

PROBABILISTIC DEPENDENCY GRAPHS AND INCONSISTENCY

HOW TO MODEL, MEASURE, AND MITIGATE INTERNAL CONFLICT

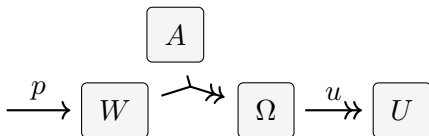
Oliver Richardson

Cornell University
Department of Computer Science

September 2021

The standard way of modeling an agent with uncertainty:

- a probability distribution $p : \Delta W$ over worlds W ,
- a utility function $u : \Omega \rightarrow \mathbb{R}$, some actions A .



Such agents cannot have internal conflict;

by construction, they have consistent beliefs and desires.

I like this slide

⟨ INCOMPLETE ⟩

Why build systems that can be inconsistent, if inconsistency is bad?

Why entertain the possibility of being wrong, if being wrong is bad?

Elsewhere, computer scientists take great care to model inconsistency:

- assertions and test cases: Text
-
- losses for training neural networks (totally separate from the probabilistic model)

I think you're agonizing too much over this. Why not just say "Because sometimes do have inconsistent beliefs. We also want to model the dynamics of making them consistent, and capture the extent to which they're inconsistent."

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

I don't think that I would remove this but modify it to say something like: not only do we want to model inconsistency, but we want to do so using a probabilistic graphical model.

Then go back to the current slides and say "In doing so, we get much more."

⟨ remove slide, once intro is finished ⟩

We introduce *probabilistic dependency graphs* (PDGs), a new class of graphical models for representing uncertainty.

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

⟨ remove slide, once intro is finished ⟩

We introduce *probabilistic dependency graphs* (PDGs), a new class of graphical models for representing uncertainty.

Why do we need another one?

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

⟨ remove slide, once intro is finished ⟩

We introduce *probabilistic dependency graphs* (PDGs), a new class of graphical models for representing uncertainty.

Why do we need another one?

- To resolve inconsistency, we must first model it.

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

⟨ remove slide, once intro is finished ⟩

We introduce *probabilistic dependency graphs* (PDGs), a new class of graphical models for representing uncertainty.

Why do we need another one?

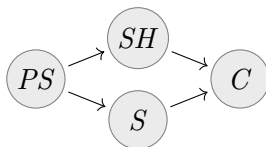
- To resolve inconsistency, we must first model it.
- In doing so, we get much more ...

TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

Qualitative BN, \mathcal{G}

an independence relation on variables

- $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid \text{Pa}(X)$, for all non-descendants Y of X



TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

Qualitative BN, \mathcal{G}

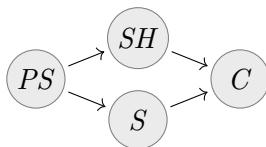
an independence relation on variables

- $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid \mathbf{Pa}(X)$, for all non-descendants Y of X

(Quantitative) BN, $\mathcal{B} = (\mathcal{G}, \mathbf{p})$

a qualitative BN (\mathcal{G}) and a cpd $p_X(X \mid \mathbf{Pa}(X))$ for each variable X .

- Defines a joint distribution $\Pr_{\mathcal{B}}$ with the independencies $\perp\!\!\!\perp_{\mathcal{G}}$.



TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

Qualitative BN, \mathcal{G}

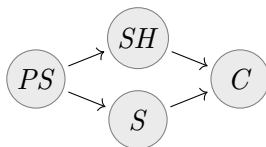
an independence relation on variables

- $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid \mathbf{Pa}(X)$, for all non-descendants Y of X

(Quantitative) BN, $\mathcal{B} = (\mathcal{G}, \mathbf{p})$

a qualitative BN (\mathcal{G}) and a cpd $p_X(X \mid \mathbf{Pa}(X))$ for each variable X .

- Defines a joint distribution $\Pr_{\mathcal{B}}$ with the independencies $\perp\!\!\!\perp_{\mathcal{G}}$.



OUTLINE FOR SECTION 1

Why is the outline coming after you've introduced BN's? What the high-level story?

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

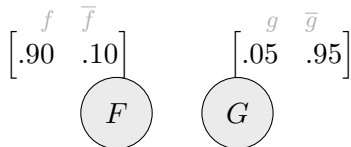
SIMPLE EXAMPLE: FLOOMPS AND GUNS

Grok thinks it likely (.95) that guns are illegal,
but that floomps (local slang) are legal (.90).

SIMPLE EXAMPLE: FLOOMPS AND GUNS

Grok thinks it likely (.95) that guns are illegal, but that floomps (local slang) are legal (.90).

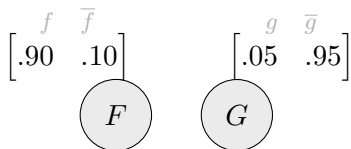
BN



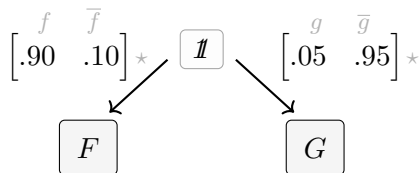
SIMPLE EXAMPLE: FLOOMPS AND GUNS

Grok thinks it likely (.95) that guns are illegal,
but that floomps (local slang) are legal (.90).

BN

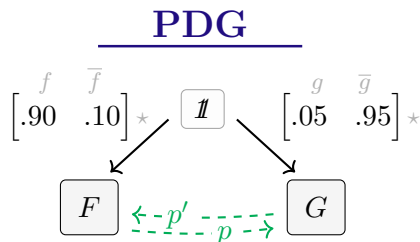
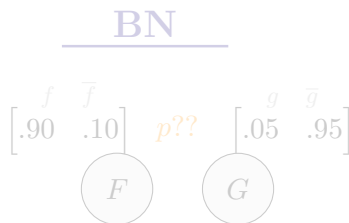


PDG



- The cpds of a PDG are attached to edges, not nodes.

SIMPLE EXAMPLE: FLOOMPS AND GUNS

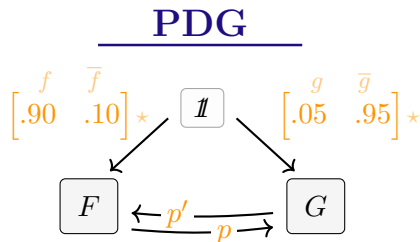
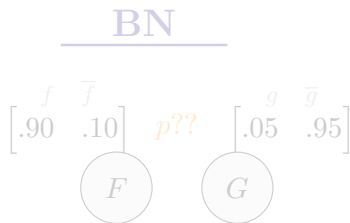


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.

Grok learns that Floomps and Guns have the same legal status (92%)

$$p(G|F) = \begin{bmatrix} g & \bar{g} \\ .92 & .08 \\ .08 & .92 \end{bmatrix} \begin{matrix} f \\ \bar{f} \end{matrix} = (p'(F|G))^T$$

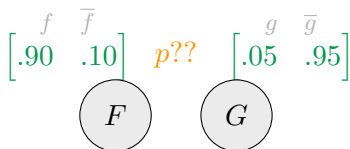
SIMPLE EXAMPLE: FLOOMPS AND GUNS



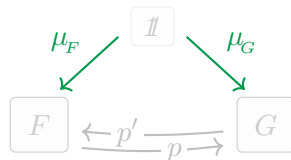
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent

SIMPLE EXAMPLE: FLOOMPS AND GUNS

BN



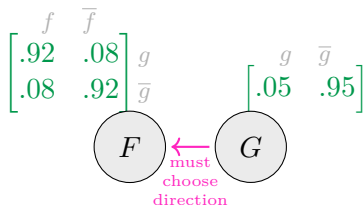
PDG



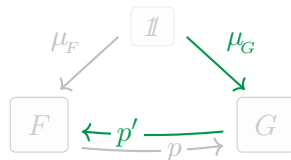
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
 - ▶ ...but BNs must resolve inconsistency first,
which may break symmetry and irrecoverably lose information.

SIMPLE EXAMPLE: FLOOMPS AND GUNS

BN



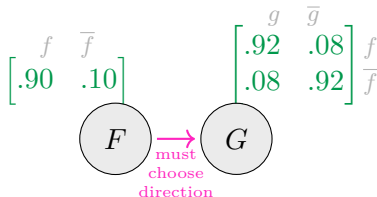
PDG



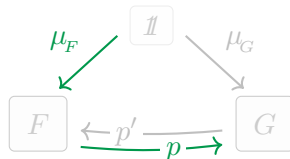
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
 - ▶ ...but BNs must resolve inconsistency first,
which may **break symmetry** and irrecoverably lose information.

SIMPLE EXAMPLE: FLOOMPS AND GUNS

BN

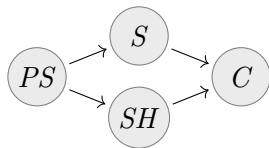


PDG

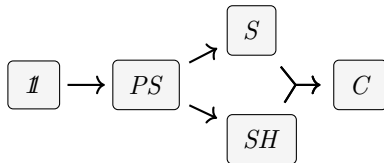
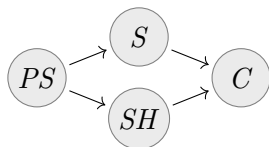


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
 - ▶ ...but BNs must resolve inconsistency first,
which may **break symmetry** and irrecoverably lose information.

