Individual Group Project 3

Jesse Cook

Class: CS 525 - Principles of Simulation Professor: Dr. Mayer

April 25, 2010



The data that I was tasked to model for our project was the percentage of the total bytes that were downloaded from the server for a given connection (HTTP and HTTPS are modeled separately). In order to determine an appropriate model, I first looked at the physical properties of several models (see Table 1). Based solely on these properties the Beta distribution would be the best fit for the data as it is both bounded and continuous and not overly simple like the triangular model. Next, I looked at the Chi-Squared, Anderson-Darling, and Kolmogorov-Smirnov statistics as well as each distribution's fit and quantile-quantile plots (see Tables 2 and 3). The Fit and Q-Q-Fit columns in the Fit Metrics tables (see page 3) use the values of Good, OK, and Poor. These values indicate how well I felt the distributions fit the data based on the graphs. This data supports the selection of the Beta distribution for the HTTPS data. However, none of the models proved to be a very good fit for the HTTP data. Because of this I split the input data into two ranges ([0.00001 – 0.74999] and [0.75000 – 0.99999]) and generated all of the statistics and graphs again for the bounded distributions (Beta, Uniform, and Triangular) over these ranges. In almost every case the use of two Beta distributions produced better test statistics and graphs than any other model. Thus, the Beta distribution will be used to model each data set using the following parameters:

HTTPS:

$$\beta_1 = 4.5809, \beta_2 = 1.9556$$

HTTP:

$$\beta_1 = \left\{ \begin{array}{ll} 14.015 & \quad 0.00001 \leq x \leq 0.74999 \\ 1.7337 & \quad 0.75000 \leq x \leq 0.99999 \end{array} \right., \\ \beta_2 = \left\{ \begin{array}{ll} 4.3029 & \quad 0.00001 \leq x \leq 0.74999 \\ 1.4296 & \quad 0.75000 \leq x \leq 0.99999 \end{array} \right.$$

Distribution	Physical Basis	Continuous	Bounded
Binomial	not modeling trials	no	no
Negative Binomial	not modeling trials	no	no
Poisson	not modeling # of independent events in a	no	no
	fixed period of time		
Normal	tails not evenly distributed; negative values	yes	no
	invalid		
Log-Normal	could be thought of as rate of return	yes	no
Exponential	not modeling time between events	yes	no
Gamma	_	yes	no
Beta	ideal range	yes	yes
Erlang	exponential not a good fit and not sum of	yes	no
	exponential processes		
Weibull	not a "time to" model	yes	no
Uniform	all outcomes not equally likely	both	yes
Triangular	too simple	yes	yes
Empirical	last resort if nothing fits	yes	no
Log-Logistic	rate increases then decreases	yes	no
Rayleigh	not two-dimensional normally distributed	yes	no
	with equal variance		
Inverse Gaussian	_	yes	no
Pearson 5	_	yes	no
Extreme Value	_	yes	no
Logistic	resembles normal distribution	yes	no
Pareto	not large amount owned by the small	yes	no
Pareto 2	not large amount owned by the small	yes	no
Error Function	derived from normal distribution	yes	no

Table 1: Distribution Properties

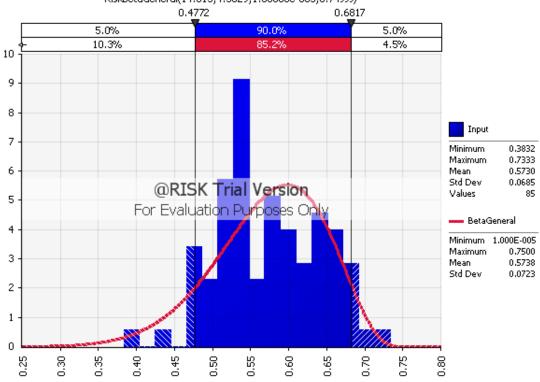
Distribution	Chi ² Statistic	A-D Statistic	K-S Statistic	Fit	Q-Q Fit
Binomial	_	_	_	_	_
Negative Binomial	_	_	_	_	_
Poisson	_	_	_	_	_
Normal	65.78	2.40	0.1268	OK	OK
Log-Normal	_	_	_	_	_
Exponential	120.22	17.60	0.3056	Poor	OK
Gamma	_	_	_	_	_
Beta	60.67	1.87	0.1422	Good	Good
Erlang	_	_	_	_	_
Weibull	_	_	_	_	_
Uniform	110.44	29.48	0.2443	Poor	OK
Triangular	54.00	6.27	0.1735	OK	Good
Empirical	_	_	_	_	_
Log-Logistic	_	_	_	_	_
Rayleigh	60.89	5.63	0.1653	OK	Poor
Inverse Gaussian	_	_	_	_	_
Pearson 5	_	_	_	_	_
Extreme Value	68.22	3.74	0.1700	OK	Poor
Logistic	68.44	2.29	0.1194	OK	OK
Pareto	286.67	$+\infty$	0.3842	Poor	Poor
Pareto 2	118.44	$+\infty$	0.2997	Poor	Poor
Error Function	447.33	83.26	0.6740	Poor	OK

