

# PROBABILISTIC DEPENDENCY GRAPHS AND INCONSISTENCY

HOW TO MODEL, MEASURE, AND MITIGATE INTERNAL CONFLICT

Oliver Richardson

Cornell University  
Department of Computer Science

September 2021

# OUTLINE FOR SECTION 1

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

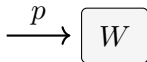
## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future 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  $u : \Omega \rightarrow \mathbb{R}$

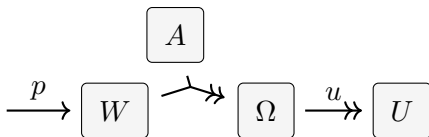
A utility function is \*not\* part of the standard way of modeling an agent with uncertainty. I would cut utility. It's a distraction.



Don't overwhelm the reader with notation. I would cut this.

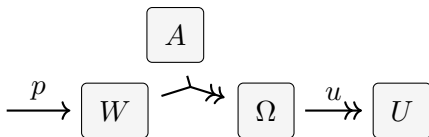
The standard way of modeling an agent with uncertainty:

- 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  $A$ .
- cut actions; it's a distraction



Such agents cannot have internal conflict;

by construction, they have consistent beliefs and desires.

# WHY INCONSISTENCY?

Why **model** inconsistency?

- Sometimes people have inconsistent beliefs;



# WHY INCONSISTENCY?

Why **model** inconsistency?

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

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

- Why entertain the possibility of being wrong?

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

- Why entertain the possibility of being wrong?
  - ▶ Assertions and test cases can fail.

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

- Why entertain the possibility of being wrong?
  - ▶ Assertions and test cases can fail.
- Why adopt beliefs without first cross-checking all knowledge?

We are *\*not\** building a system that can be inconsistent.  
I would cut all this. It's a distraction.

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

- Why entertain the possibility of being wrong?
  - ▶ Assertions and test cases can fail.
- Why adopt beliefs without first cross-checking all knowledge?
- Why consult multiple models, that can give inconsistent answers?

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

- Why entertain the possibility of being wrong?
  - ▶ Assertions and test cases can fail.
- Why adopt beliefs without first cross-checking all knowledge?
- Why consult multiple models, that can give inconsistent answers?

*A man with a watch knows what time it is.*

*A man with two watches is never sure.*

– *Segal's Law*

# WHY INCONSISTENCY?

Why **model** inconsistency?

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

Why **build a system** that can be inconsistent?

- Why entertain the possibility of being wrong?
  - ▶ Assertions and test cases can fail.
- Why adopt beliefs without first cross-checking all knowledge?
- Why consult multiple models, that can give inconsistent answers?

*A man with a watch knows what time it is.*

*A man with two watches is never sure.*

– *Segal's Law*

I would keep this and the first three lines, and cut the rest.

Freedom from perfect consistency is valuable,  
but demands that you also recognize and address internal conflict.



# YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

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

- Designed to tolerate inconsistency, so we can model it.

I don't know what it means to “tolerate” inconsistency in this context.

Why not just say that we can model inconsistent beliefs

# YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

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

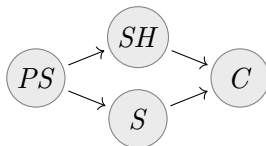
- Designed to tolerate inconsistency, so we can model it.
- In doing so, we get *much* more ...

# TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

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

an independence relation on variables

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



# TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

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

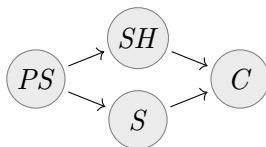
an independence relation on variables

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

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

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

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



# TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

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

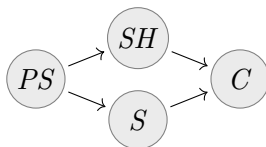
an independence relation on variables

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

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

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

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



# OUTLINE FOR SECTION 2

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

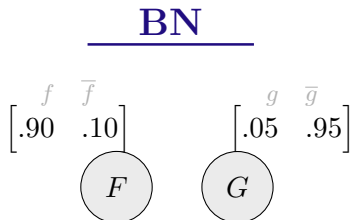
- Databases
- Open Problems + Future Work

## 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).

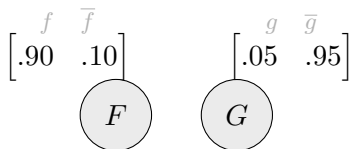




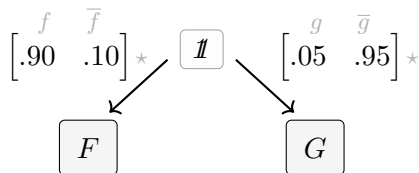
# SIMPLE EXAMPLE: FLOOMPS AND GUNS

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

**BN**

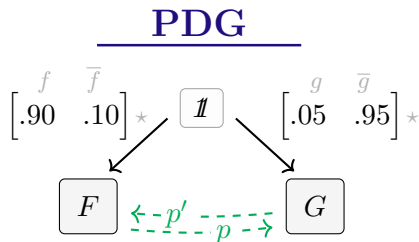
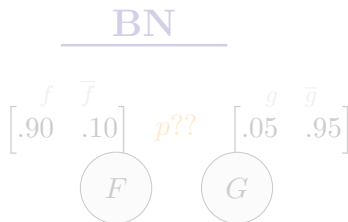


**PDG**



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

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

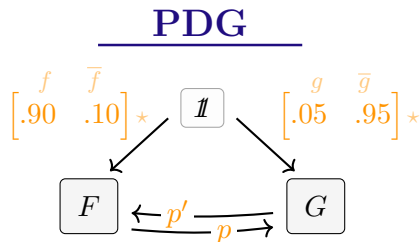
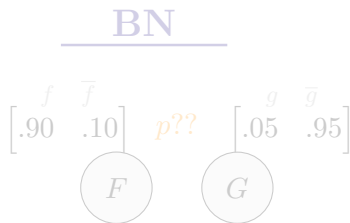


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.