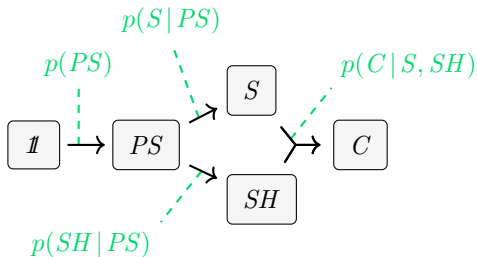
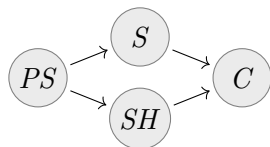
BAYESIAN NETWORKS AS PDGs



BAYESIAN NETWORKS AS PDGs



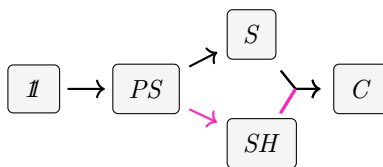
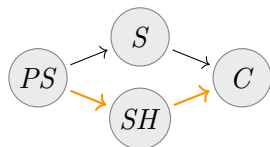
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

- edge composition has *quantitative* meaning, since edges have cpds;

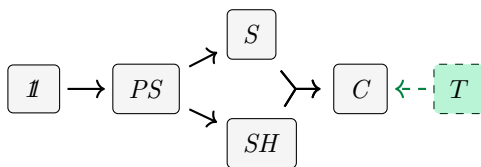
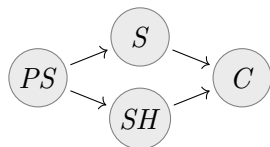
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

- edge composition has *quantitative* meaning, since edges have cpds;

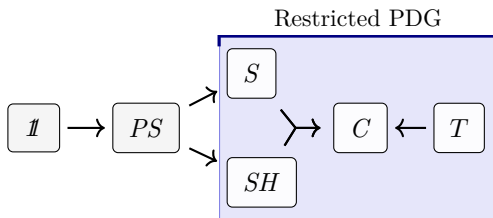
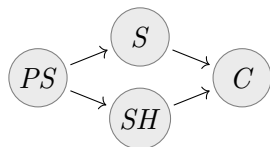
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

- edge composition has *quantitative* meaning, since edges have cpds;
- a variable can be the target of more than one cpd;

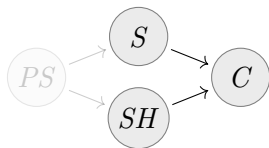
BAYESIAN NETWORKS AS PDGs



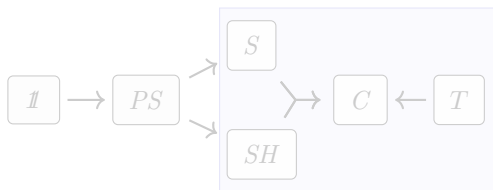
In contrast with BNs:

- edge composition has *quantitative* meaning, since edges have cpds;
- a variable can be the target of more than one cpd;
- arbitrary restrictions of PDGs are still PDGs.

BAYESIAN NETWORKS AS PDGs



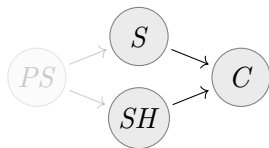
Must now give distributions on SH and S , or distinguish them as “observed” (a *conditional* BN).



In contrast with BNs:

- edge composition has *quantitative* meaning, since edges have cpds;
- a variable can be the target of more than one cpd;
- arbitrary restrictions of PDGs are still PDGs.
 - ▶ The analogue is false for BNs!

BAYESIAN NETWORKS AS PDGs



Must now give distributions on SH and S , or distinguish them as “observed” (a *conditional* BN).

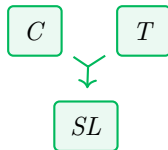
In a qualitative BN: *removing data results in new knowledge*: $A \perp\!\!\!\perp C$.



In contrast with BNs:

- edge composition has *quantitative* meaning, since edges have cpds;
- a variable can be the target of more than one cpd;
- arbitrary restrictions of PDGs are still PDGs.
 - ▶ The analogue is false for BNs!

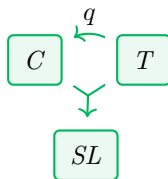
COMBINING PDGs



Grok wants to be supreme leader (SL).

- She notices that those who use tanning beds have more power, unless they get cancer

COMBINING PDGs

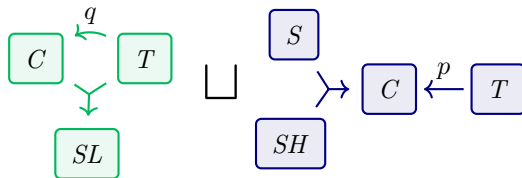


Grok wants to be supreme leader (SL).

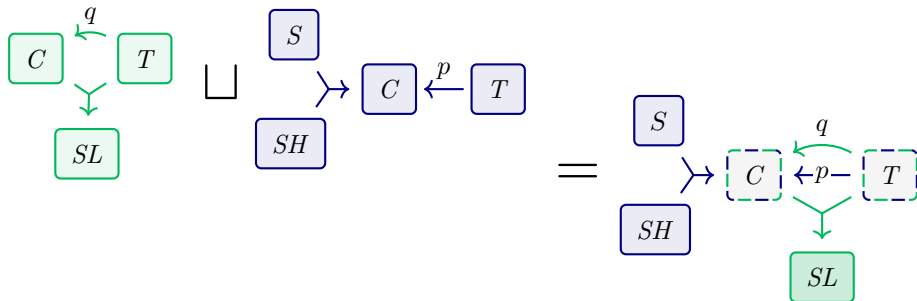
- She notices that those who use tanning beds have more power, unless they get cancer

- ...but mom says $q(C \mid T) = \begin{bmatrix} \overset{c}{.15} & \overset{\bar{c}}{.85} \\ .02 & .98 \end{bmatrix} \begin{matrix} t \\ \bar{t} \end{matrix}$.

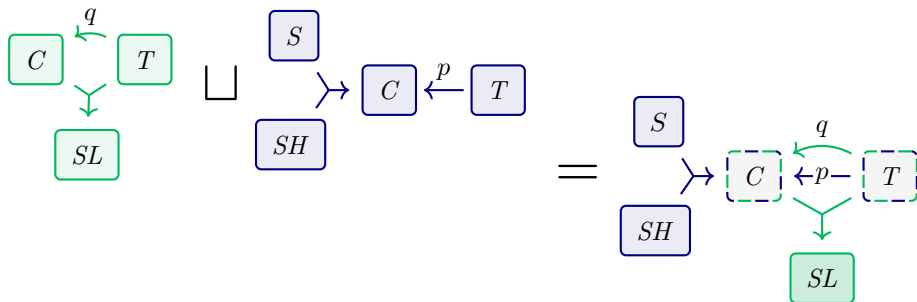
COMBINING PDGs



COMBINING PDGs

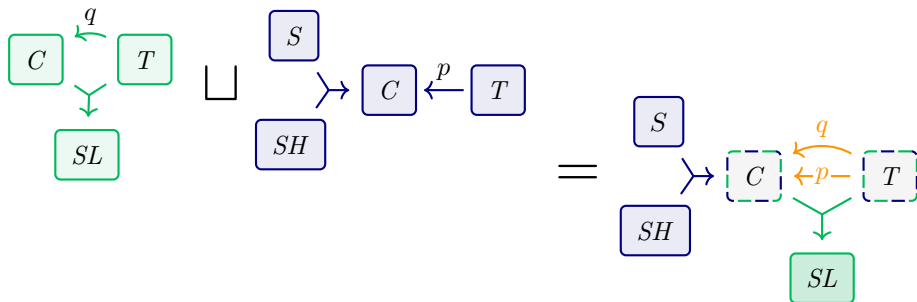


COMBINING PDGs



- Arbitrary PDGs may be combined without loss of information

COMBINING PDGs



- Arbitrary PDGs may be combined without loss of information
- They may have parallel edges which directly conflict.

OUTLINE FOR SECTION 2

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$,

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where \mathcal{N} is a finite set of nodes (variables)

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

$\mathcal{V}(\mathcal{M}) := \prod_{X \in \mathcal{N}} \mathcal{V}(X)$ is the set of possible joint variable settings.

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

(or hyper-edges)

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

(or hyper-edges)

and associated to each $X \xrightarrow{L} Y$, there is:

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

(or hyper-edges)

and associated to each $X \xrightarrow{L} Y$, there is:

\mathbf{p}_L a cpd $\mathbf{p}_L(Y \mid X)$;

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

(or hyper-edges)

and associated to each $X \xrightarrow{L} Y$, there is:

\mathbf{p}_L a cpd $\mathbf{p}_L(Y \mid X)$;

β_L a confidence in the reliability of \mathbf{p}_L .

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

(or hyper-edges)

and associated to each $X \xrightarrow{L} Y$, there is:

\mathbf{p}_L a cpd $\mathbf{p}_L(Y \mid X)$;

β_L a confidence in the reliability of \mathbf{p}_L .

write “ $p!$ ” for the limit ($\beta_p \rightarrow \infty$)
of high confidence in a cpd p .

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathcal{M} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

(or hyper-edges)

and associated to each $X \xrightarrow{L} Y$, there is:

\mathbf{p}_L a cpd $\mathbf{p}_L(Y \mid X)$;

α_L a confidence in the functional dependence $X \rightarrow Y$;

β_L a confidence in the reliability of \mathbf{p}_L .

Needless to say, I think that all the following black slides are a *terrible* idea.
This is not where you want to spend your time. It's a major distraction.

Definition (PDG)

$\mathcal{N} : \mathbf{Set}$ (node set)

$\mathcal{V} : \mathcal{N} \rightarrow \mathbf{Set}$ (node values)

$\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N} \times \text{Label}$ (edge set)

For $X \xrightarrow{L} Y \in \mathcal{E}$,

$\mathbf{p}_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$ (edge cpd)

$\alpha_L : \mathbb{R}$ (functional determination)

$\beta_L : \mathbb{R}$ (cpd confidence)

- $(\mathcal{N}, \mathcal{V})$ is a set of variables

Definition (PDG)

$\mathcal{N} : \mathbf{Set}$ (node set)

$\mathcal{V} : \mathcal{N} \rightarrow \mathbf{Set}$ (node values)

$\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N} \times \mathbf{Label}$ (edge set)

For $X \xrightarrow{L} Y \in \mathcal{E}$,

$\mathbf{p}_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$ (edge cpd)

$\alpha_L : \mathbb{R}$ (functional determination)

$\beta_L : \mathbb{R}$ (cpd confidence)

- $(\mathcal{N}, \mathcal{V})$ is a set of variables
- $(\mathcal{N}, \mathcal{E})$ is a multigraph

Definition (PDG)

$\mathcal{N} : \mathbf{Set}$ (node set)

$\mathcal{V} : \mathcal{N} \rightarrow \mathbf{Set}$ (node values)

$\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N} \times \mathbf{Label}$ (edge set)

For $X \xrightarrow{L} Y \in \mathcal{E}$,

$\mathbf{p}_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$ (edge cpd)

$\alpha_L : \mathbb{R}$ (functional determination)

$\beta_L : \mathbb{R}$ (cpd confidence)

- $(\mathcal{N}, \mathcal{V})$ is a set of variables
- $(\mathcal{N}, \mathcal{E})$ is a multigraph
- $(\mathcal{N}, \mathcal{E}, \alpha)$, the qualitative data, forms a weighted multigraph.

Definition (PDG)

$\mathcal{N} : \mathbf{Set}$ (node set)

$\mathcal{V} : \mathcal{N} \rightarrow \mathbf{Set}$ (node values)

$\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N} \times \mathit{Label}$ (edge set)

For $X \xrightarrow{L} Y \in \mathcal{E}$,

$\mathbf{p}_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$ (edge cpd)

$\alpha_L : \mathbb{R}$ (functional determination)

$\beta_L : \mathbb{R}$ (cpd confidence)

- $(\mathcal{N}, \mathcal{V})$ is a set of variables
- $(\mathcal{N}, \mathcal{E})$ is a multigraph
- $(\mathcal{N}, \mathcal{E}, \alpha)$, the qualitative data, forms a weighted multigraph.
- We call $(\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p})$ an *unweighted* PDG
 - ▶ and give it semantics as though $\alpha_L = \beta_L = 1$.

