

CS6462

Probabilistic and Explainable AI

Lesson 23

Meta-Learning Causal Structures



Causal Direction

Causality & causal inference:

- *basic goal*: learn causality from data:
 - what is the cause and what is the effect
 - correlations between variables should be enough for most of the cases
- *causal inference*: go one step further and figure out what would happen if we decide to change some of the underlying assumptions in our model

Inferring causal direction:

- the most popular tool is proper trial designs of experiments - *randomized control trial (RCT)*
- RCT are not universal, because cannot be conducted in many scenarios: can be too costly or infeasible due to the complexity of real-world systems
- the development of more causality in Machine Learning is a necessary step in building more human-like machine intelligence



Meta-Learning to Infer Causal Direction

Meta-learning in ML:

- learning algorithms learn from other learning algorithms
- use of machine learning algorithms that learn how to best combine the predictions from other machine learning algorithms
- *meta*: raising the level of abstraction one step higher + additional model information

Meta-learning for causal direction:*

- apply learned models assuming different causal directions to data with a changed transfer distribution
- meta-learning objective concerned with assumptions on data distribution and data changes due to moving from a training distribution to a transfer distribution, possibly resulting from some neuron's actions (part of the neural network we train)
- the correct causal model needs to adjust its transfer distribution only, and thus it adapts faster, which allows us to determine the underlying causal directions

*Y. Bengio et al. 2019. A Meta-Transfer Objective for Learning to Disentangle Causal Mechanisms



Learning a Causal Graph with Two Variables

Learning meta-model:

- Bayesian network of two discrete random variables **A**, **B** (each taking **n** possible values)
- training samples **(a; b)** from a pair of related distributions: **training distribution** and **transfer distribution**
- based only on samples from a single training distribution, both models **A → B** (A causes B) and **B → A** tend to perform
- determining if variable **A** causes variable **B** or vice-versa
- parametrizations of models (**A → B** and **B → A**) to estimate their joint distribution:
 $P_{A \rightarrow B}(A, B) = P_{A \rightarrow B}(A) * P_{A \rightarrow B}(B|A)$ - for training and transfer distributions
 $P_{B \rightarrow A}(A, B) = P_{B \rightarrow A}(B) * P_{B \rightarrow A}(A|B)$ - for training and transfer distributions
four distribution models (training and transfer distribution models per graph)
- **A** and **B** are completely observed: maximum likelihood estimator is used to independently train all four models



Learning a Causal Graph with Two Variables (cont.)

Distributions:

- model **$A \rightarrow B$**

- training distributions:

$$P_{A \rightarrow B}^0(A, B) = P_{A \rightarrow B}^0(A) * P_{A \rightarrow B}^0(B|A)$$

- transfer distributions

$$P_{A \rightarrow B}^1(A, B) = P_{A \rightarrow B}^1(A) * P_{A \rightarrow B}^1(B|A)$$

- model **$B \rightarrow A$**

- training distributions:

$$P_{B \rightarrow A}^0(A, B) = P_{B \rightarrow A}^0(B) * P_{B \rightarrow A}^0(A|B)$$

- transfer distributions

$$P_{B \rightarrow A}^1(A, B) = P_{B \rightarrow A}^1(B) * P_{B \rightarrow A}^1(A|B)$$

- models trained with $P_{A/B \rightarrow B/A}^0$ first and then moved to $P_{A/B \rightarrow B/A}^1$



Learning a Causal Graph with Two Variables (cont.)

Result likelihood distributions:

- both network models are meta-trained on both training and transfer distributions for T steps with resulting likelihoods:

$$L_{A \rightarrow B} = \prod_{i=1}^T P_{A \rightarrow B, i}(A_i, B_i)$$

$$L_{B \rightarrow A} = \prod_{i=1}^T P_{B \rightarrow A, i}(A_i, B_i)$$

- models are trained following the steps:
 - 1) models of $A \rightarrow B$ and $B \rightarrow A$ are trained using their training distribution
 - 2) relationship between A and B is learned using the training results
 - 3) distribution of both models is moved from a training to a transfer distribution
 - 4) both models are retrained on the new data and the resulting likelihoods are recorded
 - 5) based on the results, we evaluate the direction of the causal relationship



Learning a Causal Graph with Two Variables (cont.)

Loss function:

- the loss function of both training steps (over training and transfer distributions):

$$R(\alpha) = -\ln[\sigma(\alpha) * L_{A \rightarrow B} + (1 - \sigma(\alpha)) * L_{B \rightarrow A}]$$

- $R(\alpha)$ computed with α denoting a structural parameter defining the causal direction and $\sigma(\alpha)$ the sigmoid transformation
- α is optimized to minimize $R(\alpha)$

Loss gradient:

- $\frac{dR}{d\alpha} = \sigma(\alpha) - \sigma(\alpha + \ln(L_{A \rightarrow B}) - \ln(L_{B \rightarrow A}))$
- $\frac{dR}{d\alpha} > 0$ if $L_{A \rightarrow B} < L_{B \rightarrow A}$ - i.e., if $B \rightarrow A$ is better at explaining the transfer distribution than $A \rightarrow B$



Summary

Meta-Learning Causal Structures

Causal Directions

Meta-Learning to Infer Causal Direction

Learning a Causal Graph with Two Variables

Next Lesson:

- Structural Equation Modeling

Thank You!

Questions?