# Synthesis of Causal Discovery and Machine Learning – Questions Posed

Robert Stoddard, Principal Researcher, SEI

Mike Konrad, Principal Researcher, SEI

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA  15213

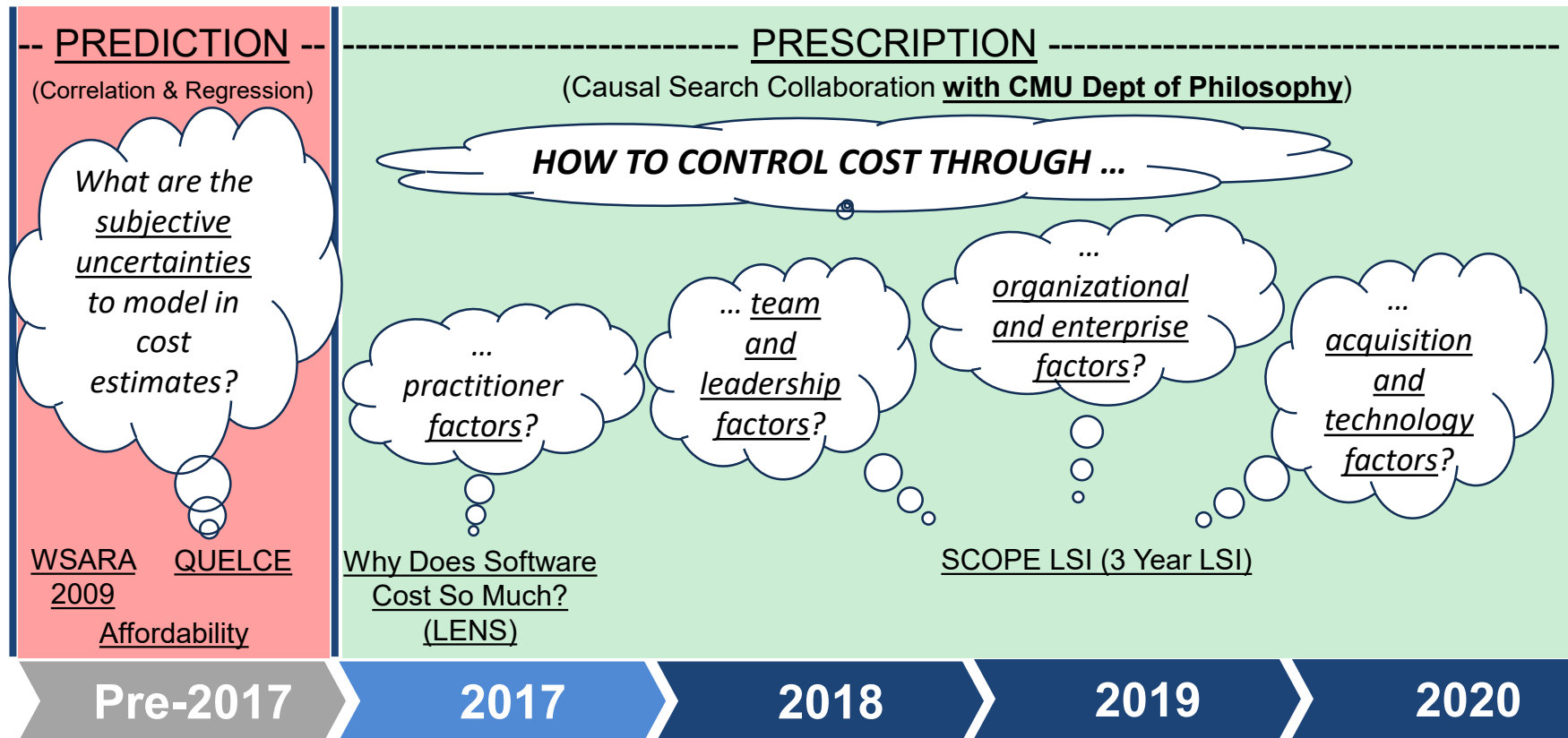# Agenda

☑ SEI SCOPE Research Focus

Use of BayesiaLab

Causal Learning

Comparison of ML and CL outputs

Questions Posed for Future Collaboration?