Definition (PDG)

$\mathcal{N} : \mathbf{Set}$ (node set)
 $\mathcal{V} : \mathcal{N} \rightarrow \mathbf{Set}$ (node values)
 $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N} \times \mathit{Label}$ (edge set)
For $X \xrightarrow{L} Y \in \mathcal{E}$,
 $\mathbf{p}_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$ (edge cpd)
 $\alpha_L : \mathbb{R}$ (functional determination)
 $\beta_L : \mathbb{R}$ (cpd confidence)

- $(\mathcal{N}, \mathcal{V})$ is a set of variables
- $(\mathcal{N}, \mathcal{E})$ is a multigraph
- $(\mathcal{N}, \mathcal{E}, \alpha)$, the qualitative data, forms a weighted multigraph.
- We call $(\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p})$ an *unweighted* PDG
 - ▶ and give it semantics as though $\alpha_L = \beta_L = 1$.

Let **Mark** be the category of measurable spaces and Markov kernels.

Definition (PDG)

$\mathcal{N} : \mathbf{Set}$ (node set)
 $\mathcal{V} : \mathcal{N} \rightarrow \mathbf{Set}$ (node values)
 $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N} \times \mathit{Label}$ (edge set)
For $X \xrightarrow{L} Y \in \mathcal{E}$,
 $\mathbf{p}_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$ (edge cpd)
 $\alpha_L : \mathbb{R}$ (functional determination)
 $\beta_L : \mathbb{R}$ (cpd confidence)

- $(\mathcal{N}, \mathcal{V})$ is a set of variables
- $(\mathcal{N}, \mathcal{E})$ is a multigraph
- $(\mathcal{N}, \mathcal{E}, \alpha)$, the qualitative data, forms a weighted multigraph.
- We call $(\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p})$ an *unweighted* PDG
 - ▶ and give it semantics as though $\alpha_L = \beta_L = 1$.

Let **Mark** be the category of measurable spaces and Markov kernels.

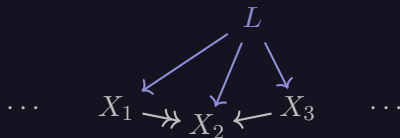
Equivalent Categorical Definition

An unweighted PDG is a functor $\langle \mathbf{p}, \mathcal{V} \rangle : \mathit{Paths}(\mathcal{N}, \mathcal{E}) \rightarrow \mathbf{Mark}$.
So a PDG is a *diagram* in **Mark**, in the usual mathematical sense.

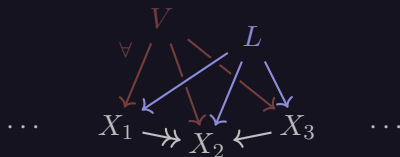
What do you do with diagrams? Take **limits** / **colimits**.

$$\cdots \quad X_1 \twoheadrightarrow X_2 \leftarrow X_3 \quad \cdots$$

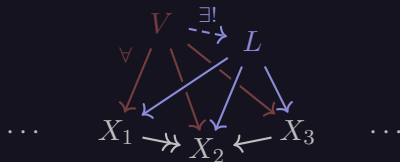
What do you do with diagrams? Take limits / colimits.



What do you do with diagrams? Take limits / colimits.



What do you do with diagrams? Take limits / colimits.



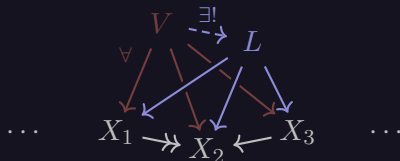
What do you do with diagrams? Take limits / colimits.



For the deterministic sub-PDG $\mathcal{m}_{\text{det}} \subseteq \mathcal{m}$:

$$\lim \mathcal{m}_{\text{det}} = \left(\begin{array}{cc} \text{natural} & \Omega, \\ \text{sample space} & \text{random variables } \left\{ \tilde{X} : \Omega \rightarrow \mathcal{V}(X) \right\}_{X \in \mathcal{N}} \end{array} \right)$$

What do you do with diagrams? Take limits / colimits.



For the deterministic sub-PDG $\mathcal{m}_{\text{det}} \subseteq \mathcal{m}$:

$$\lim \mathcal{m}_{\text{det}} = \left(\begin{array}{c} \text{natural} \\ \text{sample space} \end{array} \Omega, \begin{array}{c} \text{random} \\ \text{variables} \end{array} \left\{ \tilde{X} : \Omega \rightarrow \mathcal{V}(X) \right\}_{X \in \mathcal{N}} \right)$$

In general: $\lim \mathcal{m} = \left(\text{Verts}(\underbrace{\mathbb{L}\mathcal{m}}_{\substack{\text{Locally Consistent Polytope} \\ \text{(possible states of the Sum-Product algorithm)}}}), \{ \text{variable marginals} \} \right)$

What do you do with diagrams? Take limits / colimits.



For the deterministic sub-PDG $\mathcal{M}_{\text{det}} \subseteq \mathcal{M}$:

$$\lim \mathcal{M}_{\text{det}} = \left(\begin{array}{c} \text{natural} \\ \text{sample space} \end{array} \Omega, \begin{array}{c} \text{random} \\ \text{variables} \end{array} \left\{ \tilde{X} : \Omega \rightarrow \mathcal{V}(X) \right\}_{X \in \mathcal{N}} \right)$$

In general: $\lim \mathcal{M} = \left(\text{Verts}(\underbrace{\mathbb{L}\mathcal{M}}_{\substack{\text{Locally Consistent Polytope} \\ \text{(possible states of the Sum-Product algorithm)}}}, \{ \text{variable marginals} \} \right)$

For a BN \mathcal{B} : $\lim \mathcal{M}_{\mathcal{B}} = \left(\mathbb{I}, \left\{ \text{Pr}_{\mathcal{B}}(X) \right\}_{X \in \mathcal{N}} \right)$

OUTLINE FOR SECTION 3

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

SEMANTICS OF PDGs

$$\{\mathbf{m}\} \subseteq \Delta\mathcal{V}(\mathbf{m})$$

The set of joint distributions consistent with \mathbf{m} ;

SEMANTICS OF PDGS

$$\{\mathbf{m}\} \subseteq \Delta\mathcal{V}(\mathbf{m})$$

The set of joint distributions consistent with \mathbf{m} ;

$$[\mathbf{m}]_{\gamma} : \Delta\mathcal{V}(\mathbf{m}) \rightarrow \mathbb{R}$$

A function (parameterized by $\gamma > 0$) that scores distributions by compatibility with \mathbf{m} ;

SEMANTICS OF PDGS

$$\{\mathbf{m}\} \subseteq \Delta\mathcal{V}(\mathbf{m})$$

The set of joint distributions consistent with \mathbf{m} ;

$$[\mathbf{m}]_{\gamma} : \Delta\mathcal{V}(\mathbf{m}) \rightarrow \mathbb{R}$$

A function (parameterized by $\gamma > 0$) that scores distributions by compatibility with \mathbf{m} ;

$$[\mathbf{m}]_{\gamma}^* \subseteq \Delta\mathcal{V}(\mathbf{m})$$

The distribution(s) most compatible with \mathbf{m}
(a singleton in many cases of interest);

SEMANTICS OF PDGS

$$\{\mathbf{m}\} \subseteq \Delta\mathcal{V}(\mathbf{m})$$

The set of joint distributions consistent with \mathbf{m} ;

$$[\mathbf{m}]_{\gamma} : \Delta\mathcal{V}(\mathbf{m}) \rightarrow \mathbb{R}$$

A function (parameterized by $\gamma > 0$) that scores distributions by compatibility with \mathbf{m} ;

$$[\mathbf{m}]_{\gamma}^* \subseteq \Delta\mathcal{V}(\mathbf{m})$$

The distribution(s) most compatible with \mathbf{m}
(a singleton in many cases of interest);

$$\langle\!\langle \mathbf{m} \rangle\!\rangle_{\gamma} : \mathbb{R}$$

The best possible compatibility of \mathbf{m} with any distribution: the *inconsistency* of \mathbf{m}

SEMANTICS OF PDGs

$$\{\mathcal{m}\} \subseteq \Delta\mathcal{V}(\mathcal{m})$$

The set of joint distributions consistent with \mathcal{m} ;

$$[\mathcal{m}]_{\gamma} : \Delta\mathcal{V}(\mathcal{m}) \rightarrow \mathbb{R}$$

A function (parameterized by $\gamma > 0$) that scores distributions by compatibility with \mathcal{m} ;

$$[\mathcal{m}]_{\gamma}^* \subseteq \Delta\mathcal{V}(\mathcal{m})$$

The distribution(s) most compatible with \mathcal{m}
(a singleton in many cases of interest);

$$\langle\!\langle \mathcal{m} \rangle\!\rangle_{\gamma} : \mathbb{R}$$

The best possible compatibility of \mathcal{m} with any distribution: the *inconsistency* of \mathcal{m}

SEMANTICS OF PDGS

$$\{m\} \subseteq \Delta\mathcal{V}(m)$$

The set of joint distributions consistent with m ;

$$[[m]]_\gamma : \Delta\mathcal{V}(m) \rightarrow \mathbb{R}$$

A function (parameterized by $\gamma > 0$) that scores distributions by compatibility with m ;

$$[[m]]_\gamma^* \subseteq \Delta\mathcal{V}(m)$$

The distribution(s) most compatible with m
(a singleton in many cases of interest);

$$\langle\langle m \rangle\rangle_\gamma : \mathbb{R}$$

The best possible compatibility of m with any distribution: the *inconsistency* of m

SEMANTICS OF PDGS

$$\{m\} \subseteq \Delta\mathcal{V}(m)$$

The set of joint distributions consistent with m ;

$$[m]_{\gamma} : \Delta\mathcal{V}(m) \rightarrow \mathbb{R}$$

A function (parameterized by $\gamma > 0$) that scores distributions by compatibility with m ;

$$[m]_{\gamma}^* \subseteq \Delta\mathcal{V}(m)$$

The distribution(s) most compatible with m
(a singleton in many cases of interest);

$$\langle\!\langle m \rangle\!\rangle_{\gamma} : \mathbb{R}$$

The best possible compatibility of m with any distribution: the *inconsistency* of m

THE SCORING FUNCTION

⟨ either finish overhaul or reinstate old slides ⟩

$$\llbracket m \rrbracket_{\gamma}(\mu) := \textcolor{green}{Inc}m(\mu) + \gamma \textcolor{violet}{IDef}m(\mu)$$

tradeoff parameter $\gamma \geq 0$

⟨ Emphasize entropy balance ⟩

OUTLINE FOR SECTION 4

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

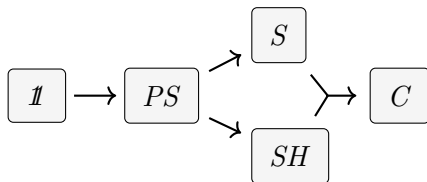
- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

CAPTURING BAYESIAN NETWORKS

For a BN \mathcal{B} with N nodes and a vector $\beta \in \mathbb{R}^N$, let $\mathbf{m}_{\mathcal{B},\beta}$ be the PDG corresponding to \mathcal{B} , with $\alpha = \mathbf{1}$, and the given vector β of confidences.



