

PROBABILISTIC DEPENDENCY GRAPHS AND INCONSISTENCY

HOW TO MODEL, MEASURE, AND MITIGATE INTERNAL CONFLICT

Oliver Richardson

Cornell University
Department of Computer Science

September 2021

OUTLINE FOR SECTION 1

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

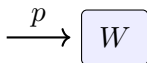
8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

The standard way of modeling an agent with uncertainty:

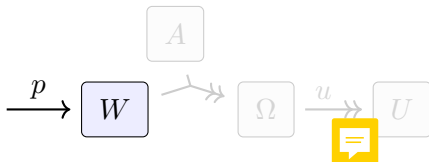
The standard way of modeling an agent with uncertainty:

- a probability distribution $p : \Delta W$ over worlds W ,



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).



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.

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

Why **build a synthetic agent** that can be inconsistent?

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

Why **build a synthetic agent** that can be inconsistent?

- Ensuring consistency is restrictive and expensive.

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

Why **build a synthetic agent** that can be inconsistent?

- Ensuring consistency is restrictive and expensive.
- Why write assertions, that could disagree with the code?

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

Why **build a synthetic agent** that can be inconsistent?

- Ensuring consistency is restrictive and expensive.
- Why write assertions, that could disagree with the code?
- Less reliant on flawless design.

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

Why **build a synthetic agent** that can be inconsistent?

- Ensuring consistency is restrictive and expensive.
- Why write assertions, that could disagree with the code?
- Less reliant on flawless design.

*A man with a watch knows what time it is;
a man with two watches is never sure.
(Segal's Law)*

WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;
 - ▶ also want to understand the process of resolving it.
- Internal conflict seems to drive important epistemic change.

Why **build a synthetic agent** that can be inconsistent?

- Ensuring consistency is restrictive and expensive.
- Why write assertions, that could disagree with the code?
- Less reliant on flawless design.



*man with a watch knows what time it is;
a man with two watches is never sure.
(Segal's Law)*



Freedom from perfect consistency frees up **a lot** of computation, but demands the ability to recognize and address internal conflict.

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

Probabilistic Dependency Graphs (PDGs),
a new class of graphical model designed to model inconsistent beliefs.

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

Probabilistic Dependency Graphs (PDGs),
a new class of graphical model designed to model inconsistent beliefs.

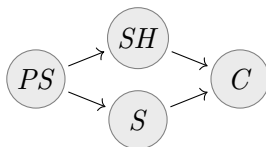
In doing so, we get much more ...

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}

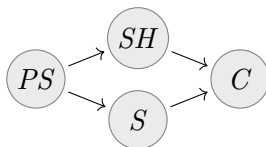
an independence relation on variables

- $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid \mathbf{Pa}(X)$, 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}

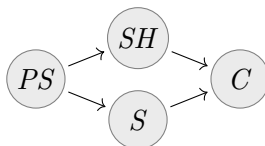
an independence relation on variables

- $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid \mathbf{Pa}(X)$, 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 2

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from 
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

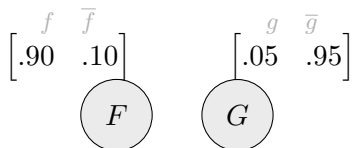
SIMPLE EXAMPLE: FLOOMPS AND GUNS

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).

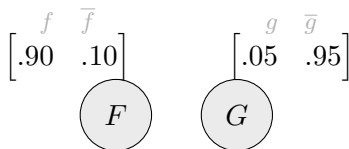
BN



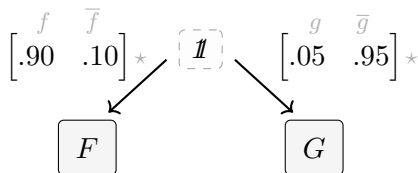
SIMPLE EXAMPLE: FLOOMPS AND GUNS

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

BN



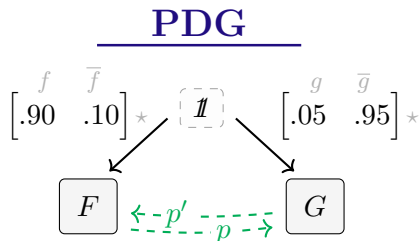
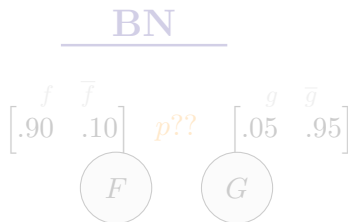
PDG



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



SIMPLE EXAMPLE: FLOOMPS AND GUNS

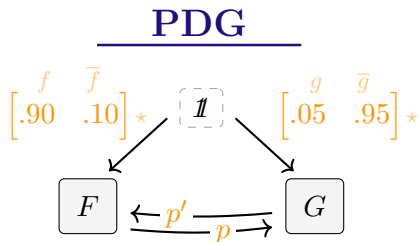
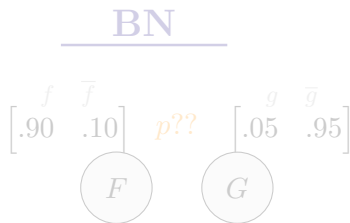


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.

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

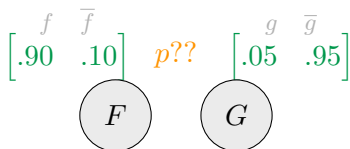


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent

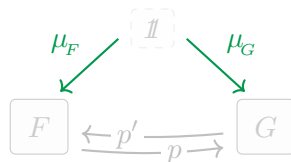


SIMPLE EXAMPLE: FLOOMPS AND GUNS

BN



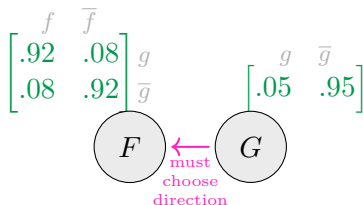
PDG



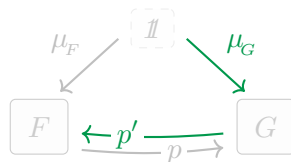
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
 - ▶ ...but BNs must resolve inconsistency first,
which may break symmetry and irrecoverably lose information.

SIMPLE EXAMPLE: FLOOMPS AND GUNS

BN



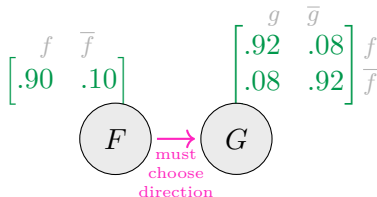
PDG



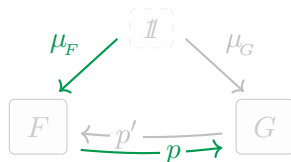
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
 - ▶ ...but BNs must resolve inconsistency first,
which may **break symmetry** and irrecoverably lose information.

SIMPLE EXAMPLE: FLOOMPS AND GUNS

BN



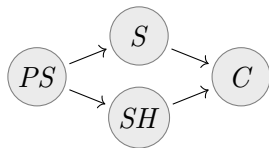
PDG



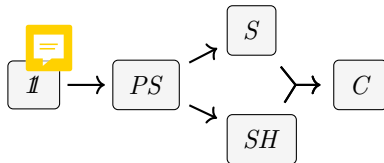
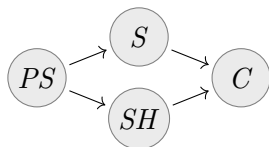
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
 - ▶ ...but BNs must resolve inconsistency first,
which may **break symmetry** and irrecoverably lose information.



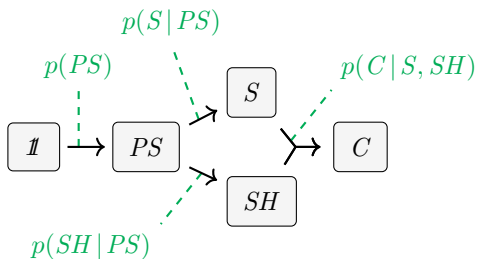
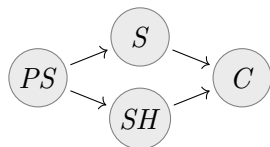
BAYESIAN NETWORKS AS PDGs



BAYESIAN NETWORKS AS PDGs



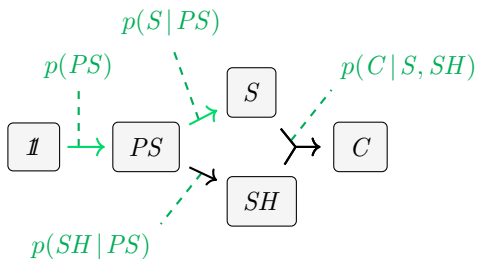
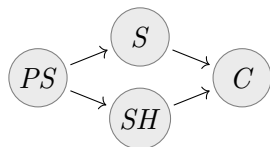
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

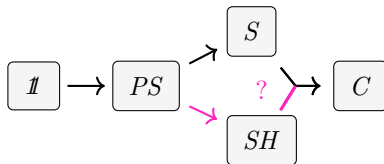
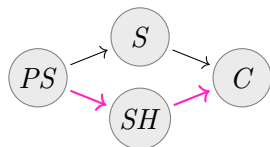
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

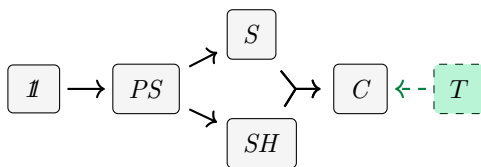
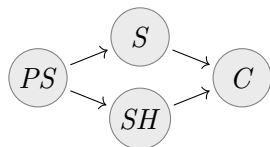
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

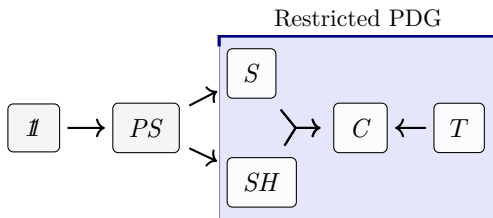
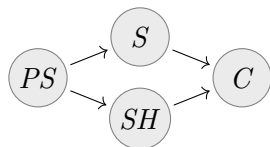
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

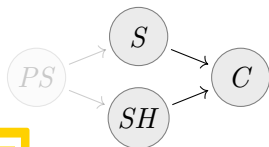
BAYESIAN NETWORKS AS PDGs



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

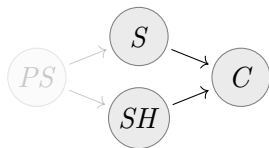


st now give distributions on SH and S , or distinguish them as “observed” (a *conditional* BN).

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).



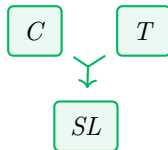
In a qualitative BN: *removing data results in new knowledge: $A \perp\!\!\!\perp C$.*



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!

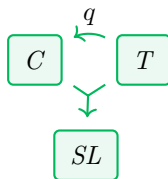
COMBINING PDGs



Grok wants to be supreme leader (SL).

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

COMBINING PDGs

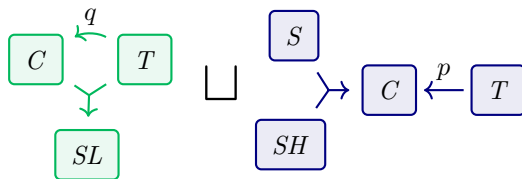


Grok wants to be supreme leader (SL).

