# NORTHROP GRUMMAN

## DEFINING THE FUTURE

# A Non-Simulation Based Method for Inducing Pearson's Correlation Between Input Random Variables

And its application on the CG(X) risk assessment

## 2008 Acquisition Research Symposium

15 May 2008

## Eric Druker

Technical/Research Lead – Northrop Grumman IT
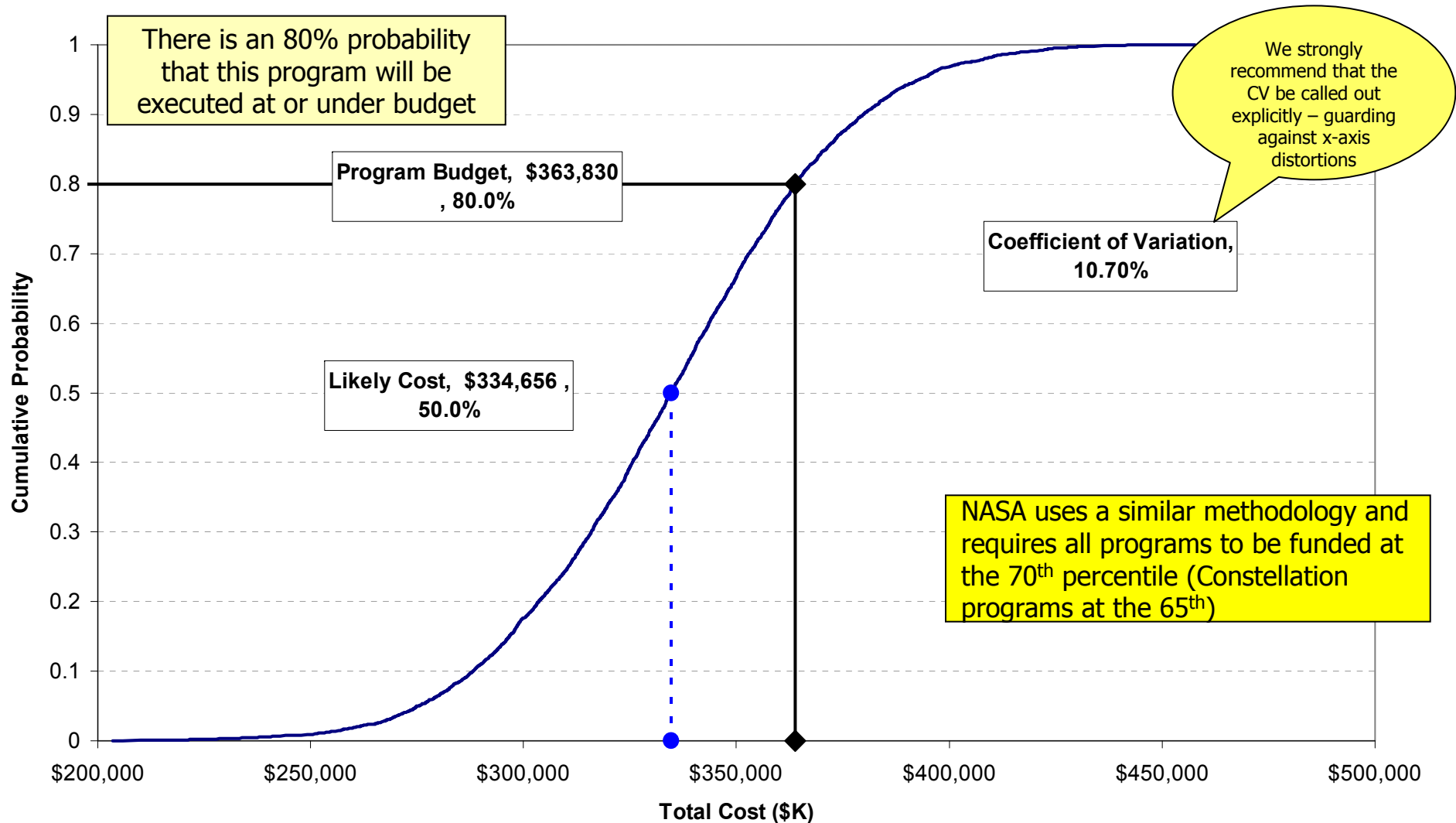
# Acknowledgements

# Introduction

- **Before moving to the main topic of the paper it is important to quickly discuss the motivation behind its development**

- **Studies have shown[1,2] that 75-85% of DoD programs experience cost overruns**
  - This suggests that as an industry, our estimates are not at the 50th percentile, but rather at about the 20th percentile

- **Recognizing this, agencies are taking the initiative to budget at higher percentiles of cost**
  - NASA requires all programs be funded at the 70th percentile
    - Constellation at the 65th
  - The Air Force (Dr. Sega) has released a memo advising that all space programs be funded at the 80th percentile
    - Rich Hartley (AFCAA) has advised against this, recommending programs be funded at the mean of the AFCAA ICE Estimate (generally between about the 55th and 60th percentiles)

- **In order to determine the appropriate funding level for programs anywhere but at the mean, it is thus imperative that the risk and uncertainty around estimates be assessed**
  - Thus S-Curves must be developed