CAPTURING BAYESIAN NETWORKS

For a BN \mathcal{B} with N nodes and a vector $\beta \in \mathbb{R}^N$, let $\mathbf{m}_{\mathcal{B},\beta}$ be the PDG corresponding to \mathcal{B} , with $\alpha = \mathbf{1}$, and the given vector β of confidences.

Theorem (*BNs are PDGs*)

If \mathcal{B} is a BN and $\text{Pr}_{\mathcal{B}}$ is the distribution it specifies, then for all $\gamma > 0$ and all vectors β ,

$$\llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket_{\gamma}^* = \{\text{Pr}_{\mathcal{B}}\}, \quad \text{and thus} \quad \llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket^* = \text{Pr}_{\mathcal{B}}.$$

CAPTURING BAYESIAN NETWORKS

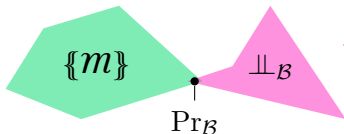
For a BN \mathcal{B} with N nodes and a vector $\beta \in \mathbb{R}^N$, let $\mathbf{m}_{\mathcal{B},\beta}$ be the PDG corresponding to \mathcal{B} , with $\alpha = \mathbf{1}$, and the given vector β of confidences.

Theorem (*BNs are PDGs*)

If \mathcal{B} is a BN and $\text{Pr}_{\mathcal{B}}$ is the distribution it specifies, then for all $\gamma > 0$ and all vectors β ,

$$\llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket_{\gamma}^* = \{\text{Pr}_{\mathcal{B}}\}, \quad \text{and thus} \quad \llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket^* = \text{Pr}_{\mathcal{B}}.$$

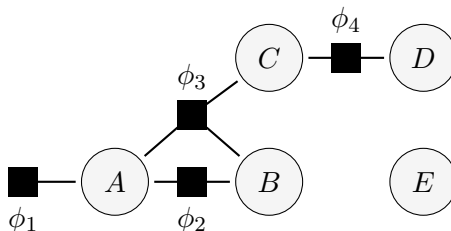
space of distributions
consistent with $\mathbf{m}_{\mathcal{B}}$
(which minimize *Inc*)



space of distributions
with independencies of \mathcal{B}
(which can be shown
to minimize *IDef*)

〈 maximum entropy result for BNs 〉

FACTOR GRAPHS



You need to add some intuition here.

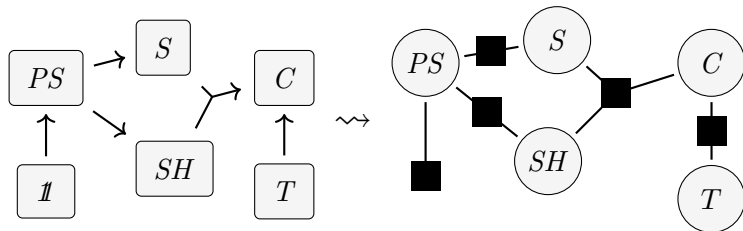
Definition

A *factor graph* Φ is a set of random variables $\mathcal{X} = \{X_i\}$ and *factors* $\{\phi_J: \mathcal{V}(X_J) \rightarrow \mathbb{R}_{\geq 0}\}_{J \in \mathcal{J}}$, where $X_J \subseteq \mathcal{X}$; define

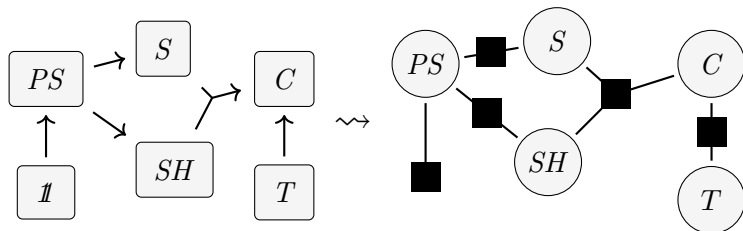
$$\Pr_{\Phi}(\vec{x}) = \frac{1}{Z_{\Phi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J),$$

where Z_{Φ} is the normalization constant.

PDGs AS FACTOR GRAPHS

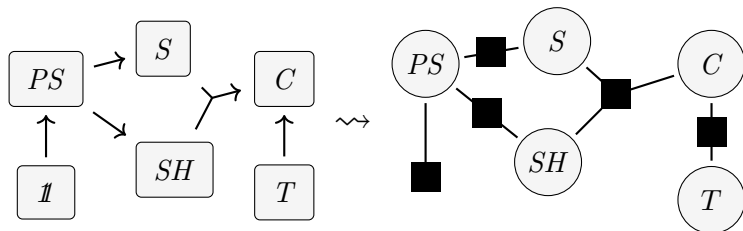


PDGs AS FACTOR GRAPHS



The cpds of a PDG are essentially factors. Are the semantics different?

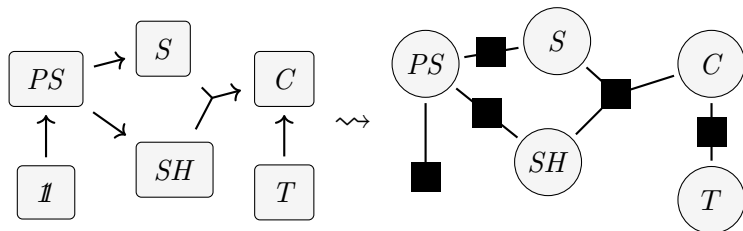
PDGs AS FACTOR GRAPHS



The cpds of a PDG are essentially factors. Are the semantics different?

Not for $\gamma = 1$.

PDGs AS FACTOR GRAPHS



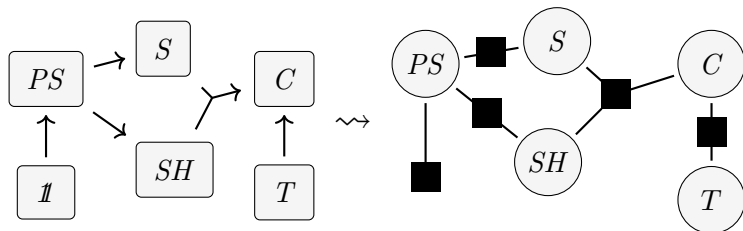
The cpds of a PDG are essentially factors. Are the semantics different?

Not for $\gamma = 1$.

Theorem

$\llbracket \mathcal{n} \rrbracket_1^* = \Pr_{\Phi_{\mathcal{n}}}$ for all unweighted PDGs \mathcal{n} .

PDGs AS FACTOR GRAPHS



The cpds of a PDG are essentially factors. Are the semantics different?

Not for $\gamma = 1$.

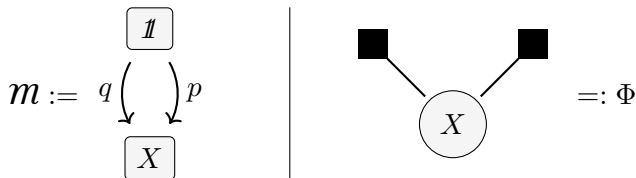
Theorem

$\llbracket \mathcal{N} \rrbracket_1^* = \text{Pr}_{\Phi_{\mathcal{N}}}$ for all unweighted PDGs \mathcal{N} .

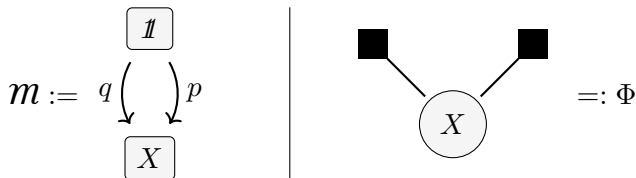
Theorem

For all unweighted PDGs \mathcal{N} and non-negative vectors \mathbf{v} over the edges of \mathcal{N} , and all $\gamma > 0$, we have that $\llbracket (\mathcal{N}, \mathbf{v}, \gamma \mathbf{v}) \rrbracket_{\gamma} = \gamma \text{GFE}_{(\Phi_{\mathcal{N}}, \mathbf{v})}$; consequently, $\llbracket (\mathcal{N}, \mathbf{v}, \gamma \mathbf{v}) \rrbracket_{\gamma}^* = \{\text{Pr}_{(\Phi_{\mathcal{N}}, \mathbf{v})}\}$.

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

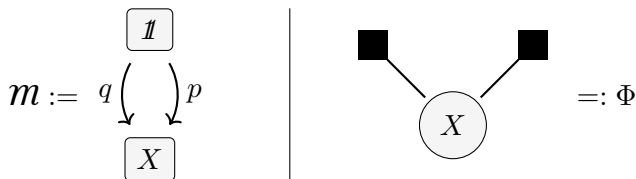


AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



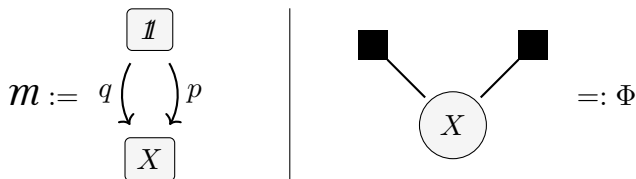
- If $p = q$, then $\llbracket m \rrbracket^* = p = q \dots$

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



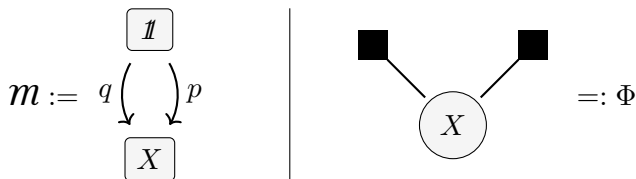
- If $p = q$, then $\llbracket m \rrbracket^* = p = q \dots$
- \dots but $\Pr_{\Phi} \propto p^2$

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



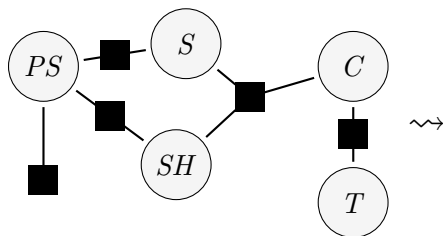
- If $p = q$, then $\llbracket m \rrbracket^* = p = q \dots$
- \dots but $\Pr_{\Phi} \propto p^2$
- Individual factors have *no probabilistic meaning*,

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

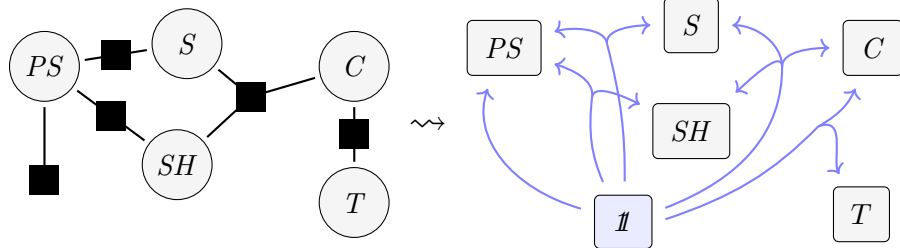


- If $p = q$, then $\llbracket m \rrbracket^* = p = q \dots$
- \dots but $\Pr_{\Phi} \propto p^2$
- Individual factors have *no probabilistic meaning*,
- a factor graph can fail to normalize, in which case it has no global semantics either.

