



CONSTRUCTING BOUNDS FOR S-CURVES



**Presented by: Christopher Mehl, Ph.D.,
Omnitec Solutions, Inc.**



What's in your P-box?





Introduction

- **Uncertainty/risk analysis in cost estimation was meant to convey impreciseness in an estimate**
- **This has morphed into an absolute probability associated with specific dollar values in the minds of decision makers.**
- **Most recipients of a cost model output don't understand fully what is presented**
- **Instead of a point estimate we now have a point estimate with a Monte Carlo Simulation making it a better point estimate.**
- **A potential tool to help diffuse some of this certainty is the P-Box.**

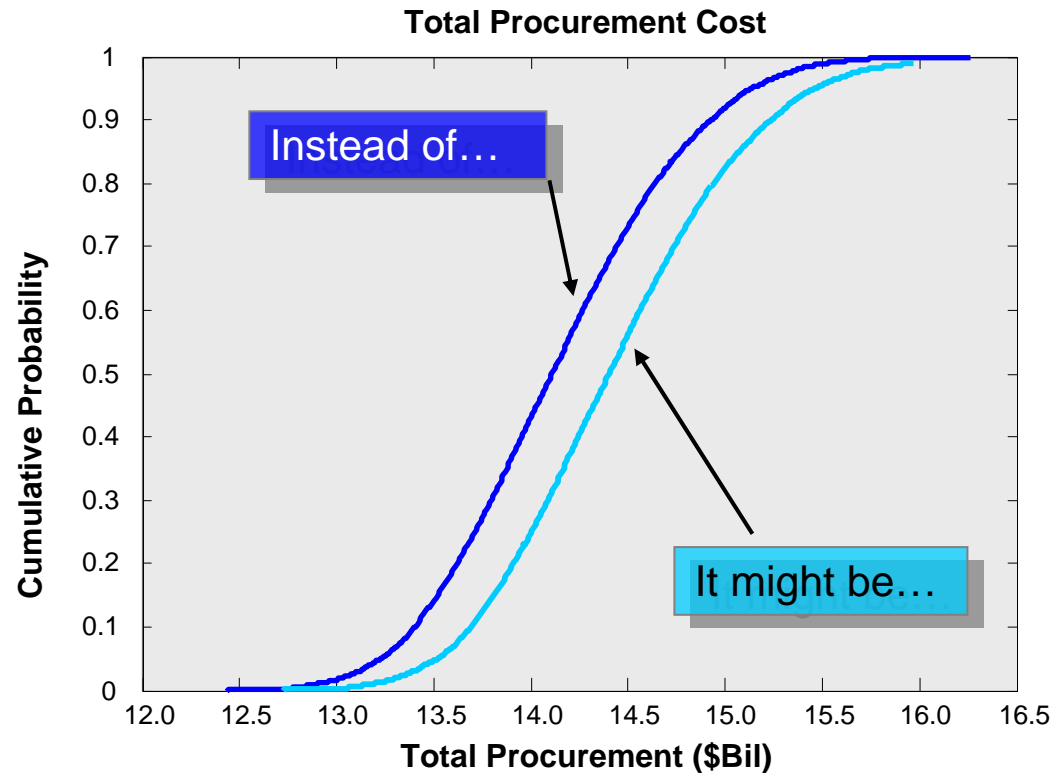




P-Boxes: Outline

- **S-curves capture only one type of uncertainty (*statistical uncertainty*) present in the cost estimate**
- **Also present is *epistemic uncertainty***
- **P-boxes are upper and lower bounds for the S-curve, showing the epistemic uncertainty**
- **Instead of the number of observations (a common parameter in statistical tools) we use the program's age**
- **Steps to construct bounds:**
 - **Kolmogorov-Smirnov bounds**
 - **Quantile bounds based on order statistics**
 - **Combine two bound types by enveloping**
- **Example**