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

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

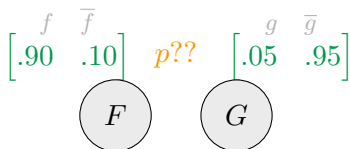
# SIMPLE EXAMPLE: FLOOMPS AND GUNS



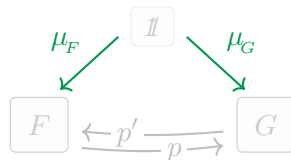
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

## BN



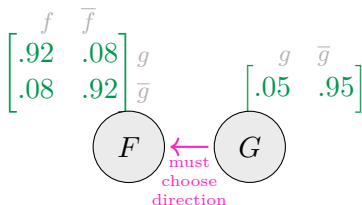
## PDG



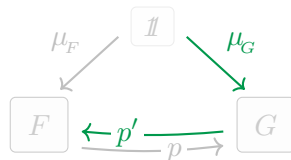
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
  - ▶ ...but BNs must resolve inconsistency first,  
which may break symmetry and irrecoverably lose information.

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

## BN



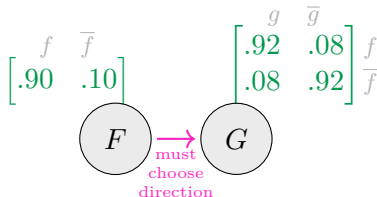
## PDG



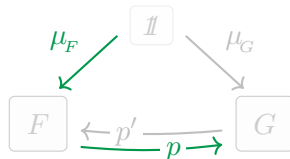
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
  - ▶ ...but BNs must resolve inconsistency first,  
which may **break symmetry** and irrecoverably lose information.

# SIMPLE EXAMPLE: FLOOMPS AND GUNS

## BN

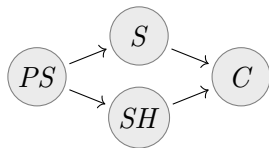


## PDG

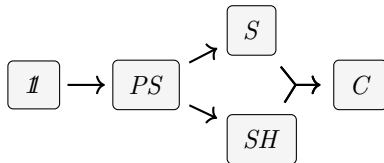
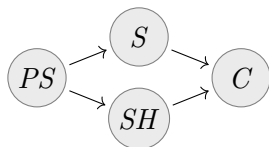


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent,
  - ▶ ...but BNs must resolve inconsistency first,  
which may **break symmetry** and irrecoverably lose information.

# BAYESIAN NETWORKS AS PDGs

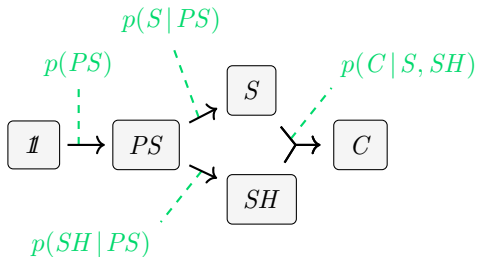
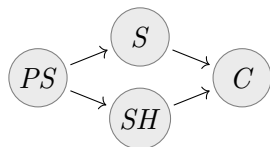


# BAYESIAN NETWORKS AS PDGs





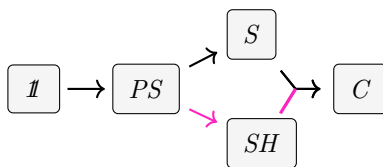
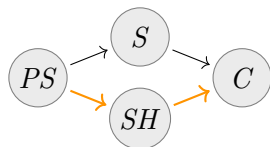
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

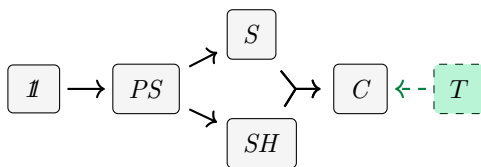
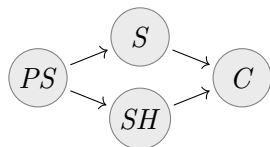
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

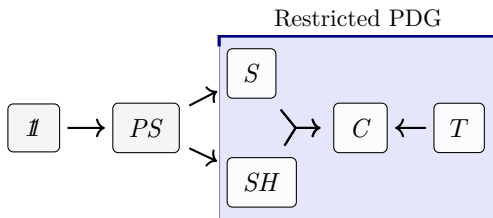
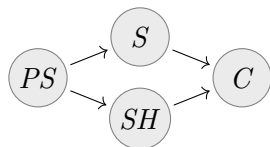
# BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

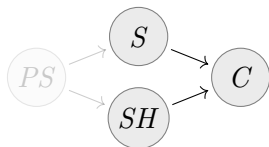
# BAYESIAN NETWORKS AS PDGs



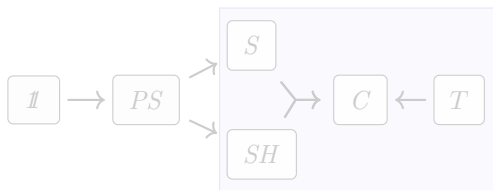
In contrast with BNs:

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

# BAYESIAN NETWORKS AS PDGs



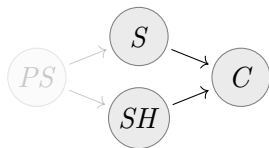
Must now give distributions on *SH* and *S*, or distinguish them as “observed” (a *conditional* BN).



In contrast with BNs:

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

# BAYESIAN NETWORKS AS PDGs



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

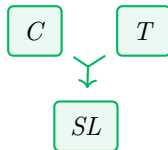
In a qualitative BN: *removing data results in new knowledge*:  $A \perp\!\!\!\perp C$ .



In contrast with BNs:

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

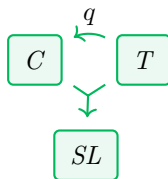
# COMBINING PDGs



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

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

# COMBINING PDGs



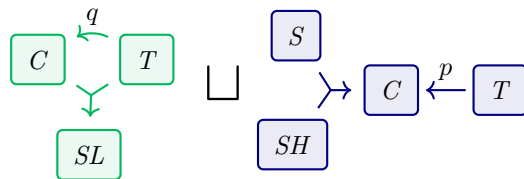
**Grok wants to be supreme leader ( $SL$ ).**

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

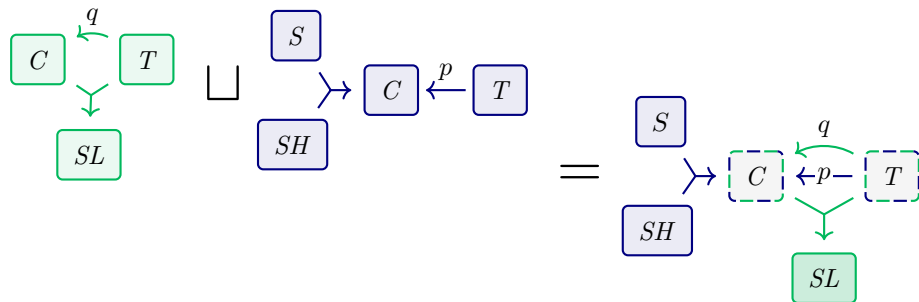
- ...but mom says  $q(C \mid T) = \begin{bmatrix} \overset{c}{.15} & \overset{\bar{c}}{.85} \\ .02 & .98 \end{bmatrix} \begin{matrix} t \\ \bar{t} \end{matrix}$ .