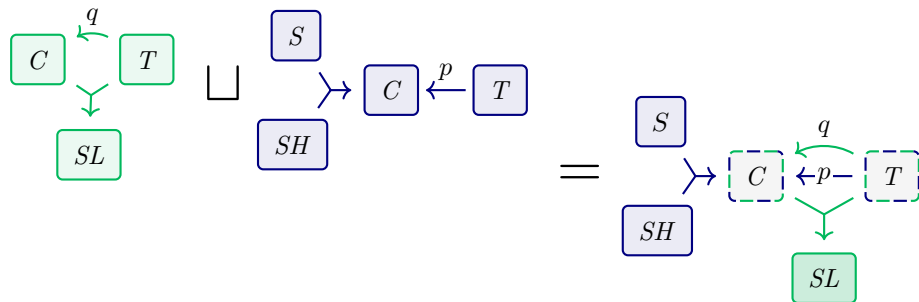
- She notices that those who use tanning beds have more power, unless they get cancer

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

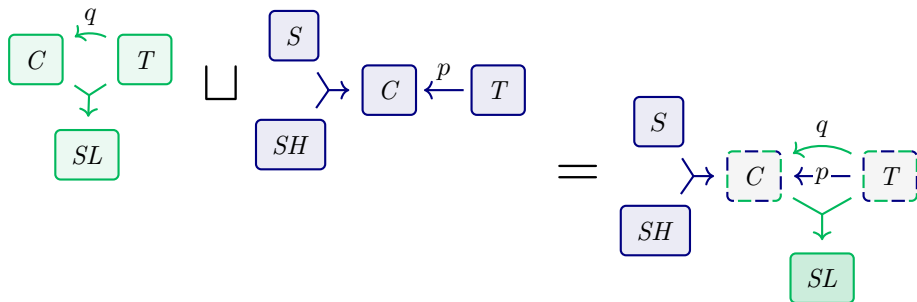
COMBINING PDGs



COMBINING PDGs

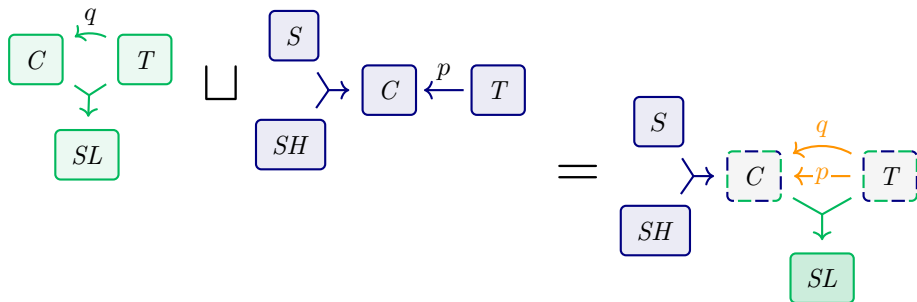


COMBINING PDGs



- 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 3

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

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 $\{X \xrightarrow{L} Y\}$,

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

► (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 $\{X \xrightarrow{L} Y\}$,

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

► (or hyper-edges)

\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 $\{X \xrightarrow{L} Y\}$,

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

► (or hyper-edges)

\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 $\{X \xrightarrow{L} Y\}$,

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

► (or hyper-edges)

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

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

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

OUTLINE FOR SECTION 4

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

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

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

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);

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

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

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}

...

► trace semantics

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

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

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}

THE SCORING FUNCTION

$$\llbracket m \rrbracket_{\gamma}(\mu) := Incm(\mu) + \gamma IDefm(\mu)$$

THE SCORING FUNCTION

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

Intuition: Measure μ 's violation of \mathbf{m} 's cpds.

THE SCORING FUNCTION

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

Definition (*Inc*)

The *incompatibility* of a joint distribution μ with \mathbf{m} is given by

$$Inc_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} D(\mu_{Y|X} \parallel \mathbf{p}_L)$$

THE SCORING FUNCTION

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

Definition (*Inc*)

The *incompatibility* of a joint distribution μ with \mathbf{m} is given by

$$\text{Inc}_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbf{D}(\mu_{Y|X} \parallel \mathbf{p}_L)$$

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_{\gamma}(\mu) := \text{Inc}_{\mathbf{m}}(\mu) + \gamma \text{IDef}_{\mathbf{m}}(\mu)$$

Definition (*Inc*)

The *incompatibility* of a joint distribution μ with \mathbf{m} is given by

$$\text{Inc}_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbf{D}(\mu_{Y|X} \parallel \mathbf{p}_L)$$

$$D(\mu \parallel \nu) = \sum_{w \in \text{Supp } \mu} \mu(w) \log \frac{\mu(w)}{\nu(w)}$$

the relative entropy
(KL Divergence)
from ν to μ .

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_{\gamma}(\mu) := \text{Inc}_{\mathbf{m}}(\mu) + \gamma \text{IDef}_{\mathbf{m}}(\mu)$$

Definition (*Inc*)

The *incompatibility* of a joint distribution μ with \mathbf{m} is given by

$$\text{Inc}_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbb{E}_{x \sim \mu_X} D\left(\mu(Y \mid X=x) \parallel \mathbf{p}_L(x)\right).$$

$$D(\mu \parallel \nu) = \sum_{w \in \text{Supp } \mu} \mu(w) \log \frac{\mu(w)}{\nu(w)}$$

the relative entropy
(KL Divergence)
from ν to μ .

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_{\gamma}(\mu) := \text{Inc}\mathbf{m}(\mu) + \gamma \text{IDef}\mathbf{m}(\mu)$$

Definition (*Inc*)

The *incompatibility* of a joint distribution μ with \mathbf{m} is given by

$$\begin{aligned} \text{Inc}\mathbf{m}(\mu) &:= \sum_{X \xrightarrow{L} Y} \beta_L \mathbb{E}_{x \sim \mu_X} D\left(\mu(Y \mid X=x) \parallel \mathbf{p}_L(x)\right). \\ &= \mathbb{E}_{\mu} \sum_{X \xrightarrow{L} Y} \beta_L \left(\underbrace{I_{\mathbf{p}_L}}_{\substack{\text{code length,} \\ \text{optimized for } \mathbf{p}_L \\ \text{to communicate} \\ Y \text{ given } X}} - \underbrace{I_{\mu}}_{\substack{\text{code length,} \\ \text{optimized for } \mu \\ \text{to communicate} \\ Y \text{ given } X}} \right) \end{aligned}$$



THE SCORING FUNCTION

$$\llbracket m \rrbracket_\gamma(\mu) := Incm(\mu) + \gamma \textcolor{red}{IDef}_m(\mu)$$

Intuition: each edge $X \xrightarrow{L} Y$ indicates that Y is determined (perhaps noisily) by X alone.

THE SCORING FUNCTION

$$\llbracket m \rrbracket_{\gamma}(\mu) := Incm(\mu) + \gamma \textcolor{red}{IDef}_m(\mu)$$

Intuition: each edge $X \xrightarrow{L} Y$ indicates that Y is determined (perhaps noisily) by X alone.

So a μ with uncertainty in Y after X is known (beyond pure noise) is qualitatively worse.

THE SCORING FUNCTION

$$\llbracket m \rrbracket_\gamma(\mu) := Incm(\mu) + \gamma \textcolor{red}{IDef}_m(\mu)$$

Intuition: each edge $X \xrightarrow{L} Y$ indicates that Y is determined (perhaps noisily) by X alone.

So a μ with $\overbrace{\text{uncertainty in } Y \text{ after } X \text{ is known}}^{\text{measured by } H(Y \mid X)}$
(beyond $\underbrace{\text{pure noise}}_{H(\mu)}$) is qualitatively worse.

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_{\gamma}(\mu) := Inc_{\mathbf{m}}(\mu) + \gamma \textcolor{red}{IDef}_{\mathbf{m}}(\mu)$$

Definition (*IDef*)

The *information deficit* of a distribution μ with respect to \mathbf{m} is

$$IDef_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \alpha_L H_{\mu}(Y | X) - H(\mu).$$

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_\gamma(\mu) := Inc_{\mathbf{m}}(\mu) + \gamma \textcolor{red}{IDef}_{\mathbf{m}}(\mu)$$

Definition (*IDef*)

The *information deficit* of a distribution μ with respect to \mathbf{m} is

$$IDef_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \alpha_L H_\mu(Y | X) - \underbrace{H(\mu)}_{\text{(a) \# bits needed to determine all variables}}.$$

(a) # bits needed to determine all variables

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_\gamma(\mu) := Inc\mathbf{m}(\mu) + \gamma \textcolor{red}{IDef}_{\mathbf{m}}(\mu)$$

Definition (*IDef*)

The *information deficit* of a distribution μ with respect to \mathbf{m} is

(b) # bits required to separately determine each target, knowing the source

$$IDef_{\mathbf{m}}(\mu) := \overbrace{\sum_{X \xrightarrow{L} Y} \alpha_L H_\mu(Y|X)} - \underbrace{H(\mu)}.$$

(a) # bits needed to determine all variables

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_\gamma(\mu) := Inc_{\mathbf{m}}(\mu) + \gamma \textcolor{red}{IDef}_{\mathbf{m}}(\mu)$$

Definition (*IDef*)

The *information deficit* of a distribution μ with respect to \mathbf{m} is

(b) # bits required to separately determine each target, knowing the source

$$IDef_{\mathbf{m}}(\mu) := \overbrace{\sum_{X \xrightarrow{L} Y} \alpha_L H_\mu(Y|X)} - \underbrace{H(\mu)}.$$

(a) # bits needed to determine all variables

THE SCORING FUNCTION

$$\llbracket \mathbf{m} \rrbracket_\gamma(\mu) := \text{Inc}_{\mathbf{m}}(\mu) + \gamma \text{IDef}_{\mathbf{m}}(\mu)$$

tradeoff parameter $\gamma \geq 0$

Definition (*Inc*)

The *incompatibility* of μ with \mathbf{m} :

$$\text{Inc}_{\mathbf{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbf{D}(\mu_{Y|X} \parallel \mathbf{p}_L)$$

Definition (*IDef*)

The \mathbf{m} -*information deficit* of μ :

bits to separately determine each target, knowing the source

$$\text{IDef}_{\mathbf{m}}(\mu) = \sum_{X \xrightarrow{L} Y} \alpha_L \mathbf{H}_\mu(Y|X) - \underbrace{\mathbf{H}(\mu)}_{\text{# bits to determine all vars}}$$

bits to determine all vars

THE SCORING FUNCTION

- A BN strictly enforces the qualitative picture (large γ)

$$\llbracket \mathcal{m} \rrbracket_{\gamma}(\mu) := \text{Inc}_{\mathcal{m}}(\mu) + \gamma \text{IDef}_{\mathcal{m}}(\mu)$$

Definition (*Inc*)

The *incompatibility* of μ with \mathcal{m} :

$$\text{Inc}_{\mathcal{m}}(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbf{D}(\mu_{Y|X} \parallel \mathbf{p}_L)$$

Definition (*IDef*)


The \mathcal{m} -*information deficit* of μ :

bits to separately determine each target, knowing the source

$$\text{IDef}_{\mathcal{m}}(\mu) = \sum_{X \xrightarrow{L} Y} \alpha_L \mathbf{H}_{\mu}(Y|X) - \underbrace{\mathbf{H}(\mu)}_{\text{# bits to determine all vars}}$$

bits to determine all vars

THE SCORING FUNCTION

- A BN strictly enforces the qualitative picture (large γ) 
- we are interested in the quantitative limit (small γ)

$$\llbracket \mathcal{M} \rrbracket_{\gamma}(\mu) := \text{Inc}_{\mathcal{M}}(\mu) + \gamma \text{IDef}_{\mathcal{M}}(\mu)$$

Definition (*Inc*)

The *incompatibility* of μ with \mathcal{M} :

$$\text{Inc}_{\mathcal{M}}(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbf{D}(\mu_{Y|X} \parallel \mathbf{p}_L)$$

Definition (*IDef*)

The \mathcal{M} -*information deficit* of μ :

bits to separately determine each target, knowing the source

$$\text{IDef}_{\mathcal{M}}(\mu) = \sum_{X \xrightarrow{L} Y} \underbrace{\alpha_L \mathbf{H}_{\mu}(Y|X)}_{\text{# bits to separately determine each target, knowing the source}} - \underbrace{\mathbf{H}(\mu)}_{\text{# bits to determine all vars}}$$

bits to determine all vars

PROPERTIES OF THE OPTIMAL DISTRIBUTION

Proposition (uniqueness for small γ)

- 1 If $0 < \gamma \leq \min_L \beta_L^m$, then $\llbracket m \rrbracket_\gamma^*$ is a singleton.
- 2 $\lim_{\gamma \rightarrow 0} \llbracket m \rrbracket_\gamma^*$ exists and is a singleton.