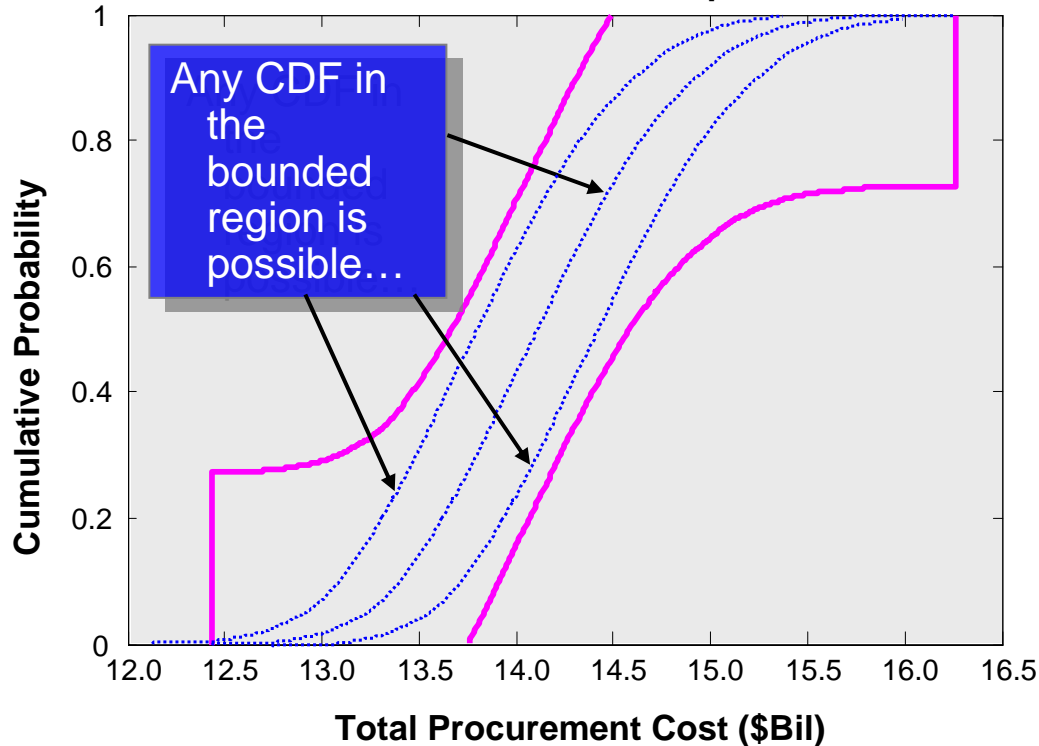
- The S-curve is a cumulative distribution function for an unknown random variable **Cost**, the statistical uncertainty inherent in the cost estimate
- The S-curve fully represents the uncertainty involved **ONLY IF** all sources of uncertainty are:
 - Known
 - Modeled correctly
 - Fixed, unchanging over time
- **BUT...**

- **...NONE** of those assumptions are true
 - Most sources of uncertainty aren't known; there are "unknown unknowns"
 - Those that are known may not be modeled correctly; different analysts choose different distributions for the same cost drivers
 - Requirements, material and labor rates, and other factors change during the life of the program



S-CURVE BOUNDS

Total Procurement, Enveloped Bounds



- By bounding the S-curve using a tool called a *p*-box, we can capture more of the uncertainty
- A *p*-box consists of an *upper* and *lower bound* for the S-curve
- With some specified probability, the “true” S-curve, which we would get if we accurately knew and modeled all sources of uncertainty, lies somewhere between these bounds
- The bounds narrow as the program progresses (and uncertainty decreases)

- Rather than a single curve representing the uncertainty, we get a region that we know contains the true S-curve with a specified probability (90% for example)
- Bounds capture both the fact that the cost is uncertain, and that the processes driving the uncertainty and how it should be modeled are unknown
- The bounds show this epistemic or type 2 uncertainty





S- CURVE AS EDF

EDF

Empirical Distribution Function

- Most often formed from a random sample
- Converges to true CDF as the number of observations increases

- **We treat the S-curve as a kind of *EDF*, an approximation to the “true” cumulative distribution function**
 - Different analysts or information could result in a different S-curve
-- a different approximation to the same “true” CDF
- **Usually, the EDF is formed from a random sample; it converges to the true CDF as the number of observations increases**
 - We use program age in lieu of number of observations
- **We construct bounds from the EDF for a region which contains the true CDF with some probability**



HOW DO WE CONSTRUCT BOUNDS?

Types of P-boxes

- 1) Kolmogorov-Smirnov bounds
- 2) Quantile bounds

- **We construct two types of bounds from the EDF for a region which contains the true CDF with some probability:**
 - Kolmogorov-Smirnov bounds: vertical bounds for the probability (risk) at fixed quantiles (cost)
 - Quantile bounds: horizontal bounds for quantiles (cost) at fixed probabilities (risk)
- **Both types of bounds are non-parametric**
- **Both bound types have number of observations as a parameter**
 - Exploiting the relationship between information, uncertainty, and observations, we derive a proxy value for observations