1 Schaffer 2004 study, referenced from *Cost Estimating Requirements to Support New Congressional Reporting Requirements*. Coonce et. Al. NASA PM Challenge, February 2008
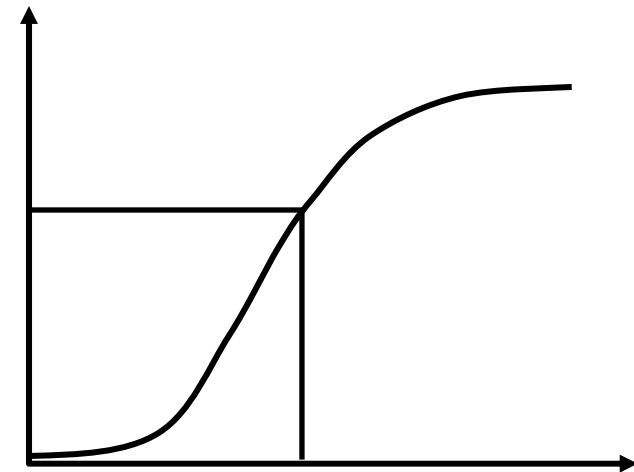
2

2 NAVAIR Cost Growth Study, R. L. Coleman, M.E. Dameron, C.L. Pullen, J.R. Summerville, D.M. Snead, 34th DoDCAS and ISPA/SCEA 2001

# Sample Program S-Curve

Program "X"
Cumulative Distribution

There is an 80% probability that this program will be executed at or under budget

Program Budget, $363,830, 80.0%

We strongly recommend that the CV be called out explicitly – guarding against x-axis distortions

Coefficient of Variation, 10.70%

Likely Cost, $334,656, 50.0%

NASA uses a similar methodology and requires all programs to be funded at the 70th percentile (Constellation programs at the 65th)

Cumulative Probability

Total Cost ($K)

Cumulative Distribution — Proposal Value — Likely Cost — Coefficient of Variation

# S-Curves

- **S-Curves are the cumulative distribution function for the cost of a system**
  - Also known as probabilistic cost estimates

- **S-Curves are generally driven by two main factors**
  - Cost Estimating Variance
    - Labor estimates
      - Data Driven
      - SME Driven
    - Escalation/Inflation Rates
    - Material Costs
    - Productivity (e.g. hrs/SLOC, hrs/ft$^2$)
  - Schedule/Technical Risks and Opportunities
    - Discrete Events
    - Continuous Events

- **Two key measures are derived from these S-Curves**
  - Confidence level of the estimate
    - What is the probability that the program will finish at or under budget?
  - Uncertainty in the estimate
    - What is the range of possibilities for the final cost of this program?
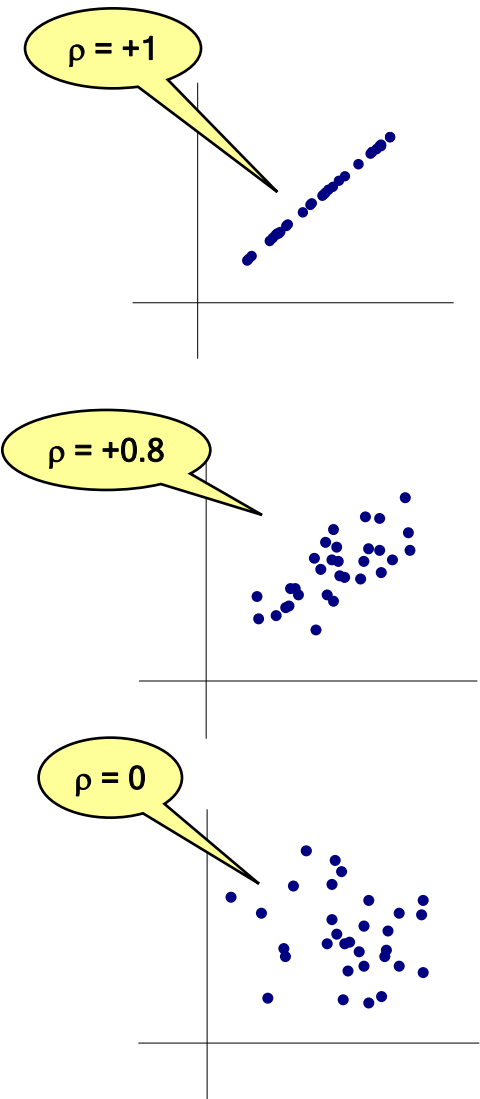
4

# Statement of Problem/Motivation

- Due to the increased focus on the reasonableness of cost estimates across the DoD community, a thorough risk assessment was conducted on the CG(X) program estimate
  - In particular, the Northrop Grumman team wanted to explore reasons that cost growth may be underestimated
  - It was determined that the treatment of correlation in risk adjusted cost estimates was one of the leading causes of this
  - Correlation directly effects the CV of the S-Curve

- In order to correctly capture program risk at a lower level, NGIT needed a way to include relational/injected correlation in our risk models
  - Without this ability the top level CV would be artificially shrunk due to the "square root of n problem"

- The following conditions lead the team away from traditional COTS models
  - The risk analysis module was to be incorporated into the CG(X) cost model
  - Both the cost and risk models were to be transitioned to a web-based platform

- In early 2006, work was begun on what would become the "Cost/Risk Correlation Module"
  - The module would have to exist entirely inside of Excel and VBA so it could be shared with any user with Office 2003 or later
  - The module would have to be *open* enough that it could be dropped quickly into most home-grown Monte Carlo models

# Outline

- Introduction to Correlation
  - Pearson's "Rho"
  - Pearson's vs. Rank Correlation

