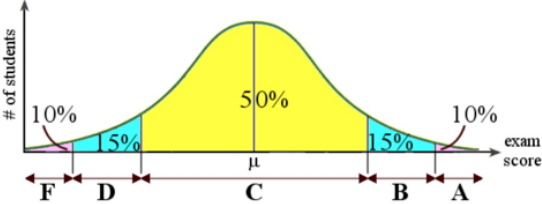


Summary	<ul style="list-style-type: none"><li>3 ways to fit distributions to data<ol style="list-style-type: none"><li>Trace-driven simulation</li><li>Theoretical distribution</li><li>Empirical distribution</li></ol></li></ul>
	<h2>Probability distributions</h2> <p>What are distributions? How to identify correct distribution of the data?</p>
	<h2>Probability distributions</h2> <ul style="list-style-type: none"><li>A statistical model that shows possible outcomes of a particular event or course of action as well as the statistical likelihood of each event.</li><li>What do we mean by “Grades in a course follows a normal distribution”?</li><li>What do we mean by “Sales for the next month may be uniformly distributed”?</li></ul>  <ul style="list-style-type: none"><li>How do you fit distributions to data?</li></ul>
<ul style="list-style-type: none"><li>In how many ways can we create models once we have the data?<ul style="list-style-type: none"><li>Name them.</li><li>Explain.</li></ul></li></ul>	<h2>How to go about this?</h2> <p>How do we use the collected business data (sales volume, loan defaulters, Salary hikes in an organization, etc.)?</p> <ol style="list-style-type: none"><li>The data values themselves are used directly in the simulation. This is called <b>trace-driven simulation</b>.</li><li>“Fit” a theoretical distribution to the data (and check whether that “fit” is good!).</li><li>The data values could be used to define an <b>empirical distribution</b> function in some way.</li></ol> <ul style="list-style-type: none"><li>In point 1, we don't need to fit any distribution. We directly use the data in our analysis.</li><li>Theoretical distributions are the distributions we've already studies, eg, normal, uniform distribution, etc.</li><li>In point 2, we first fit a distribution in the data. Check how good is the fit.</li><li>Point 3: <u>If the data doesn't fit any theoretical distribution, then instead of trying to fit already available distributions, we create our own distributions.</u> Those distributions are called <b>empirical distributions</b>. And then we use this distribution in our future analysis.</li></ul>
<ul style="list-style-type: none"><li>What are the essential building blocks of empirical distribution?</li></ul>	<h2>What are these empirical distributions?</h2> <ul style="list-style-type: none"><li>Using the data, we build our own distributions.</li><li>How does one build a distribution?</li></ul> <ul style="list-style-type: none"><li>Essential building blocks: Define the density/distribution functions. Estimate the parameters (mean, standard deviation, etc.)</li></ul> <ul style="list-style-type: none"><li>When the professor says distribution function, he means CDF.</li></ul>
	<h2>Empirical distributions</h2> <p>For ungrouped data: Let <math>X_{(i)}</math> denote the <math>i</math>th smallest of the <math>X_j</math>'s so that: <math>X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}</math>.</p>

	$F(x) = \begin{cases} 0 & \text{if } x < X_{(1)} \\ \frac{i-1}{n-1} + \frac{x - X_{(i)}}{(n-1)(X_{(i+1)} - X_{(i)})} & \text{if } X_{(i)} \leq x < X_{(i+1)} \text{ for } i = 1, 2, \dots, n-1 \\ 1 & \text{if } X_{(n)} \leq x \end{cases}$ <ul style="list-style-type: none"> <li>Given a CDF, we already know how to find the density function.</li> </ul>
	<h2>Empirical distributions</h2> <p>For grouped data:</p> <ul style="list-style-type: none"> <li>Suppose that <math>n</math> <math>X_j</math>'s are grouped in <math>k</math> adjacent intervals <math>[a_0, a_1), [a_1, a_2), \dots, [a_{k-1}, a_k)</math> so that <math>j</math>th interval contains <math>n_j</math> observations. <math>n_1 + n_2 + \dots + n_k = n</math>.</li> <li>Let a piecewise linear function <math>G</math> be such that <math>G(a_0) = 0</math>, <math>G(a_j) = (n_1 + n_2 + \dots + n_j) / n</math>, then: <math display="block">G(x) = \begin{cases} 0 &amp; \text{if } x &lt; a_0 \\ G(a_{j-1}) + \frac{x - a_{j-1}}{a_j - a_{j-1}} [G(a_j) - G(a_{j-1})] &amp; \text{if } a_{j-1} \leq x &lt; a_j, j = 1, 2, \dots, k \\ 1 &amp; \text{if } a_k \leq x. \end{cases}</math> </li> <li>We have <math>k</math> intervals, and in each interval we have <math>n_1, n_2, \dots, n_k</math> values.</li> <li><math>G(a_j)</math> is proportional to the observations up to that point/interval.</li> </ul>
<ul style="list-style-type: none"> <li>When do we use trace driven simulation? <ul style="list-style-type: none"> <li>What are its drawbacks?</li> </ul> </li> </ul> <p><input type="checkbox"/> Approach 1 is not very clear to me</p>	<h2>The three approaches...</h2> <ul style="list-style-type: none"> <li>Approach 1 is used to validate simulation model when comparing model output for an existing system with the corresponding output for the system itself.</li> <li><b>Two drawbacks of approach 1:</b> simulation can only reproduce only what happened historically; and there is seldom enough data to make all simulation runs.</li> <li>Approaches 2 and 3 avoid these shortcomings so that any value between minimum and maximum can be generated. So <b>approaches 2 and 3 are preferred over approach 1</b>.</li> <li>If theoretical distributions can be found that fits the observed data (approach 2), then <b>it is preferred over approach 3</b>.</li> <li>Approach 1: You have the output, and you want to validate if the output is correct or not. You push the already available data into your model and your model generates an output and you compare that output with the reality (the existing system, what happens in future) and check whether it matches.</li> <li>So trace-driven simulation is used to validate a model that you already may have built using some approach.</li> <li>Problem with approach 1 is you're going to <b>test the model only with the data you already have. This may not be enough.</b></li> <li>What if the data was collected in a certain circumstance, and as the circumstances change will not give you a fair values.</li> </ul>
<ul style="list-style-type: none"> <li>What are the drawbacks of empirical approach?</li> </ul>	<h2>Approach 3 v/s Approach 2</h2> <ul style="list-style-type: none"> <li>Empirical distribution may have some irregularities if small number of data points are available. Approach 2 smoothens out the data and may provide information on the overall underlying distribution.</li> <li>In <b>approach 3</b>, it is usually <b>not possible to generate values outside the range of observed data</b> in the simulation.</li> <li>If one wants to test the performance of the simulated system under extreme conditions, that can not be done using approach 3.</li> <li>There may be compelling (physical) reasons in some situations for using a particular theoretical distribution. In that case too, it is better to get empirical support for that distribution from the observed data.</li> </ul>