OBSERVATIONS AND UNCERTAINTY

Statistical Applications

- Each data point is an observation in a random sample

Cost Models: no samples

- Information from disparate sources

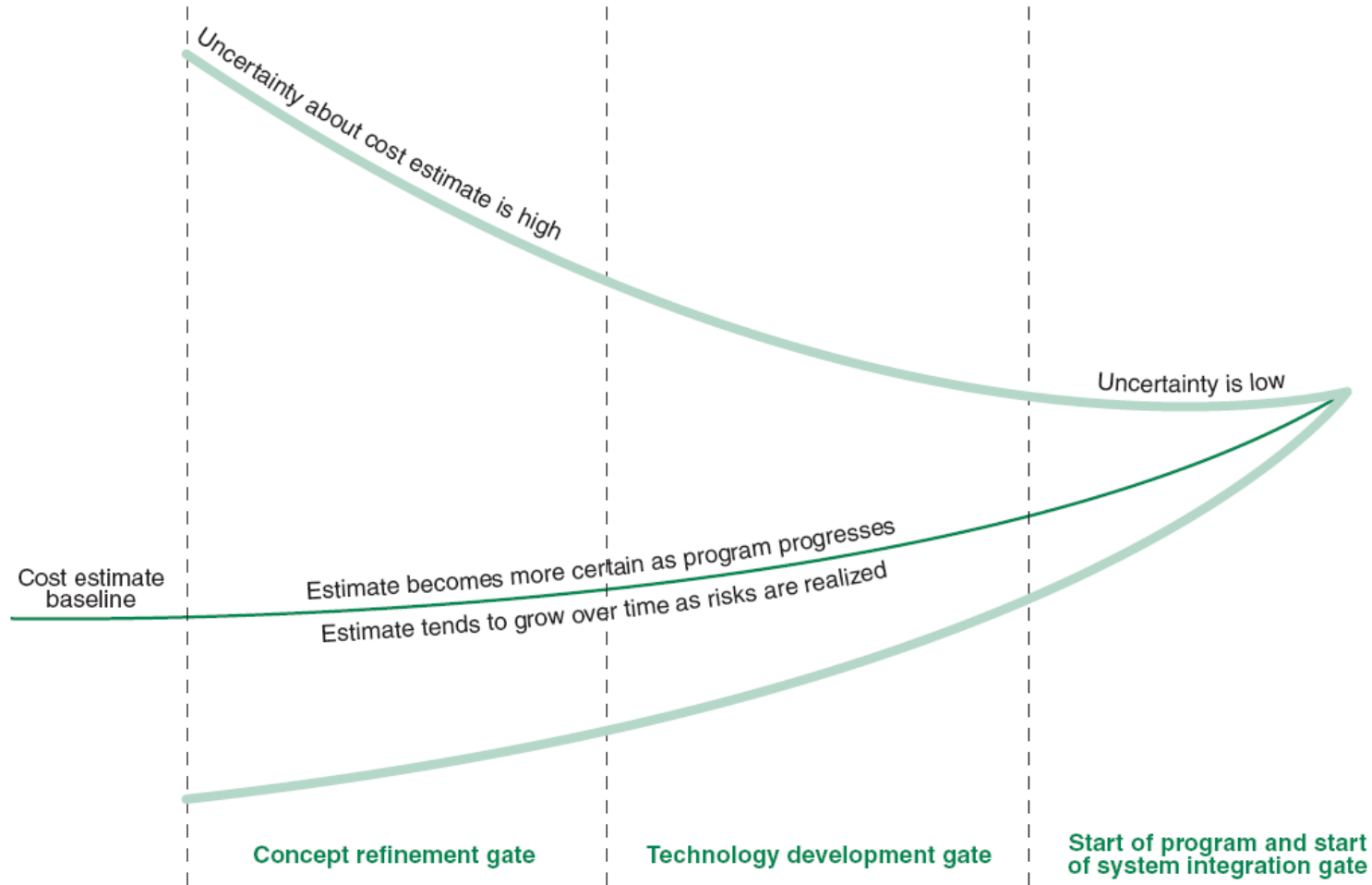
- **In most statistical applications, we have a population or process from which we collect data; each data point is an observation in a random sample**
 - The more observations, the more information
 - The more information we have, the less uncertainty
- **For a cost model, no random sample, instead information from disparate sources, like expert opinion or analogy with past programs**
 - Distributions and associated parameters chosen to model the uncertainty
 - Data becomes available and uncertainty decreases with life of program
- **In lieu of random observations, we use the program's age to determine a proxy value for observations**



Uncertainty Decreases, Information Increases



Figure 4: Cone of Uncertainty



Source: GAO.



