

# 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

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

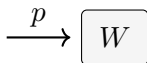
## 8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Projects Work

The standard way of modeling an agent with uncertainty:

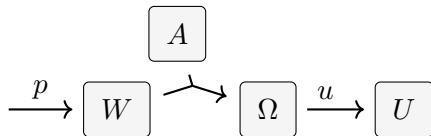
The standard way of modeling an agent with uncertainty:

- 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  $A$ ).



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?

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



# WHY INCONSISTENCY?

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

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;

# WHY INCONSISTENCY?

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

Why **model** inconsistency?

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

# WHY INCONSISTENCY?

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

Why **model** inconsistency?

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

Why **construct** an agent that can be inconsistent?

# WHY INCONSISTENCY?

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

Why **model** inconsistency?

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

Why **construct** an agent that can be inconsistent?

- Perfect consistency can be (needlessly) expensive.

# WHY INCONSISTENCY?

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

Why **model** inconsistency?

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

Why **construct** an agent that can be inconsistent?

- Perfect consistency can be (needlessly) expensive.
- Useful for identifying big problems.

# WHY INCONSISTENCY?

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

Why **model** inconsistency?

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

Why **construct** an agent that can be inconsistent?

- Perfect consistency can be (needlessly) expensive.
- Useful for identifying big problems.
  - ▶ (assertions, check-sums, paradoxes)

# WHY INCONSISTENCY?

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

Why **model** inconsistency?

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

Why **construct** an agent that can be inconsistent?

- Perfect consistency can be (needlessly) expensive.
- Useful for identifying big problems.
  - ▶ (assertions, check-sums, paradoxes)

Freedom from perfect consistency is valuable, 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.

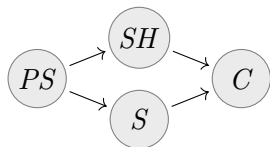


# YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

*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)



Variables:

---

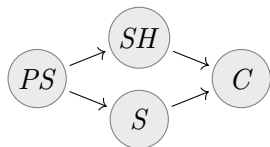
PS	Parents Smoke?
S	You smoke?
SH	Second-hand Smoke?
C	Get Cancer?

# 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  $X$



Variables:

---

PS	Parents Smoke?
S	You smoke?
SH	Second-hand Smoke?
C	Get Cancer?

# TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

## Qualitative BN, $\mathcal{G}$

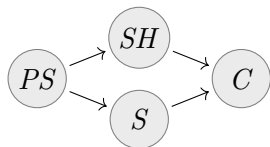
an independence relation on variables

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



Variables:

---

PS	Parents Smoke?
S	You smoke?
SH	Second-hand Smoke?
C	Get Cancer?

# OUTLINE FOR SECTION 2

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Projects Work

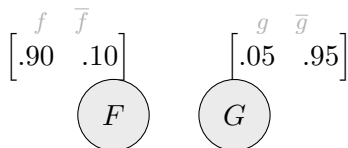
## SIMPLE EXAMPLE: FLOOMPS AND GUNS

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).

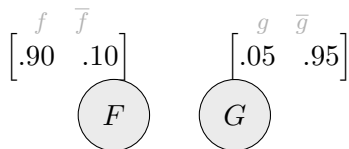
BN



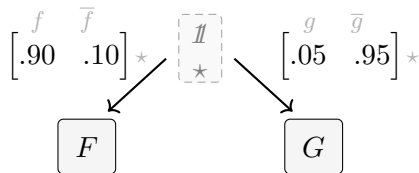
# SIMPLE EXAMPLE: FLOOMPS AND GUNS

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

**BN**



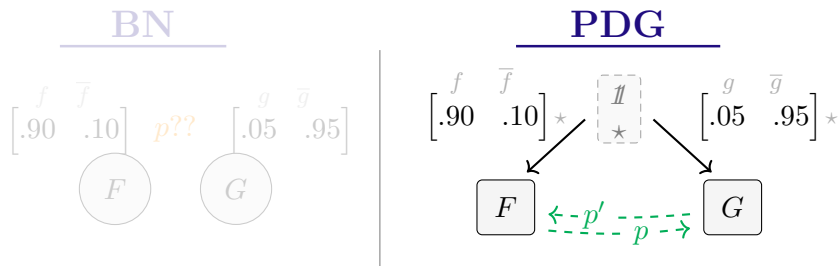
**PDG**



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



# SIMPLE EXAMPLE: FLOOMPS AND GUNS

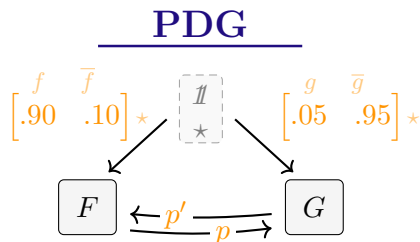
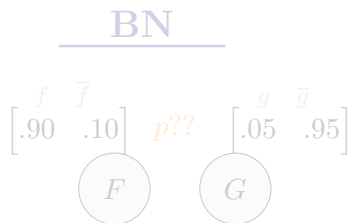


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.

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

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

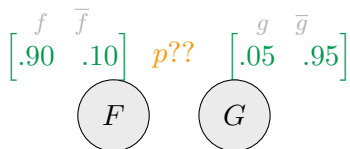
# SIMPLE EXAMPLE: FLOOMPS AND GUNS



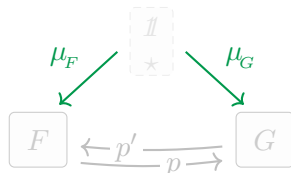
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

## BN



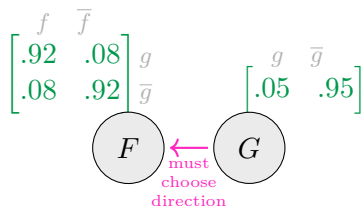
## PDG



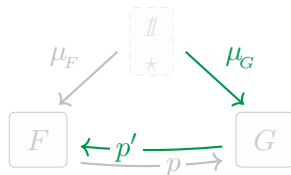
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
  - ▶ ... but BNs must resolve inconsistency first, which may break symmetry and irrecoverably lose information.

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

## BN



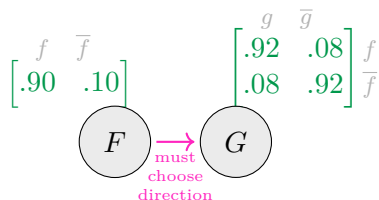
## PDG



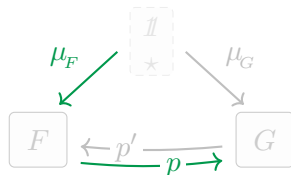
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
  - ▶ ...but BNs must resolve inconsistency first, which may **break symmetry** and irrecoverably lose information.

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

## BN

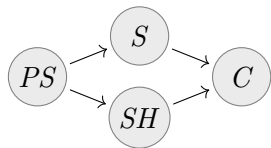


## PDG

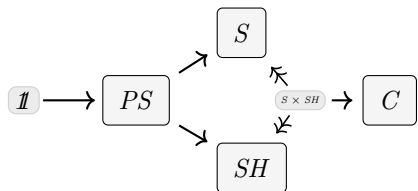
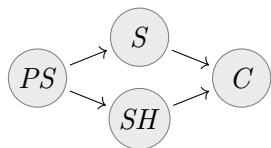


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
  - ▶ ...but BNs must resolve inconsistency first, which may **break symmetry** and irrecoverably lose information.

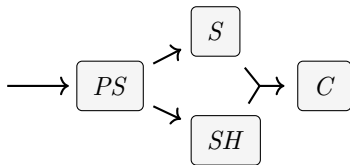
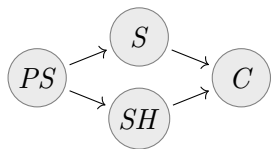
# BAYESIAN NETWORKS AS PDGs



# BAYESIAN NETWORKS AS PDGs

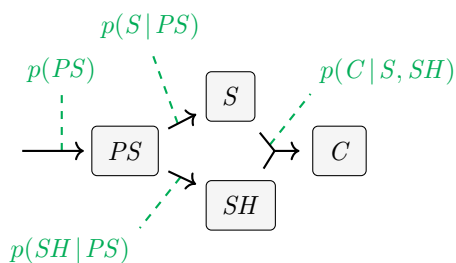
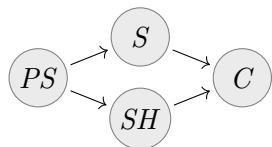


# BAYESIAN NETWORKS AS PDGs





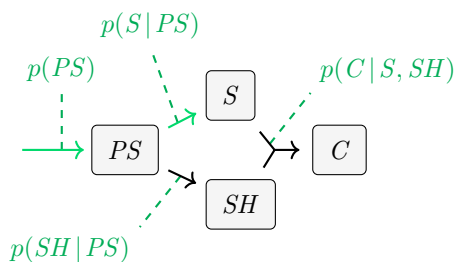
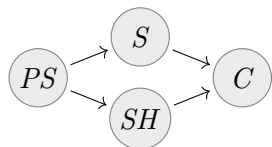
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

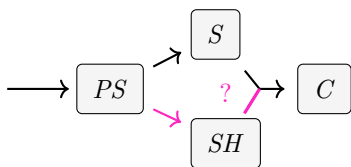
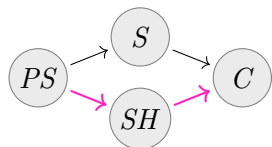
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

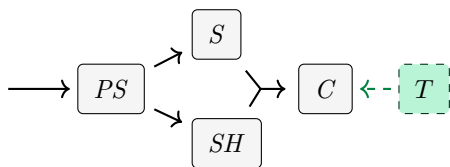
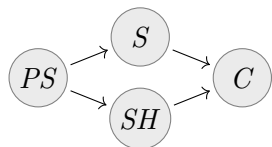
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

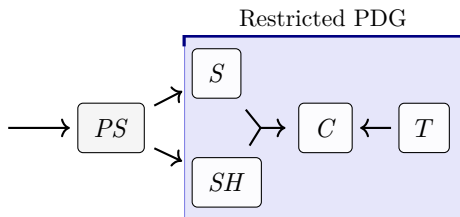
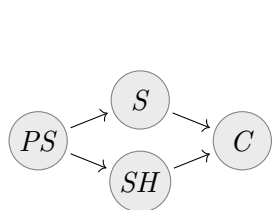
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

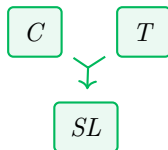
# BAYESIAN NETWORKS AS PDGs



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.