# Context of Causal Models for Software Cost Control (SCOPE)



**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

4

# Agenda

SEI SCOPE Research Focus

☑ Use of BayesiaLab

Causal Learning

Comparison of ML and CL outputs

Questions Posed for Future Collaboration?

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

**5**

# Use of BayesiaLab

1. Supervised machine learning (ML) with cost, schedule and quality as targets

2. Multi-variate outlier analysis

   a) Aid in data quality analysis

   b) Possible data segmentation strategies

3. Data imputation, when needed

4. Prediction of "what-if" scenarios of factors against outcomes

5. Classifier to assign probability of a binary outcome (e.g. good vs bad outcomes)

6. Diagnostic of most likely factors associated with a given outcome

7. All in support of DoD cost estimation and affordability analysis

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

**6**

# Agenda

SEI SCOPE Research Focus

Use of BayesiaLab

☑ Causal Learning

Comparison of ML and CL outputs

Questions Posed for Future Collaboration?

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**7**

# Why Do We Care about Causation?

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

8

# More about Misinterpreting Correlation!



Hot Temperature

Ice Cream Sales

Shark Attacks

Often, an excluded common cause results in a misinterpretation of correlation!

Does high correlation imply causation?

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

9

# Regression & ML benefit from a Structural Causal Model!

Regression and ML may be fooled by spurious association!

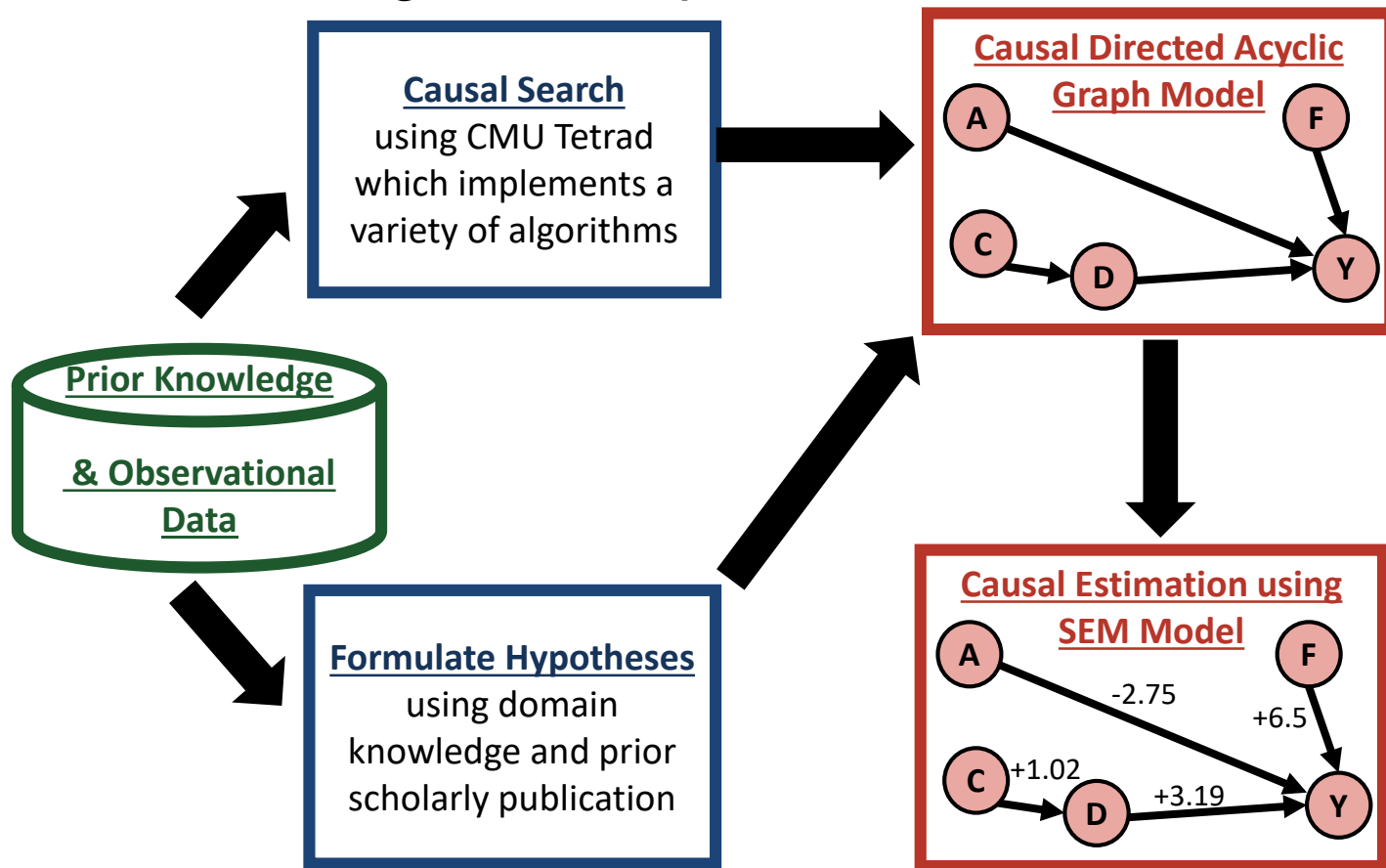Need a structural causal model (SCM) representing our theory and context