PROPERTIES OF THE OPTIMAL DISTRIBUTION



Proposition (uniqueness for small γ)

- 1 If $0 < \gamma \leq \min_L \beta_L^m$ then $\llbracket m \rrbracket_\gamma^*$ is a singleton.
- 2 $\lim_{\gamma \rightarrow 0} \llbracket m \rrbracket_\gamma^*$ exists and is a singleton.

This lets us define $\llbracket m \rrbracket^* := \text{unique element } \left(\lim_{\gamma \rightarrow 0} \llbracket m \rrbracket_\gamma^* \right)$.

PROPERTIES OF INCONSISTENCY

$$\langle\!\langle m \rangle\!\rangle_\gamma := \inf_{\mu} \llbracket m \rrbracket_\gamma$$

Nice properties for minimization:

- The function $\gamma \mapsto \langle\!\langle m \rangle\!\rangle_\gamma$ is continuous for all γ
- The function $p \mapsto \langle\!\langle m \sqcup p \rangle\!\rangle_\gamma$ is smooth and strictly convex on its interior.

OUTLINE FOR SECTION 5

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

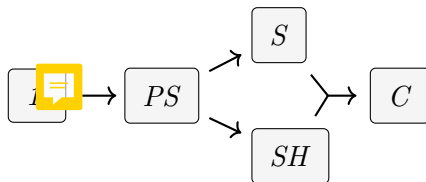
8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

► More Semantics Properties



CAPTURING BAYESIAN NETWORKS



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 $\Pr_{\mathcal{B}}$ is the distribution it specifies, then for all $\gamma > 0$ and all vectors β ,

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

CAPTURING BAYESIAN NETWORKS

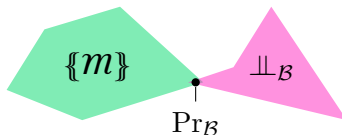
For a BN \mathcal{B} with N nodes and a vector $\beta \in \mathbb{R}^N$, let $\mathbf{m}_{\mathcal{B},\beta}$ be the PDG corresponding to \mathcal{B} , with $\alpha = \mathbf{1}$, and the given vector β of confidences.

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 ,

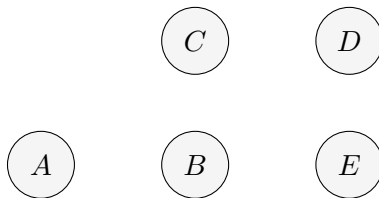
$$\llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket_{\gamma}^* = \{\text{Pr}_{\mathcal{B}}\}, \quad \text{and thus} \quad \llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket^* = \text{Pr}_{\mathcal{B}}.$$

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*)

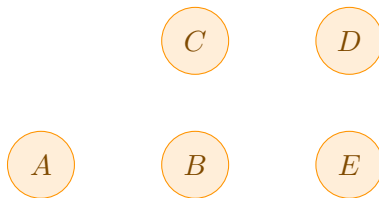
FACTOR GRAPHS



Definition

A *factor graph* Φ is

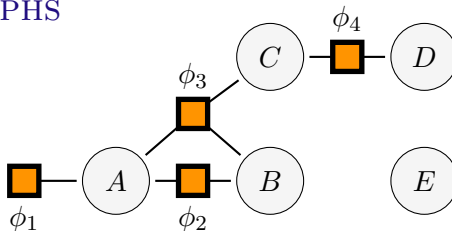
FACTOR GRAPHS



Definition

A *factor graph* Φ is a set of **variables** $\mathcal{X} = \{X_i\}$,

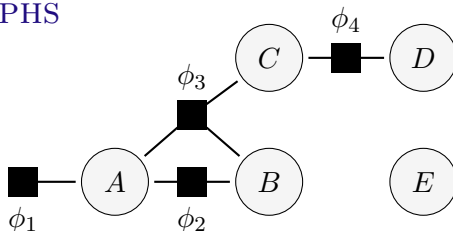
FACTOR GRAPHS



Definition

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

FACTOR GRAPHS

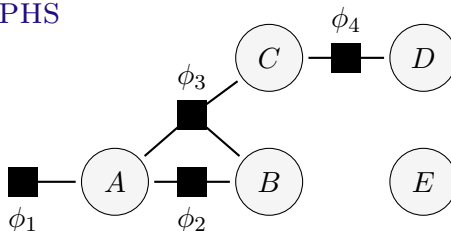


Definition

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

$$\Pr_{\Phi}(\vec{x}) := \frac{1}{Z_{\Phi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J), \quad \text{where } Z_{\Phi} \text{ is the normalization constant.}$$

FACTOR GRAPHS



Definition

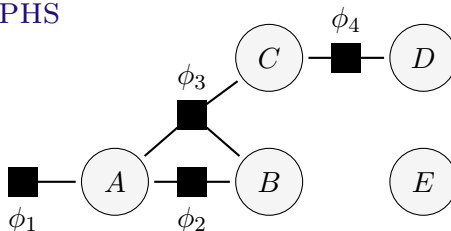
A *factor graph* Φ is a set of variables $\mathcal{X} = \{X_i\}$, and factors $\{\phi_J: \mathcal{V}(X_J) \rightarrow \mathbb{R}_{\geq 0}\}_{J \in \mathcal{J}}$, with $X_J \subseteq \mathcal{X}$; Φ defines a distribution

$$\Pr_{\Phi}(\vec{x}) := \frac{1}{Z_{\Phi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J), \quad \text{where } Z_{\Phi} \text{ is the normalization constant.}$$

Φ defines a standard scoring function

$$VFE_{\Phi}(\mu) := \mathbb{E}_{\mu} \left[- \sum_{J \in \mathcal{J}} \log \phi_J(X_J) \right] - H(\mu)$$

FACTOR GRAPHS



Definition

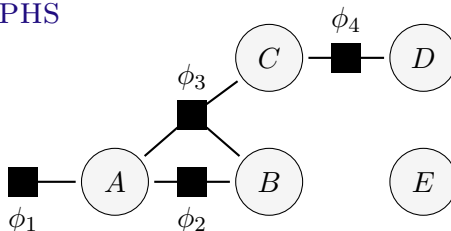
A *factor graph* Φ is a set of variables $\mathcal{X} = \{X_i\}$, and factors $\{\phi_J: \mathcal{V}(X_J) \rightarrow \mathbb{R}_{\geq 0}\}_{J \in \mathcal{J}}$, with $X_J \subseteq \mathcal{X}$; Φ defines a distribution

$$\Pr_{\Phi}(\vec{x}) := \frac{1}{Z_{\Phi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J), \quad \text{where } Z_{\Phi} \text{ is the normalization constant.}$$

Φ defines a standard scoring function

$$VFE_{\Phi}(\mu) := \mathbb{E}_{\mu} \left[- \sum_{J \in \mathcal{J}} \log \phi_J(X_J) \right] - H(\mu)$$

FACTOR GRAPHS



Definition

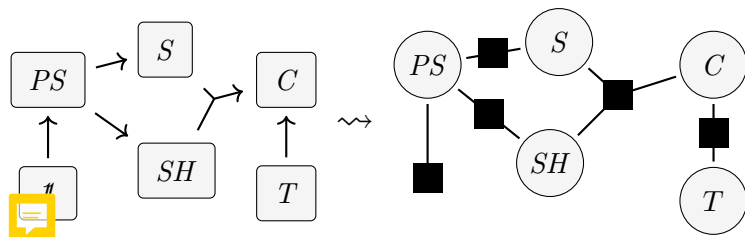
A **weighted factor graph** $\Psi = (\Phi, \theta)$ is a set of variables $\mathcal{X} = \{X_i\}$, factors $\{\phi_J: \mathcal{V}(X_J) \rightarrow \mathbb{R}_{\geq 0}\}_{J \in \mathcal{J}}$, and weights $(\theta_J)_{J \in \mathcal{J}}$ with $X_J \subseteq \mathcal{X}$; Ψ defines a distribution

$$\Pr_{\Psi}(\vec{x}) := \frac{1}{Z_{\Psi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J)^{\theta_J}, \quad \text{where } Z_{\Psi} \text{ is the normalization constant.}$$

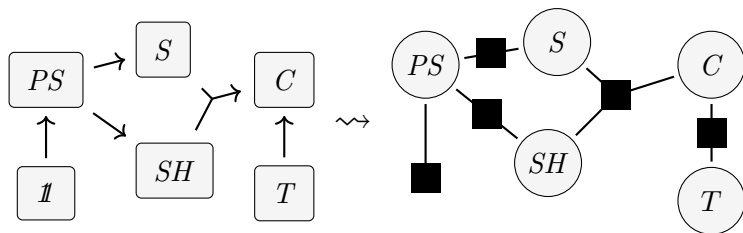
Ψ defines a standard scoring function

$$VFE_{\Psi}(\mu) := \mathbb{E}_{\mu} \left[- \sum_{J \in \mathcal{J}} \theta_J \log \phi_J(X_J) \right] - H(\mu)$$

PDGs AS FACTOR GRAPHS

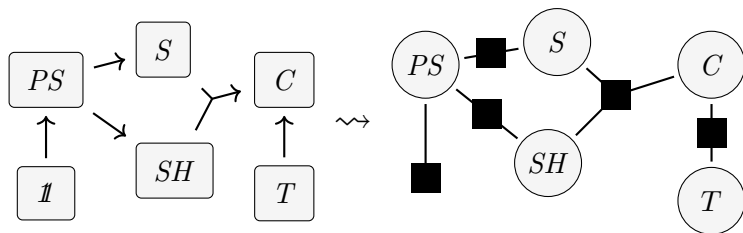


PDGs AS FACTOR GRAPHS



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

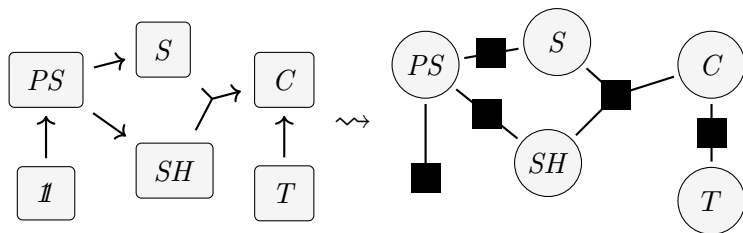
PDGs AS FACTOR GRAPHS



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

Not for $\gamma = 1$.

PDGs AS FACTOR GRAPHS



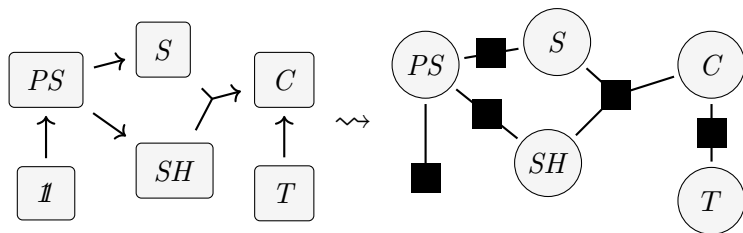
The cpds of a PDG are essentially factors. Are the semantics different?

Not for $\gamma = 1$.

Theorem

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

PDGs AS FACTOR GRAPHS



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

Not for $\gamma = 1$.

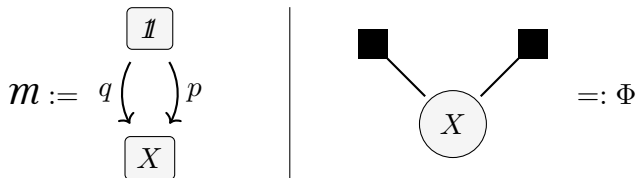
Theorem

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

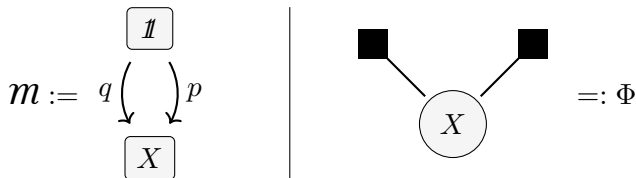
Theorem

If $\beta = \gamma\alpha$, then $\llbracket \mathcal{m} \rrbracket_{\gamma}^* = \text{Pr}_{(\Phi_{\mathcal{m}}, \beta)}$.