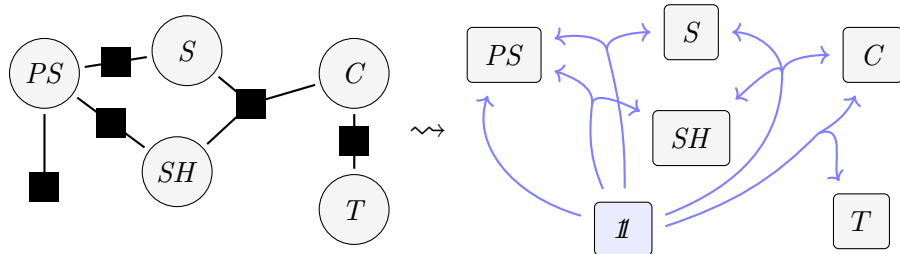
FACTOR GRAPHS AS PDGs



FACTOR GRAPHS AS PDGs



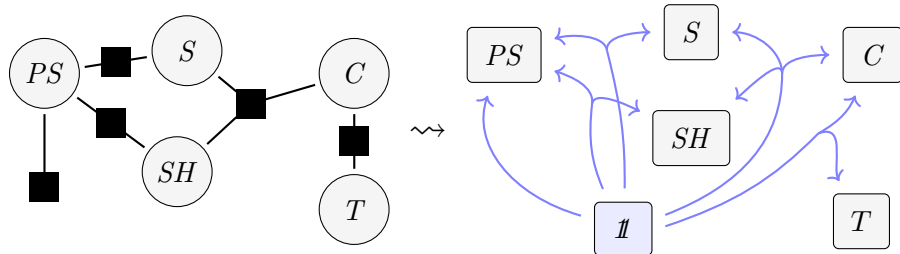
FACTOR GRAPHS AS PDGs



Theorem

$\Pr_{\Phi} = \llbracket n_{\Phi} \rrbracket_1^*$ for all factor graphs Φ .

FACTOR GRAPHS AS PDGs



Theorem

$\Pr_{\Phi} = \llbracket \mathbf{n}_{\Phi} \rrbracket_1^*$ for all factor graphs Φ .

Theorem

For all weighted factor graphs $\Psi = (\Phi, \theta)$ and all $\gamma > 0$, we have that $GFE_{\Psi} = 1/\gamma \llbracket \mathbf{m}_{\Psi, \gamma} \rrbracket_{\gamma} + C$ for some constant C , so \Pr_{Ψ} is the unique element of $\llbracket \mathbf{m}_{\Psi, \gamma} \rrbracket_{\gamma}^*$.

I suspect that
you'll lose your
audience with
this

< Add theorem: $\log Z_{\Phi} = \llbracket \mathbf{n}_{\Phi} \rrbracket$ >

Letting $x^{\mathbf{w}}$ and $y^{\mathbf{w}}$ denote the values of X and Y , respectively, in $\mathbf{w} \in \mathcal{V}(\mathcal{M})$, we have

$$\llbracket \mathcal{M} \rrbracket(\mu) = \mathbb{E}_{\mathbf{w} \sim \mu} \left\{ \sum_{X \xrightarrow{L} Y} \left[\overbrace{\beta_L \log \frac{1}{\mathbf{p}_L(y^{\mathbf{w}}|x^{\mathbf{w}})}}^{\text{log likelihood / cross entropy}} + \underbrace{(\alpha_L \gamma - \beta_L) \log \frac{1}{\mu(y^{\mathbf{w}}|x^{\mathbf{w}})}}_{\text{local regularization } (\beta_L > \alpha_L \gamma)} \right] - \underbrace{\gamma \log \frac{1}{\mu(\mathbf{w})}}_{\text{global regularization}} \right\}.$$

OUTLINE FOR SECTION 5

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

INFERENCE VIA INCONSISTENCY REDUCTION

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{I} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $\llbracket \mathcal{M}_{Y=y} \rrbracket^* = \llbracket \mathcal{M} \rrbracket^* \mid (Y=y)$.

Add some words for what this is saying (e.g., remind the reader of what $\llbracket \mathcal{M} \rrbracket^*$ is). And

INFERENCE VIA INCONSISTENCY REDUCTION

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{m} , simply add the edge $\mathbb{I} \xrightarrow{\delta_y} Y$ to get $\mathcal{m}_{Y=y}$.
Then $\llbracket \mathcal{m}_{Y=y} \rrbracket^* = \llbracket \mathcal{m} \rrbracket^* \mid (Y=y)$.

Querying $\Pr(Y \mid X)$ in a PDG \mathcal{m} .

- We can add $X \xrightarrow{p} Y$ to \mathcal{m} with a cpt p , to get \mathcal{m}^{+p} .

Text

\langle Convexity Result \rangle

\langle Hardness Result \rangle

INFERENCE VIA INCONSISTENCY REDUCTION

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{I} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $\llbracket \mathcal{M}_{Y=y} \rrbracket^* = \llbracket \mathcal{M} \rrbracket^* \mid (Y=y)$.

Querying $\Pr(Y \mid X)$ in a PDG \mathcal{M} .

- We can add $X \xrightarrow{p} Y$ to \mathcal{M} with a cpt p , to get \mathcal{M}^{+p} .
- The choice of cpd p that minimizes the inconsistency of \mathcal{M}^{+p} (which is strongly convex and smooth in p) is $\llbracket \mathcal{M} \rrbracket^*(Y \mid X)$,

⟨ Convexity Result ⟩

⟨ Hardness Result ⟩

INFERENCE VIA INCONSISTENCY REDUCTION

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{I} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $\llbracket \mathcal{M}_{Y=y} \rrbracket^* = \llbracket \mathcal{M} \rrbracket^* \mid (Y=y)$.

Querying $\Pr(Y \mid X)$ in a PDG \mathcal{M} .

- We can add $X \xrightarrow{p} Y$ to \mathcal{M} with a cpt p , to get \mathcal{M}^{+p} .
 - The choice of cpd p that minimizes the inconsistency of \mathcal{M}^{+p} (which is strongly convex and smooth in p) is $\llbracket \mathcal{M} \rrbracket^*(Y \mid X)$,
 - so oracle access to inconsistency yields fast inference by gradient descent.
- Make this a sub-bullet

⟨ Convexity Result ⟩

⟨ Hardness Result ⟩

OUTLINE FOR SECTION 6

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...
 - ▶ choice is made by instinct, tradition, or pragmatics, which makes results difficult to motivate, and vulnerable to overfitting.

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...
 - ▶ choice is made by instinct, tradition, or pragmatics, which makes results difficult to motivate, and vulnerable to overfitting.
- Choice of *model* admits more principled discussion.

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...
 - ▶ choice is made by instinct, tradition, or pragmatics, which makes results difficult to motivate, and vulnerable to overfitting.
- Choice of *model* admits more principled discussion.
 - ▶ and a possibly-inconsistent model comes with a natural objective.

This needs more intuition. What's the model doing? Why might it be inconsistent?

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...
 - ▶ choice is made by instinct, tradition, or pragmatics, which makes results difficult to motivate, and vulnerable to overfitting.
- Choice of *model* admits more principled discussion.
 - ▶ and a possibly-inconsistent model comes with a natural objective.

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...
 - ▶ choice is made by instinct, tradition, or pragmatics, which makes results difficult to motivate, and vulnerable to overfitting.
- Choice of *model* admits more principled discussion.
 - ▶ and a possibly-inconsistent model comes with a natural objective.

Surprising Result

Most standard objectives arise naturally as the inconsistency of the **obvious** PDG describing the situation.

obvious -> natural

INCONSISTENCY: THE UNIVERSAL LOSS

- Fruitful to cast AI as optimization. But what to optimize?
 - ▶ Cross Entropy, Square Loss, Accuracy, ...
 - ▶ choice is made by instinct, tradition, or pragmatics, which makes results difficult to motivate, and vulnerable to overfitting.
- Choice of *model* admits more principled discussion.
 - ▶ and a possibly-inconsistent model comes with a natural objective.

Surprising Result

Most standard objectives arise naturally as the inconsistency of the obvious PDG describing the situation.

Pedagogical Bonus

An intuitive visual language for reasoning about inequalities

This seems misplaced. “Inequalities” comes out of the blue.

SURPRISE AS INCONSISTENCY

Proposition

Consider a distribution over X with mass function $p(X)$. The surprise (or information content) $I_p(x) := -\log p(X=x)$ at seeing a sample x is the inconsistency of the pdg containing p and the event $X = x$, i.e.,

$$I_p(x) = \log \frac{1}{p(X=x)} = \langle\langle \xrightarrow{p} X \xleftarrow{X=x} \rangle\rangle.$$

Where did the $\langle\langle \rangle\rangle$ notation come from?

SURPRISE AS INCONSISTENCY

Proposition

Consider a distribution over X with mass function $p(X)$. The surprise (or information content) $I_p(x) := -\log p(X=x)$ at seeing a sample x is the inconsistency of the pdg containing p and the event $X = x$, i.e.,

$$I_p(x) = \log \frac{1}{p(X=x)} = \left\langle\!\left\langle \xrightarrow{p} \boxed{X} \xleftarrow{X=x} \right\rangle\!\right\rangle.$$

- PDG semantics just so happen to give the standard measure of compatibility between a sample and distribution.

SURPRISE AS INCONSISTENCY

Proposition

Consider a distribution over X with mass function $p(X)$. The surprise (or information content) $I_p(x) := -\log p(X=x)$ at seeing a sample x is the inconsistency of the pdg containing p and the event $X = x$, i.e.,

$$I_p(x) = \log \frac{1}{p(X=x)} = \left\langle\left\langle \xrightarrow{p} \boxed{X} \xleftarrow{X=x} \right\rangle\right\rangle.$$

- PDG semantics just so happen to give the standard measure of compatibility between a sample and distribution.
- Known as “surprise”, a particular kind of internal conflict.