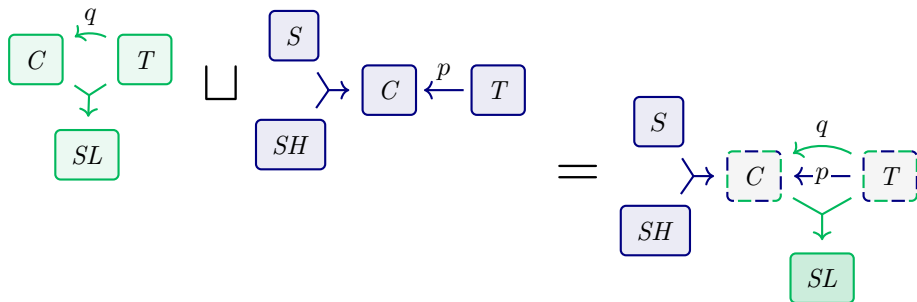
# COMBINING PDGs



# COMBINING PDGs

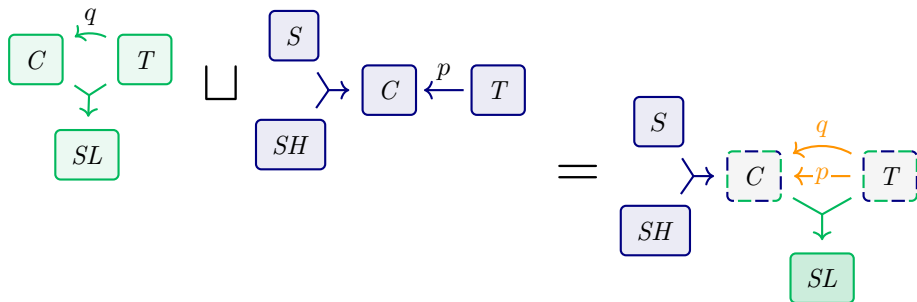


# COMBINING PDGs



- Arbitrary PDGs may be combined without loss of information

# COMBINING PDGs



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

# OUTLINE FOR SECTION 3

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future Work

## Definition (Probabilistic Dependency Graph)

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

## Definition (Probabilistic Dependency Graph)

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

## Definition (Probabilistic Dependency Graph)

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

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

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



## Definition (Probabilistic Dependency Graph)

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

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

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

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

## Definition (Probabilistic Dependency Graph)

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

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

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

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

(or hyper-edges)

## Definition (Probabilistic Dependency Graph)

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

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

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

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

(or hyper-edges)

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

## Definition (Probabilistic Dependency Graph)

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

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

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

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

(or hyper-edges)

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

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

## Definition (Probabilistic Dependency Graph)

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

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

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

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

(or hyper-edges)

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

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

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

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

$\alpha_L$  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

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future Work



# SEMANTICS OF PDGs

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

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

# SEMANTICS OF PDGS

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

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

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

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

# SEMANTICS OF PDGs

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

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

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

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

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

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

# SEMANTICS OF PDGS

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

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

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

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

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

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

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

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

# SEMANTICS OF PDGS

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

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

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

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

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

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

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

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

...

(other possibilities as well)

# SEMANTICS OF PDGs

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

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

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

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

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

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

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

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

# SEMANTICS OF PDGS

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

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

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

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

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

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

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

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

# THE SCORING FUNCTION

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

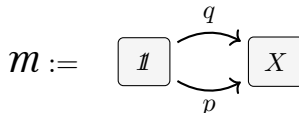


# THE SCORING FUNCTION

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

**Intuition:** Measure  $\mu$ 's violation of  $\mathbf{m}$ 's cpds.

MOTIVATING EXAMPLES.



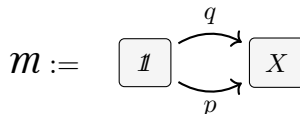
Suppose  $p = [.4, .6]$ .

# 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.

MOTIVATING EXAMPLES.



Suppose  $p = [.4, .6]$ .

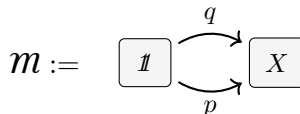
- If  $p = q$ , then  $\mathcal{m}$  is clearly consistent, and compatible with the joint distribution  $\mu(X) = p = q$ , so  $\text{Inc}_{\mathcal{m}}(p) = 0$ .

# 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.

MOTIVATING EXAMPLES.



Suppose  $p = [.4, .6]$ .

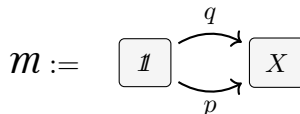
- If  $p = q$ , then  $\mathcal{m}$  is clearly consistent, and compatible with the joint distribution  $\mu(X) = p = q$ , so  $\text{Inc}_{\mathcal{m}}(p) = 0$ .
- If  $q = [.5, .5]$  then  $\mathcal{m}$  is not consistent, but  $\mu = [.45, .55]$  matches better than  $\mu = [.9, .1]$

# 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.

## MOTIVATING EXAMPLES.



Suppose  $p = [.4, .6]$ .

- If  $p = q$ , then  $\mathcal{M}$  is clearly consistent, and compatible with the joint distribution  $\mu(X) = p = q$ , so  $\text{Inc}_{\mathcal{M}}(p) = 0$ .
- If  $q = [.5, .5]$  then  $\mathcal{M}$  is not consistent, but  $\mu = [.45, .55]$  matches better than  $\mu = [.9, .1]$
- If  $q = [0, 1]$ , then  $\mathcal{M}$  is much more inconsistent than before, even though  $\llbracket \mathcal{M} \rrbracket = \emptyset$  in both cases.