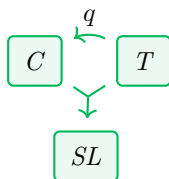
# COMBINING PDGs



Grok wants to be supreme leader ( $SL$ ).

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

# COMBINING PDGs

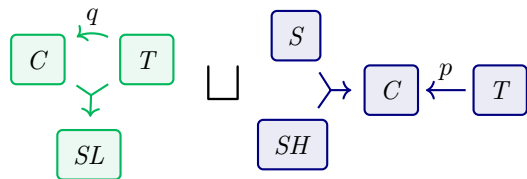


Grok wants to be supreme leader ( $SL$ ).

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

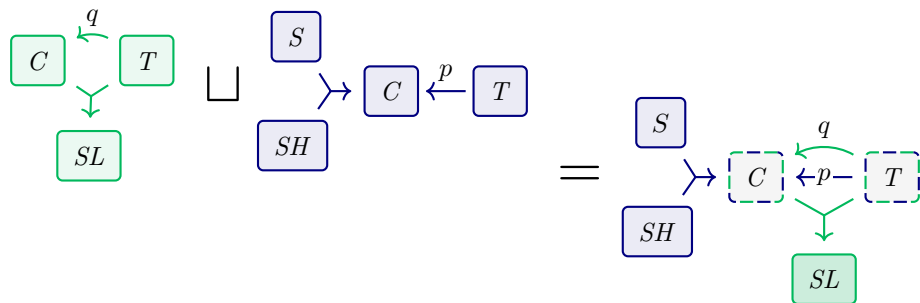
- ...but mom says  $q(C | T) = \begin{bmatrix} c & \bar{c} \\ .15 & .85 \\ .02 & .98 \end{bmatrix} \begin{matrix} t \\ \bar{t} \end{matrix}$ .

# COMBINING PDGs

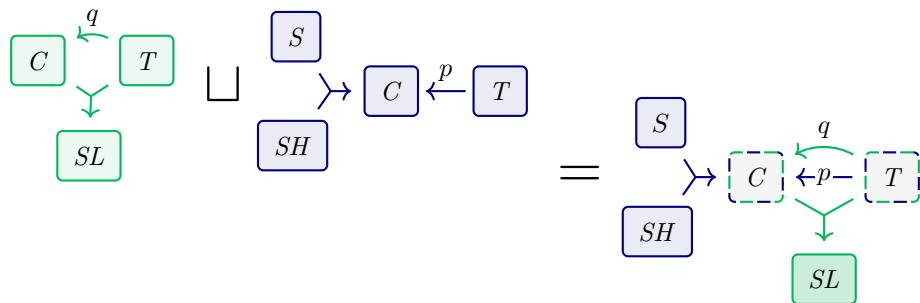




# COMBINING PDGs

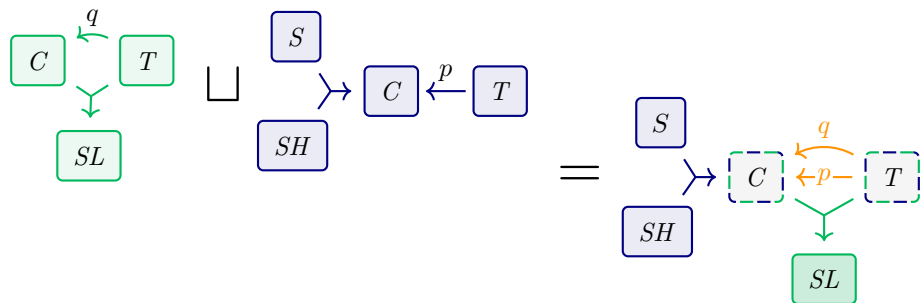


# COMBINING PDGs



- 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

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Projects 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  $\mathbf{m} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$ , where

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

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

$\mathcal{V}(\mathbf{m}) := \prod_{X \in \mathcal{N}} \mathcal{V}(X)$  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:

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

► (or hyper-edges)

## Definition (Probabilistic Dependency Graph)

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

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

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

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

► (or hyper-edges)

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

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

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

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

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

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

# OUTLINE FOR SECTION 4

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Projects Work

# SEMANTICS OF PDGS

How to stitch cpds together?

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

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



# SEMANTICS OF PDGS

How to stitch cpds together?

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

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  $\mathbf{m}$ ;

# THE SCORING FUNCTION

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

# THE SCORING FUNCTION

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

**Intuition:** Measure  $\mu$ 's violation of  $\mathcal{M}$ 's cpds.

# THE SCORING FUNCTION

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

## Definition (*Inc*)

The *incompatibility* of a joint distribution  $\mu$  with  $\mathcal{m}$  is given by

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

# 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  $\mu$  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

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

## Definition (*Inc*)

The *incompatibility* of a joint distribution  $\mu$  with  $\mathcal{m}$  is given by

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

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

the relative entropy  
(KL Divergence)  
from  $\nu$  to  $\mu$ .

# THE SCORING FUNCTION

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

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

# THE SCORING FUNCTION

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

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

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



# THE SCORING FUNCTION

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

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

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

$H(\mu)$

# THE SCORING FUNCTION

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

## Definition (*IDef*)

The *information deficiency* of a distribution  $\mu$  with respect to  $\mathcal{M}$  is

$$\text{IDef}_{\mathcal{M}}(\mu) := \sum_{X \xrightarrow{L} Y} \alpha_L \text{H}_{\mu}(Y | X) - \text{H}(\mu).$$

# THE SCORING FUNCTION

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

## Definition (*IDef*)

The *information deficiency* of a distribution  $\mu$  with respect to  $\mathcal{M}$  is

$$\text{IDef}_{\mathcal{M}}(\mu) := \sum_{X \xrightarrow{L} Y} \alpha_L \text{H}_{\mu}(Y | X) - \underbrace{\text{H}(\mu)}.$$

(a) # bits needed to determine all variables

# THE SCORING FUNCTION

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

## Definition (*IDef*)

The *information deficiency* of a distribution  $\mu$  with respect to  $\mathcal{M}$  is

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

$$\text{IDef}_{\mathcal{M}}(\mu) := \underbrace{\sum_{X \xrightarrow{L} Y} \alpha_L \text{H}_{\mu}(Y|X)}_{\text{(b)}} - \underbrace{\text{H}(\mu)}_{\text{(a)}}$$

(a) # bits needed to determine all variables

# THE SCORING FUNCTION

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

## Definition (*IDef*)

The *information deficiency* of a distribution  $\mu$  with respect to  $\mathcal{M}$  is

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

$$\text{IDef}_{\mathcal{M}}(\mu) := \sum_{X \xrightarrow{L} Y} \alpha_L \text{H}_{\mu}(Y | X) - \text{H}(\mu).$$

(a) # bits needed to determine all variables

# THE SCORING FUNCTION

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

tradeoff parameter  $\gamma \geq 0$

## Definition (*Inc*)

The *incompatibility* of  $\mu$  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 deficiency* of  $\mu$ :

$$\text{IDef}_{\mathcal{m}}(\mu) = \sum_{X \xrightarrow{L} Y} \alpha_L \mathbf{H}_{\mu}(Y|X) - \mathbf{H}(\mu)$$

# THE SCORING FUNCTION

- We are interested in the quantitative limit (small  $\gamma$ )

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

## Definition (*Inc*)

The *incompatibility* of  $\mu$  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 deficiency* of  $\mu$ :

$$\text{IDef}_{\mathcal{m}}(\mu) = \sum_{X \xrightarrow{L} Y} \alpha_L \mathbf{H}_{\mu}(Y|X) - \mathbf{H}(\mu)$$

# THE OPTIMAL DISTRIBUTION(S)

We have a scoring function  $[[\mathbf{m}]]_\gamma : \Delta\mathcal{V}(\mathbf{m}) \rightarrow \mathbb{R}$ .



# THE OPTIMAL DISTRIBUTION(S)

We have a scoring function  $[[\mathcal{m}]]_{\gamma} : \Delta\mathcal{V}(\mathcal{m}) \rightarrow \mathbb{R}$ .

Let  $[[\mathcal{m}]]_{\gamma}^*$  be the set of best-scoring distributions.

# THE OPTIMAL DISTRIBUTION(S)

We have a scoring function  $[[\mathcal{m}]]_\gamma : \Delta\mathcal{V}(\mathcal{m}) \rightarrow \mathbb{R}$ .

Let  $[[\mathcal{m}]]_\gamma^*$  be the set of best-scoring distributions.

**Proposition** (uniqueness for small  $\gamma$ , informal)

*As  $\gamma \rightarrow 0$ , there is a unique optimal distribution, which we call  $[[\mathcal{m}]]^*$ .*

# OUTLINE FOR SECTION 5

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

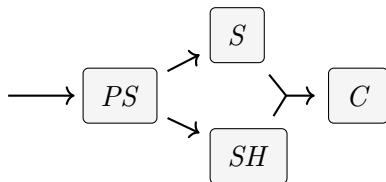
## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

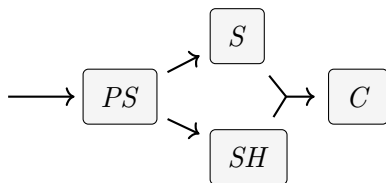
- Category Theory
- Databases
- Other Projects Work

# CAPTURING BAYESIAN NETWORKS



# CAPTURING BAYESIAN NETWORKS

For any  $\beta$ , let  $\mathbf{m}_{\mathcal{B},\beta}$  be the PDG corresponding to  $\mathcal{B}$ , with  $\alpha = \mathbf{1}$ , and confidences  $\beta$ .



# CAPTURING BAYESIAN NETWORKS

For any  $\beta$ , let  $\mathbf{m}_{\mathcal{B},\beta}$  be the PDG corresponding to  $\mathcal{B}$ , with  $\alpha = \mathbf{1}$ , and confidences  $\beta$ .