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

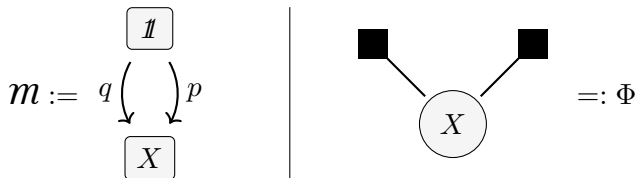


AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



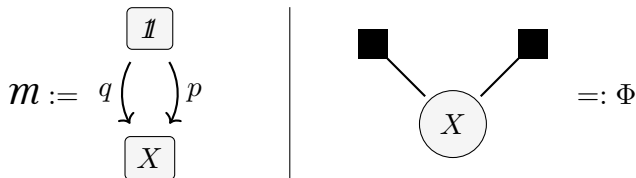
- If $p = q$, then $[[m]]^* = p = q \dots$

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



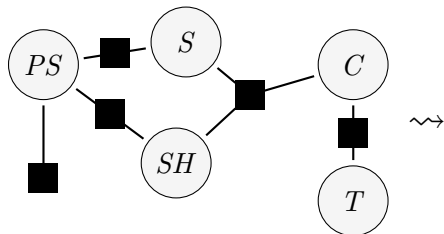
- If $p = q$, then $\llbracket m \rrbracket^* = p = q \dots$
- ... but $\Pr_{\Phi} \propto p^2$

AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

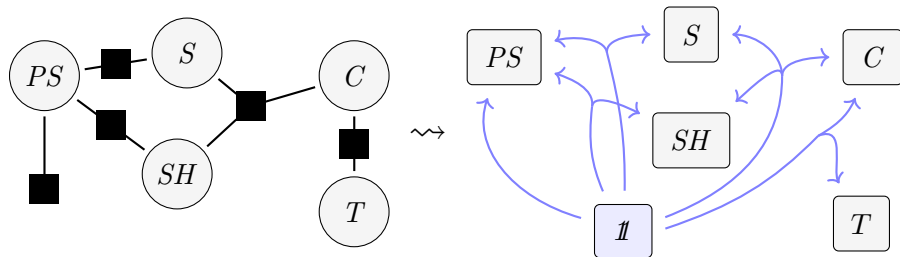


- If $p = q$, then $\llbracket m \rrbracket^* = p = q \dots$
- \dots but $\Pr_{\Phi} \propto p^2$
- Individual factors have *no probabilistic meaning*.

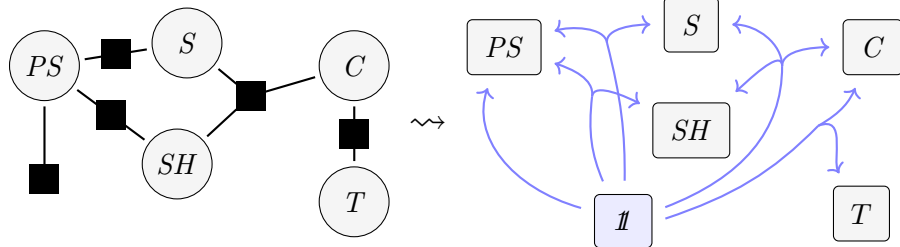
FACTOR GRAPHS AS PDGs



FACTOR GRAPHS AS PDGs



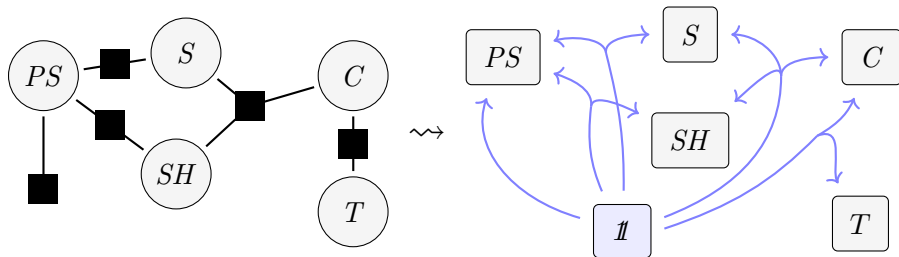
FACTOR GRAPHS AS PDGs



Theorem

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

FACTOR GRAPHS AS PDGs



Theorem

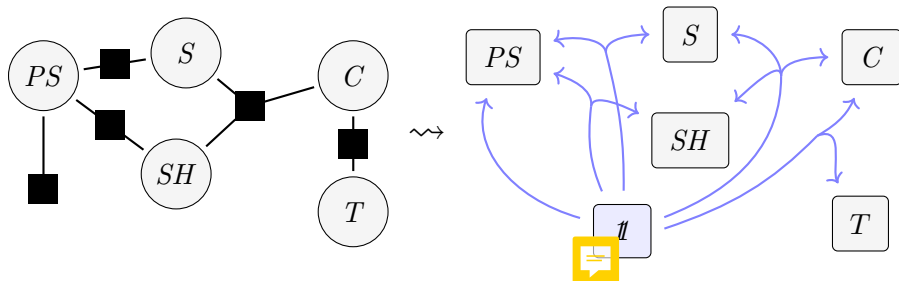
$\Pr_{\Phi} = \llbracket \mathbf{n}_{\Phi} \rrbracket_1^*$ for all factor graphs Φ .

Theorem

Ψ can be translated to a PDG \mathbf{m}_{Ψ} where $\beta = k\alpha$, and

$$\Pr_{\Psi} = \llbracket \mathbf{m}_{\Psi} \rrbracket_k^*$$

FACTOR GRAPHS AS PDGs



Theorem

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

Theorem

Ψ can be translated to a PDG \mathbf{m}_{Ψ} where $\beta = k\alpha$, and

$$\Pr_{\Psi} = \llbracket \mathbf{m}_{\Psi} \rrbracket_k^*$$

Also: $\log Z_{\Phi} = \llbracket \mathbf{n}_{\Phi} \rrbracket_1$.

PDG SEMANTICS, COMPARED TO THAT OF FACTOR GRAPHS

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

PDG SEMANTICS, COMPARED TO THAT OF FACTOR GRAPHS

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

And recall that

$$VFE_\Psi(\mu) := \mathbb{E}_\mu \left[\sum_{J \in \mathcal{J}} \theta_J \log \frac{1}{\phi_J(X_J)} \right] - H(\mu)$$

OUTLINE FOR SECTION 6

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

Conditioning as inconsistency resolution.

To condition on an event $(Y=y)$, simply add it to the PDG. Then the new best distribution is the old one, conditioned on $(Y=y)$. That is,

$$\llbracket m \sqcup (Y=y) \rrbracket^* = \llbracket m \rrbracket^* \mid (Y=y).$$

INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

Conditioning as inconsistency resolution.

To condition on an event $(Y=y)$, simply add it to the PDG. Then the new best distribution is the old one, conditioned on $(Y=y)$. That is,

$$\llbracket m \sqcup (Y=y) \rrbracket^* = \llbracket m \rrbracket^* \mid (Y=y).$$

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

- We can add $X \xrightarrow{p} Y$ to m , to get $m \sqcup p$.

INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

Conditioning as inconsistency resolution.

To condition on an event $(Y=y)$, simply add it to the PDG. Then the new best distribution is the old one, conditioned on $(Y=y)$. That is,

$$\llbracket \mathcal{M} \sqcup (Y=y) \rrbracket^* = \llbracket \mathcal{M} \rrbracket^* \mid (Y=y).$$

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

- We can add $X \xrightarrow{p} Y$ to \mathcal{M} , to get $\mathcal{M} \sqcup p$.
- The choice of cpd p that minimizes the inconsistency $\llbracket \mathcal{M} \sqcup p \rrbracket$ is $\llbracket \mathcal{M} \rrbracket^*(Y \mid X)$,

INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

Conditioning as inconsistency resolution.

To condition on an event $(Y=y)$, simply add it to the PDG. Then the new best distribution is the old one, conditioned on $(Y=y)$. That is,

$$\llbracket \mathcal{M} \sqcup (Y=y) \rrbracket^* = \llbracket \mathcal{M} \rrbracket^* \mid (Y=y).$$

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

- We can add $X \xrightarrow{p} Y$ to \mathcal{M} , to get $\mathcal{M} \sqcup p$.
- The choice of cpd p that minimizes the inconsistency $\langle\langle \mathcal{M} \sqcup p \rangle\rangle$ is $\llbracket \mathcal{M} \rrbracket^*(Y \mid X)$,
 - ▶ so an inconsistency oracle yields fast inference by gradient descent.

INFERENCE VIA INCONSISTENCY REDUCTION


Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

Conditioning as inconsistency resolution.

To condition on an event $(Y=y)$, simply add it to the PDG. Then the new best distribution is the old one, conditioned on $(Y=y)$. That is,

$$\llbracket \mathbf{m} \sqcup (Y=y) \rrbracket^* = \llbracket \mathbf{m} \rrbracket^* \mid (Y=y).$$

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

- We can add $X \xrightarrow{p} Y$ to \mathbf{m} , to get $\mathbf{m} \sqcup p$.
- The choice of cpd p that minimizes the inconsistency $\llbracket \mathbf{m} \sqcup p \rrbracket$ is $\llbracket \mathbf{m} \rrbracket^*(Y \mid X)$,

 - ▶ so an inconsistency oracle yields fast inference by gradient descent.

(Theorem): Unfortunately,

- ❶ Deciding if \mathbf{m} is consistent is NP-hard.
- ❷ Computing $\llbracket \mathbf{m} \rrbracket_\gamma$ is #P-hard, for $\gamma > 0$.

INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event $Y=y$ with the cpd $\mathbb{I} \xrightarrow{\delta_y} Y$.

Conditioning as inconsistency resolution.

To condition on an event $(Y=y)$, simply add it to the PDG. Then the new best distribution is the old one, conditioned on $(Y=y)$. That is,

$$\llbracket \mathbf{m} \sqcup (Y=y) \rrbracket^* = \llbracket \mathbf{m} \rrbracket^* \mid (Y=y).$$

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

- We can add $X \xrightarrow{p} Y$ to \mathbf{m} , to get $\mathbf{m} \sqcup p$.
- The choice of cpd p that minimizes the inconsistency $\llbracket \mathbf{m} \sqcup p \rrbracket$ is $\llbracket \mathbf{m} \rrbracket^*(Y \mid X)$,
 - ▶ so an inconsistency oracle yields fast inference by gradient descent.

(Theorem): Unfortunately,

- ❶ Deciding if \mathbf{m} is consistent is NP-hard.
- ❷ Computing $\llbracket \mathbf{m} \rrbracket_\gamma$ is #P-hard, for $\gamma > 0$.

...just like for BDDs and Factor Graphs.

OUTLINE FOR SECTION 7

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

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.
 - ▶ Model makes claims about reality

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.
 - ▶ Model makes claims about reality
 - ▶ For instance: priors correspond to (but are easier to criticize than) regularizers.

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.
 - ▶ Model makes claims about reality
 - ▶ For instance: priors correspond to (but are easier to criticize than) regularizers.



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.
 - ▶ Model makes claims about reality
 - ▶ For instance: priors correspond to (but  easier to criticize than) regularizers.

Surprising Result

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

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.
 - ▶ Model makes claims about reality
 - ▶ For instance: priors correspond to (but are easier to criticize than) regularizers.

Surprising Result

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

Bonus

An visual language for reasoning about relationships between objective functions.

SURPRISE AS INCONSISTENCY

Consider a distribution $p(X)$.

SURPRISE AS INCONSISTENCY

Consider a distribution $p(X)$.