# 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

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

# 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

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

# 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

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

$$D(\mu \parallel \nu) = \sum_{w \in \text{Supp}(\mu)} \mu(w) \log \frac{\mu(w)}{\nu(w)} \text{ is the relative entropy from } \nu \text{ to } \mu.$$

# 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

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

$$D(\mu \parallel \nu) = \sum_{w \in \text{Supp}(\mu)} \mu(w) \log \frac{\mu(w)}{\nu(w)} \text{ is the relative entropy from } \nu \text{ to } \mu.$$



# THE SCORING FUNCTION

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

## Definition (*Inc*)

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

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

# THE SCORING FUNCTION

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

## Definition (*IDef*)

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

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

# THE SCORING FUNCTION

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

## Definition (*IDef*)

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

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

(a) # bits needed to determine all variables

# THE SCORING FUNCTION

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

## Definition (*IDef*)

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

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

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

(a) # bits needed to determine all variables

# THE SCORING FUNCTION

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

## Definition ( $IDef$ )

The  $\mathcal{m}$ -information deficit of  $\mu$ :

# bits to separately determine  
each target, knowing the source

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

# bits to determine all vars

## EXAMPLES

•  $\mathcal{m}_0 =$  X Y

$$IDef_{\mathcal{m}_0}(\mu) = -H_{\mu}(X, Y)$$

(optimal  $\mu$  maximizes entropy of  $X, Y$ )

People won't know this notation. You can't assume that they do.

# THE SCORING FUNCTION

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

## Definition ( $IDef$ )

The  $\mathbf{m}$ -information deficit of  $\mu$ :

# bits to separately determine each target, knowing the source

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

# bits to determine all vars

## EXAMPLES

•  $\mathbf{m}_0 =$  X Y

$$IDef_{\mathbf{m}_0}(\mu) = -H_{\mu}(X, Y)$$

(optimal  $\mu$  maximizes entropy of  $X, Y$ )

•  $\mathbf{m}_1 =$  X  $\longrightarrow$  Y

$$IDef_{\mathbf{m}_1}(\mu) = -H_{\mu}(X)$$

(optimal  $\mu$  maximizes entropy of  $X$ )

# THE SCORING FUNCTION

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

## Definition (*IDef*)

The  $\mathbf{m}$ -information deficit of  $\mu$ :

# bits to separately determine each target, knowing the source

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

# bits to determine all vars

## EXAMPLES

$$\bullet \mathbf{m}_0 = \boxed{X} \quad \boxed{Y}$$

$$\text{IDef}_{\mathbf{m}_0}(\mu) = -H_{\mu}(X, Y)$$

(optimal  $\mu$  maximizes entropy of  $X, Y$ )

$$\bullet \mathbf{m}_1 = \boxed{X} \longrightarrow \boxed{Y}$$

$$\text{IDef}_{\mathbf{m}_1}(\mu) = -H_{\mu}(X)$$

(optimal  $\mu$  maximizes entropy of  $X$ )

$$\bullet \mathbf{m}_2 = \boxed{X} \rightleftarrows \boxed{Y}$$

$$\text{IDef}_{\mathbf{m}_2}(\mu) = -H_{\mu}(X) + H_{\mu}(Y | X)$$

(optimal  $\mu$  maximizes entropy for  $X$ , and makes  $Y$  a function of  $X$ )

# THE SCORING FUNCTION

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

## Definition ( $IDef$ )

The  $\mathbf{m}$ -information deficit of  $\mu$ :

# bits to separately determine each target, knowing the source

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

# bits to determine all vars

## EXAMPLES

•  $\mathbf{m}_0 = \boxed{X} \quad \boxed{Y}$

$$IDef_{\mathbf{m}_0}(\mu) = -H_{\mu}(X, Y)$$

(optimal  $\mu$  maximizes entropy of  $X, Y$ )

•  $\mathbf{m}_1 = \boxed{X} \rightarrow \boxed{Y}$

$$IDef_{\mathbf{m}_1}(\mu) = -H_{\mu}(X)$$

(optimal  $\mu$  maximizes entropy of  $X$ )

•  $\mathbf{m}_2 = \boxed{X} \rightleftarrows \boxed{Y}$

$$IDef_{\mathbf{m}_2}(\mu) = -H_{\mu}(X) + H_{\mu}(Y | X)$$

(optimal  $\mu$  maximizes entropy for  $X$ , and makes  $Y$  a function of  $X$ )

•  $\mathbf{m}_3 = \boxed{X} \rightleftarrows \boxed{Y}$

$$IDef_{\mathbf{m}_3}(\mu) = -I_{\mu}(X; Y)$$

(opt.  $\mu$  makes  $X, Y$  share information)



# THE SCORING FUNCTION

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

## Definition ( $IDef$ )

The  $m$ -information deficit of  $\mu$ :

# bits to separately determine each target, knowing the source

$$IDef_m(\mu) = \sum_{X \xrightarrow{L} Y} \alpha_L H_\mu(Y | X) - H(\mu)$$

# bits to determine all vars

You've lost at least half of your audience at this point, using notation and concepts that they won't understand and you won't have time to explain.

I would cut all this.

## EXAMPLES

$$\bullet m_0 = \boxed{X} \quad \boxed{Y}$$

$$IDef_{m_0}(\mu) = -H_\mu(X, Y)$$

(optimal  $\mu$  maximizes entropy of  $X, Y$ )

$$\bullet m_1 = \boxed{X} \longrightarrow \boxed{Y}$$

$$IDef_{m_1}(\mu) = -H_\mu(X)$$

(optimal  $\mu$  maximizes entropy of  $X$ )

$$\bullet m_2 = \boxed{X} \rightleftarrows \boxed{Y}$$

$$IDef_{m_2}(\mu) = -H_\mu(X) + H_\mu(Y | X)$$

(optimal  $\mu$  maximizes entropy for  $X$ , and makes  $Y$  a function of  $X$ )

$$\bullet m_3 = \boxed{X} \rightleftarrows \boxed{Y}$$

$$IDef_{m_3}(\mu) = -I_\mu(X; Y)$$

(opt.  $\mu$  makes  $X, Y$  share information)

Information Diagrams

# THE SCORING FUNCTION

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

tradeoff parameter  $\gamma \geq 0$

## Definition (*Inc*)

The *incompatibility* of  $\mu$  with  $\mathbf{m}$ :

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

## Definition (*IDef*)

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

# bits to separately determine each target, knowing the source

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

# bits to determine all vars

# THE SCORING FUNCTION

- A BN strictly enforces the qualitative picture (large  $\gamma$ )

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

## Definition (*Inc*)

The *incompatibility* of  $\mu$  with  $\mathbf{m}$ :