VARIATIONS: SURPRISE AS INCONSISTENCY

Proposition (marginal information as inconsistency)

If $p(X, Z)$ is a joint distribution, the (marginal) information of the (partial) observation $X = x$ is given by

$$I_p(x) = \log \frac{1}{p(x)} = \left\langle\!\!\left\langle \begin{array}{c} \swarrow \quad \searrow \\ Z \quad X \\ \nwarrow \quad \nearrow \end{array} \right\rangle\!\!\right\rangle.$$

While you may be used to all these funny symbols, the listener won't be. You have to be careful not to go overboard here.

VARIATIONS: SURPRISE AS INCONSISTENCY

Proposition (marginal information as inconsistency)

If $p(X, Z)$ is a joint distribution, the (marginal) information of the (partial) observation $X = x$ is given by

$$I_p(x) = \log \frac{1}{p(x)} = \left\langle\!\left\langle \begin{array}{c} \swarrow \quad \searrow \\ Z \quad X \\ \nwarrow \quad \nearrow \end{array} \begin{array}{c} p \\ x \end{array} \right\rangle\!\right\rangle.$$

Proposition (cross entropy as inconsistency)

Given a dataset $\underline{\mathbf{x}}$, the cross entropy $\text{CE}(p, \underline{\mathbf{x}}) := -\frac{1}{|\underline{\mathbf{x}}|} \sum_{x \in \underline{\mathbf{x}}} \log p(x)$ is the inconsistency of the PDG containing p and the data distribution $\text{Pr}_{\underline{\mathbf{x}}}$, plus the entropy of the data distribution (constant in p). That is,

I still object to the !, and I assure you the listener will have no clue of what it means.

$$\text{CE}(p; \underline{\mathbf{x}}) = \left\langle\!\left\langle \begin{array}{c} \swarrow \quad \searrow \\ Z \quad X \\ \nwarrow \quad \nearrow \end{array} \begin{array}{c} p \\ \text{Pr}_{\underline{\mathbf{x}}}! \end{array} \right\rangle\!\right\rangle + H(\text{Pr}_{\underline{\mathbf{x}}}).$$

You're now officially overboard.

Proposition (Accuracy as Inconsistency)

Consider a predictor $h : X \rightarrow Y$ for true labels $f : X \rightarrow Y$, and a distribution $D(X)$. The inconsistency of believing all three is

What is accuracy?

$$\left\langle \begin{array}{c} \xrightarrow{D^{(\beta)}} \\ X \end{array} \begin{array}{c} \xrightarrow{h} \\ \xrightarrow{f} \\ Y \end{array} \right\rangle = -\beta \log \left(\text{accuracy}_{f,D}(h) \right) = \beta I_D[f = h].$$

Proposition (Mean Square Error as Inconsistency)

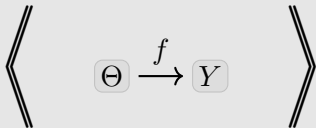
$$\begin{aligned}
 \left\langle\!\!\left\langle \begin{array}{ccc} & \mathcal{N}(f(x), 1) & \\ D! \nearrow & X & \searrow \\ & \mathcal{N}(g(x), 1) & \end{array} \right\rangle\!\!\right\rangle &= \left\langle\!\!\left\langle \begin{array}{ccccc} & & f & \mu_f & \mathcal{N}_1 \\ D! \nearrow & X & \nearrow & & \searrow \\ & & h & \mu_h & \nearrow \\ & & & & \mathcal{N}_1 \end{array} \right\rangle\!\!\right\rangle \\
 &= \mathbb{E}_D \left(f(X) - h(X) \right)^2 =: \text{MSE}(f, h)
 \end{aligned}$$

where $\mathcal{N}_1 = \mathcal{N}(-, 1)$ is the normal distribution with unit variance, and mean equal to its argument.

Proposition (Regularizers as priors)

Suppose you believe $Y \sim f_\theta(Y)$,

That is,

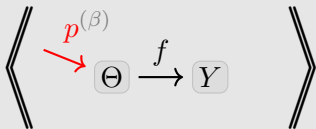


What does this have to do with belief?
I'm lost.

Proposition (Regularizers as priors)

Suppose you believe $Y \sim f_{\theta}(Y)$, have a prior $p(\theta)$,

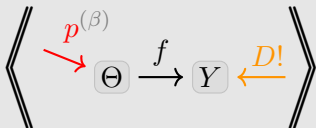
That is,



Proposition (Regularizers as priors)

Suppose you believe $Y \sim f_{\theta}(Y)$, have a prior $p(\theta)$, and have an empirical distribution $D(Y)$ which you trust.

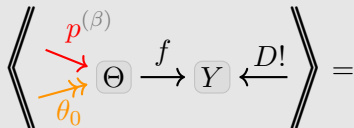
That is,



Proposition (Regularizers as priors)

Suppose you believe $Y \sim f_{\theta}(Y)$, have a prior $p(\theta)$, and have an empirical distribution $D(Y)$ which you trust. Then the inconsistency of also believing $\Theta = \theta_0$ is

That is,



Proposition (Regularizers as priors)

Suppose you believe $Y \sim f_\theta(Y)$, have a prior $p(\theta)$, and have an empirical distribution $D(Y)$ which you trust. Then the inconsistency of also believing $\Theta = \theta_0$ is the *regularized-cross entropy loss*, and controlled by the strength β_p of the prior. That is,

$$\left\langle \begin{array}{c} \text{red arrow } p^{(\beta)} \\ \text{black arrow } \theta_0 \end{array} \rightarrow \Theta \xrightarrow{f} Y \xleftarrow{D!} \right\rangle = \mathbb{E}_{y \sim D} \left[\log \frac{1}{f(y \mid \theta_0)} \right] + \beta \log \frac{1}{p(\theta_0)} - H(D)$$

Proposition (Regularizers as priors)

Suppose you believe $Y \sim f_\theta(Y)$, have a prior $p(\theta)$, and have an empirical distribution $D(Y)$ which you trust. Then the inconsistency of also believing $\Theta = \theta_0$ is the **regularized**-cross entropy loss, and controlled by the strength β_p of the prior. That is,

$$\left\langle \begin{array}{c} \text{red arrow } p^{(\beta)} \\ \text{black arrow } \theta_0 \end{array} \rightarrow \Theta \xrightarrow{f} Y \xleftarrow{D!} \right\rangle = \mathbb{E}_{y \sim D} \left[\log \frac{1}{f(y | \theta_0)} \right] + \beta \log \frac{1}{p(\theta_0)} - H(D)$$

Using a (discretized) unit gaussian as a prior, $p(\theta) = \frac{1}{k} \exp(-\frac{1}{2}\theta^2)$ for a normalization constant k , the RHS becomes

$$\underbrace{\mathbb{E}_D \left[\log \frac{1}{f(Y | \theta_0)} \right]}_{\text{Cross entropy loss of } f_\theta \text{ w.r.t. } D \text{ (data-fit cost of } \theta_0)} + \underbrace{\frac{\beta}{2} \theta_0^2}_{\text{red } \ell_2 \text{ regularizer (complexity cost of } \theta_0)} + \underbrace{\beta \log k - H(D)}_{\text{constant in } f \text{ and } \theta_0}.$$

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:

$e(Z | X)$: encodes X 's in a (small) latent space Z ;

Encodes in what way? What's a latent space. I'm lost.



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:

$e(Z | X)$: encodes X 's in a (small) latent space Z ;

$d(X | Z)$: generate samples of X from Z .



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:
 - $e(Z | X)$: encodes X 's in a (small) latent space Z ;
 - $d(X | Z)$: generate samples of X from Z .
- Objective:



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:
 - $e(Z | X)$: encodes X 's in a (small) latent space Z ;
 - $d(X | Z)$: generate samples of X from Z .
- Objective:
 - ▶ Want to minimize “reconstruction error” for each x

$$\text{Rec}(x) = - \mathbb{E}_{z \sim e|x} \log d(x | z)$$



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:

$e(Z | X)$: encodes X 's in a (small) latent space Z ;

$d(X | Z)$: generate samples of X from Z .

- Objective:

- ▶ Want to minimize “reconstruction error” for each x

$$\text{Rec}(x) = - \mathbb{E}_{z \sim e|x} \log d(x | z)$$

- ▶ you also have a prior $p(Z)$, which you want to match.



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:

$e(Z | X)$: encodes X 's in a (small) latent space Z ;

$d(X | Z)$: generate samples of X from Z .

- Objective:

- ▶ Want to minimize “reconstruction error” for each x

$$\text{Rec}(x) = - \mathbb{E}_{z \sim e|x} \log d(x | z)$$

- ▶ you also have a prior $p(Z)$, which you want to match.
- ▶ Together, maximize $\text{ELBO}_{p,e,d}(x) :=$

$$-D(e(Z|x) \parallel p(Z)) - \text{Rec}(x) = \mathbb{E}_{z \sim e|x} \left[\log \frac{p(z)d(x | z)}{e(z | x)} \right]$$

What does this expression represent? Why should you want to maximize it? I continue to be lost.



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:

$e(Z | X)$: encodes X 's in a (small) latent space Z ;

$d(X | Z)$: generate samples of X from Z .

- Objective:

- ▶ Want to minimize “reconstruction error” for each x

$$\text{Rec}(x) = - \mathbb{E}_{z \sim e|x} \log d(x | z)$$

- ▶ you also have a prior $p(Z)$, which you want to match.
- ▶ Together, maximize $\text{ELBO}_{p,e,d}(x) :=$

$$-D(e(Z|x) \parallel p(Z)) - \text{Rec}(x) = \mathbb{E}_{z \sim e|x} \left[\log \frac{p(z)d(x | z)}{e(z | x)} \right] \leq \log \Pr_{pd}(x)$$



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:

$e(Z | X)$: encodes X 's in a (small) latent space Z ;

$d(X | Z)$: generate samples of X from Z .

- Objective:

- ▶ Want to minimize “reconstruction error” for each x

$$\text{Rec}(x) = - \mathbb{E}_{z \sim e|x} \log d(x | z)$$

- ▶ you also have a prior $p(Z)$, which you want to match.
- ▶ Together, maximize $\text{ELBO}_{p,e,d}(x) :=$

$$-D(e(Z|x) \parallel p(Z)) - \text{Rec}(x) = \mathbb{E}_{z \sim e|x} \left[\log \frac{p(z)d(x | z)}{e(z | x)} \right] \leq \log \Pr_{pd}(x)$$

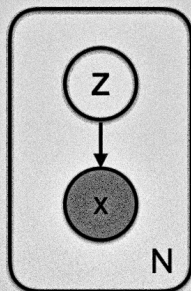


Note: not a graphical model (despite shoe-horning attempts)

I have no idea why I should think it's a graphical model or what this note is referring to.