FINDING THE EFFECTIVE OBSERVATIONS VALUE



Effective Observations

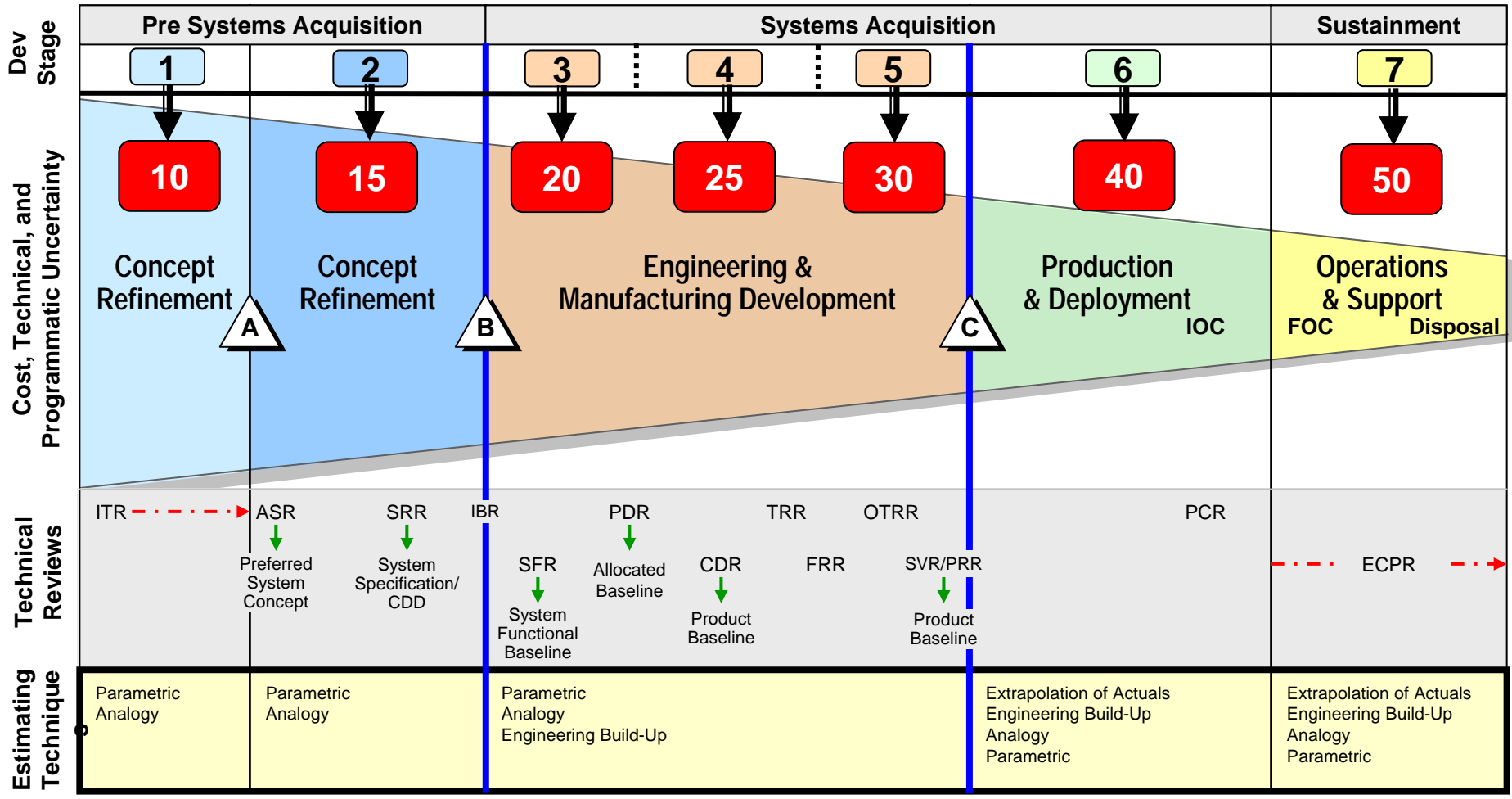
- “Stabilization point” around 25
- Magnitude effect diminished over time
- Baseline of 25 effective observations in stage 4 (near CDR)
- Obs>25 in later stages; Obs<25 in earlier stages

- **Used cost data from several existing programs**
 - Ranging from pre-milestone B to programs in production
 - Program life cycle divided into 7 stages, pinned to technical reviews and decision milestones
- **We calculate the percent range (bound width), ranging effective observations from 10 to 50 to quantify the effect**
- **Analysis revealed following key points about the effect of number of observations:**
 - At all stages of program development, there was a “stabilization point” around 25 beyond which the effect on bounds width is small, relative to effect before this point
 - Magnitude of effect diminished as programs progress; choice of observations less important for older programs
 - Most data for programs in stage 4 (near CDR), giving baseline of 25 observations
 - Programs in later stages have Obs>25, those in earlier stages have Obs<25
- **The resulting equivalence between program life cycle and number of effective observations is shown on the next slide**