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

## Definition (*IDef*)

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

# bits to separately determine each target, knowing the source

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

# THE SCORING FUNCTION

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

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

## 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 deficit* of  $\mu$ :

# bits to separately determine each target, knowing the source

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

# bits to determine all vars

# PROPERTIES OF THE OPTIMAL DISTRIBUTION

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

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

# PROPERTIES OF THE OPTIMAL DISTRIBUTION

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

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

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

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

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

## Proposition (If there are distributions consistent with $m$ , the best distribution is one of them.)

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

# PROPERTIES OF INCONSISTENCY

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

Nice properties for minimization:

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

# OUTLINE FOR SECTION 5

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

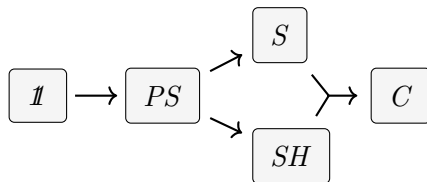
## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future Work



# CAPTURING BAYESIAN NETWORKS

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



# CAPTURING BAYESIAN NETWORKS

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

## Theorem (*BNs are PDGs*)

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

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

# CAPTURING BAYESIAN NETWORKS

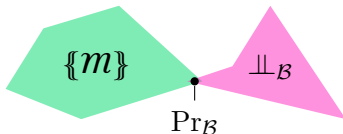
For a BN  $\mathcal{B}$  with  $N$  nodes and a vector  $\beta \in \mathbb{R}^N$ , let  $\mathbf{m}_{\mathcal{B},\beta}$  be the PDG corresponding to  $\mathcal{B}$ , with  $\alpha = \mathbf{1}$ , and the given vector  $\beta$  of confidences.

## Theorem (*BNs are PDGs*)

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

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

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*)

# 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?

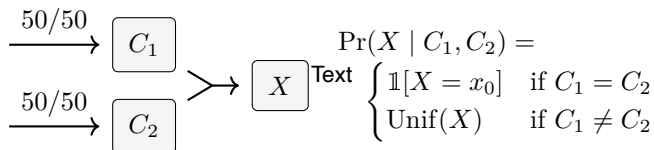
Common distributions  
tend to maximize  
entropy subject to  
natural constraints.

<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 + ???

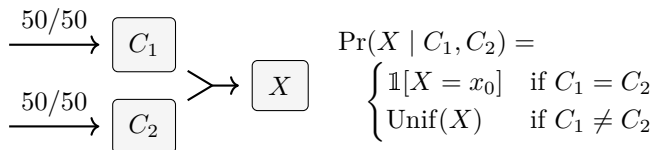


You're going into way too much technical detail here. Your audience is not just Ziv! I would cut all this.

# 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 + ???

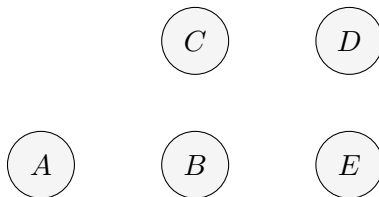


## Corollary

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

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

# FACTOR GRAPHS

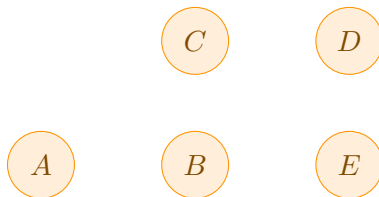


## Definition

A *factor graph*  $\Phi$  is



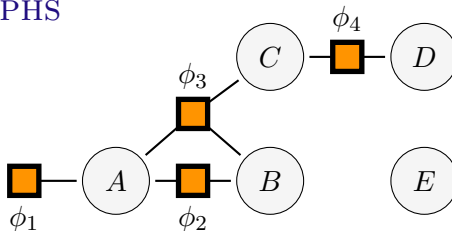
# FACTOR GRAPHS



## Definition

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

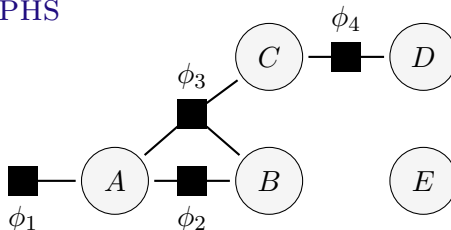
# FACTOR GRAPHS



## 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



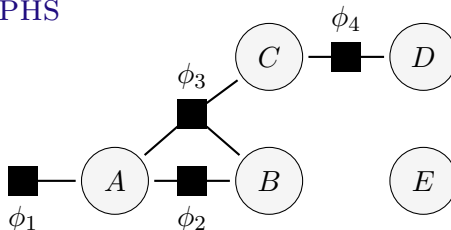
## Definition

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

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

Give some intuition for the  $\phi_J$ ? What do they represent?

# 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

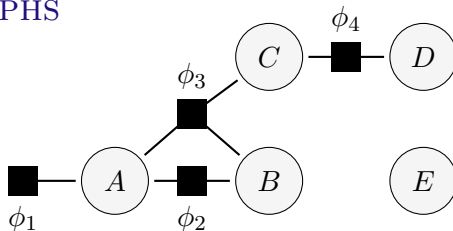
$$\Pr_{\Phi}(\vec{x}) := \frac{1}{Z_{\Phi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J), \quad \text{where } Z_{\Phi} \text{ is the normalization constant.}$$

$\Phi$  defines a “variational free energy”

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

Why do we need this? You’re overwhelming the poor listener.

# FACTOR GRAPHS



## Definition

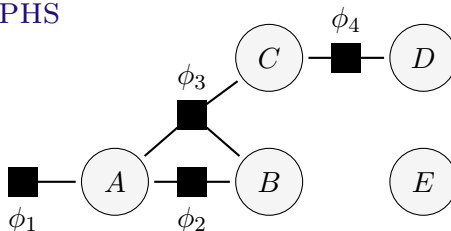
A *factor graph*  $\Phi$  is a set of variables  $\mathcal{X} = \{X_i\}$ , and *factors*  $\{\phi_J: \mathcal{V}(X_J) \rightarrow \mathbb{R}_{\geq 0}\}_{J \in \mathcal{J}}$ , with  $X_J \subseteq \mathcal{X}$ ;  $\Phi$  defines a distribution

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

$\Phi$  defines a “variational free energy”