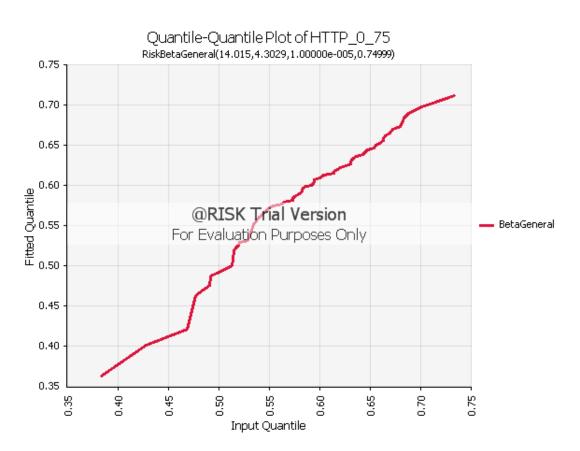
Table 2: HTTPS Fit Metrics

Distribution	Chi ² Statistic	A-D Statistic	K-S Statistic	Fit	Q-Q Fit
Binomial	_	_	_	_	_
Negative Binomial	_	_	_	_	_
Poisson	_	_	_	_	_
Normal	75.21	5.32	0.1433	OK	OK
Log-Normal	61.62	2.95	0.1392	OK	Poor
Exponential	108.11	12.63	0.2635	Poor	Poor
Gamma	_	_	_	_	_
Beta	89.00	6.43	0.1913	OK	OK
Beta $(0 - 75)$	22.54	1.02	0.1204	Good	Good
Beta $[75 - 100)$	15.81	1.31	0.1769	Good	OK
Erlang	_	_	_	_	_
Weibull	70.68	3.44	0.1374	OK	Poor
Uniform	172.12	3.74	0.4545	Poor	OK
Triangular	130.56	6.65	0.2089	OK	OK
Empirical	_	_	_	_	_
Log-Logistic	99.24	2.74	0.1313	OK	Poor
Rayleigh	70.68	3.39	0.1371	OK	Poor
Inverse Gaussian	61.62	2.90	0.1389	OK	Poor
Pearson 5	61.82	3.03	0.1402	OK	Poor
Extreme Value	60.64	3.48	0.1469	OK	Poor
Logistic	86.64	5.09	0.1524	OK	OK
Pareto	149.86	$+\infty$	0.3194	Poor	Poor
Pareto 2	106.53	$+\infty$	0.2586	Poor	Poor
Error Function	537.70	101.95	0.7314	Poor	OK

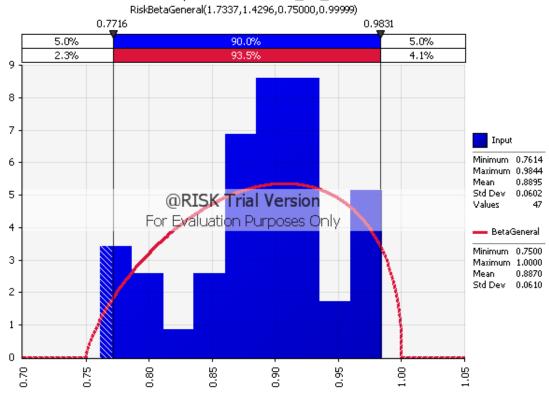
Table 3: HTTP Fit Metrics

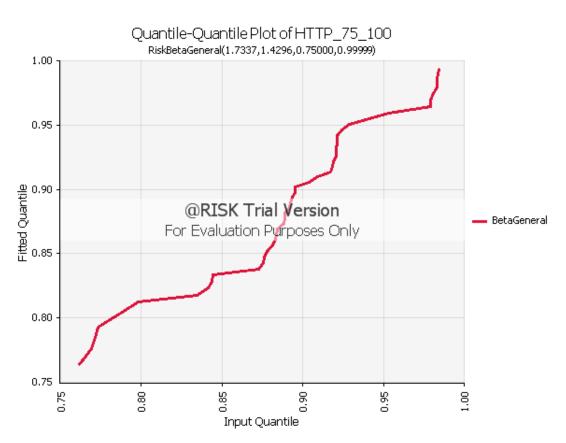
Fit Comparison for HTTP_0_75 RiskBetaGeneral(14.015,4.3029,1.00000e-005,0.74999)





Fit Comparison for HTTP_75_100





Fit Comparison for HTTPS