PROGRAM UNCERTAINTY VS. TIME

$$O_{eff} = f(\text{DEVELOPMENT STAGE})$$



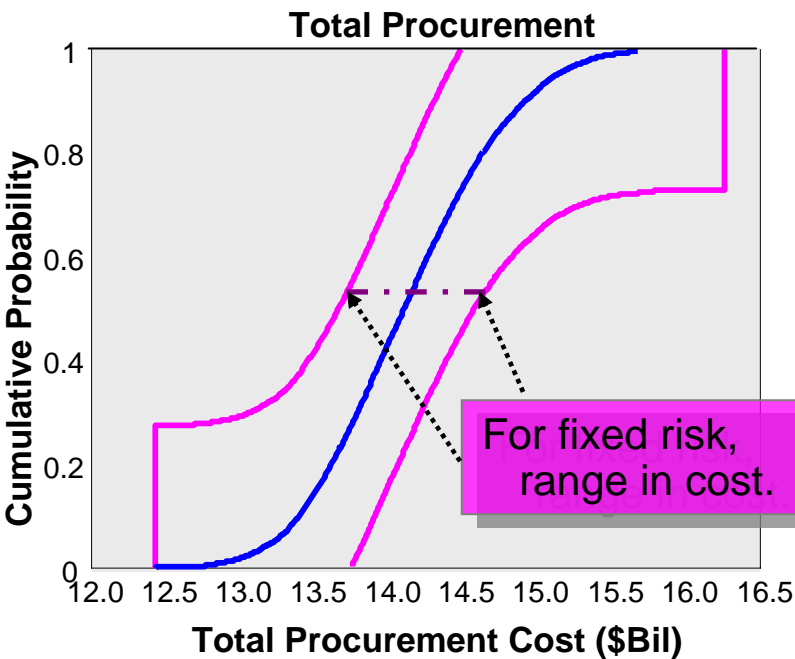
UNCERTAINTY LESSENS AND INFORMATION INCREASES OVER TIME, O_{eff} INCREASES



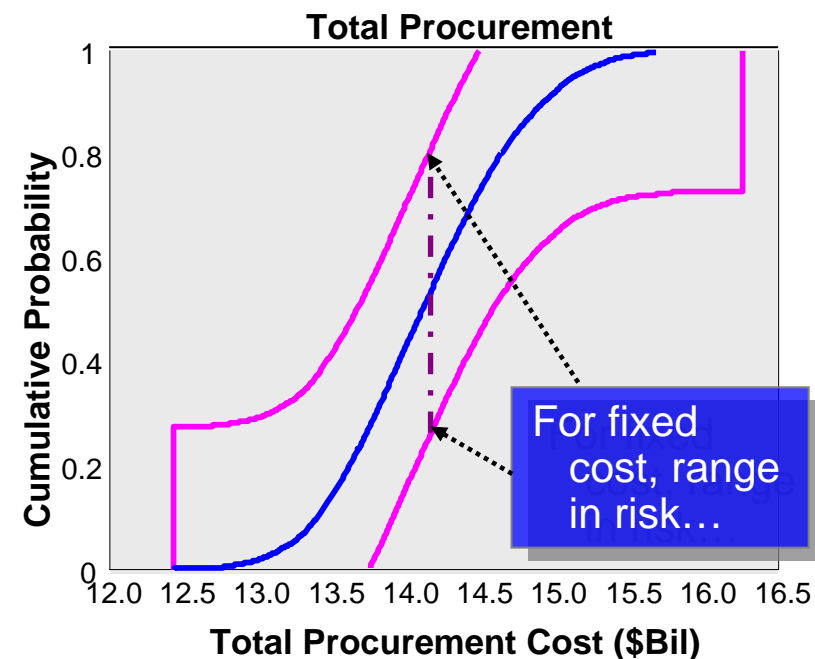
Constructing Bounds for S-Curves



- We construct bounds from the EDF for a region which contains the true CDF with some probability. Two ways of constructing non-parametric bounds:
 - Kolmogorov-Smirnov bounds: vertical bounds for the probability (risk) at fixed quantiles (cost)
 - Quantile bounds: horizontal bounds for quantiles (cost) at fixed probabilities (risk)
- Bound width is a function of a program's stage in it's life cycle.



Quantile Bounds



Kolmogorov-Smirnov Bounds





Kolmogorov-Smirnov P-Box

- Treating the S-curve as an empirical distribution function (EDF) enables use of the EDF's convergence properties
- In particular we use convergence for continuous CDFs $F(x)$ to the Kolmogorov Distribution of:

$$\sqrt{n} \left\| \hat{F}_n(x) - F(x) \right\|_{\infty}$$

- Using the γ variate from the Kolmogorov Distribution we can construct Kolmogorov-Smirnov Bounds:

$$\Pr\left[\max\left(0, \hat{F}_n(x) - \frac{k_{\gamma}}{\sqrt{n}}\right) \leq F(x) \leq \min\left(1, \hat{F}_n(x) + \frac{k_{\gamma}}{\sqrt{n}}\right)\right] \approx \gamma$$