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

# FACTOR GRAPHS



## Definition

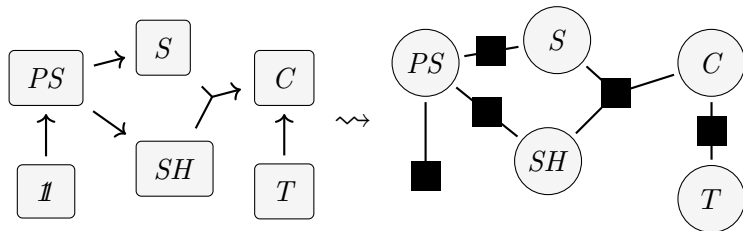
A **weighted factor graph**  $\Psi = (\Phi, \theta)$  is a set of variables  $\mathcal{X} = \{X_i\}$ , factors  $\{\phi_J: \mathcal{V}(X_J) \rightarrow \mathbb{R}_{\geq 0}\}_{J \in \mathcal{J}}$ , **and weights**  $(\theta_J)_{J \in \mathcal{J}}$  with  $X_J \subseteq \mathcal{X}$ ;  $\Psi$  defines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

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

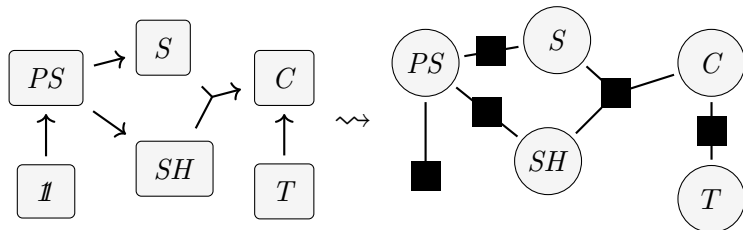
$\Psi$  defines a “variational free energy”

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

# PDGs AS FACTOR GRAPHS



# PDGs AS FACTOR GRAPHS

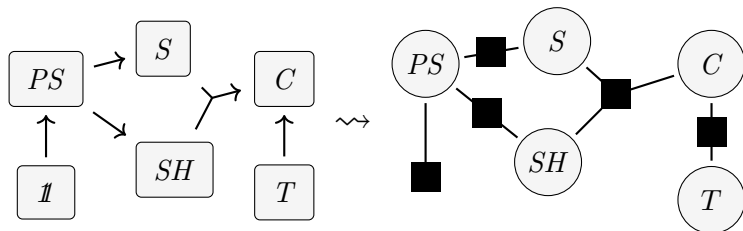


The cpds of a PDG are essentially factors. Are the semantics different?

What, intuitively, is a factor?



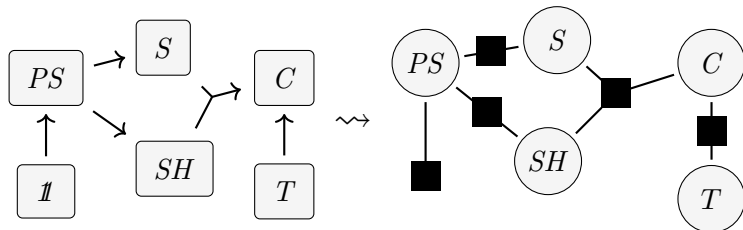
# PDGs AS FACTOR GRAPHS



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

Not for  $\gamma = 1$ .

# PDGs AS FACTOR GRAPHS



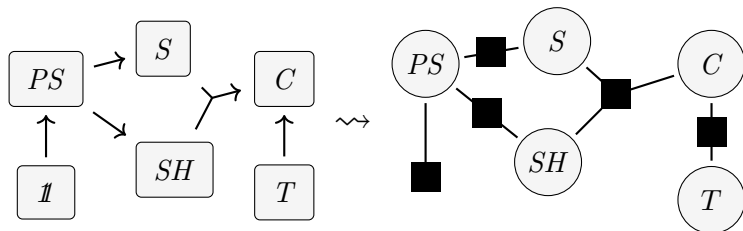
The cpds of a PDG are essentially factors. Are the semantics different?

Not for  $\gamma = 1$ .

## Theorem

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

# PDGs AS FACTOR GRAPHS



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

Not for  $\gamma = 1$ .

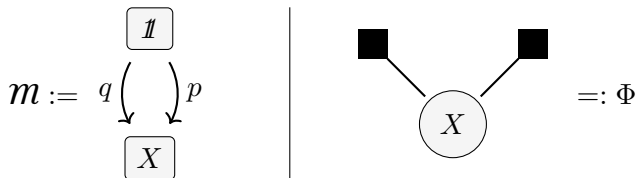
## Theorem

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

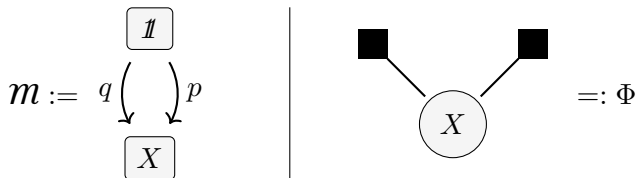
## Theorem

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

# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS

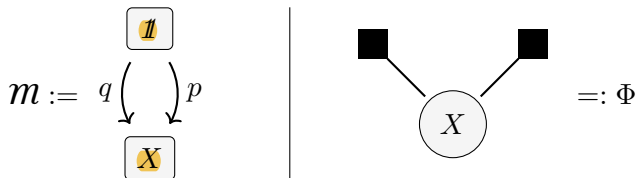


# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



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

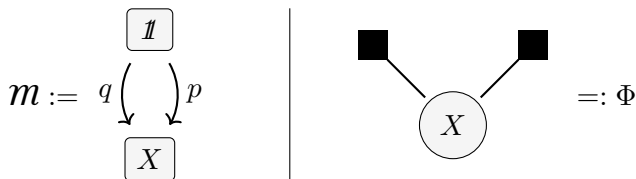
# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



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

Cut this. It will be of interest only to experts in factor graphs (which I believe is none of your audience).

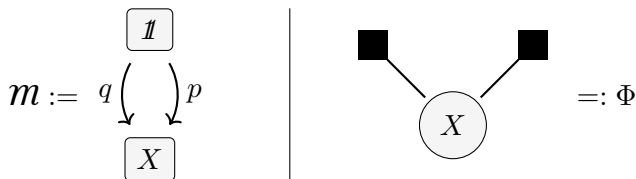
# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



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

So what do they mean? You have to say something about this.