- The Problem

- Correlation Matrix Definitions

- Correlation in Risk Models

- Cost/Risk Correlation Algorithm
  - Correcting the user-input matrix
  - Correlating the random variables
  - Optimizing the Applied Matrix
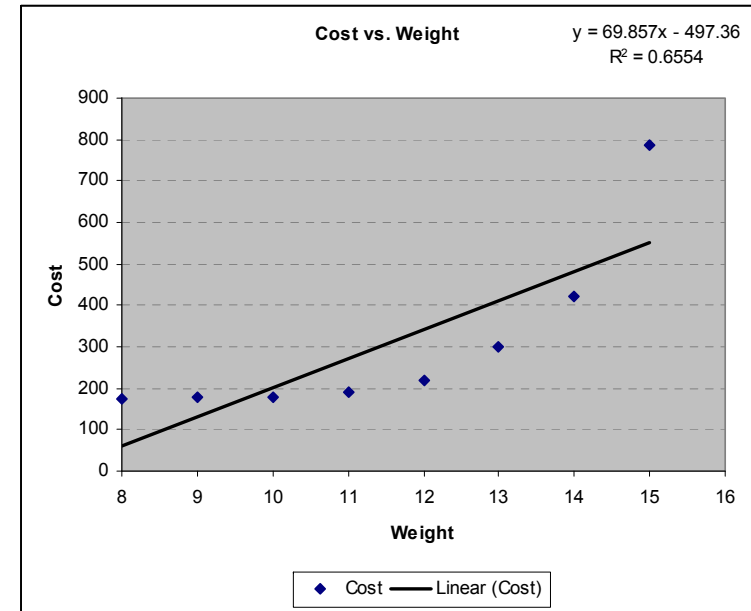
# Correlation (Pearson's)

- Although this paper is not about correlation itself, it's important to briefly review the two most common measures
  - Pearson's Product-Moment Correlation
  - Spearman's Rank Correlation

- When correlation is discussed in terms of cost estimating, Pearson's correlation is generally described

- Pearson's Correlation is a measure of the linear relationship between two or more variables
  - This is as opposed to Rank Correlation, which will be discussed on the next slide

- It is identified using the Greek symbol $\rho$ and is always between [-1,1]

- The correlation of a linear regression is the square root of $r^2$

- The examples on the right show representative data sets for three values of $\rho$

$\rho = +1$

$\rho = +0.8$

$\rho = 0$

# Pearson's Correlation vs. Rank Correlation

- Most commercial risk programs (e.g. Crystal Ball & @Risk) use Spearman's rank correlation rather than Pearson's correlation because it is easier to simulate

- Spearman's rank correlation is used to detect correlation between two variables, without assuming a linear relationship
  - It is concerned with whether or not the function is monotonic

- Some other differences include
  - Pearson's is parametric, Spearman's is not
  - Spearman's is not to be used for predictive purposes

- In the example to the right, rank correlation and Pearson's correlation yield very different answers

- Although it is important to distinguish between these two types of correlation, past research has shown that in cost risk simulations using the most common families of distributions, the two yield fairly similar results[1]
  - The aim of the authors is to "commit no avoidable errors"

**Cost vs. Weight**

$y = 69.857x - 497.36$
$R^2 = 0.6554$

| | |
|---|---|
| **Pearson's Rho** | **0.81** |
| **Spearman's Rho** | **1.00** |

[1] Robinson, M and Salls, W. *More on Correlation Accuracy in Crystal Ball Simulations (or What We've Now Learned about Spearman's R in Cost Risk Analyses)*. Presented at the 2004 SCEA Conference, Manhattan Beach, CA, June 2004

# Correlation in Risk Models

- In risk analysis, correlations are critical to successful simulations used to find distributions of cost
  - Correlations are thought to be widely present among elements of cost, but little data exists to determine them, principally because to determine correlations among any set of variables, data points must contain those variables in common, and this is rarely the case
  - Without accounting for correlation, summing multiple independent risk distributions will lead to an artificial degradation in the CV
    - This is known as the "Square Root of N" problem

- Lacking discernable correlations, risk analysts are forced to rely on Subject Matter Experts to estimate correlations
  - These correlations are subtle and difficult to estimate
  - Estimated correlations, to be usable, must be "coherent", as discussed later

- Once the desired correlation between all cost elements is determined, the next problem is to build these correlations into the risk model

- The following slides will lay out the algorithms used in the correlation module and demonstrate how they were applied to the CG(X) program

# Definitions: Matrices

- Before proceeding, it is important to define several matrices that will be used in the algorithm

- Input Correlation Matrix:
  - The correlation matrix inputted by the user, may or may not be a consistent correlation matrix

| User-Input Matrix | | |
|---|---|---|
| 1.0000 | 0.8000 | 0.1000 |
| 0.8000 | 1.0000 | 0.8000 |
| 0.1000 | 0.8000 | 1.0000 |

- Adjusted Correlation Matrix:
  - The consistent correlation matrix found by the model that is as close as possible to the Input Correlation Matrix
    - This matrix is positive semidefinite
    - It is also coherent given the distributions being correlated

| Adjusted Matrix | | |
|---|---|---|
| 1.0000 | 0.7522 | 0.1322 |
| 0.7522 | 1.0000 | 0.7522 |
| 0.1322 | 0.7522 | 1.0000 |

- Applied Correlation Matrix:
  - The correlation matrix utilized by the algorithm to generate correlated random number draws

| Applied Matrix | | |
|---|---|---|
| 1.0000 | 0.7915 | 0.2263 |
| 0.7915 | 1.0000 | 0.7744 |
| 0.2263 | 0.7744 | 1.0000 |

- Outcome Correlation Matrix
  - The correlation matrix of the simulation variables after the simulation is run
  - Ideally it is identical to the Adjusted Correlation Matrix

| Outcome | | |
|---|---|---|
| 1.0000 | 0.7522 | 0.1316 |
| 0.7522 | 1.0000 | 0.7521 |
| 0.1316 | 0.7521 | 1.0000 |

# Definitions: Eigenvalues/Eigenvectors

- An eigenvector is a vector $v$ such that for a square matrix A and a scalar $\lambda$, $Av = \lambda v$