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:



THE GRAPHICAL MODEL REPRESENTATION OF THE MODEL IN THE VARIATIONAL AUTOENCODER.
 THE LATENT VARIABLE Z IS A STANDARD NORMAL, AND THE DATA ARE DRAWN FROM $P(x|z)$. THE
 SHADED NODE FOR x DENOTES OBSERVED DATA. FOR BLACK AND WHITE IMAGES OF
 HANDWRITTEN DIGITS, THIS DATA LIKELIHOOD IS BERNOULLI DISTRIBUTED.

l) latent space Z ;
 from Z .

on error” for each x
 $\log d(x | z)$

h you want to match.
 $c) :=$

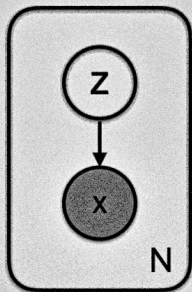
$$\mathbb{E}_{e|x} \left[\log \frac{p(z)d(x | z)}{e(z | x)} \right] \leq \log \Pr_{pd}(x)$$



Note: not a graphical model (despite shoehorning attempts)

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:



THE GRAPHICAL MODEL REPRESENTATION OF THE MODEL IN THE VARIATIONAL AUTO-ENCODER. THE LATENT VARIABLE Z IS A STANDARD NORMAL, AND THE DATA ARE DRAWN FROM $P(X|Z)$. THE SHADED NODE FOR X DENOTES OBSERVED DATA. FOR BLACK AND WHITE IMAGES OF HANDWRITTEN DIGITS, THIS DATA LIKELIHOOD IS BERNOLLI DISTRIBUTED.

1) latent space Z ;

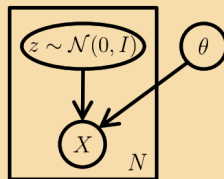
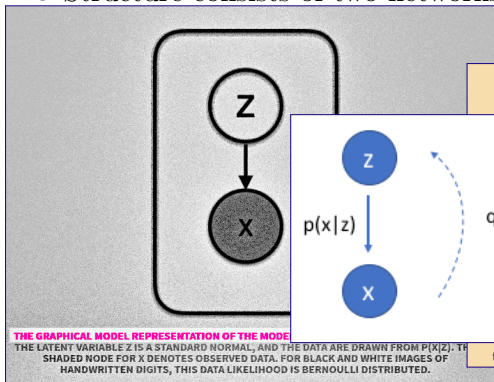


Figure 1: The standard VAE model represented as a graphical model. Note the conspicuous lack of any structure or even an “encoder” pathway: it is

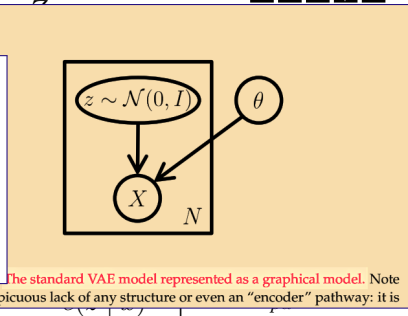
Note: not a graphical model (despite shoehorning attempts)

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:



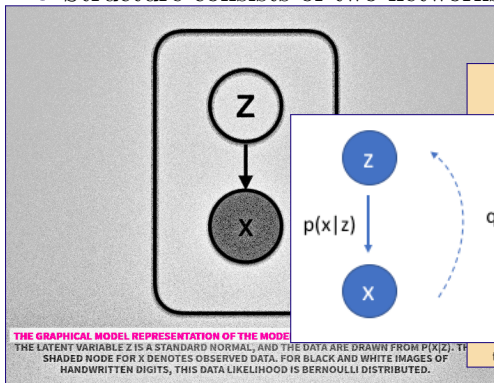
1) latent space Z ;



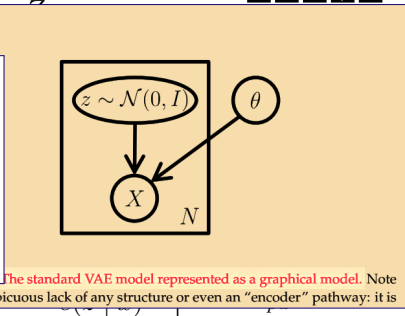
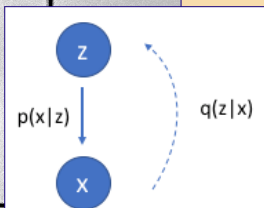
Note: not a graphical model (despite shoehorning attempts)

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:



1) latent space Z ;

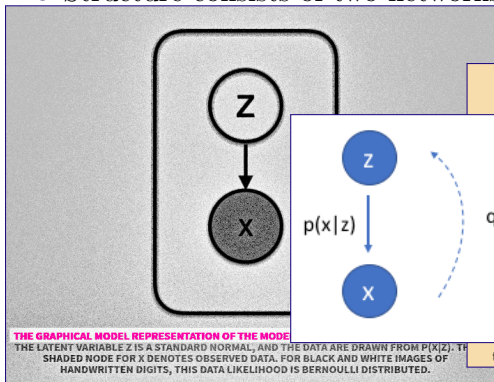


Note: not a graphical model (despite shoehorning attempts)

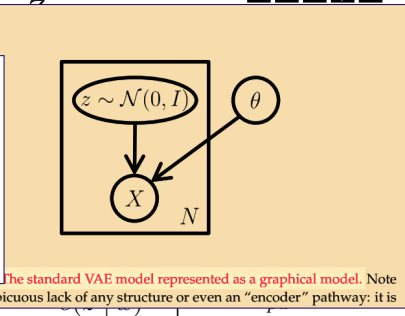
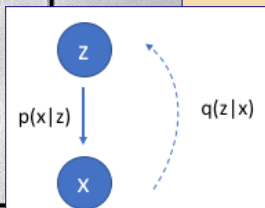
- $e(Z | X)$ has same target as $p(Z)$, so BN doesn't work

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two networks:



1) latent space Z ;



Note: not a graphical model (despite shoehorning attempts)

- $e(Z | X)$ has same target as $p(Z)$, so BN doesn't work
- The heart of the VaE is not its structure, but its objective.

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$$\langle \boxed{Z} \quad \boxed{X} \rangle = \text{ELBO}_{p,e,d}(x)$$

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encodes X in a latent space Z ;

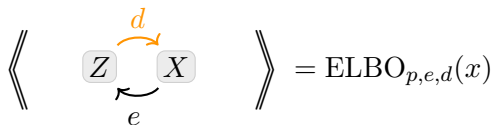
$$\left\langle \begin{array}{c} \boxed{Z} \quad \boxed{X} \\ \quad \quad \quad \curvearrowright \\ \quad \quad \quad e \end{array} \right\rangle = \text{ELBO}_{p,e,d}(x)$$

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encodes X in a latent space Z ;

$d(X | Z)$: generate samples of X from Z .



The diagram illustrates the structure of a Variational Auto-Encoder. It features two gray rectangular nodes labeled Z and X . A curved black arrow labeled e points from X to Z , representing the encoding process. A curved orange arrow labeled d points from Z to X , representing the decoding process. The entire diagram is enclosed in large double-angle brackets $\langle \rangle$. To the right of the brackets is the expression $= \text{ELBO}_{p,e,d}(x)$.

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encodes X in a latent space Z ;

$d(X | Z)$: generate samples of X from Z .

$p(Z)$: a prior over from Z .

The diagram illustrates the structure of a Variational Auto-Encoder (VAE). It features two nodes, Z and X , represented as light gray rounded rectangles. An orange arrow labeled p points from the left into node Z , representing the prior distribution. A curved arrow labeled d points from node Z to node X , representing the decoder. A curved arrow labeled e points from node X back to node Z , representing the encoder. The entire diagram is enclosed in large double-angle brackets $\langle\!\langle$ and $\rangle\!\rangle$. To the right of the brackets is the expression $= \text{ELBO}_{p,e,d}(x)$.

$$\langle\!\langle \begin{array}{c} \xrightarrow{p} Z \\ \xrightarrow{d} X \\ \xleftarrow{e} \end{array} \rangle\!\rangle = \text{ELBO}_{p,e,d}(x)$$

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encodes X in a latent space Z ;

$d(X | Z)$: generate samples of X from Z .

$p(Z)$: a prior over from Z .

- want to do a gradient step for a specific x .

I'm totally lost.

$$\left\langle \left\langle \begin{array}{c} \xrightarrow{p} Z \xrightarrow{d} X \xleftarrow{x} \\ \xleftarrow{e!} \end{array} \right\rangle \right\rangle = \text{ELBO}_{p,e,d}(x)$$

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encodes X in a latent space Z ;

$d(X | Z)$: generate samples of X from Z .

$p(Z)$: a prior over from Z .

- want to do a gradient step for a specific x .

Objective function is free:

$$\left\langle\left\langle \begin{array}{c} \xrightarrow{p} \boxed{Z} \begin{array}{c} \xrightarrow{d} \\ \xleftarrow{e!} \end{array} \boxed{X} \xleftarrow{x} \end{array} \right\rangle\right\rangle = \text{ELBO}_{p,e,d}(x)$$

A VERY USEFUL FACT

Believing more things can't make you any less inconsistent.

Lemma (monotonicity of inconsistency)

For all pdgs \mathcal{M} , \mathcal{M}' , and all $\gamma > 0$,

- 1 $\langle\langle \mathcal{M} \sqcup \mathcal{M}' \rangle\rangle_{\gamma} \geq \langle\langle \mathcal{M} \rangle\rangle_{\gamma}$.
- 2 If \mathcal{M} and \mathcal{M}' have respective confidence vectors β and β' , and $\beta \succeq \beta'$ (that is, $\beta_L \geq \beta'_L$ for all $L \in \mathcal{E}$), then $\langle\langle \mathcal{M} \rangle\rangle_{\gamma} \geq \langle\langle \mathcal{M}' \rangle\rangle_{\gamma}$.