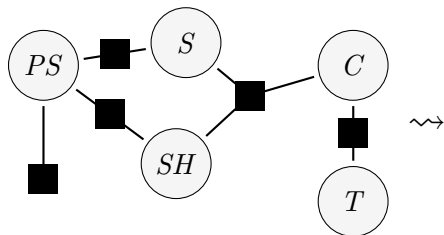
# AN IMPORTANT DIFFERENCE BETWEEN PDGs AND FACTOR GRAPHS



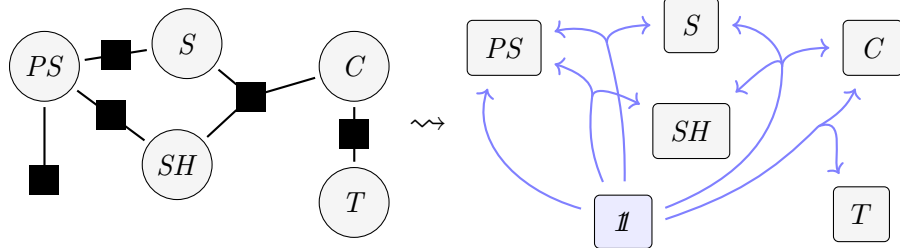
- If  $p = q$ , then  $\llbracket m \rrbracket^* = p = q \dots$
- $\dots$  but  $\Pr_{\Phi} \propto p^2$
- Individual factors have *no probabilistic meaning*,
- a factor graph can fail to normalize, in which case it has no global semantics either.



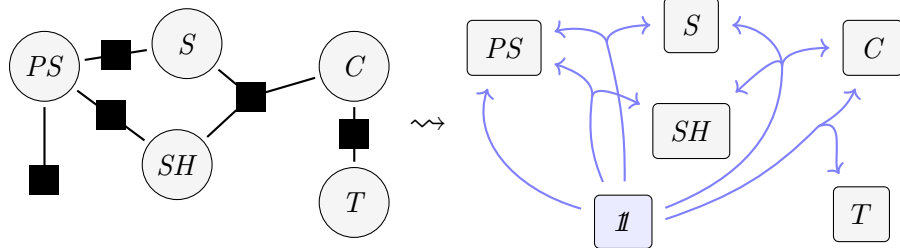
# FACTOR GRAPHS AS PDGs



# FACTOR GRAPHS AS PDGs



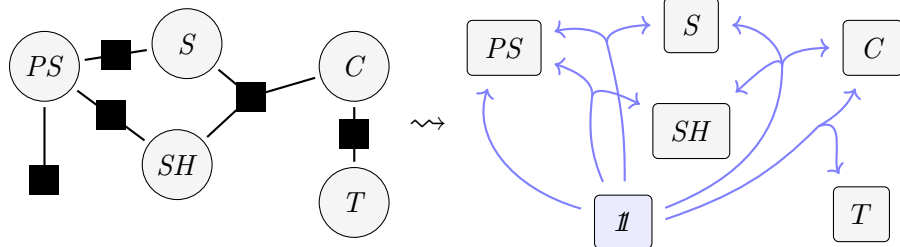
# FACTOR GRAPHS AS PDGs



## Theorem

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

# FACTOR GRAPHS AS PDGs



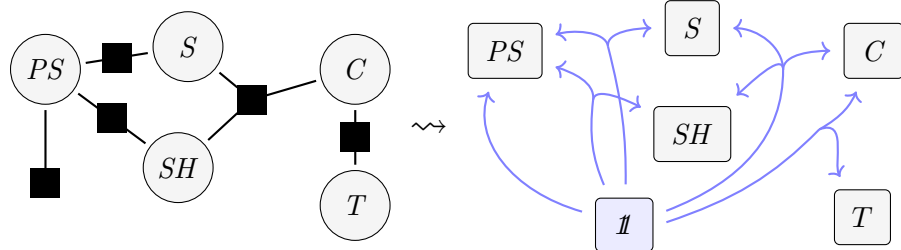
## Theorem

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

## Theorem

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

# FACTOR GRAPHS AS PDGs



## Theorem

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

## Theorem

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

Also:  $\log Z_{\Phi} = \langle \mathbf{n}_{\Phi} \rangle_1$ .

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

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

# OUTLINE FOR SECTION 6

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future Work

# INFERENCE VIA INCONSISTENCY REDUCTION

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



# INFERENCE VIA INCONSISTENCY REDUCTION

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

## Conditioning as inconsistency resolution.

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

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

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## Conditioning as inconsistency resolution.

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

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

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

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

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## Conditioning as inconsistency resolution.

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

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

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

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

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## Conditioning as inconsistency resolution.

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

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

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

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

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## Conditioning as inconsistency resolution.

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

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

## 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 of  $\mathcal{M} \sqcup p$  (which is strictly convex and smooth in  $p$ ) is  $\llbracket \mathcal{M} \rrbracket^*(Y \mid X)$ ,
  - ▶ so an inconsistency oracle yields fast inference by gradient descent.

**(Theorem):** Unfortunately,

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

# INFERENCE VIA INCONSISTENCY REDUCTION

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

## Conditioning as inconsistency resolution.

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

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

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

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

**(Theorem):** Unfortunately,

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

...just like for BNs and Factor Graphs.

# OUTLINE FOR SECTION 7

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future 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.

This is important, and it needs \*much\* more discussion.  
How do you use the model? Where might it come from?

# 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.

## Surprising Result

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

# INCONSISTENCY: THE UNIVERSAL LOSS

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

## Surprising Result

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

## Bonus

An intuitive visual language for reasoning about relationships between objective functions.

# SURPRISE AS INCONSISTENCY

This can't come out of the blue. You need to introduce surprise, and explain why people think it's a reasonable objective to minimize.

## Proposition

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

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

# SURPRISE AS INCONSISTENCY

## Proposition

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

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

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



# SURPRISE AS INCONSISTENCY

## Proposition

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

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

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

It's too late to say this here. “Surprise” has already been used in the proposition. You need to introduce surprise at the beginning, not the end!

## VARIATIONS: SURPRISE AS INCONSISTENCY

### Proposition (marginal information as inconsistency)

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

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

# 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

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

## Proposition (cross entropy as inconsistency)

What's f?