Quantile Interval P-box

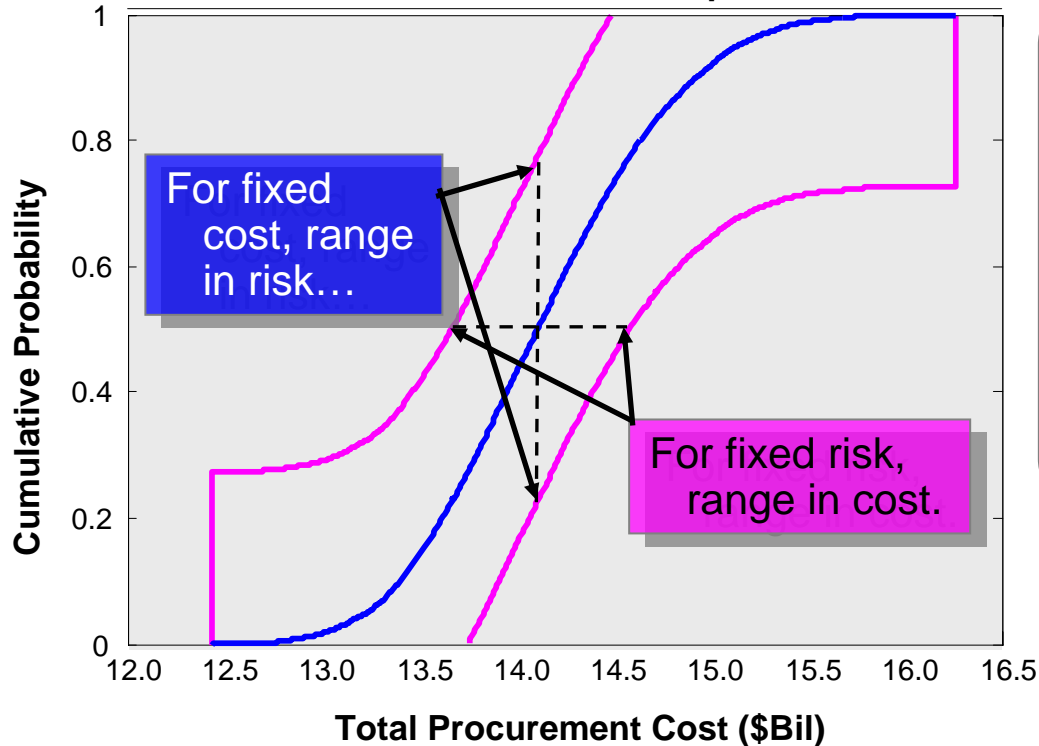
- **Quantile interval based p-boxes are created by dividing the S-curve into O_{eff} equiprobable regions, pulling a cost value from each region to create a pseudo-sample, and using the order statistics to find intervals for each quantile.**
- **For O_{eff} effective observations, the number of values from the pseudo-sample less than the q^{th} quantile will follow a binomial distribution with parameters O_{eff} and q .**



THE BENEFITS OF S-CURVE BOUNDS



Total Procurement, Enveloped Bounds



- Recent acquisition reforms (e.g. ACQ Reform 2009) call for disclosing confidence levels with the cost estimate.
- In particular, there is a focus on determining a “confidence level” for cost estimates of 80%.
- Bounds convey the uncertainty in both cost and risk.

- The bounds help us answer two distinct, key questions about uncertainty:
 - For a fixed dollar value (cost) what is the range in probability (or risk)?
 - For a fixed level of risk (probability) what is the range in cost?
- The bounds and derived ranges show the upside risk to the cost; i.e., how bad do we think it could get?