The surprise (information content) at seeing a sample x is:

$$I_p(x) := \log \frac{1}{p(X=x)}.$$

Proposition

*Surprise is the inconsistency of simultaneously believing p and $X = x$.
That is,*

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

SURPRISE AS INCONSISTENCY

Consider a distribution $p(X)$.

The surprise (information content) at seeing a sample x is:

$$I_p(x) := \log \frac{1}{p(X=x)}.$$

Proposition

Surprise is the inconsistency of simultaneously believing p and $X = x$. That is,

$$I_p(x) = \left\langle\!\left\langle \overset{p}{\longrightarrow} \boxed{X} \overset{X=x}{\longleftarrow} \right\rangle\!\right\rangle.$$

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

SURPRISE AS INCONSISTENCY

Consider a distribution $p(X)$.

The surprise (information content) at seeing a sample x is:

$$I_p(x) := \log \frac{1}{p(X=x)}.$$

Proposition

Surprise is the inconsistency of simultaneously believing p and $X = x$. That is,

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

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

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

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 (supervised setting: conditional cross entropy)



The inconsistency of the PDG containing $f(Y | X)$ and a high-confidence empirical distribution $\text{Pr}_{\underline{\mathbf{xy}}}$ of samples $\underline{\mathbf{xy}} = \{(x_i, y_i)\}$ is equal to the cross entropy (plus $H(Y | X)$, a constant that depends only on the data $\text{Pr}_{\underline{\mathbf{xy}}}$). That is,

$$\left\langle\left\langle \begin{array}{c} \text{Pr}_{\underline{\mathbf{xy}}} \quad (\beta:\infty) \\ \swarrow \quad \searrow \\ X \quad Y \\ \xrightarrow{f} \end{array} \right\rangle\right\rangle = \frac{1}{|\underline{\mathbf{xy}}|} \sum_{(x,y) \in \underline{\mathbf{xy}}} \left[\log \frac{1}{f(y | x)} \right] - H_{\text{Pr}_{\underline{\mathbf{xy}}}}(Y | X).$$

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

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

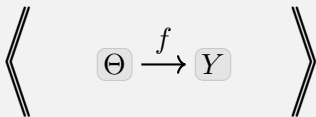
Proposition (Mean Square Error as Inconsistency)

$$\left\langle \begin{array}{c} \xrightarrow[\substack{D \\ (\beta:\infty)}]{} \\ \begin{array}{ccc} & \mathcal{N}(f(x), 1) & \\ X & \begin{array}{c} \curvearrowright \\ \curvearrowleft \end{array} & Y \\ & \mathcal{N}(g(x), 1) & \end{array} \end{array} \right\rangle = \mathbb{E}_D \left(f(X) - h(X) \right)^2 =: \text{MSE}(f, h)$$

Proposition (Regularizers as priors)

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

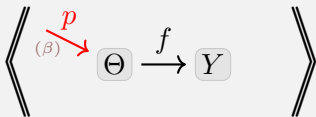
That is,



Proposition (Regularizers as priors)

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

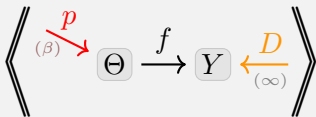
That is,



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,

$$\left\langle \begin{array}{c} \xrightarrow[p(\theta)]{(\beta)} \\ \xrightarrow[\theta_0]{} \end{array} \Theta \xrightarrow{f} Y \xleftarrow[D_{(\infty)}]{} \right\rangle =$$

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,

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

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,

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

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,

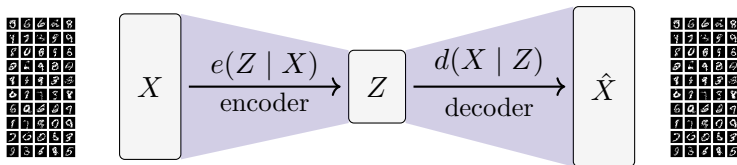
$$\left\langle \begin{array}{c} \xrightarrow[\theta_0]{(\beta)} \\ \Theta \end{array} \xrightarrow{f} Y \xleftarrow{(\infty)} D \right\rangle = \mathbb{E}_{y \sim D} \left[\log \frac{1}{f(y \mid \theta_0)} \right] + \beta \log \frac{1}{p(\theta_0)} - H(D)$$

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

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

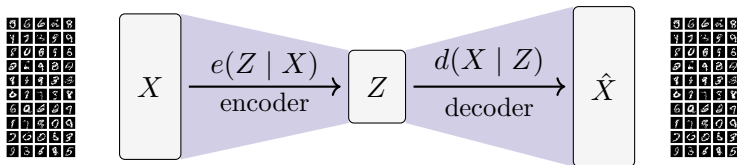
VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



VARIATIONAL AUTO-ENCODERS, TAKE 1

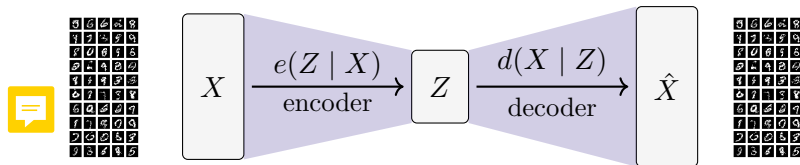
- Structure consists of two neural networks (cpds):



- Objective:

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

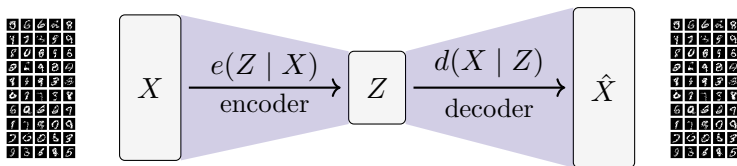


- Objective:

- ▶ For each x , want to minimize $\text{Rec}(x) := - \overset{\text{"reconstruction error"}}{\mathbb{E}_{z \sim e|x}} \log d(x | z)$

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



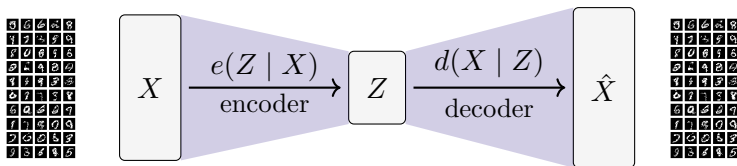
- Objective:

- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ ^{“reconstruction error”}
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

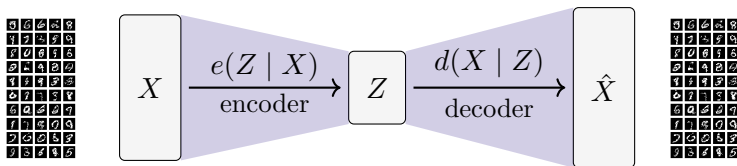


- Objective:

- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ “reconstruction error”
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.
- ▶ Combine to get VaE objective:
 $\text{ELBO}_{p,e,d}(x) :=$

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



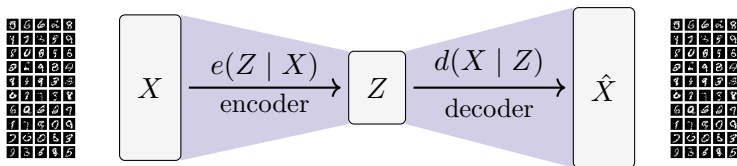
- Objective:

- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ “reconstruction error”
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.
- ▶ Combine to get VaE objective:

$$\text{ELBO}_{p,e,d}(x) :=$$
$$\underbrace{-D(e(Z|x) \parallel p(Z))}_{\text{divergence from prior}}$$

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

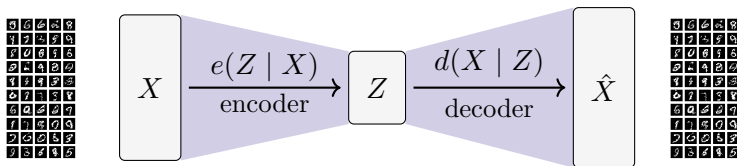
- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ “reconstruction error”
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.
- ▶ Combine to get VaE objective:

$$\text{ELBO}_{p,e,d}(x) :=$$

$$\underbrace{-D(e(Z|x) \parallel p(Z))}_{\text{divergence from prior}} - \text{Rec}(x)$$

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

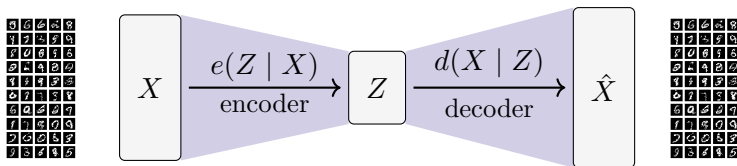
- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ “reconstruction error”
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.
- ▶ Combine to get VaE objective:

$$\text{ELBO}_{p,e,d}(x) :=$$

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

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



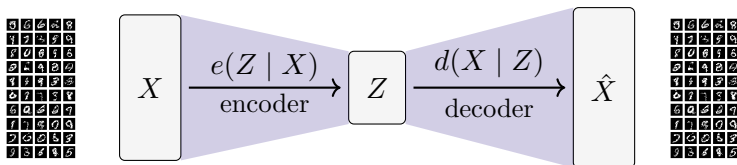
- Objective:

- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ “reconstruction error”
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.
- ▶ Combine to get VaE objective:

$$\text{ELBO}_{p,e,d}(x) := \underbrace{-D(e(Z|x) \parallel p(Z))}_{\text{divergence from prior}} - \text{Rec}(x) = \mathbb{E}_{z \sim e|x} \left[\log \frac{p(z)d(x | z)}{e(z | x)} \right] \leq \overset{\text{“evidence”}}{\log \text{Pr}_{pd}(x)}$$

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

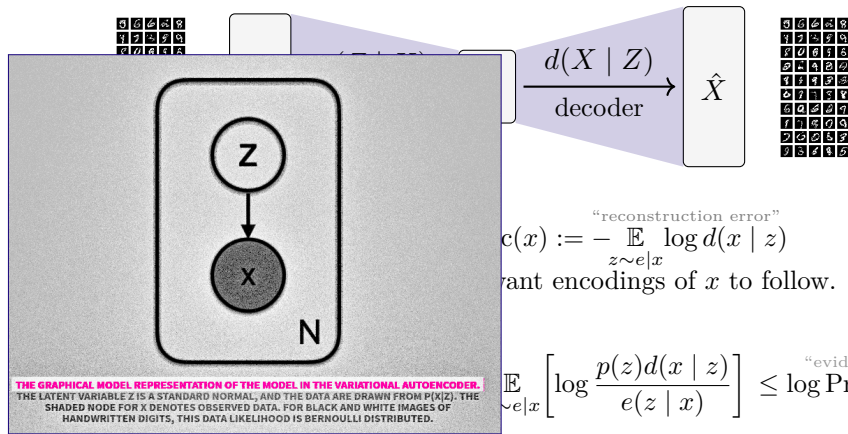
- ▶ For each x , want to minimize $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$ (“reconstruction error”)
- ▶ Also have a prior $p(Z)$ that we want encodings of x to follow.
- ▶ Combine to get VaE objective:

$$\text{ELBO}_{p,e,d}(x) := \underbrace{-D(e(Z|x) \parallel p(Z))}_{\text{divergence from prior}} - \text{Rec}(x) = \mathbb{E}_{z \sim e|x} \left[\log \frac{p(z)d(x | z)}{e(z | x)} \right] \leq \log \text{Pr}_{pd}(x) \quad \text{“evidence”}$$

Urge to use graphical models (even if can’t quite capture *entire* VaE)

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



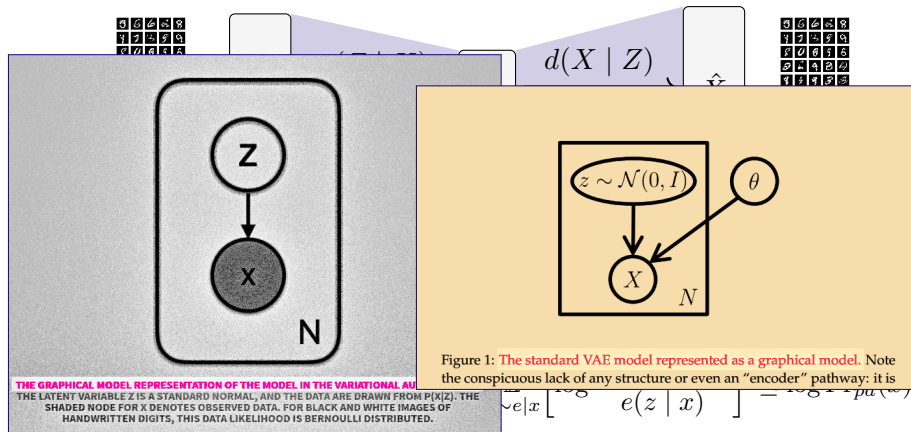
“reconstruction error”
 $c(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$
 want encodings of x to follow.