## 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  $\beta$ ,*

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

# CAPTURING BAYESIAN NETWORKS

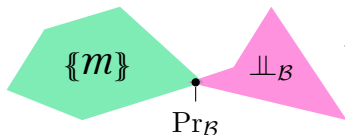
For any  $\beta$ , let  $\mathbf{m}_{\mathcal{B},\beta}$  be the PDG corresponding to  $\mathcal{B}$ , with  $\alpha = \mathbf{1}$ , and confidences  $\beta$ .

## 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  $\beta$ ,

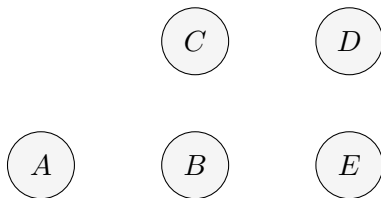
$$\llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket_{\gamma}^* = \{\text{Pr}_{\mathcal{B}}\}, \quad \text{and thus} \quad \llbracket \mathbf{m}_{\mathcal{B},\beta} \rrbracket^* = \text{Pr}_{\mathcal{B}}.$$

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

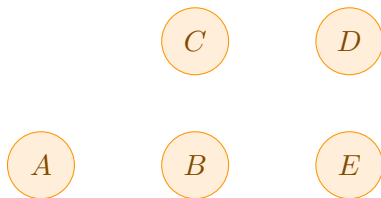


## Definition

A *factor graph*  $\Phi$  is



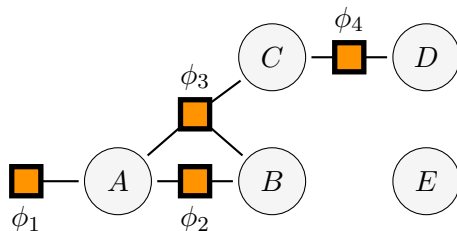
# FACTOR GRAPHS



## Definition

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

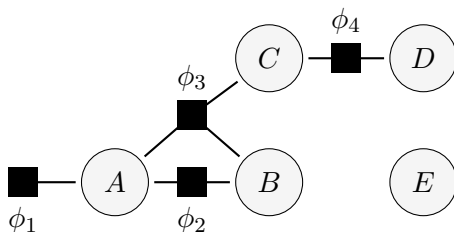
# FACTOR GRAPHS



## Definition

A *factor graph*  $\Phi$  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



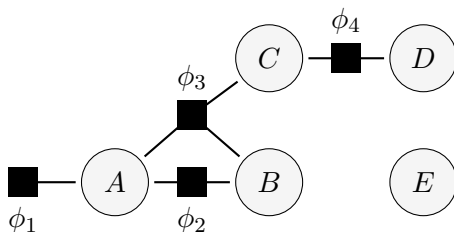
## Definition

A *factor graph*  $\Phi$  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  $X_J \subseteq \mathcal{X}$ ;  $\Phi$  defines a **distribution**

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

where  $Z_{\Phi}$  is the normalization constant.

# FACTOR GRAPHS



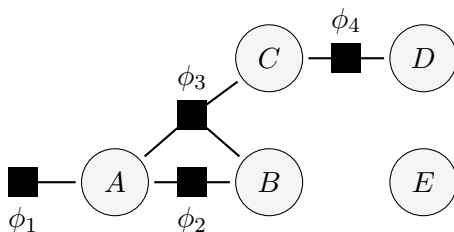
## Definition

A *factor graph*  $\Phi$  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  $X_J \subseteq \mathcal{X}$ ;  $\Phi$  defines a distribution

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

where  $Z_{\Phi}$  is the normalization constant.

# FACTOR GRAPHS



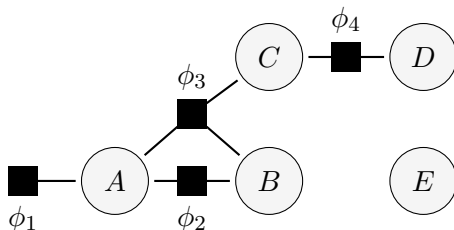
## Definition

A *factor graph*  $\Phi$  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  $X_J \subseteq \mathcal{X}$ ;  $\Phi$  defines a distribution

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

where  $Z_{\Phi}$  is the normalization constant.

# FACTOR GRAPHS



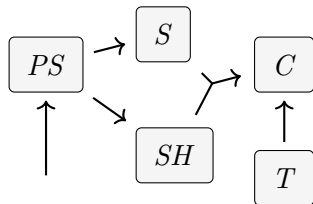
## Definition

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  $X_J \subseteq \mathcal{X}$ ;  $\Psi$  defines a distribution

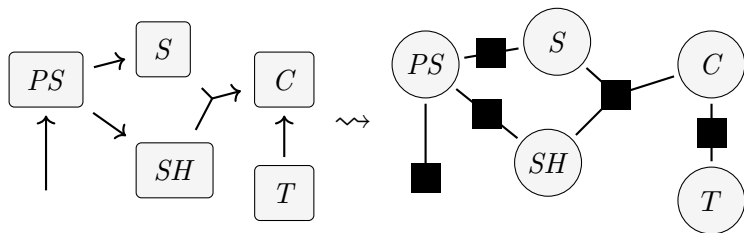
$$\Pr_{\Psi}(\vec{x}) := \frac{1}{Z_{\Psi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J)^{\theta_J},$$

where  $Z_{\Psi}$  is the normalization constant.

# PDGs AS FACTOR GRAPHS

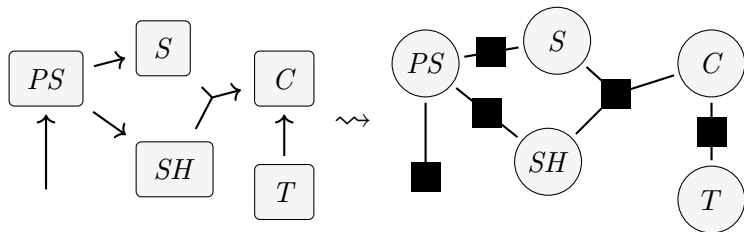


# PDGs AS FACTOR GRAPHS



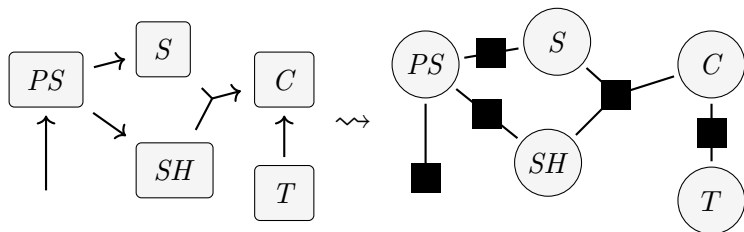


## PDGs AS FACTOR GRAPHS



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

# PDGs AS FACTOR GRAPHS

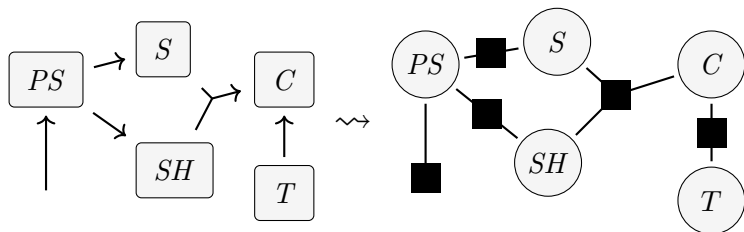


The cpds of a PDG are essentially factors. Are the semantics the same?

**Theorem (Yes, for  $\gamma = 1$ )**

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

# PDGs AS FACTOR GRAPHS



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

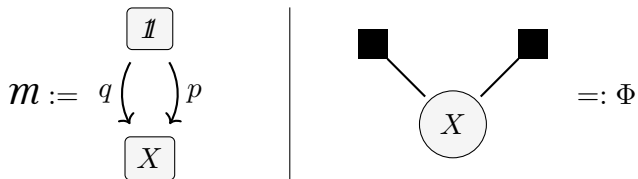
**Theorem (Yes, for  $\gamma = 1$ )**

$[[\mathcal{N}]]_1^* = \text{Pr}_{\Phi_{\mathcal{N}}}$  for all unweighted PDGs  $\mathcal{N}$ .

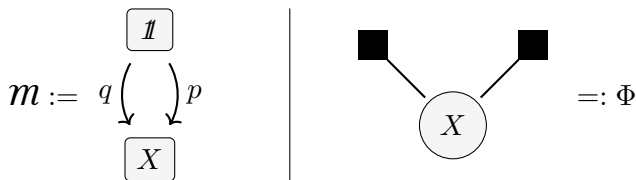
**Theorem (generalization to weighted factor graphs)**

*Semantics match (for specific  $\gamma$ ) if  $\beta \propto \alpha$ .*

# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

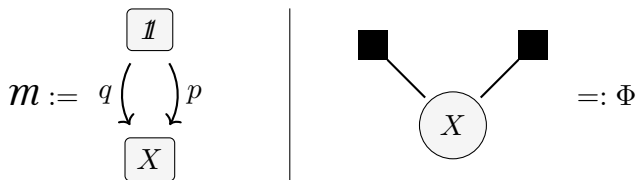


# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



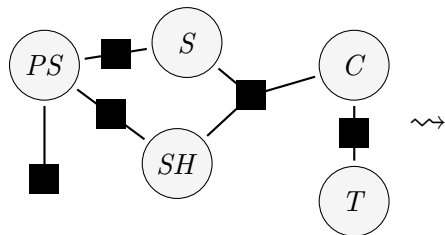
- If  $p = q$ , then  $[[m]]^* = p = q \dots$

# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

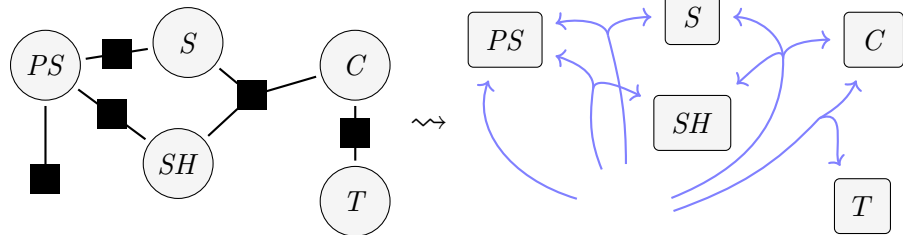


- If  $p = q$ , then  $[[\mathcal{m}]]^* = p = q \dots$
- $\dots$  but  $\Pr_{\Phi} \propto p^2$