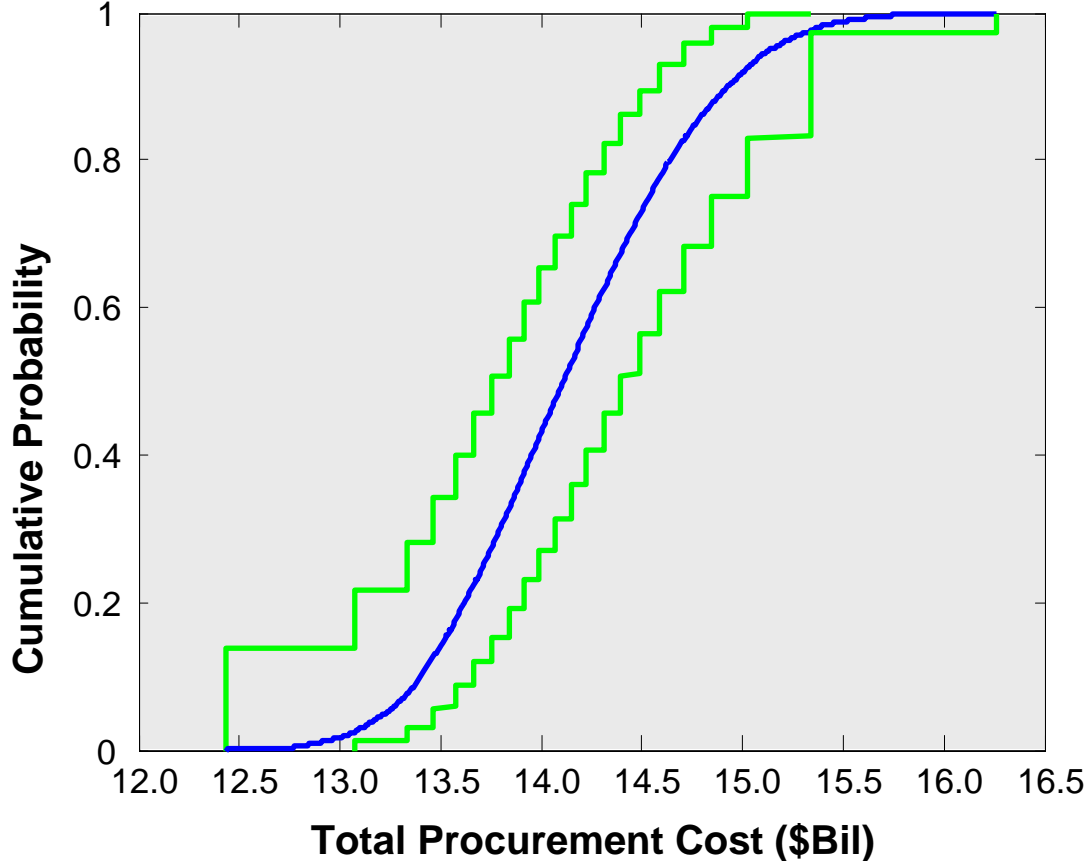
BOUNDS APPLIED TO AN EXAMPLE





QUANTILE BOUND P-BOX TOTAL PROCUREMENT COST (\$TY)

Total Procurement, Quantile Bounds



- Confidence Level = 90%
- Effective Observations = 20

- Cost Ranges (for fixed P, range in cost):

Bound/ Percentile	Lower (\$Bil)	Upper (\$Bil)
20th	13.08	13.91
50th	13.75	14.40
80th	14.31	15.03

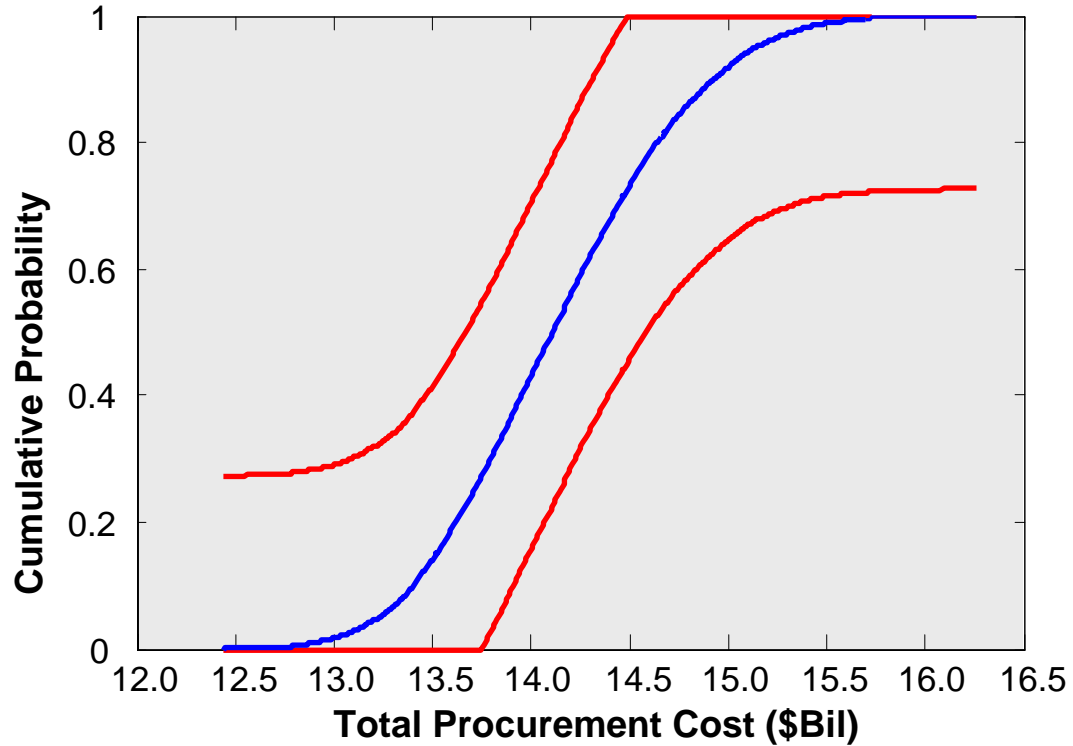
For a fixed level of risk,
what is the range in cost?