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

Urge to use graphical models (even if can't quite capture *entire* VaE)

VARIATIONAL AUTO-ENCODERS, TAKE 1

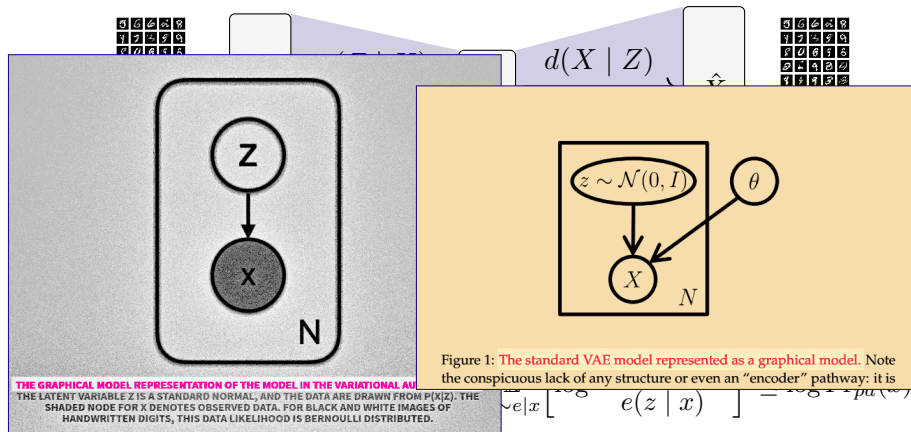
- Structure consists of two neural networks (cpds):



Urge to use graphical models (even if can't quite capture *entire* VaE)

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

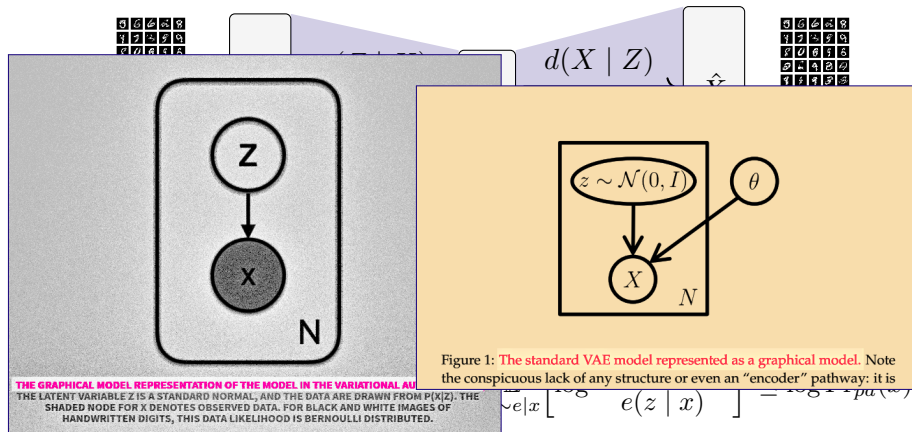


Urge to use graphical models (even if can't quite capture *entire* VaE)

- $e(Z | X)$ has same target as $p(Z)$, so can't put in BN;

VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



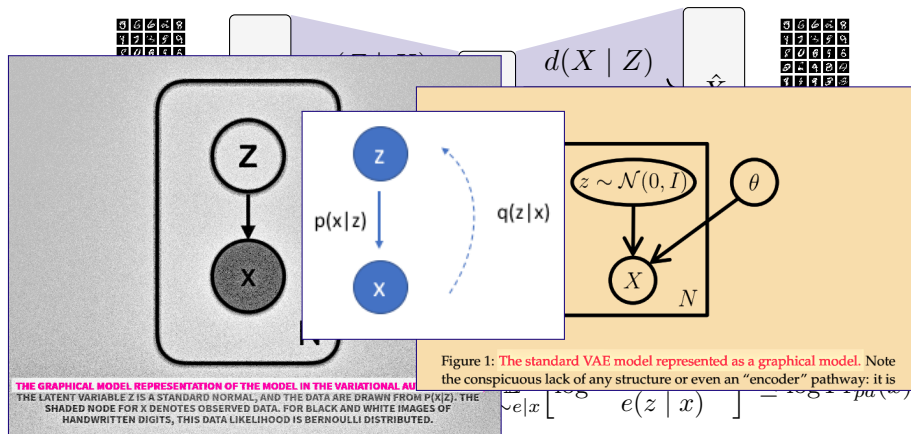
Urge to use graphical models (even if can't quite capture *entire* VaE)

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



VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



Urge to use graphical models (even if can't quite capture *entire* VaE)

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

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

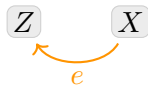
Z

X

VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encoder

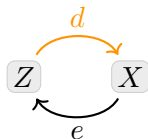


VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encoder

$d(X | Z)$: decoder



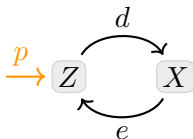
VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$: encoder

$d(X | Z)$: decoder

$p(Z)$: prior



VARIATIONAL AUTO-ENCODERS, TAKE 2

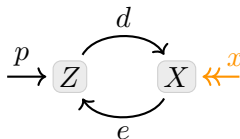
- Structure:

$e(Z | X)$: encoder

$d(X | Z)$: decoder

$p(Z)$: prior

- observe a sample x



VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

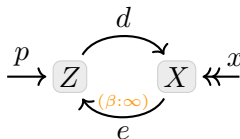
$e(Z | X)$: encoder

$d(X | Z)$: decoder

$p(Z)$: prior

- observe a sample x

▶ and trust encoding



VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

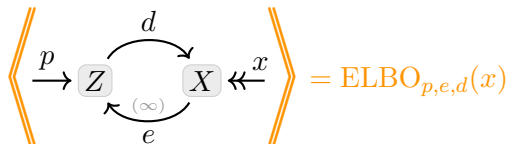
$e(Z | X)$: encoder

$d(X | Z)$: decoder

$p(Z)$: prior

- observe a sample x
 - ▶ and trust encoding

Objective function is free:



VARIATIONAL AUTO-ENCODERS, TAKE 2

Objective function is free:

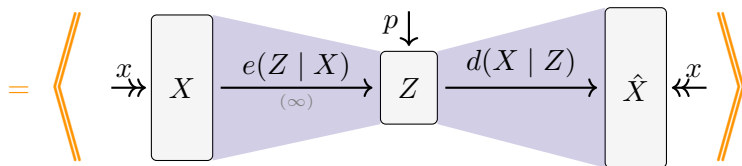
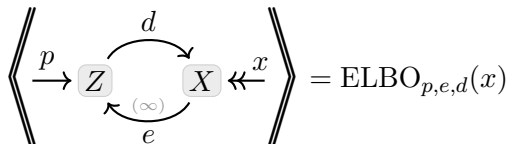
- Structure:

$e(Z | X)$: encoder

$d(X | Z)$: decoder

$p(Z)$: prior

- observe a sample x
 - and trust encoding



VARIATIONAL AUTO-ENCODERS, TAKE 2

Objective function is free:

- Structure:

$e(Z | X)$: encoder

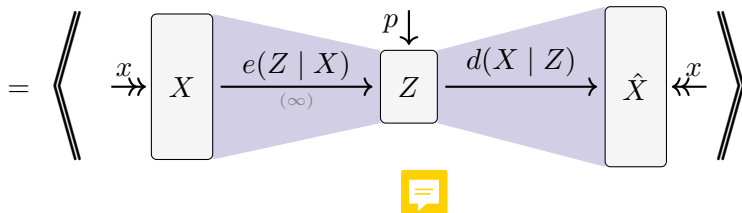
$d(X | Z)$: decoder

$p(Z)$: prior

- observe a sample x

► and trust encoding

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



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

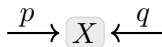
$$-\log \Pr_{p,d}(X=x) = \left\langle \begin{array}{c} p \rightarrow Z \xrightarrow{d} X \leftarrow x \end{array} \right\rangle \left\langle \begin{array}{c} p \rightarrow Z \xrightleftharpoons[e!]{d} X \leftarrow x \end{array} \right\rangle = -\text{ELBO}_{p,e,d}(x).$$

VISUAL PROOF: THE VARIATIONAL BOUND

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

INCONSISTENCY AS A DIVERGENCE

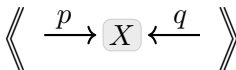
You believe both $p(X)$ and $q(X)$.



INCONSISTENCY AS A DIVERGENCE

You believe both $p(X)$ and $q(X)$.

Your inconsistency: a divergence between p and q ?



INCONSISTENCY AS A DIVERGENCE

You believe both $p(X)$ and $q(X)$.

Your inconsistency: a divergence between p and q ?

$$\left\langle \begin{array}{c} \xrightarrow[p(\beta:r)]{} \\ \end{array} X \begin{array}{c} \xleftarrow[q(\beta:s)]{} \\ \end{array} \right\rangle$$

INCONSISTENCY AS A DIVERGENCE

You believe both $p(X)$ and $q(X)$.

Your inconsistency: a divergence between p and q ?

$$\text{Let } D_{(r,s)}^{\text{PDG}}(p, q) := \left\langle \frac{p}{(\beta:r)} \rightarrow X \leftarrow \frac{q}{(\beta:s)} \right\rangle$$

INCONSISTENCY AS A DIVERGENCE

You believe both $p(X)$ and $q(X)$.

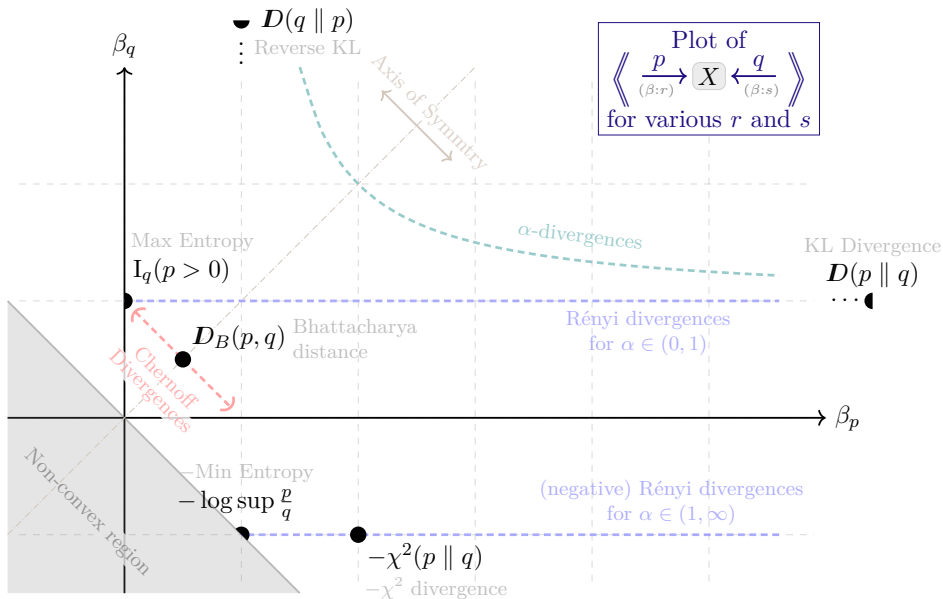
Your inconsistency: a divergence between p and q ?