# FACTOR GRAPHS AS PDGs

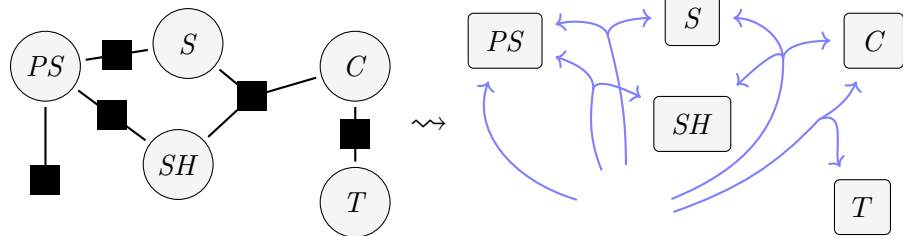


# FACTOR GRAPHS AS PDGs





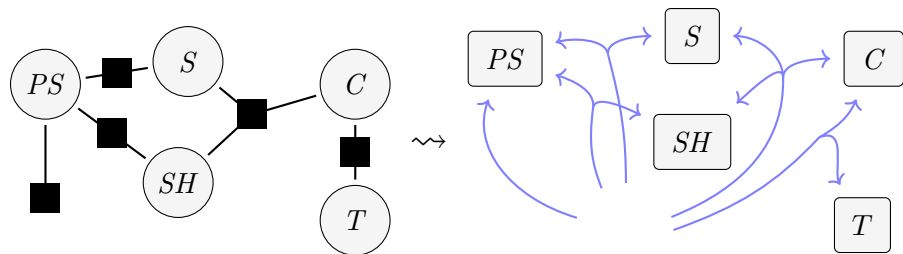
# FACTOR GRAPHS AS PDGs



## Theorem

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

# FACTOR GRAPHS AS PDGs



## Theorem

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

## Theorem

A similar result holds for weighted factor graphs.

# OUTLINE FOR SECTION 6

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Projects Work

# INCONSISTENCY

Inconsistency: the optimal value of the scoring function.

$$\langle\langle \mathcal{m} \rangle\rangle_{\gamma} := \inf_{\mu} \llbracket \mathcal{m} \rrbracket_{\gamma}(\mu)$$

# INCONSISTENCY

Inconsistency: the optimal value of the scoring function.

$$\langle\langle \mathbf{m} \rangle\rangle_{\gamma} := \inf_{\mu} \llbracket \mathbf{m} \rrbracket_{\gamma}(\mu)$$

Nice property for minimization:

- $\langle\langle \mathbf{m} \rangle\rangle$  is strictly convex and smooth in cpds (on the interior)

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## INFERENCE VIA INCONSISTENCY REDUCTION

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

### 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,

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

## INFERENCE VIA INCONSISTENCY REDUCTION

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

### 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  $m$ .**



## INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event  $Y=y$  with the cpd  $\frac{\delta_y}{\rightarrow} Y$ .

### 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  $\mathcal{m} \sqcup p$ .

## INFERENCE VIA INCONSISTENCY REDUCTION

Identify the event  $Y=y$  with the cpd  $\delta_{y \rightarrow Y}$ .

### 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  $\delta_{y \rightarrow Y}$ .

### 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  $\llbracket \mathcal{m} \rrbracket^*(Y \mid X)$ ,
  - ▶ so an inconsistency oracle yields fast inference by gradient descent.

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## 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,

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

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## 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,

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

...just like for BNs and Factor Graphs.

# OUTLINE FOR SECTION 7

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Category Theory
- Databases
- Other Projects 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 regularizers, but can be wrong.

# 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 regularizers, but can be wrong.

Bayes Rule: posterior  $\propto$  likelihood  $\cdot$  prior

# 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 regularizers, but can be wrong.

Bayes Rule: posterior  $\propto$  likelihood  $\cdot$  prior

$$\log \text{posterior} = \log \text{likelihood} + \log \text{prior} + C$$

# 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 regularizers, but can be wrong.

Bayes Rule: posterior  $\propto$  likelihood  $\cdot$  prior

$$\begin{array}{ccc} \log \text{posterior} & = & \log \text{likelihood} + \log \text{prior} + C \\ \text{(new objective)} & & \text{(old objective)} \quad \text{(regularizer)} \end{array}$$

# 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 regularizers, but can be wrong.

## Surprising Result

Most standard objectives arise as the inconsistency of the natural 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 regularizers, but can be wrong.

## Surprising Result

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

## Bonus

A visual language for reasoning about relationships between objectives.

# 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:

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

# 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,*

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

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

# SURPRISE 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,*

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

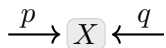
- PDG semantics just so happen to give the standard measure of compatibility between a sample and distribution.
- “surprise”: a particular kind of internal conflict.

# BIG TABLE OF OBJECTIVES

Objective	PDG	Equation
Marginal Information		$-\log p(X=x)$
Cross Entropy (Supervised)		$\frac{1}{m} \sum_{i=1}^m \left[ \log \frac{1}{f(y^i x^i)} \right] - H_{\text{data}}(Y X)$
Accuracy		$-\beta \log (\text{accuracy}_{f,D}(h))$
Square Loss		$\mathbb{E}_D \left( f(X) - h(X) \right)^2$

# INCONSISTENCY AS A DIVERGENCE

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

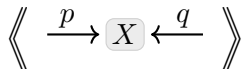




# INCONSISTENCY AS A DIVERGENCE

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  $q$ ?

$$\left\langle \left\langle \frac{p}{(\beta:r)} \rightarrow X \leftarrow \frac{q}{(\beta:s)} \right\rangle \right\rangle$$

# INCONSISTENCY AS A DIVERGENCE

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

Your inconsistency: a divergence between  $p$  and  $q$ ?

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

# INCONSISTENCY AS A DIVERGENCE

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

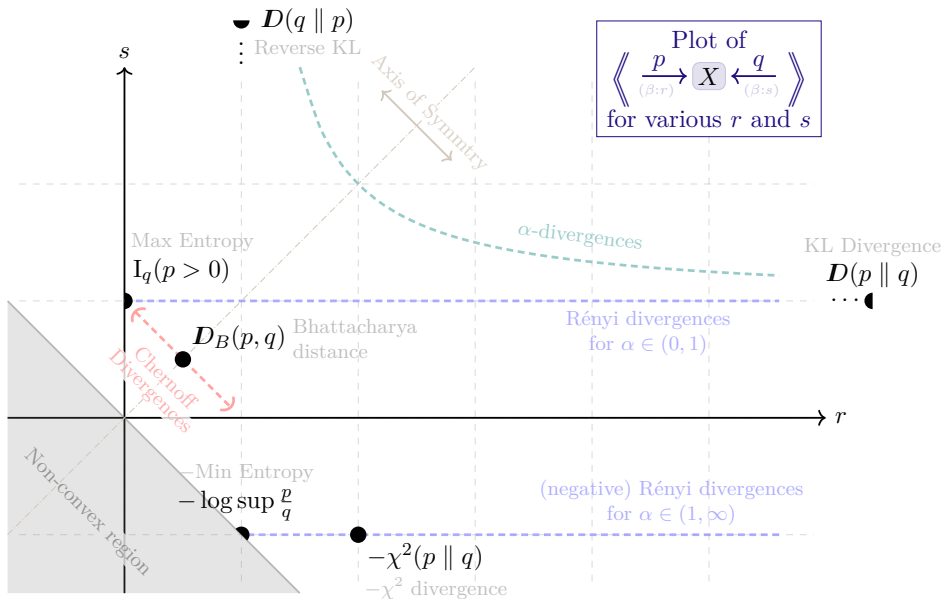
Your inconsistency: a divergence between  $p$  and  $q$ ?

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

## Lemma (Closed Form)

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

# DIVERGENCES AS INCONSISTENCIES



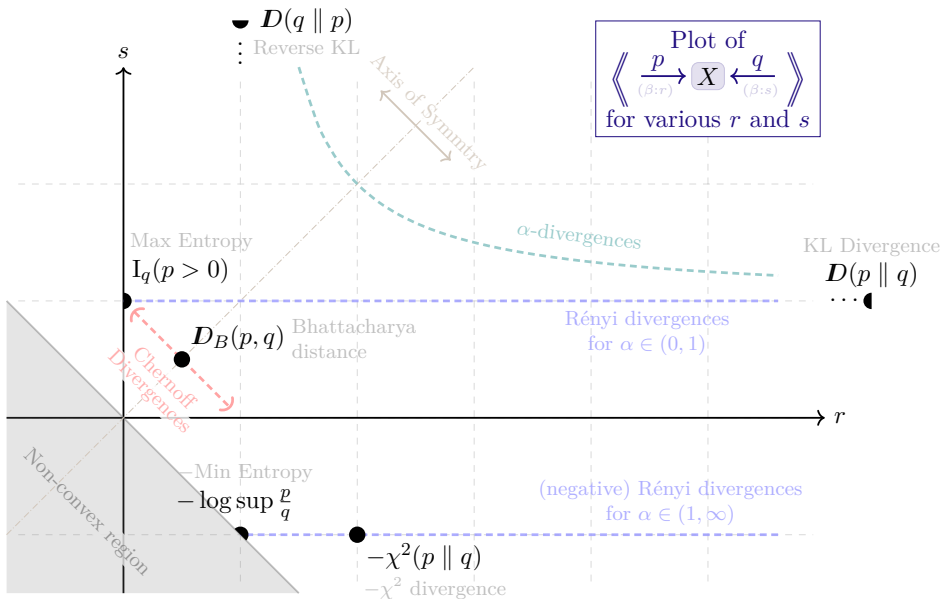
# A VERY USEFUL FACT

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

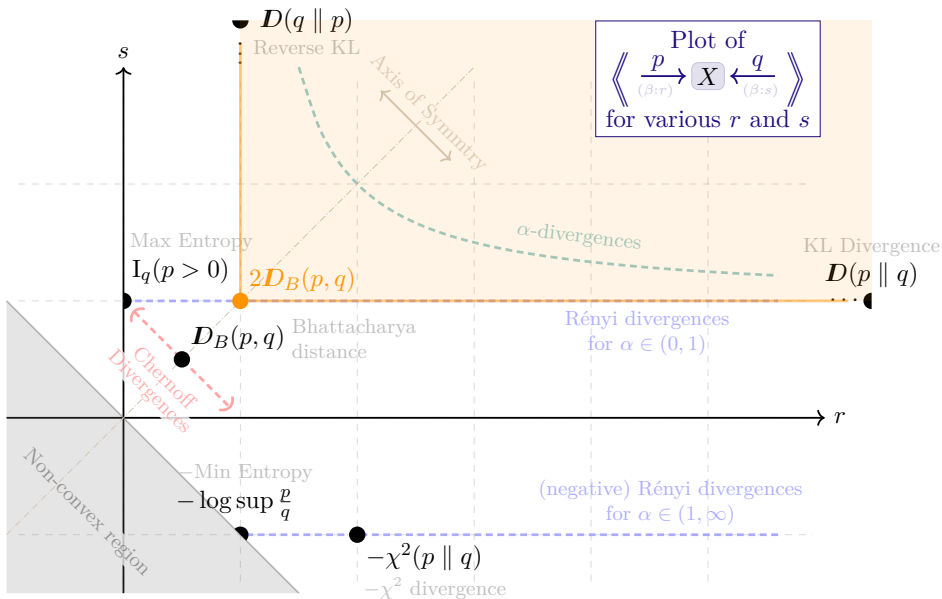
## Lemma (monotonicity of inconsistency)

- 1  $\langle\langle m \sqcup m' \rangle\rangle \geq \langle\langle m \rangle\rangle$ .
- 2 If  $\beta > \beta'$ , then  $\langle\langle m \rangle\rangle \geq \langle\langle m' \rangle\rangle$ .