KOLMOGOROV-SMIRNOV P-BOX TOTAL PROCUREMENT COST (\$TY)

Total Procurement, Kolmogorov-Smirnov Bounds



- Confidence Level = 90%
- Effective Observations = 20

For a fixed cost, what is the range in risk?

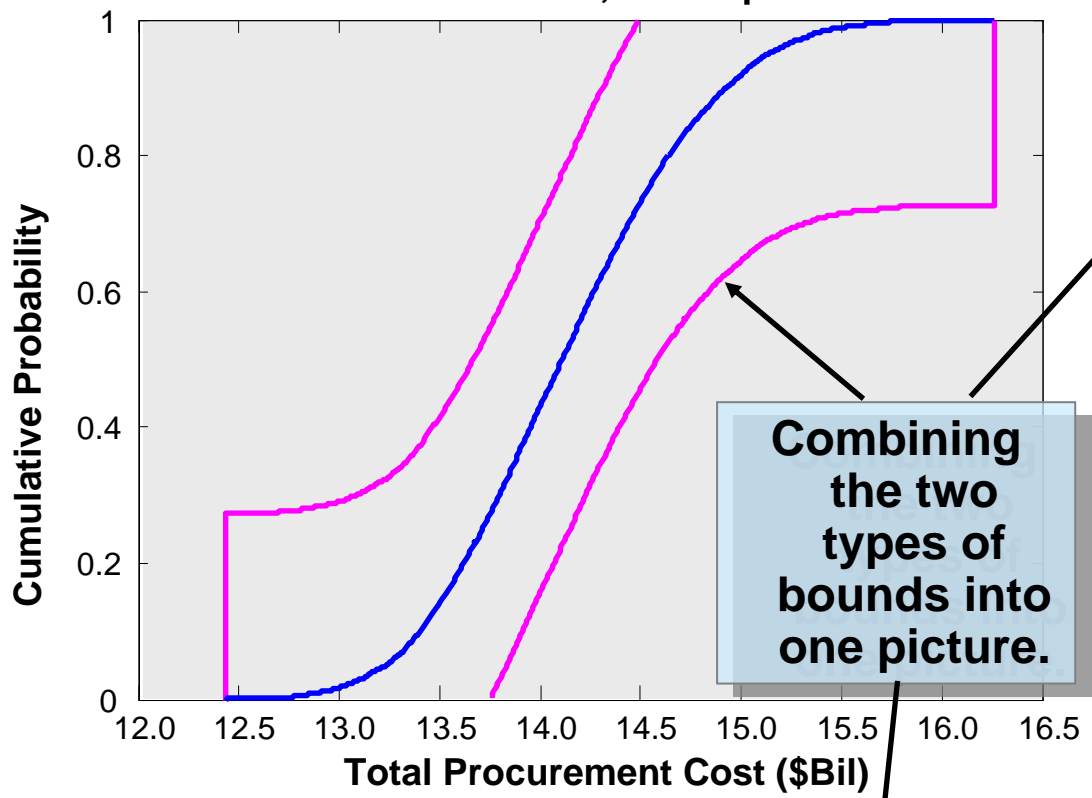
Cost (\$Bill)\Range	Lower	Upper	Risk Percentile
13.62	0%	47.37%	20%
14.11	22.63%	77.37%	50%
14.64	52.63	100%	80%





ENVELOPED P-BOX: COMBINING KS AND QUANTILE

Total Procurement, Enveloped Bounds



- Confidence Level = 90%
- Effective Observations = 20

Bound/ Percentile	Lower (\$Bil)	Upper (\$Bill)
20th	12.43	14.06
50th	13.66	14.58
80th	14.15	16.25

**Combining
the two
types of
bounds into
one picture.**

Range/Cost (\$Bil)	Lower	Upper	Risk Percentile
13.62	0%	47.37%	20%
14.11	22.63%	77.37%	50%
14.64	52.63	100%	80%





Conclusion

- **S-curves alone do not capture all of the uncertainty present in the cost estimate**
- **P-boxes help capture the *epistemic uncertainty* that is also present**
- **Changes in the techniques used for developing cost models, i.e. producing “better” S-curves, will not eliminate utility of P-boxes**
 - **Periodic updating of the effective observations analysis will ensure more accurate bounds as cost estimating techniques progress**
 - **Epistemic uncertainty will always be present**