- It follows that if Q is an indexed set of linearly independent eigenvectors for matrix A and $\Lambda$ is the diagonal matrix containing the corresponding eigenvalues of A as its diagonal entries then:

    $A = Q\Lambda Q^{-1}$

- By altering $\Lambda$, the diagonal matrix consisting of A's eigenvalues, we eventually arrive at a positive definite correlation matrix that is close to the user input matrix

- The Jacobi Eigenvalue algorithm is used to find both the eigenvalues and eigenvectors of the user input correlation matrix

# The Cost Risk Correlation Algorithm

Correcting the User Input Matrix

Correlating the Uniform Random Number Draws

Optimizing the Applied Matrix

# Correcting the User Input Matrix

- As a rule, correlation matrices must be positive semidefinite
  - Positive semidefinite matrices have all non-negative eigenvalues

- When using data to generate correlation matrices, they will necessarily be positive definite

- Unfortunately, when generating matrices based on SME judgment, this condition may not be met

- To correct these matrices, an algorithm developed by Iman and Davenport[1] was used
  - The criteria for "closest matrix" that comes out of this algorithm is unknown to the authors but it is computationally efficient and relatively simple to implement
  - Because the generation of the "closest viable correlation matrix" is so critical in finance, there are several more robust algorithms available[2]

- The following slide will outline the algorithm used in the Cost-Risk Correlation Module

[1] Iman, R and Davenport J. *An Intterative Algorithm to Produce a Positive Definite Matrix from an "Approximated Correlation Matrix" (With a Program User's Guide)* Sandia National Laboratories for the US DoE, June 1982
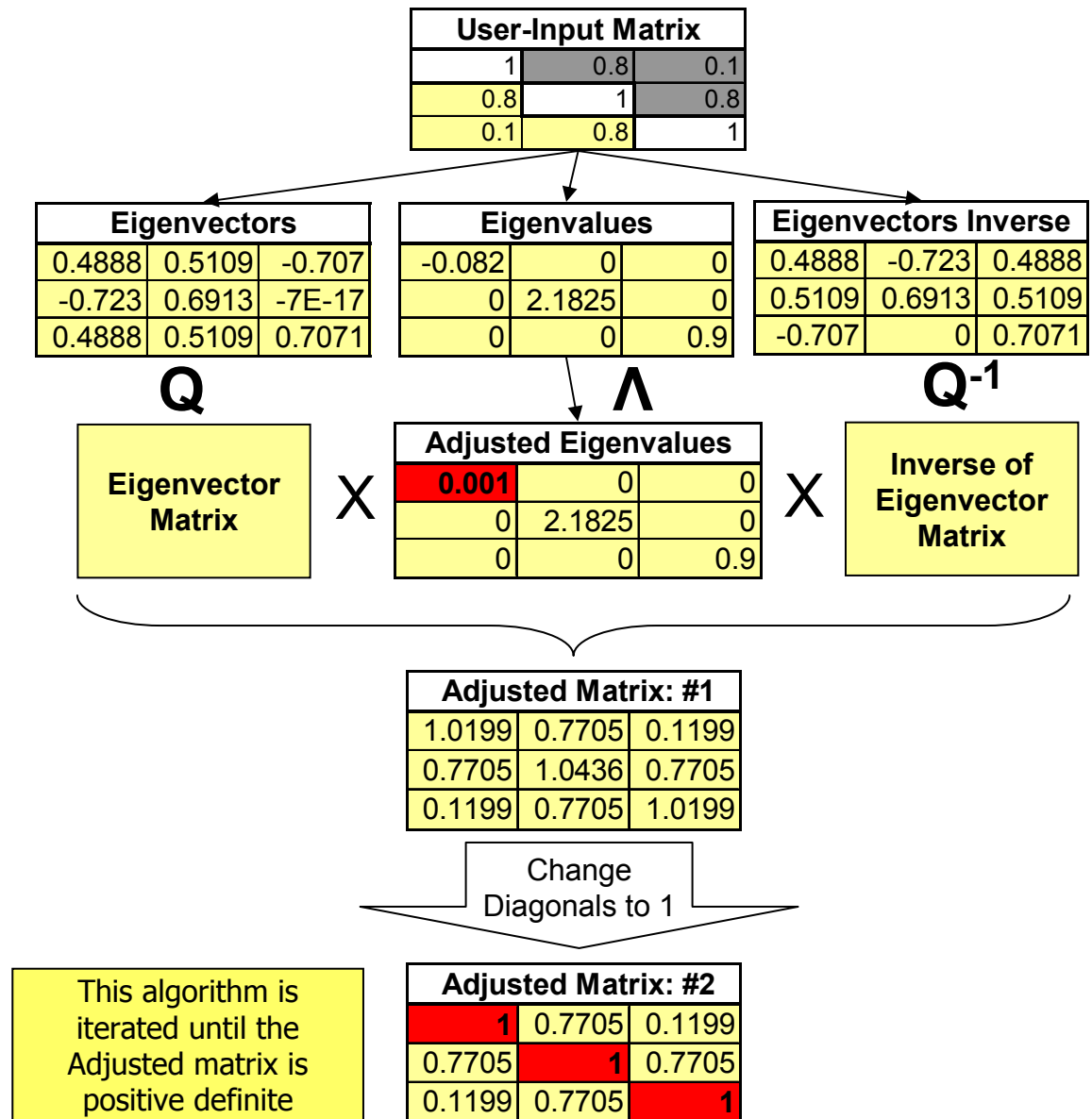
[2] Higham, N. *Computing the Nearest Correlation Matrix – A Problem from Finance.* IMA Journal of Numerical Analysis. 2002

# Correcting the User Input Matrix - Hurdles

- Two hurdles existed in implementing the algorithm
  - Excel doesn't have a function that finds Eigenvalues and Eigenvectors for the correlation matrices
  - Excel doesn't have a function to compute the Cholesky Decomposition matrix