# DIVERGENCES AS INCONSISTENCIES



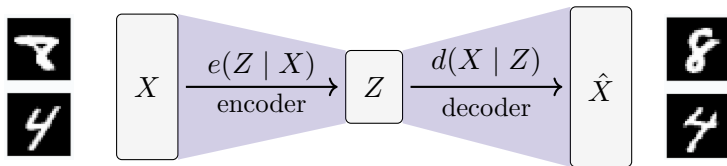
# DIVERGENCES AS INCONSISTENCIES





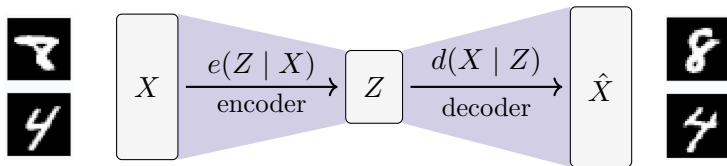
# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



# VARIATIONAL AUTO-ENCODERS, TAKE 1

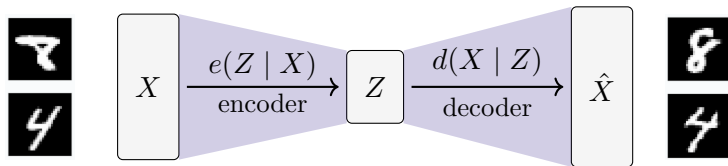
- Structure consists of two neural networks (cpds):



- Objective:

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

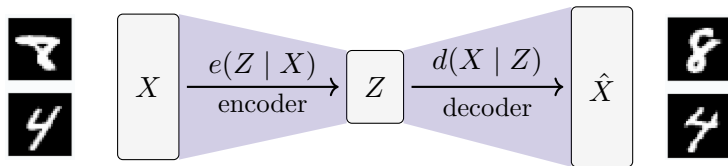


- Objective:

- ▶ 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):

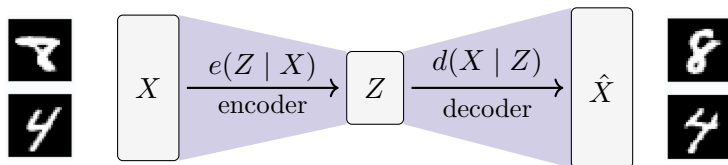


- Objective:

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



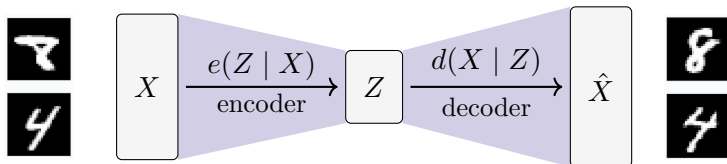
- Objective:

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



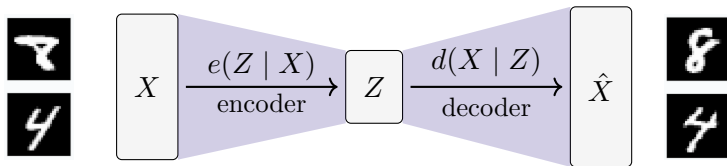
- Objective:

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



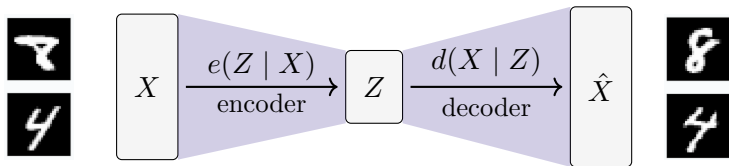
- Objective:

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

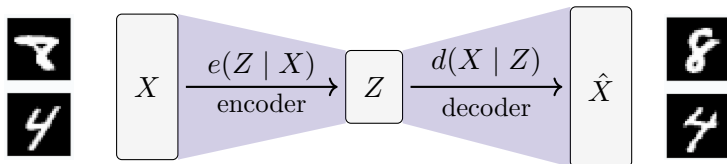
- ▶ For each  $x$ , want to minimize  $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$  <sup>“reconstruction error”</sup>
- ▶ Also have a distribution  $p(Z)$  that we want encodings to match.
- ▶ Combine to get VaE objective:

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



# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



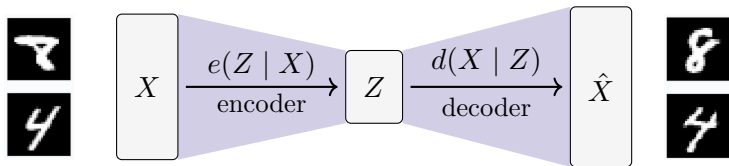
- Objective:

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

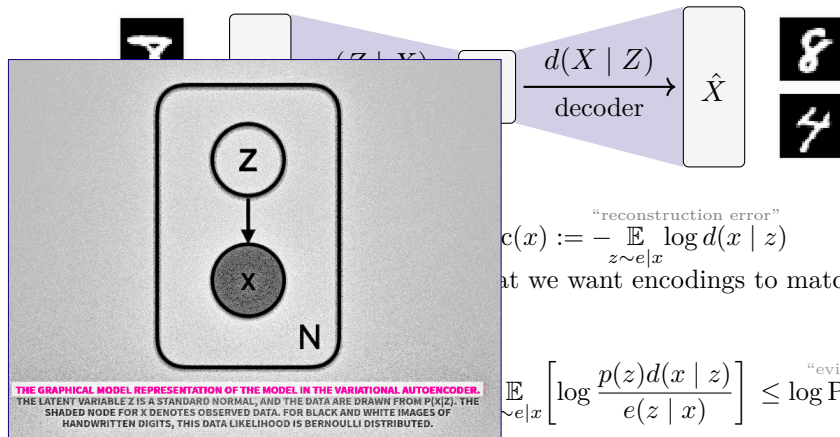
- ▶ For each  $x$ , want to minimize  $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$  “reconstruction error”
- ▶ Also have a distribution  $p(Z)$  that we want encodings to match.
- ▶ Combine to get VaE objective:

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

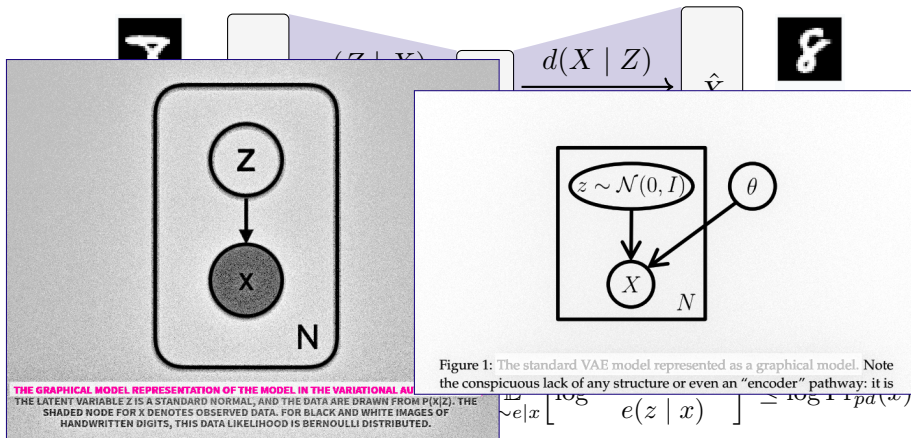
- Structure consists of two neural networks (cpds):



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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

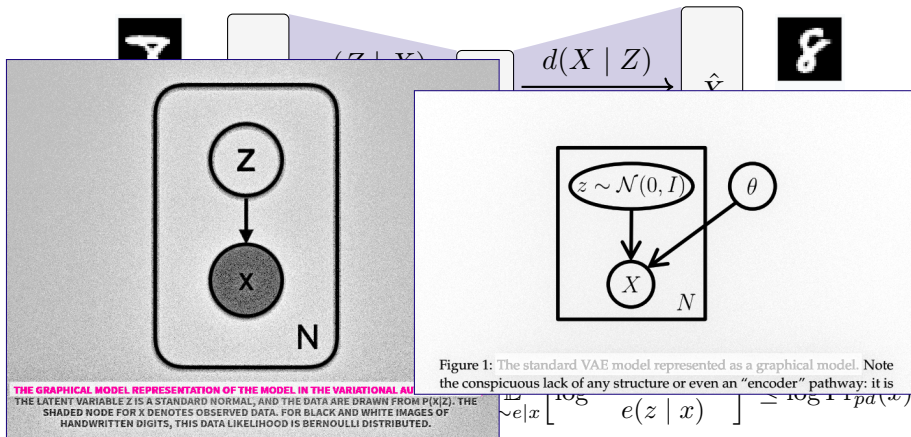
- Structure consists of two neural networks (cpds):



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):