VISUAL PROOF: THE VARIATIONAL BOUND

VISUAL PROOF: THE VARIATIONAL BOUND

$$\left\langle \begin{array}{c} \xrightarrow{p} \boxed{Z} \begin{array}{c} \xrightarrow{d} \boxed{X} \\ \xleftarrow{e!} \end{array} \xleftarrow{x} \end{array} \right\rangle = -\text{ELBO}_{p,e,d}(x).$$

VISUAL PROOF: THE VARIATIONAL BOUND

$$-\log \Pr_{p,d}(X=x) = \left\langle \begin{array}{c} p \rightarrow Z \xrightarrow{d} X \leftarrow x \end{array} \right\rangle \left\langle \begin{array}{c} p \rightarrow Z \xrightleftharpoons[e!]{d} X \leftarrow x \end{array} \right\rangle = -\text{ELBO}_{p,e,d}(x).$$

VISUAL PROOF: THE VARIATIONAL BOUND

There's no way you'll have the time for this.

$$\begin{aligned} -\log \Pr_{p,d}(X=x) &= \\ &\left\langle \begin{array}{c} p \rightarrow Z \xrightarrow{d} X \leftarrow x \end{array} \right\rangle \leq \left\langle \begin{array}{c} p \rightarrow Z \xrightleftharpoons[e!]{d} X \leftarrow x \end{array} \right\rangle \\ &= -\text{ELBO}_{p,e,d}(x). \end{aligned}$$

DIVERGENCES AND INCONSISTENCY

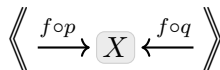
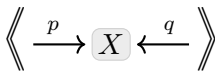
Lemma (Divergences and PDGs)

The PDG divergence $D_{(r,s)}^{\text{PDG}}(p, q)$, the inconsistency of a PDG containing $p(X)$ with confidence r and $q(X)$ with confidence s , is given by

$$D_{(r,s)}^{\text{PDG}}(p, q) := \left\langle\left\langle \frac{p}{(\beta:r)} \rightarrow X \leftarrow \frac{q}{(\beta:s)} \right\rangle\right\rangle = -(r+s) \log \sum_x \left(p(x)^r q(x)^s \right)^{\frac{1}{r+s}}.$$

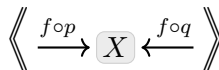
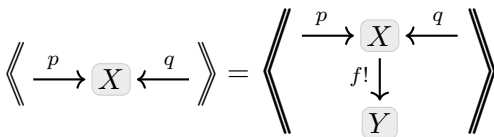
VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D(p \parallel q) \geq D(f \circ p \parallel f \circ q)$$



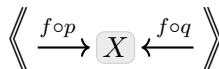
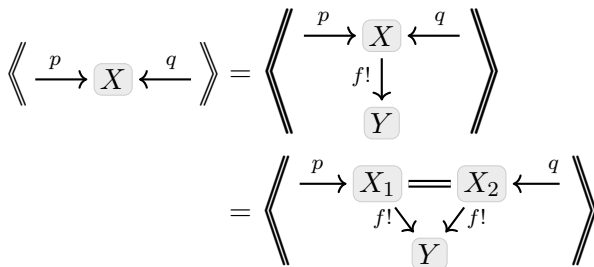
VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D(p \parallel q) \geq D(f \circ p \parallel f \circ q)$$



VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D(p \parallel q) \geq D(f \circ p \parallel f \circ q)$$

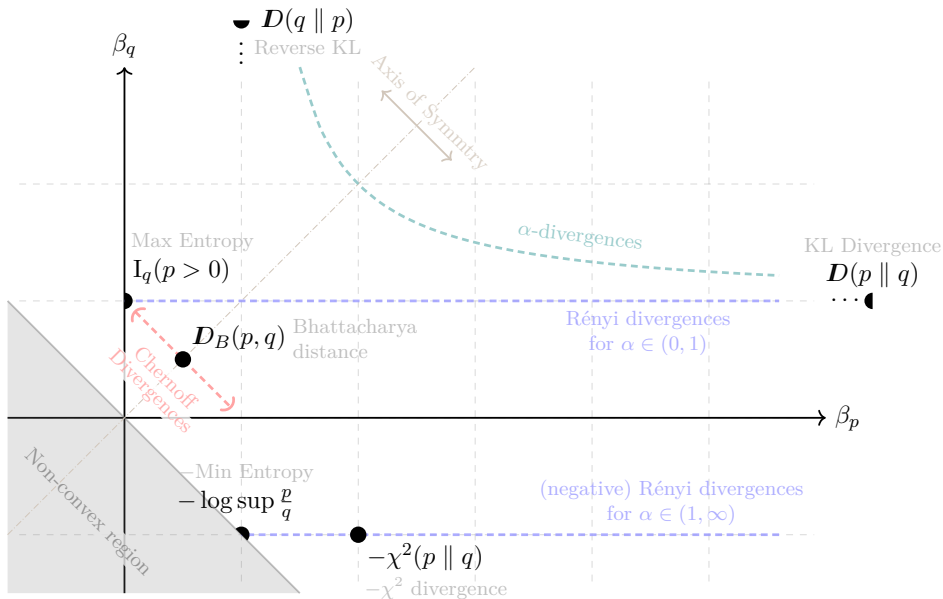


VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D(p \parallel q) \geq D(f \circ p \parallel f \circ q)$$

$$\begin{aligned}
 \langle\!\langle p \rightarrow X \leftarrow q \rangle\!\rangle &= \langle\!\langle \begin{array}{c} p \rightarrow X \leftarrow q \\ f! \downarrow \\ Y \end{array} \rangle\!\rangle \\
 &= \langle\!\langle \begin{array}{c} p \rightarrow X_1 = X_2 \leftarrow q \\ f! \searrow \swarrow f! \\ Y \end{array} \rangle\!\rangle \\
 &\geq \langle\!\langle \begin{array}{c} p \rightarrow X_1 \quad X_2 \leftarrow q \\ f! \searrow \swarrow f! \\ Y \end{array} \rangle\!\rangle = \langle\!\langle f \circ p \rightarrow X \leftarrow f \circ q \rangle\!\rangle
 \end{aligned}$$

DIVERGENCES AS INCONSISTENCIES



OUTLINE FOR SECTION 7

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

⟨ Schema Picture ⟩
⟨ Add Universal Relation Theorem ⟩

OUTLINE FOR SECTION 8

1 MODELING EXAMPLES

- The Simplest Inconsistency
- Differences from BNs
- PDG Union and Restriction

2 SYNTAX

- Formal Definitions of PDGs
- PDGs as diagrams of the Markov Category

3 SEMANTICS OF PDGs

4 PDGs AND OTHER

GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

5 INFERENCE

6 INCONSISTENCY AS LOSS

- Motivation
- Standard Metrics as Inconsistency
- Variational AutoEncoders
- Inconsistency and Statistical Divergences

7 DATABASES

8 OPEN PROBLEMS

OPEN PROBLEMS AND FUTURE WORK

⟨ INCOMPLETE ⟩

- - ▶ Trace Semantics
 - ▶ Composition

SUMMARY

⟨ update with second half ⟩

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.

SUMMARY

⟨ update with second half ⟩

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular; to combine info from two sources, simply take a PDG union. This incorporates new data (edge cpds) and concepts (nodes) without affecting previous information.

SUMMARY

⟨ update with second half ⟩

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular; to combine info from two sources, simply take a PDG union. This incorporates new data (edge cpds) and concepts (nodes) without affecting previous information.
- cleanly separate quantitative info (the cpds) from qualitative info (the edges), with variable confidence in both (the weights β and α). This is captured by terms *Inc* and *IDef* in our scoring function.

SUMMARY

⟨ update with second half ⟩

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular; to combine info from two sources, simply take a PDG union. This incorporates new data (edge cpds) and concepts (nodes) without affecting previous information.
- cleanly separate quantitative info (the cpds) from qualitative info (the edges), with variable confidence in both (the weights β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- have (several) natural semantics; one of them allows us to pick out a unique distribution. Using this distribution, PDGs can capture BNs and factor graphs.

SUMMARY

⟨ update with second half ⟩

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular; to combine info from two sources, simply take a PDG union. This incorporates new data (edge cpds) and concepts (nodes) without affecting previous information.
- cleanly separate quantitative info (the cpds) from qualitative info (the edges), with variable confidence in both (the weights β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- have (several) natural semantics; one of them allows us to pick out a unique distribution. Using this distribution, PDGs can capture BNs and factor graphs.

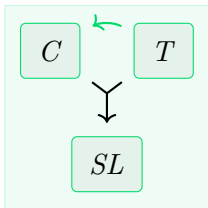
But there is much more to be done!

⟨ return to initial slide, but with more
conflicts ⟩

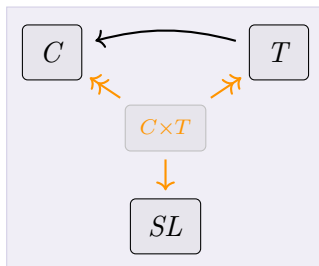
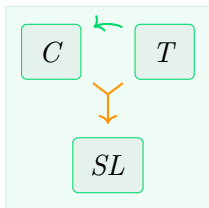
OUTLINE FOR SECTION 9

9 APPENDIX

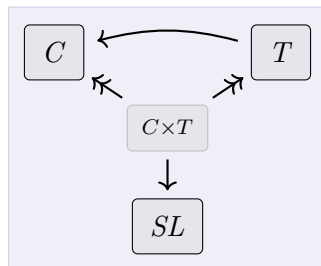
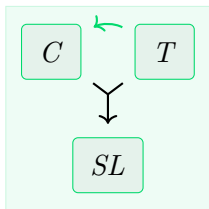
HYPER-GRAPHS? OR MERELY GRAPHS?



HYPER-GRAPHS? OR MERELY GRAPHS?

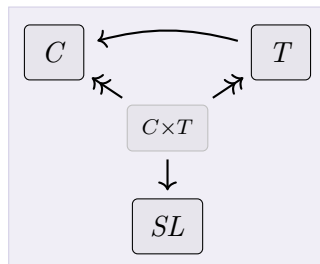
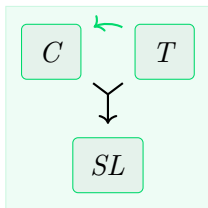


HYPER-GRAPHS? OR MERELY GRAPHS?



- This widget expands state space, but graphs are simpler.

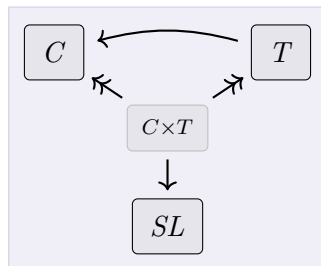
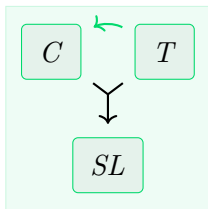
HYPER-GRAPHS? OR MERELY GRAPHS?



- This widget expands state space, but graphs are simpler.
- There is a natural correspondence

joint distributions \iff expanded joint distributions
satisfying coherence constraints

HYPER-GRAPHS? OR MERELY GRAPHS?



- This widget expands state space, but graphs are simpler.
- There is a natural correspondence

joint distributions \iff expanded joint distributions
satisfying coherence constraints

(working directly with hypergraphs is also possible)

ILLUSTRATIONS OF $IDef$