Need to determine which paths are causal versus non-causal
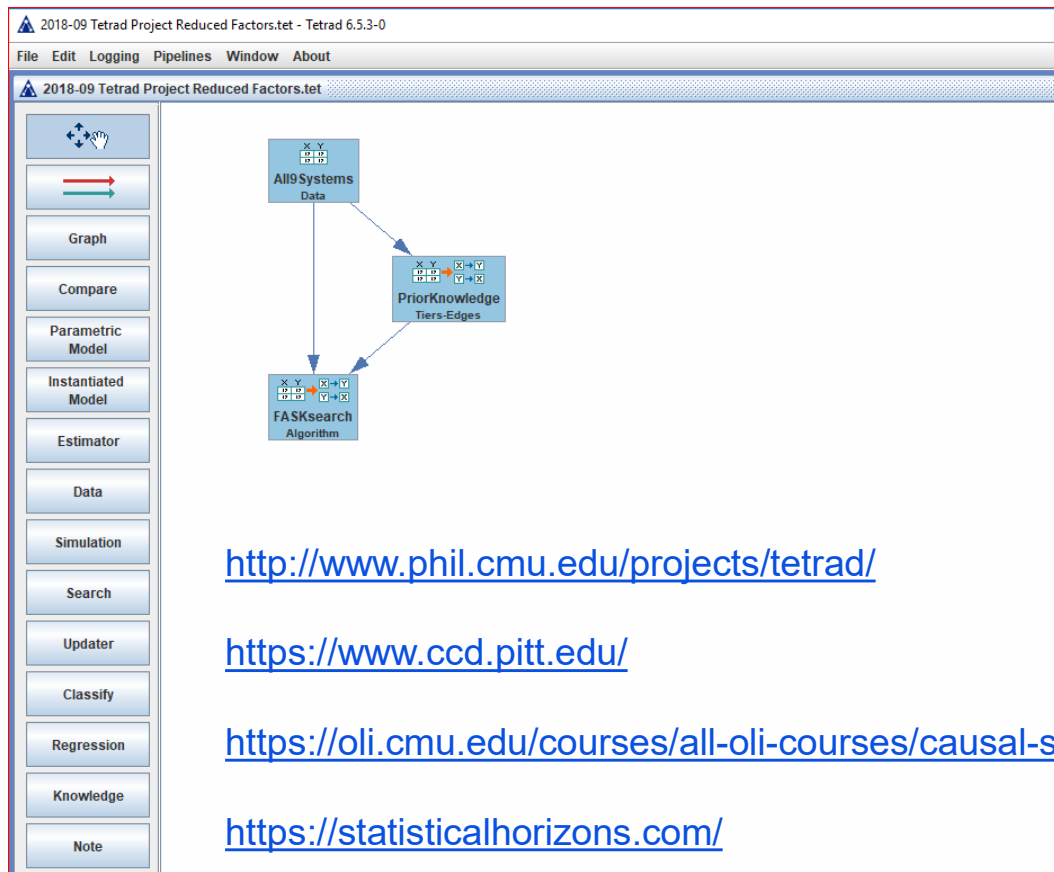
Must block non-causal paths

Then conduct regression and ML with the correct set of factors!