$$\text{Let } D_{(r,s)}^{\text{PDG}}(p, q) := \left\langle \frac{p}{(\beta:r)} \rightarrow X \leftarrow \frac{q}{(\beta:s)} \right\rangle$$

Lemma

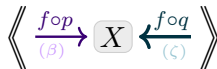
$$D_{(r,s)}^{\text{PDG}}(p, q) = -(r+s) \log \sum_x \left(p(x)^r q(x)^s \right)^{\frac{1}{r+s}}.$$

DIVERGENCES AS INCONSISTENCIES



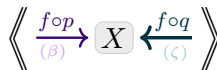
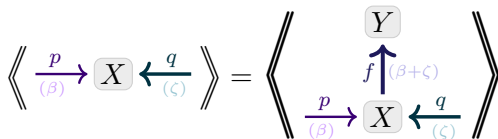
VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D_{(\beta, \zeta)}^{\text{PDG}}(p \parallel q) \geq D_{(\beta, \zeta)}^{\text{PDG}}(f \circ p \parallel f \circ q)$$



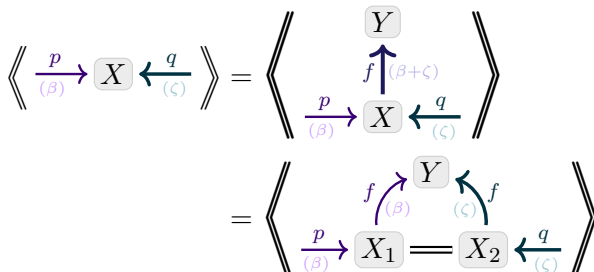
VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D_{(\beta, \zeta)}^{\text{PDG}}(p \parallel q) \geq D_{(\beta, \zeta)}^{\text{PDG}}(f \circ p \parallel f \circ q)$$



VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D_{(\beta, \zeta)}^{\text{PDG}}(p \parallel q) \geq D_{(\beta, \zeta)}^{\text{PDG}}(f \circ p \parallel f \circ q)$$



$$\left\langle \left\langle \frac{f \circ p}{(\beta)} \rightarrow X \leftarrow \frac{f \circ q}{(\zeta)} \right\rangle \right\rangle$$

VISUAL PROOF: DATA PROCESSING INEQUALITY

$$D_{(\beta, \zeta)}^{\text{PDG}}(p \parallel q) \geq D_{(\beta, \zeta)}^{\text{PDG}}(f \circ p \parallel f \circ q)$$

$$\begin{aligned}
 \left\langle \left\langle \frac{p}{(\beta)} \rightarrow X \leftarrow \frac{q}{(\zeta)} \right\rangle \right\rangle &= \left\langle \left\langle \begin{array}{c} Y \\ \uparrow f_{(\beta+\zeta)} \\ X \end{array} \begin{array}{c} \xleftarrow{q_{(\zeta)}} \\ \end{array} \end{array} \right\rangle \right\rangle \\
 &= \left\langle \left\langle \begin{array}{c} Y \\ \nearrow f_{(\beta)} \quad \nwarrow f_{(\zeta)} \\ X_1 = X_2 \end{array} \begin{array}{c} \xleftarrow{q_{(\zeta)}} \\ \end{array} \end{array} \right\rangle \right\rangle \\
 &\geq \left\langle \left\langle \begin{array}{c} Y \\ \nearrow f_{(\beta)} \quad \nwarrow f_{(\zeta)} \\ X_1 \quad X_2 \end{array} \begin{array}{c} \xleftarrow{q_{(\zeta)}} \\ \end{array} \end{array} \right\rangle \right\rangle = \left\langle \left\langle \frac{f \circ p}{(\beta)} \rightarrow X \leftarrow \frac{f \circ q}{(\zeta)} \right\rangle \right\rangle
 \end{aligned}$$

RECAP

PDGs...

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

RECAP

PDGs...

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

RECAP

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular.
- cleanly separate quantitative from qualitative information, and can express confidence in each (β and α). This is captured by terms *Inc* and *IDef* in our scoring function.

RECAP

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular.
- cleanly separate quantitative from qualitative information, and can express confidence in each (β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- naturally capture BNs and factor graphs, with the best-scoring distribution.

RECAP

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular.
- cleanly separate quantitative from qualitative information, and can express confidence in each (β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- naturally capture BNs and factor graphs, with the best-scoring distribution.
- simultaneously capture many standard loss functions and divergences with the value of the scoring function.

RECAP

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular.
- cleanly separate quantitative from qualitative information, and can express confidence in each (β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- naturally capture BNs and factor graphs, with the best-scoring distribution.
- simultaneously capture many standard loss functions and divergences with the value of the scoring function.
- give us a clean visual language for reasoning about the relationships between objectives.

RECAP

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular.
- cleanly separate quantitative from qualitative information, and can express confidence in each (β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- naturally capture BNs and factor graphs, with the best-scoring distribution.
- simultaneously capture many standard loss functions and divergences with the value of the scoring function.
- give us a clean visual language for reasoning about the relationships between objectives.

But there is much more to be done!

OUTLINE FOR SECTION 8

1 INTRODUCTION

2 MODELING EXAMPLES

- A Simple Example: What are Floomps?
- Differences from BNs
- PDG Union and Restriction

3 SYNTAX

- Formal Definitions of PDGs

4 SEMANTICS

5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

6 INFERENCE

7 INCONSISTENCY AS LOSS

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

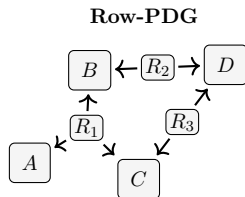
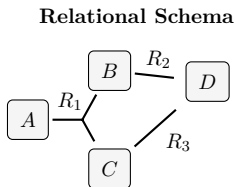
8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Ongoing Work

$$\begin{array}{ccc}
 A & \xrightarrow{f} & B \\
 g \downarrow & & \downarrow k \\
 C & \xrightarrow{h} & D
 \end{array}$$

$$\llbracket m \rrbracket = \mathbb{I}[kfa = hga] = -\log \#\{a : kfa = hga\}$$

Database						
R_1				R_2		
	A	B	C		B	D
	a_1	b_1	c_1		b_2	d_1
	a_2	b_2	c_2	b_3	d_2	
				b_4	d_3	
		R_3	C	D		
			c_2	d_1		
			c_1	d_3		



Proposition

If \mathcal{D} is a database and μ is a joint distribution over $\mathcal{M}_{\mathcal{D}}$, then $\mu \in \{\mathcal{M}_{\mathcal{D}}\}$ iff $\text{Supp}(\mu)$ is a universal relation for \mathcal{D} .

Corollary

$\mathcal{M}_{\mathcal{D}}$ is consistent iff \mathcal{D} is join consistent.

OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction

OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: incomplete cpds

OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: incomplete cpds
- Encoding preferences, and understanding preference changes

OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: incomplete cpds
- Encoding preferences, and understanding preference changes
- Trace Semantics and Composition

OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: incomplete cpds
- Encoding preferences, and understanding preference changes
- Trace Semantics and Composition
 - ▶ Extend semantics to score other objects, not just joint distributions.

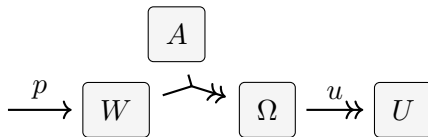
OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: incomplete cpds
- Encoding preferences, and understanding preference changes
- Trace Semantics and Composition
 - ▶ Extend semantics to score other objects, not just joint distributions.
 - ▶ Regarding PDGs as probabilistic automata.

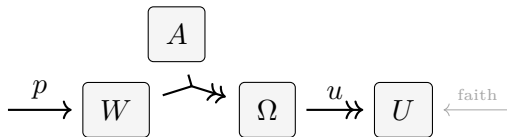
OPEN PROBLEMS AND FUTURE WORK

- Fleshing out the details of belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: incomplete cpds
- Encoding preferences, and understanding preference changes
- Trace Semantics and Composition
 - ▶ Extend semantics to score other objects, not just joint distributions.
 - ▶ Regarding PDGs as probabilistic automata.
 - ▶ Open Question: Do PDGs capture Dependency Networks? *

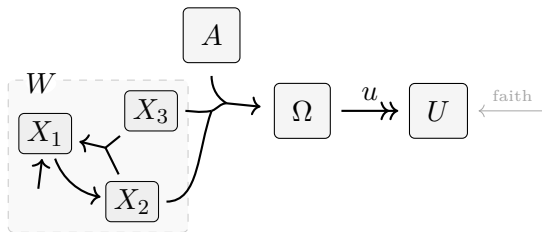
A DIFFERENT PICTURE OF AGENCY



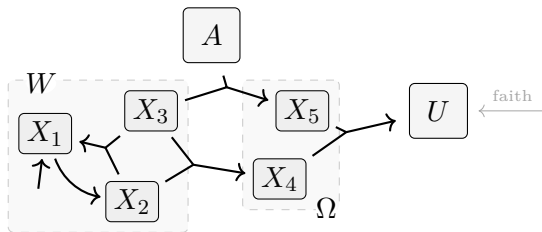
A DIFFERENT PICTURE OF AGENCY



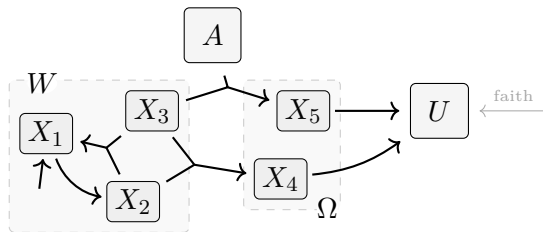
A DIFFERENT PICTURE OF AGENCY



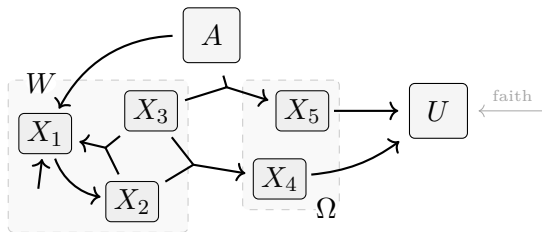
A DIFFERENT PICTURE OF AGENCY



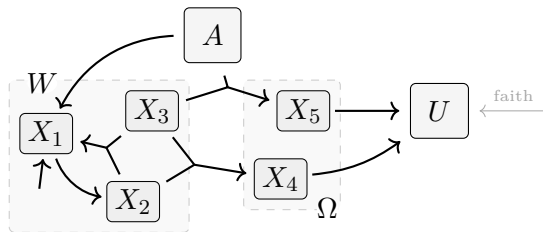
A DIFFERENT PICTURE OF AGENCY



A DIFFERENT PICTURE OF AGENCY

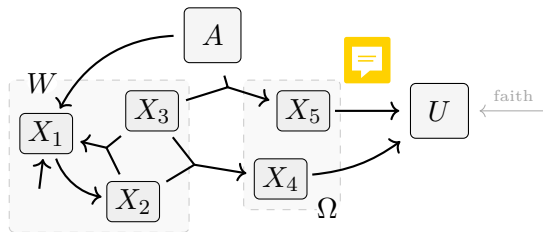


A DIFFERENT PICTURE OF AGENCY



- driven by pursuit of coherent identity;
not necessarily “favorite number go up”.

A DIFFERENT PICTURE OF AGENCY



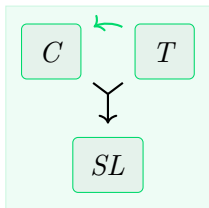
- driven by pursuit of coherent identity;
not necessarily “favorite number go up”.