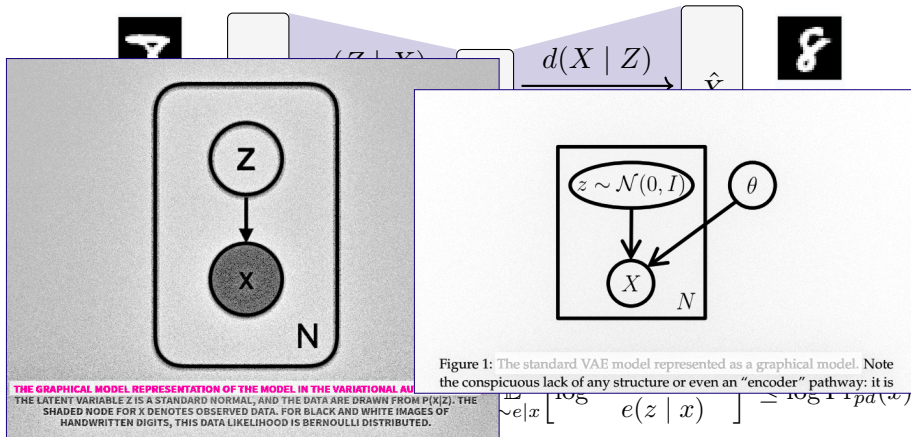


Urge to use graphical models (even if can't quite capture *entire* VaE)

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

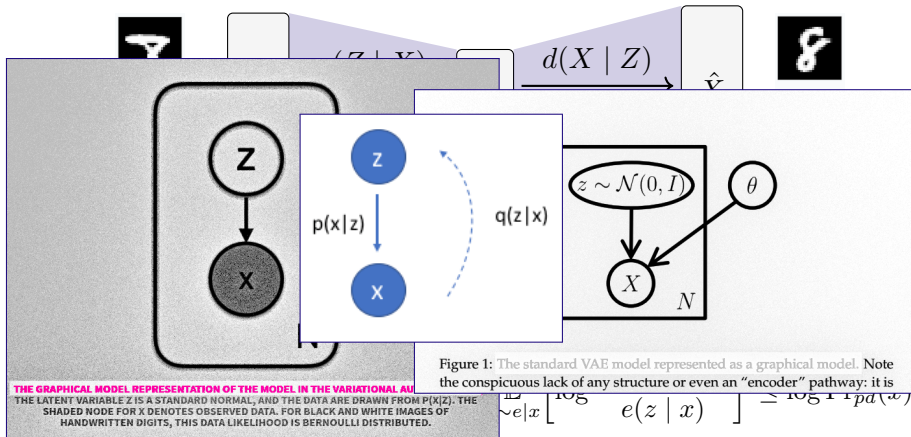


Urge to use graphical models (even if can't quite capture *entire* VaE)

- $e(Z | X)$  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:

$Z$

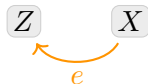
$X$



# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$  : encoder

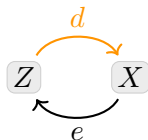


# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder



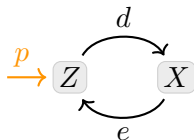
# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior



# VARIATIONAL AUTO-ENCODERS, TAKE 2

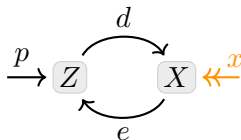
- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$



# VARIATIONAL AUTO-ENCODERS, TAKE 2

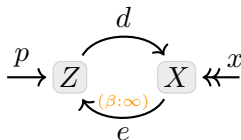
- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

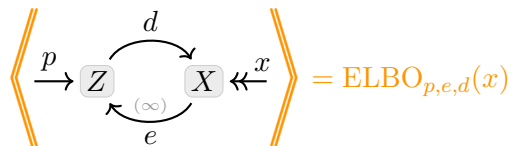
$p(Z)$  : prior

- observe a sample  $x$



# VARIATIONAL AUTO-ENCODERS, TAKE 2

Objective function is free:



- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$

# VARIATIONAL AUTO-ENCODERS, TAKE 2

Objective function is free:

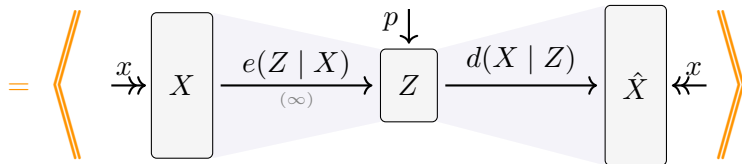
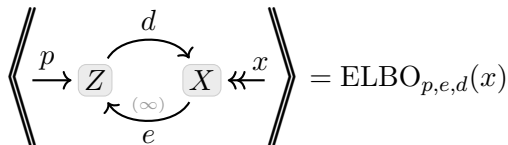
- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$



# VARIATIONAL AUTO-ENCODERS, TAKE 2

Objective function is free:

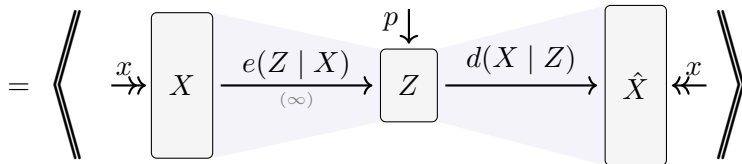
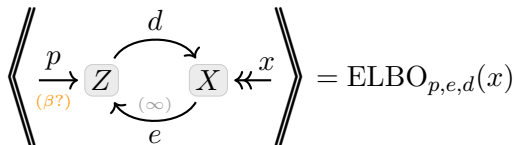
- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$





# VISUAL PROOF: THE VARIATIONAL BOUND

# VISUAL PROOF: THE VARIATIONAL BOUND

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

# VISUAL PROOF: THE VARIATIONAL BOUND

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



# 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 variable confidence in each ( $\beta$  and  $\alpha$ ).

# 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 variable confidence in each ( $\beta$  and  $\alpha$ ).
- 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 variable confidence in each ( $\beta$  and  $\alpha$ ).
- 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 variable confidence in each ( $\beta$  and  $\alpha$ ).
- 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 variable confidence in each ( $\beta$  and  $\alpha$ ).
- 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

- What are Floomps?
- Smoking BN Manipulations
- Union and Restriction

## 3 SYNTAX

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

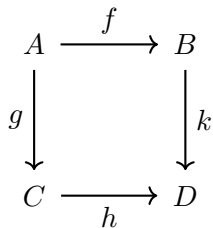
- Motivation
- Standard Metrics
- Inconsistency and Statistical Divergences
- Variational AutoEncoders

## 8 OTHER ASPECTS OF PDGs

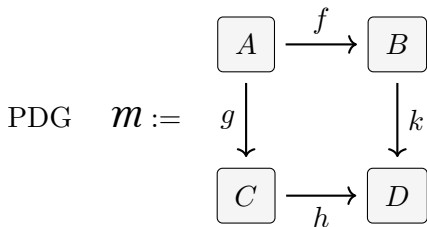
- Category Theory
- Databases
- Other Projects Work

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

Commutative diagram



says  $\forall a. kfa = hga.$



PDG inconsistency measures how close the diagram is to commuting

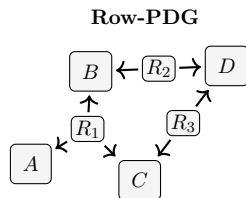
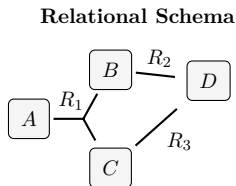
$$\exp\left(-\frac{1}{\gamma}\llbracket m \rrbracket_{\gamma}\right) = \#\{a : kfa = hga\}$$

# DATABASES

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  $\mu$  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.



# ONGOING AND FUTURE WORK

- Belief propagation as local inconsistency reduction

# ONGOING AND FUTURE WORK

- Belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: semantics for incomplete cpds

# ONGOING AND FUTURE WORK

- Belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: semantics for incomplete cpds
- Trace Semantics and Composition

# ONGOING AND FUTURE WORK

- Belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: semantics for incomplete cpds
- Trace Semantics and Composition
  - ▶ Extend semantics to score other objects, not just joint distributions.

# ONGOING AND FUTURE WORK

- Belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: semantics for incomplete cpds
- Trace Semantics and Composition
  - ▶ Extend semantics to score other objects, not just joint distributions.
  - ▶ Regarding PDGs as probabilistic automata.

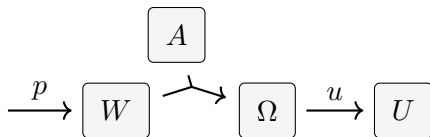
# ONGOING AND FUTURE WORK

- Belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: semantics for incomplete cpds
- 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? \*

# ONGOING AND FUTURE WORK

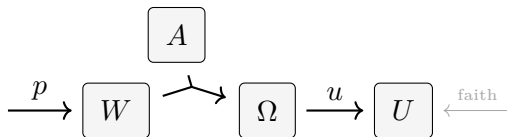
- Belief propagation as local inconsistency reduction
- Properties of “sub-stochastic” PDGs: semantics for incomplete cpds
- Trace Semantics and Composition
  - ▶ Extend semantics to score other objects, not just joint distributions.
  - ▶ Regarding PDGs as probabilistic automata.
  - ▶ Open Question: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