### *Suitability of the model depends on the SCM!*

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

10

# The Causal Learning Landscape

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

11

# Conduct Causal Search using Tetrad



http://www.phil.cmu.edu/projects/tetrad/

https://www.ccd.pitt.edu/

https://oli.cmu.edu/courses/all-oli-courses/causal-statistical-reasoning/

https://statisticalhorizons.com/

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

12

# A View of the Data File Loaded into Tetrad



All9Systems (Data)

File   Edit   Tools

**All 9 for Tetrad-v010.csv**

|    | C1 AgeMonths | C2 NumDev | C3 LOC | C4 NumBugs | C5 BugChurn | C6 NumCyclic... | C7 NumModul... | C8 NumUnsta... | C9 NumImpro... |
|----|----------|---------|----------|----------|----------|----------|---------|---------|---------|
| 1  | 71.0000  | 8.0000  | 491.0000 | 18.0000  | 241.0000 | 8.0000   | 2.0000  | 3.0000  | 1.0000  |
| 2  | 35.0000  | 5.0000  | 270.0000 | 10.0000  | 329.0000 | 167.0000 | 1.0000  | 1.0000  | 4.0000  |
| 3  | 52.0000  | 2.0000  | 58.0000  | 0.0000   | 0.0000   | 0.0000   | 0.0000  | 0.0000  | 0.0000  |
| 4  | 42.0000  | 1.0000  | 47.0000  | 2.0000   | 13.0000  | 0.0000   | 0.0000  | 0.0000  | 0.0000  |
| 5  | 49.0000  | 1.0000  | 10.0000  | 0.0000   | 0.0000   | 0.0000   | 0.0000  | 1.0000  | 0.0000  |
| 6  | 36.0000  | 2.0000  | 103.0000 | 0.0000   | 0.0000   | 0.0000   | 1.0000  | 0.0000  | 0.0000  |
| 7  | 54.0000  | 2.0000  | 29.0000  | 2.0000   | 0.0000   | 0.0000   | 0.0000  | 0.0000  | 0.0000  |
| 8  | 75.0000  | 8.0000  | 163.0000 | 13.0000  | 134.0000 | 0.0000   | 1.0000  | 3.0000  | 0.0000  |
| 9  | 74.0000  | 2.0000  | 15.0000  | 0.0000   | 0.0000   | 0.0000   | 1.0000  | 0.0000  | 0.0000  |
| 10 | 57.0000  | 2.0000  | 26.0000  | 1.0000   | 16.0000  | 22.0000  | 0.0000  | 0.0000  | 0.0000  |
| 11 | 48.0000  | 4.0000  | 81.0000  | 2.0000   | 6.0000   | 0.0000   | 1.0000  | 0.0000  | 0.0000  |
| 12 | 39.0000  | 1.0000  | 30.0000  | 0.0000   | 0.0000   | 0.0000   | 0.0000  | 0.0000  | 0.0000  |
| 13 | 49.0000  | 2.0000  | 46.0000  | 3.0000   | 36.0000  | 0.0000   | 0.0000  | 0.0000  | 0.0000  |
| 14 | 46.0000  | 3.0000  | 34.0000  | 0.0000   | 0.0000   | 0.0000   | 0.0000  | 0.0000  | 1.0000  |
| 15 | 75.0000  | 0.0000  | 0.0000   | 0.0000   | 0.0000   | 0.0000   | 0.0000  | 0.0000  | 1.0000  |

Done

# Prior Knowledge Entered into Tetrad

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**14**

# Causal Learning Algorithms

**<u>Constraint-based:</u>** Calculate independences in the data and do "backwards inference"; used to minimize the degree of false negative edges

