve with $cv > 1$:: Lognormal Distribution• In some cases, cv may not be defined. Take an example of standard normal distribution where $\mu = 0$.• Here, $cv = \frac{0.38093}{0.4012} = 0.9495 \approx 1$
<ul style="list-style-type: none">• Lexis ratio• Skewness (v)<ul style="list-style-type: none">• $v = 0$ for?• $v > 0$?<ul style="list-style-type: none">◦ $v = 2$ for?• $v < 0$?	<h3>Clues from summary statistics</h3> <ul style="list-style-type: none">• Lexis ratio: same as cv for discrete distributions.• Skewness (v): measure of symmetry of a distribution. For normal dist. $v = 0$. For $v > 0$, the distribution is skewed towards right (exponential dist, $v = 2$). And for $v < 0$, the distribution is skewed towards left.• If Skewness,<ul style="list-style-type: none">◦ $v = 0$:: Normal distribution◦ $v > 0$:: Distribution is skewed towards right<ul style="list-style-type: none">• $v = 2$:: Exponential Distribution◦ $v < 0$:: Distribution is skewed towards left• We'll try to fit exponential distribution here.
	<h3>Parameter estimation</h3> <ul style="list-style-type: none">• Once distribution is guessed, the next step is estimating the parameters of the distribution.• Each distribution has a set of parameters.<ul style="list-style-type: none">✓ Normal distribution has mean and standard deviation✓ Exponential distribution has a “λ”.• Most common method of parameter estimation: MLE (What is this?)
<ul style="list-style-type: none">• How can we check the goodness-of-fit of the fitted distribution?<ul style="list-style-type: none">• Name the methods• Explain	<h3>Goodness-of-fit</h3> <ul style="list-style-type: none">• For the input data we have, we have assumed a probability distribution.• We also have estimated the parameters for the same.• How do we know this fitted distribution is “good enough?”• It can be checked by several methods:<ol style="list-style-type: none">1. Frequency comparison (a bit technical)2. Probability plots (visual tool)3. Goodness-of-fit tests (statistical test of goodness. Very widely used).• We're trying to fit exponential distribution here.1. In frequency comparison, we can compare the frequency that comes from exponential distribution with the frequency you have observed from the dataset. If the frequencies match, we say that exponential distribution is good fit.2. Probability plots are visual tools that tell you if the observed frequency/quartiles/percentiles matches with the frequency/quartiles/percentiles of the distribution. If it fits, you'll get a line , if it doesn't fit, you'll be far away from that line.3. Goodness of fit tests: Many of them uses Chi-Square tests.