- Research was conducted and algorithms (and the associated VBA source code) that conquered both hurdles were found
  - Both were part of the MATRIX and LINEAR ALGEBRA Package For EXCEL developed by The Foxes team in Italy
  - The Cholesky Decomposition, Eigenvalues and Eigenvectors functions were taken from this package and added into the tool

# Correcting the User Input Matrix - Algorithm

**NORTHROP GRUMMAN**

- The algorithm iteratively adjusts the eigenvalues of user-inputted correlation matrices until the resulting matrix has all non-negative

- During each iteration of the algorithm, there are two adjustments

  1. Adjustment of the negative eigenvalues to small, positive values

  2. Adjustment of the first adjusted matrix's diagonal entities to values of 1

- Once the adjusted matrix (#2) is found to have all non-negative Eigenvalues, the algorithm has found its solution

**User-Input Matrix**

| 1 | 0.8 | 0.1 |
|---|-----|-----|
| 0.8 | 1 | 0.8 |
| 0.1 | 0.8 | 1 |

**Eigenvectors**

| 0.4888 | 0.5109 | -0.707 |
|--------|--------|--------|
| -0.723 | 0.6913 | -7E-17 |
| 0.4888 | 0.5109 | 0.7071 |

**Q**

**Eigenvalues**

| -0.082 | 0 | 0 |
|--------|---|---|
| 0 | 2.1825 | 0 |
| 0 | 0 | 0.9 |

**Λ**

**Eigenvectors Inverse**

| 0.4888 | -0.723 | 0.4888 |
|--------|--------|--------|
| 0.5109 | 0.6913 | 0.5109 |
| -0.707 | 0 | 0.7071 |

**Q⁻¹**

**Eigenvector Matrix**   X

**Adjusted Eigenvalues**

| 0.001 | 0 | 0 |
|-------|---|---|
| 0 | 2.1825 | 0 |
| 0 | 0 | 0.9 |

X   **Inverse of Eigenvector Matrix**

**Adjusted Matrix: #1**

| 1.0199 | 0.7705 | 0.1199 |
|--------|--------|--------|
| 0.7705 | 1.0436 | 0.7705 |
| 0.1199 | 0.7705 | 1.0199 |

Change Diagonals to 1

This algorithm is iterated until the Adjusted matrix is positive definite

**Adjusted Matrix: #2**

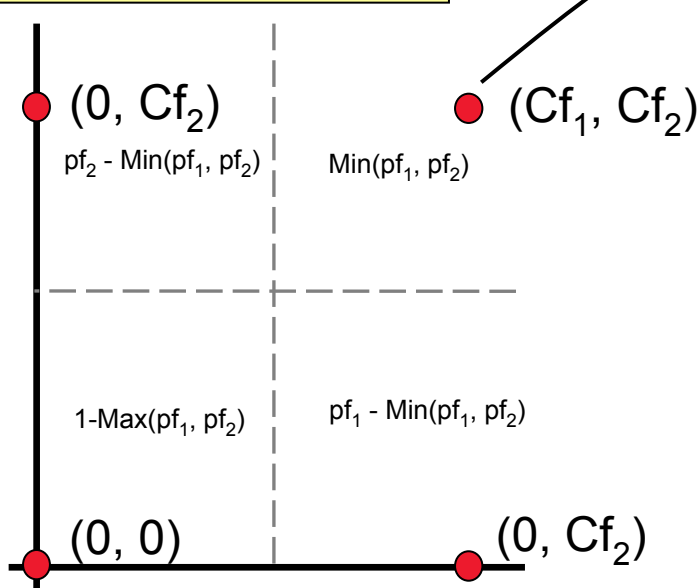| 1 | 0.7705 | 0.1199 |
|---|--------|--------|
| 0.7705 | 1 | 0.7705 |
| 0.1199 | 0.7705 | 1 |

# Correcting the User Input Matrix – Other Complications

- Although the matrix produced using the algorithm on the preceding slides is a consistent correlation matrix, depending on the random variables being correlated it may or may not be feasible
    - At least if the marginal distributions are to be preserved

- The best way to illustrate this is to examine the maximum possible correlation between two Bernoulli risks
    - As shown below, unless the probabilities of the two risks are equal, there is a maximum possible correlation between them

- The final step to correcting the User Input Matrix is to adjust the matrix so that all correlations are feasible based on the distributions being correlated

> Although this example seems odd, this is an efficient way of inducing conditional probabilities between Bernoulli random variables

> The only case in which XY ≠ 0 is when both risks occur, it follows that E(XY) simplifies down to $Cf_1 \times Cf_2$ times the probability that both risks occur. The highest this probability can possibly be is the minimum of the two probabilities

$(0, Cf_2)$

$pf_2 - Min(pf_1, pf_2)$     $Min(pf_1, pf_2)$

$(Cf_1, Cf_2)$

$$\rho_{X,Y} = \frac{E(XY) - E(X)E(Y)}{\sqrt{E(X^2) - E^2(X)}\sqrt{E(Y^2) - E^2(Y)}}$$

$$\rho(max)_{X,Y} = \frac{Min(Pf_1, Pf_2) \times Cf_1 \times Cf_2 - (Pf_1 \times Cf_1) \times (Pf_2 \times Cf_2)}{\sigma_X \times \sigma_Y}$$

$1-Max(pf_1, pf_2)$     $pf_1 - Min(pf_1, pf_2)$

$$\rho(max)_{X,Y} = \frac{Min(Pf_1, Pf_2) \times Cf_1 \times Cf_2 - (Pf_1 \times Cf_1) \times (Pf_2 \times Cf_2)}{Cf_1\sqrt{Pf_1 \times Qf_1} \times Cf_2\sqrt{Pf_2 \times Qf_2}}$$