**<u>Score-based (Bayesian):</u>** Calculate the likelihood of different DAGs given the data; used to minimize the degree of false positive edges
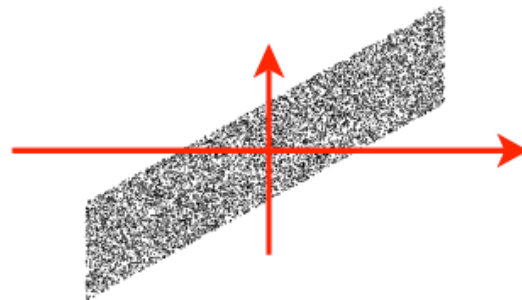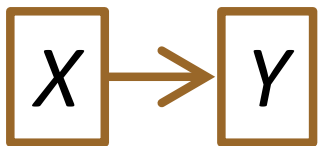
**<u>Hybrid:</u>** Use constraint-based to get "close," then Bayesian search around neighborhood

| | |
|---|---|
| A      B | No evidence of a causal link |
| A ⟶ B | Evidence of a causal link from A to B |
| A ⟵ B | Evidence of a causal link from B to A |
| A ⟷ B | Evidence of an unmeasured confounder |

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**15**

# Some Algorithms Exploit Non-Gaussianality



Linear Gaussian    Linear non-Gaussian

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

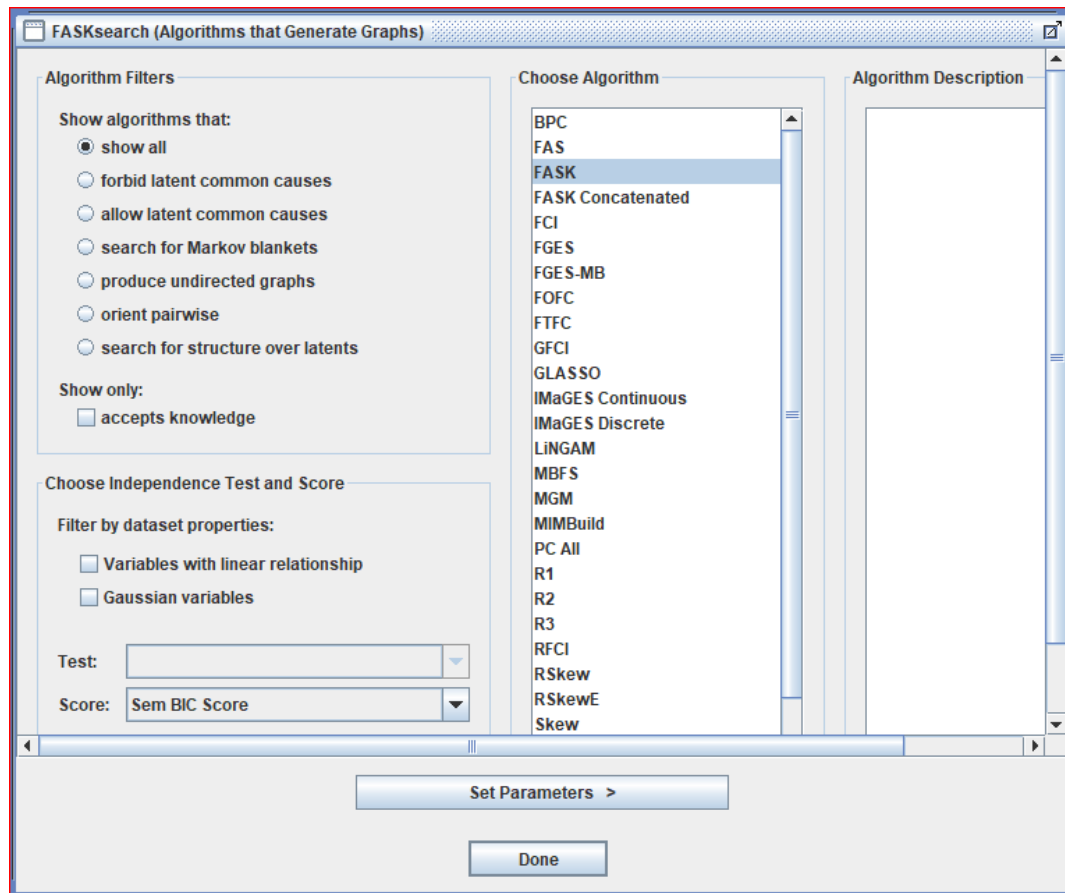[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