Python library available at <https://orichardson.github.io/pdg/>

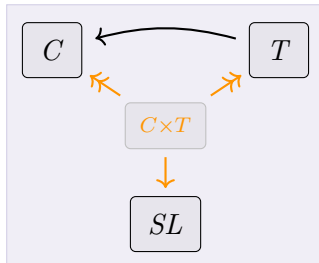
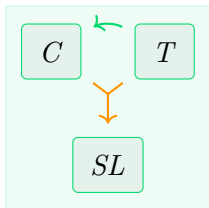
OUTLINE FOR SECTION 9

- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS
 - BNs as MaxEnt
- 13 MORE LOSSES
- 14 MORE CATEGORY THEORY
 - PDGs as diagrams of the Markov Category

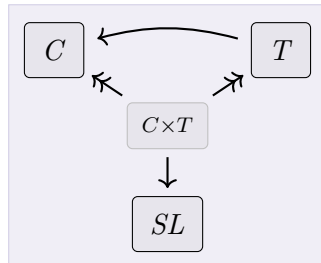
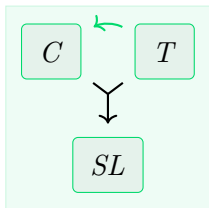
HYPER-GRAPHS? OR MERELY GRAPHS?



HYPER-GRAPHS? OR MERELY GRAPHS?

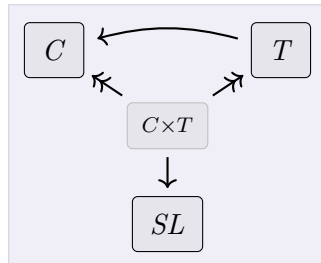
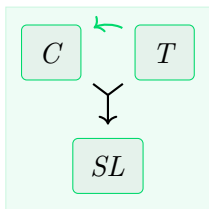


HYPER-GRAPHS? OR MERELY GRAPHS?



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

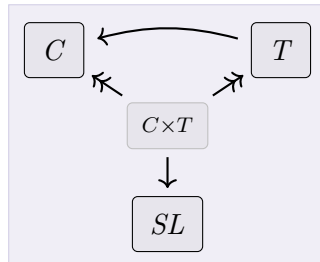
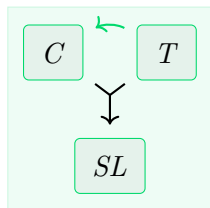
HYPER-GRAPHS? OR MERELY GRAPHS?



- 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

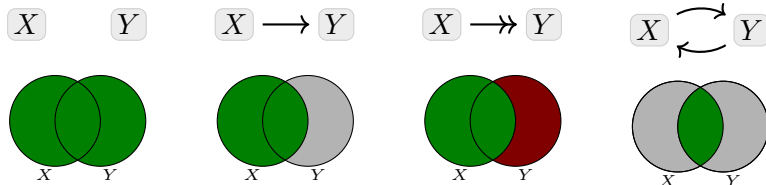
joint distributions \iff expanded joint distributions
satisfying coherence constraints

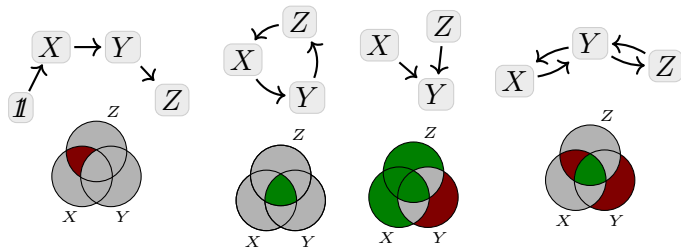
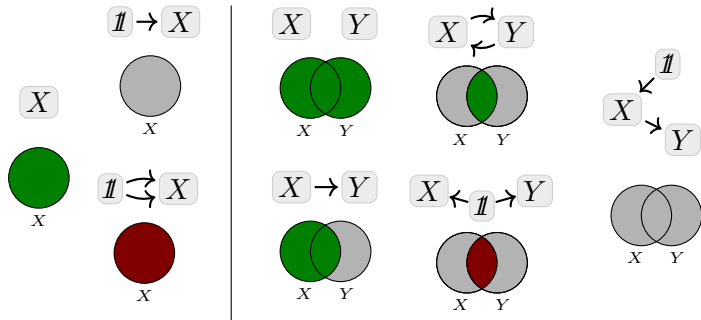
(working directly with hypergraphs is also possible)

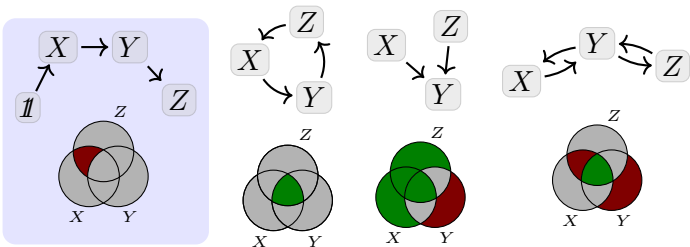
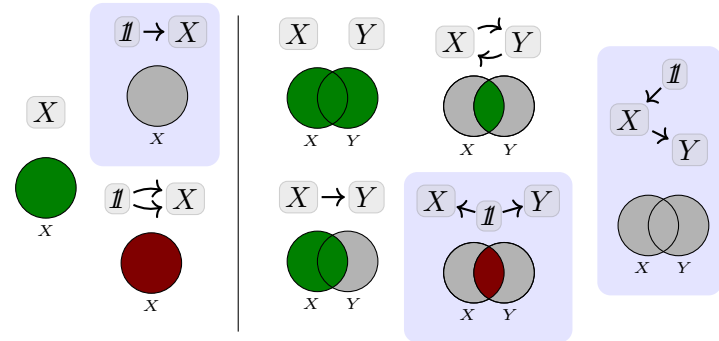
OUTLINE FOR SECTION 10

- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS
 - BNs as MaxEnt
- 13 MORE LOSSES
- 14 MORE CATEGORY THEORY
 - PDGs as diagrams of the Markov Category

ILLUSTRATIONS OF *IDef*







OUTLINE FOR SECTION 11

9 HYPER-GRAPHS

10 THE INFORMATION
DEFICIENCY

11 MORE ON SEMANTICS

12 MORE ON GRAPHICAL
MODELS
• BNs as MaxEnt

13 MORE LOSSES

14 MORE CATEGORY THEORY
• PDGs as diagrams of the
Markov Category

RELATIONSHIPS BETWEEN SEMANTICS

Proposition (*the set of consistent distributions is the zero set of the scoring function*)

$$\{\mathcal{m}\} = \{\mu : \llbracket \mathcal{m} \rrbracket_0(\mu) = 0\}.$$

Proposition (*If there are distributions consistent with \mathcal{m} , the best distribution is one of them.*)

$$\llbracket \mathcal{m} \rrbracket^* \in \llbracket \mathcal{m} \rrbracket_0^*, \text{ so if } \mathcal{m} \text{ is consistent, then } \llbracket \mathcal{m} \rrbracket^* \in \{\mathcal{m}\}.$$

$$\llbracket m \rrbracket_\gamma(\mu) = \mathbb{E}_\mu \log \prod_{X \xrightarrow{L} Y} \left(\frac{\mu(Y \mid X)}{\mathbf{p}_L(Y \mid X)} \right)^{\beta_L} \left(\frac{\mu(\mathcal{N})}{\prod_{X \xrightarrow{L} Y} \mu(Y \mid X)^{\alpha_L}} \right)^\gamma$$

OUTLINE FOR SECTION 12

- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS**
 - BNs as MaxEnt
- 13 MORE LOSSES
- 14 MORE CATEGORY THEORY
 - PDGs as diagrams of the Markov Category

BAYESIAN NETWORKS: MAXIMUM ENTROPY?

Common distributions
tend to maximize
entropy subject to
natural constraints.

distribution	constraints
Gaussian $\mathcal{N}(\mu, \sigma^2)$	mean μ , variance σ^2
Exponential $\text{Exp}(\lambda)$	positive support, mean λ
Factor graphs	moment matching.
...	...

BAYESIAN NETWORKS: MAXIMUM ENTROPY?

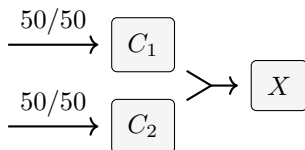
Common distributions
tend to maximize
entropy subject to
natural constraints.

distribution	constraints
Gaussian $\mathcal{N}(\mu, \sigma^2)$	mean μ , variance σ^2
Exponential $\text{Exp}(\lambda)$	positive support, mean λ
Factor graphs	moment matching.
...	...
Bayesian Networks	cpds + ???

BAYESIAN NETWORKS: MAXIMUM ENTROPY?

Common distributions
tend to maximize
entropy subject to
natural constraints.

distribution	constraints
Gaussian $\mathcal{N}(\mu, \sigma^2)$	mean μ , variance σ^2
Exponential $\text{Exp}(\lambda)$	positive support, mean λ
Factor graphs	moment matching.
...	...
Bayesian Networks	cpds + ???

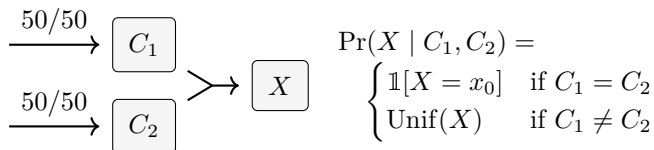


$$\Pr(X \mid C_1, C_2) = \begin{cases} \mathbb{1}[X = x_0] & \text{if } C_1 = C_2 \\ \text{Unif}(X) & \text{if } C_1 \neq C_2 \end{cases}$$

BAYESIAN NETWORKS: MAXIMUM ENTROPY?

Common distributions
tend to maximize
entropy subject to
natural constraints.

distribution	constraints
Gaussian $\mathcal{N}(\mu, \sigma^2)$	mean μ , variance σ^2
Exponential $\text{Exp}(\lambda)$	positive support, mean λ
Factor graphs	moment matching.
...	...
Bayesian Networks	cpds + ???



Corollary

Among the distributions in $\{\mathcal{B}\}$, $\Pr_{\mathcal{B}}$ has the maximum entropy, beyond the entropy of the given cpds.

$$\text{IDef says maximize: } H(\mu) - \sum_{X \in \mathcal{N}} H_{\mu}(X \mid \mathbf{Pa} X)$$

FULL FACTOR GRAPH RESULTS

Theorem (PDGs are WFGs)

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 $\llbracket(\mathcal{N}, \mathbf{v}, \gamma \mathbf{v})\rrbracket_\gamma = \gamma \text{VFE}_{(\Phi_{\mathcal{N}}, \mathbf{v})}$; consequently, $\llbracket(\mathcal{N}, \mathbf{v}, \gamma \mathbf{v})\rrbracket_\gamma^ = \{\text{Pr}_{(\Phi_{\mathcal{N}}, \mathbf{v})}\}$.*

Theorem (WFGs are PDGs)

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

OUTLINE FOR SECTION 13

- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS
 - BNs as MaxEnt
- 13 MORE LOSSES**
- 14 MORE CATEGORY THEORY
 - PDGs as diagrams of the Markov Category

OUTLINE FOR SECTION 14

- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS
 - BNs as MaxEnt
- 13 MORE LOSSES
- 14 MORE CATEGORY THEORY
 - PDGs as diagrams of the Markov Category

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.

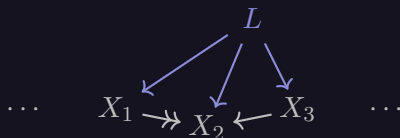
Equivalent Categorical Definition

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.

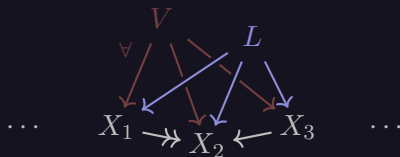
What do you do with diagrams? Take **limits** / **colimits**.

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

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.



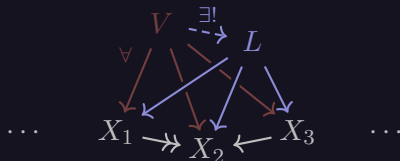
What do you do with diagrams? Take limits / colimits.



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} :

$$\lim \mathcal{m}_{\mathcal{B}} = \left(\mathbb{I}, \left\{ \text{Pr}_{\mathcal{B}}(X) \right\}_{X \in \mathcal{N}} \right)$$