$(0, 0)$            $(0, Cf_2)$

$$\rho(max)_{X,Y} = \frac{Min(Pf_1, Pf_2) - Pf_1 \times Pf_2}{\sqrt{Pf_1 \times Qf_1}\sqrt{Pf_2 \times Qf_2}}$$

16

# Correlating Random Variables:
## An Introduction to the Lurie-Goldberg Method[1]

- The only method the authors were aware of for inducing Pearson's correlation between input random variables is the Lurie-Goldberg Algorithm
  - The Lurie-Goldberg Algorithm aims to find an applied correlation matrix such that the input correlation and output correlation are as close as possible

- Find matrix **L** such that series of transformations

$$X \xrightarrow{\mathbf{L}} Y \xrightarrow{\Phi} U \xrightarrow{F^{-1}} V$$

  indep. normal → mult. normal → uniform → desired

  lead to random variables with desired correlations and marginal distributions
  - **L**: Cholesky factor transforms independent normals to correlated normals
  - $\Phi$: normal c.d.f. transforms correlated normals to correlated uniforms
  - $F^{-1}$: transforms correlated uniforms to correlated random variables with desired marginal distributions $F$

  - Unfortunately, the authors could not find a method for finding this optimal matrix (L... referenced as A' in this paper)
  - One obvious solution is to optimize the matrix by examining the post-simulation correlations
    - Given the computing power needed to complete each simulation, this could be a time consuming endeavor

[1]Goldberg, Matthew S, Lurie, Phillip M. *Correlating Random Variables*, 32nd DoDCAS, Williamsburg, VA. February 1999
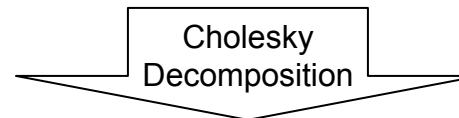
# Correlating the Uniform Draws:
# The Lurie-Goldberg Method

- Once a viable correlation matrix exists Uniform (0,1) correlated random numbers must be generated which in turn are used to generate the desired random variables

- To accomplish this, the Cholesky Decomposition Matrix of the adjusted matrix is found
  - L is the Cholesky Decomposition of A iff L is a lower triangular matrix such that:

$$A = LL^T$$

- After the Cholesky Decomposition Matrix is found, the algorithm at right is run to produce correlated Uniform (0,1) random numbers

- These random numbers, vice the originals, are used in the risk model to generate points off of distributions

| Adjusted Matrix | | |
|---|---|---|
| 1.0000 | 0.7522 | 0.1322 |
| 0.7522 | 1.0000 | 0.7522 |
| 0.1322 | 0.7522 | 1.0000 |

Cholesky Decomposition

| Cholesky Decomposition | | |
|---|---|---|
| 1.0000 | 0.0000 | 0.0000 |
| 0.7522 | 0.6589 | 0.0000 |
| 0.1322 | 0.9907 | 0.0321 |

**X**

| U(0,1) Random Draws |
|---|
| 0.26271853333989800 |
| 0.79616660202169400 |
| 0.15362541632109700 |

Inverse CDF Technique

| Random N(0,1) |
|---|
| (0.63498673467686800) |
| 0.82800654029771300 |
| (1.02100761130346000) |

Multiply N(0,1) by Cholesky

| Correlated Random N(0,1) |
|---|
| (0.63498673467686800) |
| 0.06794090908429620 |
| 0.70360627328862900 |

Multiply N(0,1) by Cholesky

Note: The resulting correlation between the Correlated Random U(0,1) random numbers will not be exactly the same as the adjusted correlation matrix… **more on this soon**

| Outcome Correlation | | |
|---|---|---|
| 1.0000 | 0.7386 | 0.1367 |
| 0.7386 | 1.0000 | 0.7395 |
| 0.1367 | 0.7395 | 1.0000 |

Resulting Correlation

| Correlated Random U(0,1) |
|---|
| 0.26271853333989800 |
| 0.52708366338494000 |
| 0.75916099823068700 |

# Optimizing the Applied Correlation Matrix

**NORTHROP GRUMMAN**

- Non-linear transformations are used to correlate random variables in the model
  - Because of this, the outcome correlation may be different from the intended correlation

- The biggest hurdle this module faced was in the correction of this discrepancy

- Northrop Grumman has developed a method that can find the outcome correlation matrix for any applied correlation matrix prior to the simulation being run
  - In other words, the algorithm can determine $\rho_{Output}$ given $\rho_{Applied}$
  - The applied correlation matrix can then be optimized so that the outcome correlation matrix is equal to the adjusted correlation matrix

- Additionally, it follows from mathematical proofs that the optimal applied correlation matrix will induce the desired correlation
  - This infers that any variation in $\rho$ in the simulation runs is due solely to Monte Carlo sampling error

Find:

| Applied Correlation Matrix | | |
|---|---|---|
| 1.0000 | 0.7915 | 0.2263 |
| 0.7915 | 1.0000 | 0.7744 |
| 0.2263 | 0.7744 | 1.0000 |

Such that after the Lurie-Goldberg method takes place:

| Outcome Correlation Matrix | | |
|---|---|---|
| 1.0000 | 0.7522 | 0.1322 |
| 0.7522 | 1.0000 | 0.7522 |
| 0.1322 | 0.7522 | 1.0000 |

=

| Adjusted Correlation Matrix | | |
|---|---|---|
| 1.0000 | 0.7522 | 0.1322 |
| 0.7522 | 1.0000 | 0.7522 |
| 0.1322 | 0.7522 | 1.0000 |