16

# Causal Search Capable with Small Data

**Challenge**:  Which genes regulate flowering time in Arabidopsis thaliana?

Using only 47 observations, causal search identified 9 out of 21,326 genes as causal on gene activation

Subsequent greenhouse study, that used knockout variants, confirmed that 4 of the 9 were actual regulators

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**17**

# Using FASK Search with Associated Parameters



**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

18

# Additional FASK Search Parameter Settings

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

**19**

# Causal Structure Graph Result

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

# Markov Blanket of the NumBugs Factor



**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

21

# Traditional SEM Results from Tetrad

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

22

# Additional Causal Learning Topics

1. Algorithms operating on the Structural Causal Model (see Judea Pearl, 2018, "The Book of Why")

2. Propensity Scoring (see Shenyang Guo and Mark W. Fraser, 2014, "Propensity Score Analysis")

3. Instrumental Variables (see Felix Elwert, publications on Instrumental Variables)

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**23**

# Agenda

SEI SCOPE Research Focus

Use of BayesiaLab

Causal Learning

☑ Comparison of ML and CL outputs

Questions Posed for Future Collaboration?

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

24

# ML and CL Graph Structures May Be Different

## CL Markov Blanket



## ML Markov Blanket

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

25

# Agenda

SEI SCOPE Research Focus

Use of BayesiaLab

Causal Learning

Comparison of ML and CL outputs

☑ Questions Posed for Future Collaboration?

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

**26**

# When ML and CL Graph Structure Results Differ?

1. Choose to instantiate the Tetrad causal structure in BayesiaLab as a PSEM?

2. Use BayesiaLab to conduct Pearl graph surgery or Jouffe's likelihood matching for causal modeling?

3. Pursue metrics such as Average Causal Effect (ACE) and Total Causal Effect (TCE)?

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

27

# Opportunities to Integrate ML & CL? - 01

1. Can a ML association graph structure result inform a CL causal search?

2. For extremely large datasets and # variables, would ML require significantly less computer time than a CL causal search? If so, could ML serve as a pre-screen of a CL causal search?

3. Could ML graph structure results inform opportunities for research into new CL causal search algorithms?

4. Could/should CL causal search be combined with ML graphical results for a new, superior output?

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**28**

# Opportunities to Integrate ML & CL? - 02

5. Is there a possible superior understanding obtainable from graphical structural results of both ML and CL?

   a) Can differences between the two graphs provide insight?

   b) Can commonality across the two graphs provide insight?

   c) More generally, is there greater knowledge of combining Shannon Information Theory with Causal Theory?

6. Can combined use of ML and CL graphical structures enable an improved method of "stitching together" separate, but overlapping results towards a more holistic result?

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**29**

# Conclusion

We are seeking research collaboration in two ways:

1. Collaboration and data access for software project cost estimation and control, and

2. Collaboration to gain insight and answer the questions posed in this presentation

**Carnegie Mellon University**
Software Engineering Institute

Synthesis of Causal Discovery and Machine Learning – Questions Posed
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.]

**30**

# Contact Information

**Presenter Contact Information**



Dr. Mike Konrad
Principal Researcher,
SEI / CMU
mdk@sei.cmu.edu
1-412-268-5813



Robert Stoddard
Principal Researcher,
SEI / CMU
rws@sei.cmu.edu
1-412-268-1121

**Carnegie Mellon University**
Software Engineering Institute

**Synthesis of Causal Discovery and Machine Learning – Questions Posed**
© 2018 Carnegie Mellon University

[DISTRIBUTION STATEMENT A] Approved for
public release and unlimited distribution.]

31