# A DIFFERENT FLAVOR OF AGENT

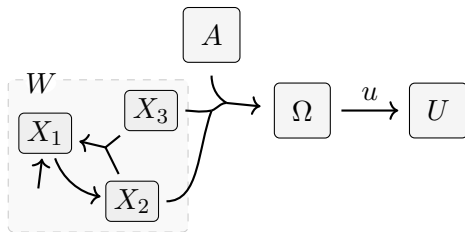




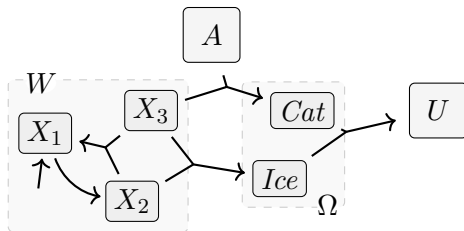
# A DIFFERENT FLAVOR OF AGENT



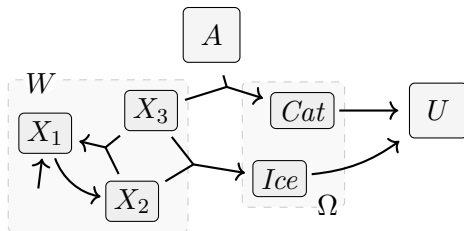
# A DIFFERENT FLAVOR OF AGENT



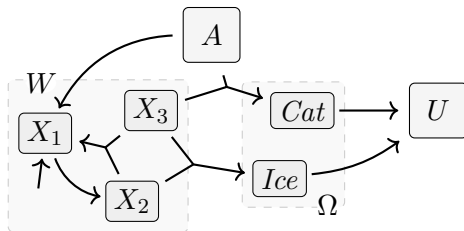
# A DIFFERENT FLAVOR OF AGENT



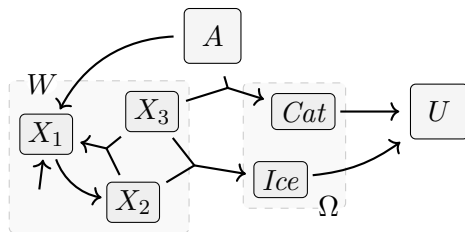
# A DIFFERENT FLAVOR OF AGENT



# A DIFFERENT FLAVOR OF AGENT

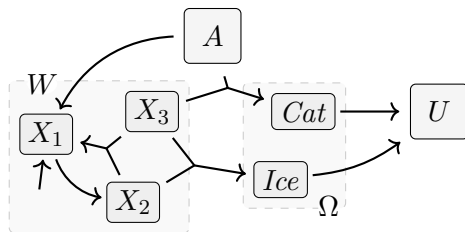


# A DIFFERENT FLAVOR OF AGENT



- Driven by pursuit of coherence; not (necessarily) maximization.

# A DIFFERENT FLAVOR OF AGENT



- Driven by pursuit of coherence; not (necessarily) maximization.

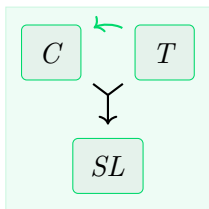
PDG Python library available at  
<https://orichardson.github.io/pdg/>

# OUTLINE FOR SECTION 9

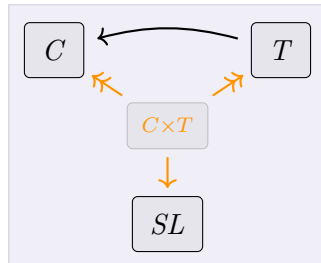
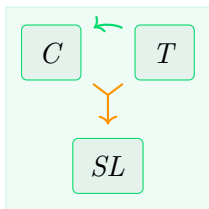
- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS
  - BNs as MaxEnt
- 13 MORE LOSSES
  - Regularizers
- 14 MORE VISUAL PROOFS
- 15 MORE CATEGORY THEORY
  - PDGs as diagrams of the Markov Category



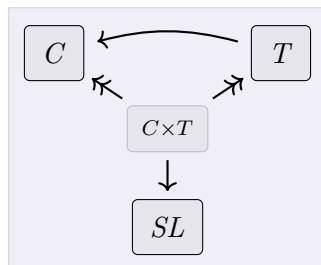
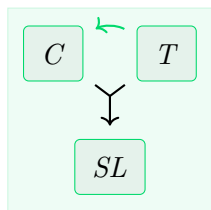
# HYPER-GRAPHS? OR MERELY GRAPHS?



# HYPER-GRAPHS? OR MERELY GRAPHS?

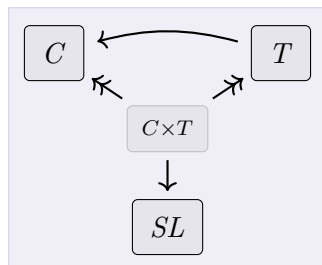
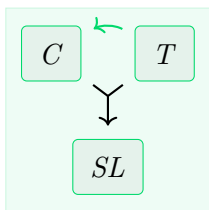


# HYPER-GRAPHS? OR MERELY GRAPHS?



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

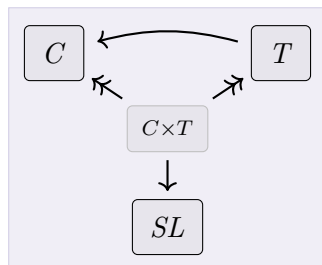
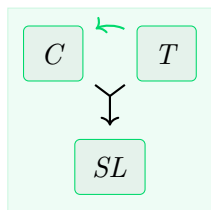
# HYPER-GRAPHS? OR MERELY GRAPHS?



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

joint distributions  $\Leftrightarrow$  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

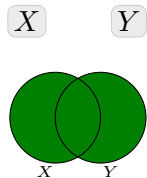
joint distributions  $\Leftrightarrow$  expanded joint distributions  
satisfying coherence constraints

(working directly with hypergraphs is also possible)

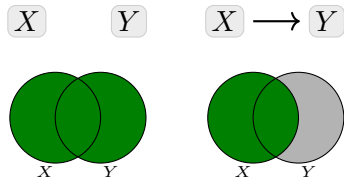
# OUTLINE FOR SECTION 10

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

# ILLUSTRATIONS OF *IDef*

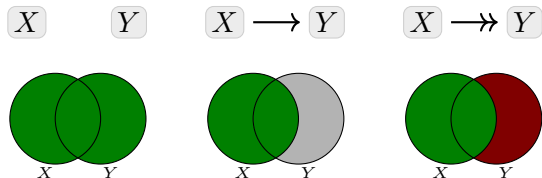


# ILLUSTRATIONS OF $IDef$



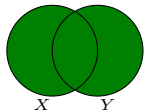


# ILLUSTRATIONS OF $IDef$

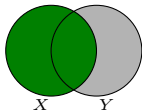


# ILLUSTRATIONS OF $IDef$

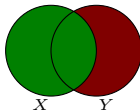
$X$     $Y$



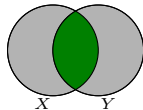
$X \rightarrow Y$

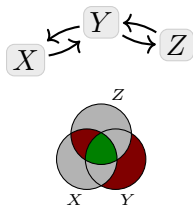
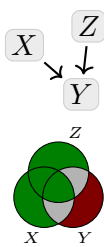
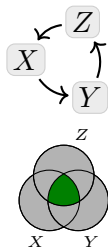
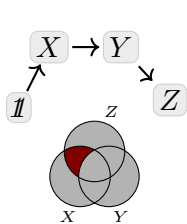
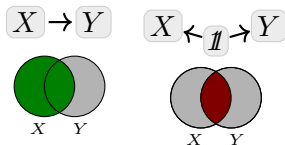
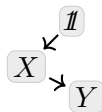
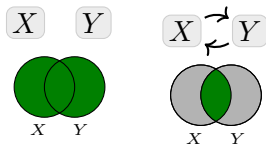
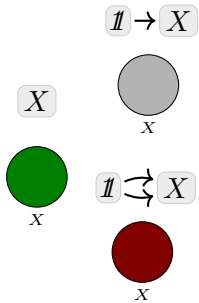


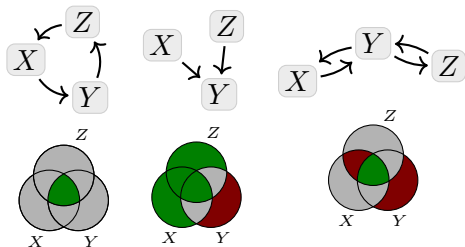
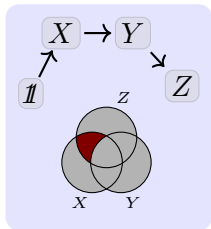
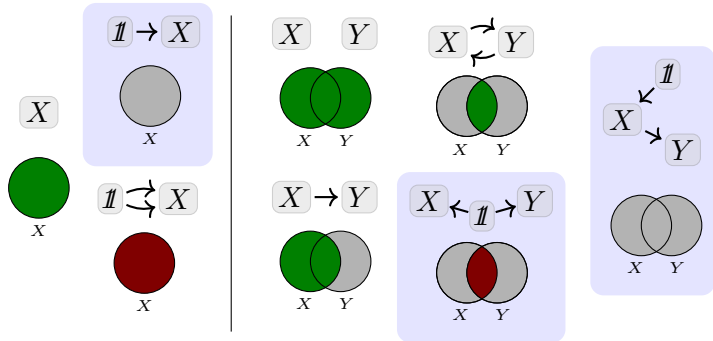
$X \twoheadrightarrow Y$



$X \rightleftarrows Y$







# 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
  - Regularizers
- 14 MORE VISUAL PROOFS
- 15 MORE CATEGORY THEORY
  - PDGs as diagrams of the Markov Category

## Proposition (uniqueness for small $\gamma$ )

- 1 If  $0 < \gamma \leq \min_L \beta_L^m$ , then  $[[\mathcal{M}]]_\gamma^*$  is a singleton.
- 2  $\lim_{\gamma \rightarrow 0} [[\mathcal{M}]]_\gamma^*$  exists and is a singleton.

## Proposition (uniqueness for small $\gamma$ )

- 1 If  $0 < \gamma \leq \min_L \beta_L^m$ , then  $[[\mathcal{M}]]_\gamma^*$  is a singleton.
- 2  $\lim_{\gamma \rightarrow 0} [[\mathcal{M}]]_\gamma^*$  exists and is a singleton.

This lets us define  $[[\mathcal{M}]]^* := \text{unique element } \left( \lim_{\gamma \rightarrow 0} [[\mathcal{M}]]_\gamma^* \right)$ .

# 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.*)

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



# ANOTHER VIEW OF PDG SEMANTICS

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

# COMPARING PDG TO FACTOR GRAPH SEMANTICS

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

# COMPARING PDG TO FACTOR GRAPH SEMANTICS

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

And the weighted factor graph's canonical scoring function:

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

# PROPERTIES OF INCONSISTENCY, FOR MINIMIZATION

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

Nice properties for minimization:

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

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

# 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
  - Regularizers
- 14 MORE VISUAL PROOFS
- 15 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 $\mu$ , variance $\sigma^2$
Exponential $\text{Exp}(\lambda)$	positive support, mean $\lambda$
Factor graphs	moment matching.
...	...

# BAYESIAN NETWORKS: MAXIMUM ENTROPY?

Common distributions  
tend to maximize  
entropy subject to  
natural constraints.