# Optimizing the Applied Correlation Matrix

**NORTHROP GRUMMAN**

- The algorithm developed by Northrop Grumman finds the optimal applied correlation matrix given:
    1. The parent distributions being correlated
    2. The adjusted correlation matrix

- The algorithm runs prior to the simulation being executed and once performed, only needs to be re-ran as variables are added or changed
    - And in those cases, only for the new/modified distributions

- Although the algorithm was originally developed for cost risk analysis, it has applications wherever a user needs to account for correlation between independent random variables
    - For example: the modeling of mutual fund performance given it is made up of a group of correlated stocks and bonds

- In fact, the algorithm's first use is in the modeling of conditional probabilities between Bernoulli independent random variables
    - The customer needed an efficient way to model the conditional probabilities they found between parameters in their data while preserving the marginal probabilities
    - It can be shown using the same general methodology on slide 14 that Pearson's correlation between two Bernoulli random variables equates to a conditional probability between them
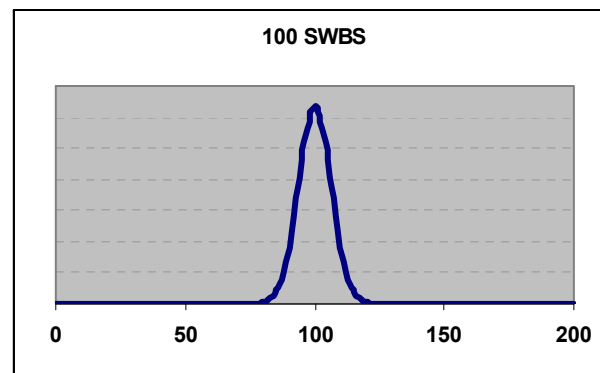
Application to the CG(X) Program Risk Assessment
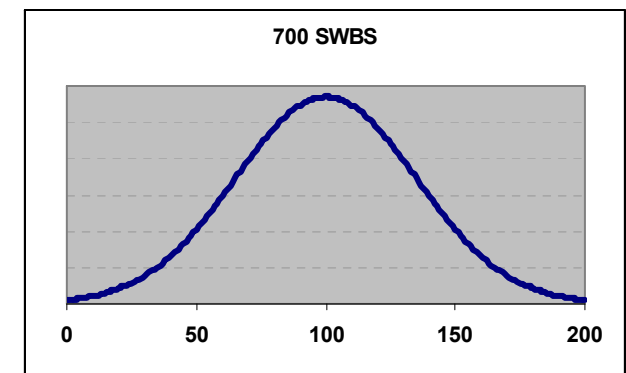
# Correlation Data

- One of the most difficult steps in the risk assessment process is in determining the correlation between the elements

- In this assessment, correlation is currently being measured using the relationship between the SWBS hours for three classes of surface combatants

- Just recently, data was obtained showing estimates vs. actuals, by SWBS, for various ships
    - The plan is to switch to correlations using this data once the analysis is complete

- Once uncertainty was evaluated for each lower level SWBS, correlation was applied between them to produce the top level risk adjusted estimate

# Estimating Variance (diagram)

- This simplified example shows only the 100 and 700 SWBS
  - 100 may have a lower level of uncertainty around its estimate than 700

- Using the correlation algorithm, accurate distributions can be generated for the lower level SWBSs that, when added together, still produce the known historical distribution
  - This allows decision makers to see what areas of the ship contain the greatest variance
  - It also allows risk to be applied at the 1-digit-level (see next slides)

**Whole Ship Cost**

| | | | | |
|100|150|200|250|300|

The bigger ratio of new to repeat work in the 700 SWBS is reflected in its larger CV (wider curve)

**100 SWBS**

| | | | | |
|0|50|100|150|200|

**+**

$\rho_{100,700}$

**700 SWBS**

| | | | | |
|0|50|100|150|200|

# Schedule/Technical Risks & Opportunities

- The next step in the risk assessment was adding in schedule and technical risks
  - Opportunities are just risks with a negative cost impact (cost is decreased)
  - From this point on, risks &opportunities will be referred to simply as risks

- Technical experts involved in CG(X) from across the corporation were interviewed to produce the schedule/technical risks associated with their area of the ship

- The following information was collected:
  - Description of the risk
  - Probability of occurrence
  - Description of the impact
    - This is the consequence of the risk occurring
  - Mitigation plans
    - Description of the mitigation plan
    - Cost of the mitigation plan
    - Probability and impact if the risk is mitigated
    - Whether or not the mitigation plan is included in the cost baseline
  - Other areas of the ship affected if the risk were to occur
    - If a schedule/technical risk increased the probability of occurrence for another risk, this was captured using the previously described correlation algorithm

# Schedule/Technical Risk Template

| Risk ID: | An ID used to identify the risk. Label Sequentially |
|---|---|
| Risk Description: | The risk description is a basic description of what the risk is. In particular, what could go wrong. |
| Probability of Occurrence: | The probability that the risk will occur. |
| Impact Description: | The impact description is all the information that would be needed from the SME in order to estimate the cost impact of the risk independently. Wherever, possible, please include schedule impacts as well |
| Mitigation Plans(s): | The mitigation plan(s) are all activities that would lower the expected value of the risk. These activities do not have to completely eliminate the risks, they could just lower either the probability of occurrence or cost impact. Information to be included: 1. Cost of Mitigation Plan (both schedule and $) 2. Affect mitigation plan has on the risk (what is the decrease in probability or cost/schedule impact) |
| Other Areas Affected: | Are there any other areas of the ship that could be impacted if this risk were to occur (or if the mitigation plans are put into motion)? If so, describe the impact and the area it would affect. Then, interview the owner of that area to determine if there are anymore residual impacts not forseen originally. |

# Schedule/Technical Risk Modeling

- Once the risks are collected, they were input into the model

- For risks with mitigation strategies, whether or not the mitigation strategy is implemented was selected using a drop-down menu
  - Mitigated risks (whose cost of mitigation is not included in the cost baseline) will add cost to the baseline cost
  - Mitigated risks will use mitigated probabilities and consequences

- Each risk is assigned to a 1-digit-level SWBS
  - This, along with the fact that cost estimating variability is also assessed at the 1-digit-level, allows cost distributions to be produced accurately at the 1-digit-level

- Risks can also be inputted as continuous risks (as appropriate):
  - Triangular Distributions
  - Normal Distributions
  - Log-Normal Distributions
  - All of these distributions can have probabilities assigned to them as well
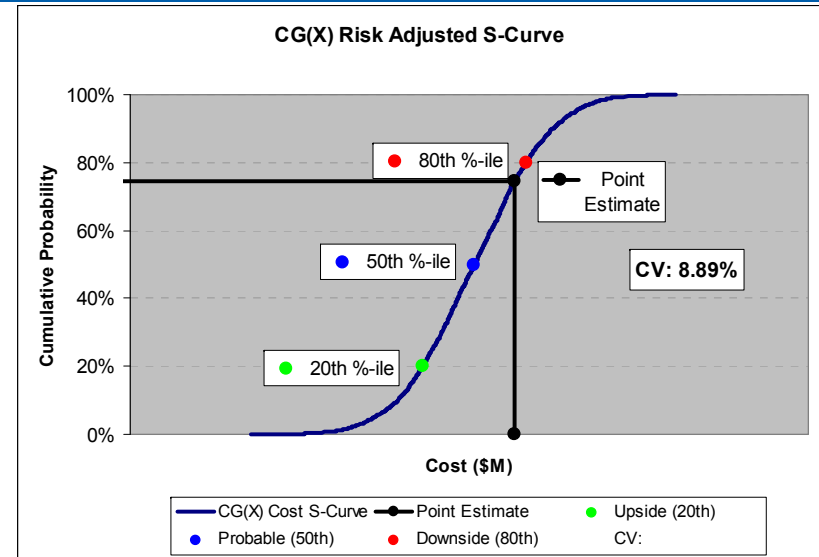
# Schedule/Technical Risks

| Risk ID | SWBS | Description | Probability of Occurrence | Bernoulli Cost Impact | Mitigation Plan Description | Mitigated Probability | Mitigated Cost Impact | Cost of Mitigation | Mitigation Plan Implemented? |
|---------|------|-------------|---------------------------|-----------------------|----------------------------|-----------------------|-----------------------|--------------------|------------------------------|
| 1 | 000 | Sample Risk 1 | 90% | $ 25,000,000 | Mitigation Plan 1 | 30% | $ 10,000,000 | $ 7,500,000 | Yes |
| 2 | 200 | Sample Risk 2 | 52% | | | | | | No |
| 3 | 300 | Sample Risk 3 | 75% | | | | | | No |
| 4 | 400 | Sample Risk 4 | 100% | | | | | | No |
| 5 | 500 | Sample Risk 5 | 10% | $ 100,000,000 | Mitigation Plan 5 | 1% | $ 50,000,000 | $ 10,000,000 | Yes |
| 6 | 600 | Sample Risk 6 | 25% | $ 13,000,000 | | | | | No |
| 7 | 700 | Sample Risk 7 | 90% | $ 9,000,000 | | | | | No |
| 8 | 800 | Sample Risk 8 | 100% | | | | | | No |
| 9 | 900 | Sample Risk 9 | 100% | | | | | | No |

Model Also Accepts Triangular, Normal and Lognormal Risk Distributions

- Several sets of results are produced automatically by the simulation when the "Run Simulation" button is hit

- CG(X) Risk Adjusted S-Curve
  - Shows the whole-ship cost distribution with the point estimate and its confidence on the graph

- CG(X) Risk Adjusted Estimate by 1-digit-level SWBS
  - Shows upside (20th Percentile), Probable (50th Percentile) and Downside (80th Percentile) by 1-digit-level SWBS

- CG(X) Risk by SWBS
  - Shows upside, probable and downside risk $'s by SWBS
  - These are the $'s due entirely to the risks, not estimating variation



CG(X) Risk Adjusted S-Curve

CG(X) Risk Adjusted Estimate

| SWBS | Description | Upside | Probable | Downside |
|------|-------------|--------|----------|----------|
| 000 | Administration | | | |
| 100 | Hull | | | |
| 200 | Propulsion | | | |
| 300 | Electric Plant | | | |
| 400 | Electonics Systems | | | |
| 500 | Auxillary Systems | | | |
| 600 | Outfit & Furnishings | | | |
| 700 | Weapons | | | |
| 800 | Integration & Engineering | | | |
| 900 | Ship Assembly & Support | | | |
| | Total | | | |

CG(X) Risk by SWBS

| SWBS | Description | Upside | Probable | Downside |
|------|-------------|--------|----------|----------|
| 000 | Administration | | | |
| 100 | Hull | | | |
| 200 | Propulsion | | | |
| 300 | Electric Plant | | | |
| 400 | Electonics Systems | | | |
| 500 | Auxillary Systems | | | |
| 600 | Outfit & Furnishings | | | |
| 700 | Weapons | | | |
| 800 | Integration & Engineering | | | |
| 900 | Ship Assembly & Support | | | |
| | Total | | | |

# Conclusion

NORTHROP GRUMMAN

- The previously discussed method is an attempt at producing a risk adjusted estimate for the CG(X) program that is also accurate at the SWBS level

- This analysis would not have been possible were it not for the creation of the cost/risk correlation module