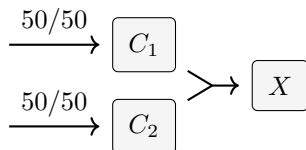
distribution	constraints
Gaussian $\mathcal{N}(\mu, \sigma^2)$	mean $\mu$ , variance $\sigma^2$
Exponential $\text{Exp}(\lambda)$	positive support, mean $\lambda$
Factor graphs	moment matching.
...	...
Bayesian Networks	cpds + ???



# BAYESIAN NETWORKS: MAXIMUM ENTROPY?

Common distributions  
tend to maximize  
entropy subject to  
natural constraints.

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

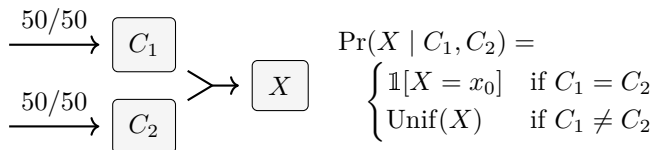


$$\Pr(X | C_1, C_2) = \begin{cases} \mathbb{1}[X = x_0] & \text{if } C_1 = C_2 \\ \text{Unif}(X) & \text{if } C_1 \neq C_2 \end{cases}$$

# BAYESIAN NETWORKS: MAXIMUM ENTROPY?

Common distributions  
tend to maximize  
entropy subject to  
natural constraints.

distribution	constraints
Gaussian $\mathcal{N}(\mu, \sigma^2)$	mean $\mu$ , variance $\sigma^2$
Exponential $\text{Exp}(\lambda)$	positive support, mean $\lambda$
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.*

*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)

If  $\beta = \gamma\alpha$ , then  $[[\mathcal{M}]]_\gamma^* = \text{Pr}_{(\Phi_m, \beta)}$ .

Concretely, 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  $[[\langle \mathcal{N}, \mathbf{v}, \gamma\mathbf{v} \rangle]]_\gamma = \gamma \text{VFE}_{(\Phi_n, \mathbf{v})}$ ; consequently,  $[[\langle \mathcal{N}, \mathbf{v}, \gamma\mathbf{v} \rangle]]_\gamma^* = \{\text{Pr}_{(\Phi_n, \mathbf{v})}\}$ .

## Theorem (WFGs are PDGs)

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

# 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
  - Regularizers
- 14 MORE VISUAL PROOFS
- 15 MORE CATEGORY THEORY
  - PDGs as diagrams of the Markov Category

## 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

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

# 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

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

## Proposition (supervised setting: conditional cross entropy)

The inconsistency of the PDG containing  $f(Y | X)$  and a high-confidence empirical distribution  $\Pr_{\underline{\mathbf{xy}}}$  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  $\Pr_{\underline{\mathbf{xy}}}$ ). That is,

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

## 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

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

## 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

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

- Thought of as a feature of  $h$ , but as a PDG, symmetry between  $f, h$  is clear.



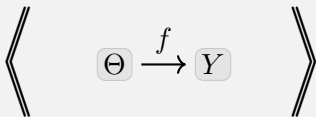
## Proposition (Mean Square Error as Inconsistency)

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

## Proposition (Regularizers as priors)

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

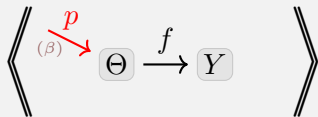
That is,



## Proposition (Regularizers as priors)

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

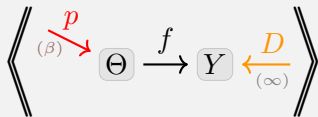
That is,



## 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,

$$\left\langle \begin{array}{c} \xrightarrow{p} \\ (\beta) \\ \xrightarrow{\theta_0} \end{array} \Theta \xrightarrow{f} Y \xleftarrow{D} \right\rangle_{(\infty)} =$$

## 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  $\beta_p$  of the prior. That is,

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

## 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  $\beta_p$  of the prior. That is,

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

## 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  $\beta_p$  of the prior. That is,

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

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

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



# 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

*Average Surprise is the inconsistency of simultaneously believing  $p$  and an empirical distribution  $\text{Pr}_{\underline{x}}$ , with high confidence (plus  $H(Y | X)$ , a constant that depends only on the data  $\text{Pr}_{\underline{xy}}$ ) That is,*

$$I_p(x) = \left\langle \begin{array}{c} \xrightarrow{p} \\ \text{X} \\ \xleftarrow[\text{(\beta:\infty)}]{\text{Pr}_{\underline{x}}} \end{array} \right\rangle + H(\text{Pr}_{\underline{x}}).$$

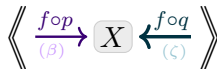
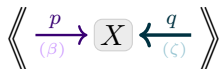
- PDG semantics just so happen to give the standard measure of compatibility between a sample and distribution.
- “surprise”: a particular kind of internal conflict.

# 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
  - Regularizers
- 14 MORE VISUAL PROOFS
- 15 MORE CATEGORY THEORY
  - PDGs as diagrams of the Markov Category

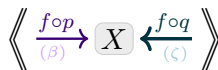
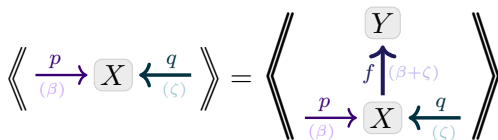
# VISUAL PROOF: DATA PROCESSING INEQUALITY

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



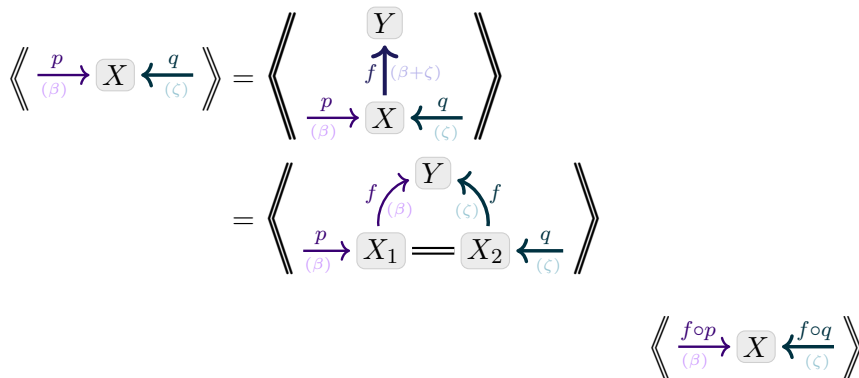
# VISUAL PROOF: DATA PROCESSING INEQUALITY

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



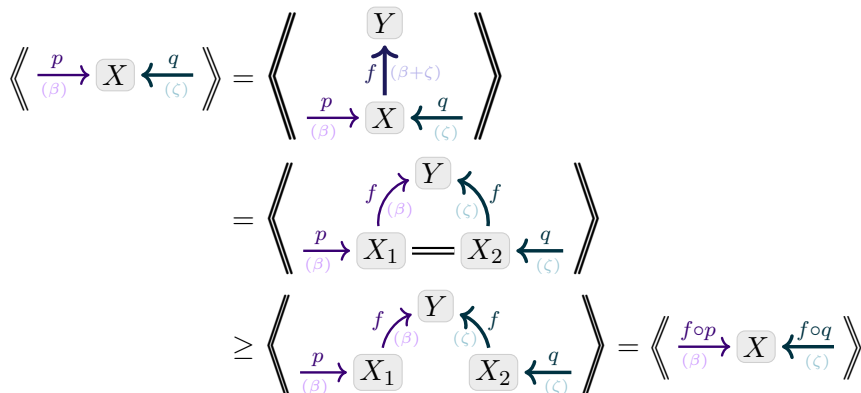
# VISUAL PROOF: DATA PROCESSING INEQUALITY

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



# VISUAL PROOF: DATA PROCESSING INEQUALITY

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



# OUTLINE FOR SECTION 15

- 9 HYPER-GRAPHS
- 10 THE INFORMATION DEFICIENCY
- 11 MORE ON SEMANTICS
- 12 MORE ON GRAPHICAL MODELS
  - BNs as MaxEnt
- 13 MORE LOSSES
  - Regularizers
- 14 MORE VISUAL PROOFS
- 15 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)

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



## Definition (PDG)

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

$\mathcal{V} : \mathcal{N} \rightarrow \text{Set}$  (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
- $(\mathcal{N}, \mathcal{E})$  is a multigraph

## Definition (PDG)

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

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

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

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

$p_L : \mathcal{V}(X) \rightarrow \Delta \mathcal{V}(Y)$  (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)

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

## Definition (PDG)

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

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

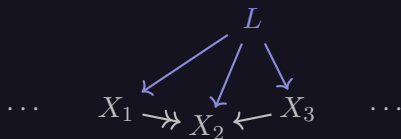
### Equivalent Categorical Definition

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

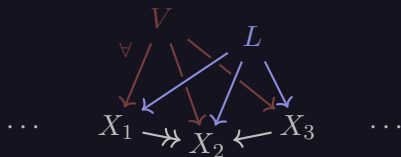
What do you do with diagrams? Take limits / colimits.

$$\dots \quad X_1 \twoheadrightarrow X_2 \longleftarrow X_3 \quad \dots$$

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

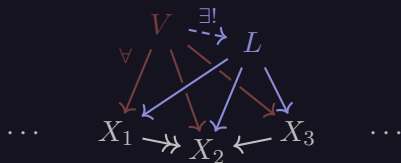


What do you do with diagrams? Take limits / colimits.





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



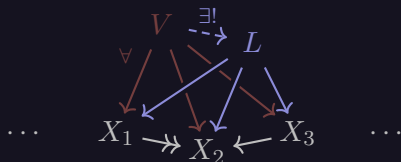
What do you do with diagrams? Take limits / colimits.



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

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

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

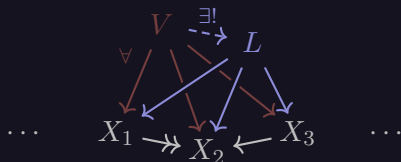


For the deterministic sub-PDG  $\mathcal{m}_{\text{det}} \subseteq \mathcal{m}$ :

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

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

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



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

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

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  $\mathcal{B}$ :  $\lim \mathcal{m}_{\mathcal{B}} = \left( \mathbb{1}, \left\{ \text{Pr}_{\mathcal{B}}(X) \right\}_{X \in \mathcal{N}} \right)$