The inconsistency of the PDG containing  $f(Y | X)$  and a **high-confidence** empirical distribution  $\text{Pr}_{\underline{\mathbf{xy}}}$  of samples  $\underline{\mathbf{xy}} = \{(x_i, y_i)\}$  is equal to the cross entropy (plus  $H(Y | X)$ , a constant that depends only on the data  $\text{Pr}_{\underline{\mathbf{xy}}}$ ). That is, Where did “high confidence” come from? This needs much more discussion and intuition. You're far too focused on listing results, rather than explaining them

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

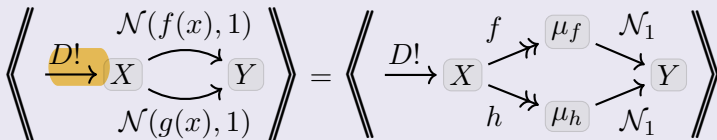
## Proposition (Accuracy as Inconsistency)

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

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

I have no clue what the notation  $D^{(\beta)}$  means. Showing me with notation like this is a \*bad\* idea.

## Proposition (Mean Square Error as Inconsistency)



Same comment as above  
for  $D^{\beta}$

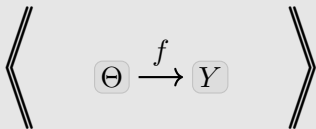
$$= \mathbb{E}_D \left( f(X) - h(X) \right)^2 =: \text{MSE}(f, h)$$

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

## Proposition (Regularizers as priors)

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

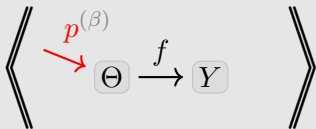
That is,



## Proposition (Regularizers as priors)

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

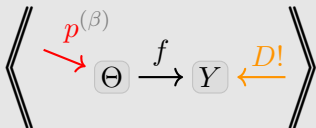
That is,



## Proposition (Regularizers as priors)

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

That is,

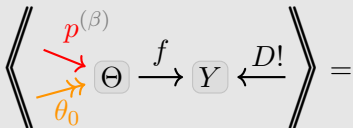




## Proposition (Regularizers as priors)

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

That is,



## Proposition (Regularizers as priors)

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

I'm feeling overwhelmed by notation.

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

## 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,

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

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

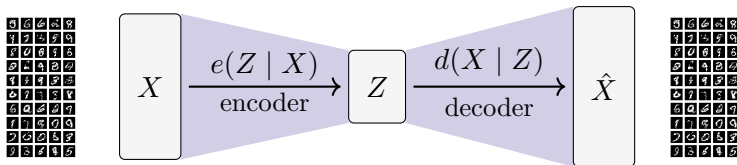
$$\underbrace{\mathbb{E}_D \left[ \log \frac{1}{f(Y | \theta_0)} \right]}_{\text{Cross entropy loss of } f_\theta \text{ w.r.t. } D} + \underbrace{\frac{\beta}{2} \theta_0^2}_{\ell_2 \text{ regularizer}} + \underbrace{+\beta \log k - H(D)}_{\text{constant in } f \text{ and } \theta_0}.$$

(data-fit cost of  $\theta_0$ )  
 (complexity cost of  $\theta_0$ )

At this point, I guarantee you've lost your audience, unless you're prepared to spend 10 minutes explaining this.

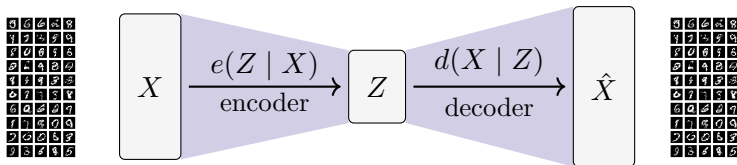
# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

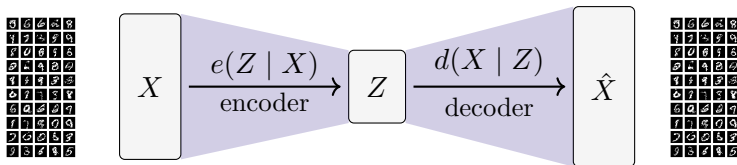


- Objective:

What's  $Z$ ? Where did it come from?

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

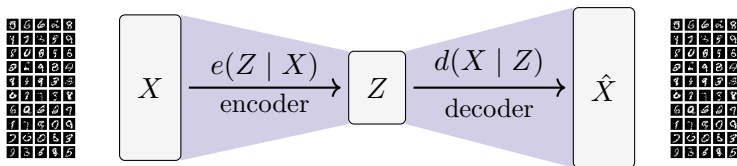


- Objective:

- ▶ 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):



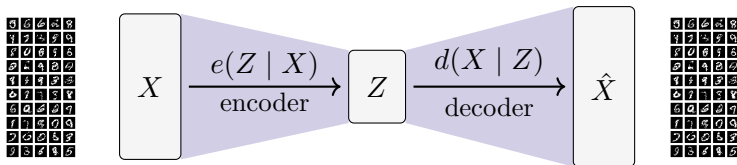
- Objective:

- ▶ 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 prior  $p(Z)$  that we want encodings of  $x$  to follow.

Why? Where did this come from? (I still don't know what  $Z$  is.)

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



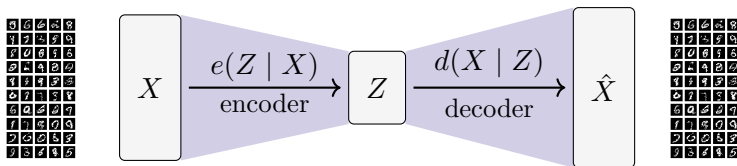
- Objective:

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



# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



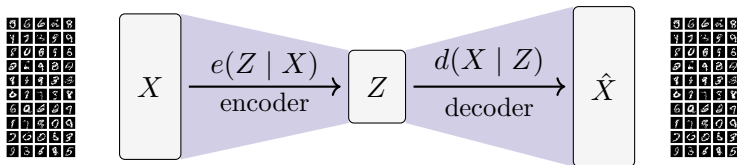
- Objective:

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

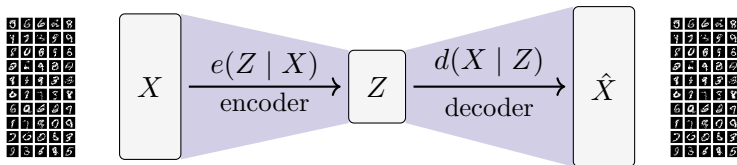


- Objective:

- For each  $x$ , want to minimize  $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$  “reconstruction error”
  - Also have a prior  $p(Z)$  that we want encodings of  $x$  to follow.
  - Add terms & negate to get
- $$\text{ELBO}_{p,e,d}(x) :=$$
- $$-\underbrace{D(e(Z|x) \parallel p(Z))}_{\text{divergence from prior}} - \text{Rec}(x)$$

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

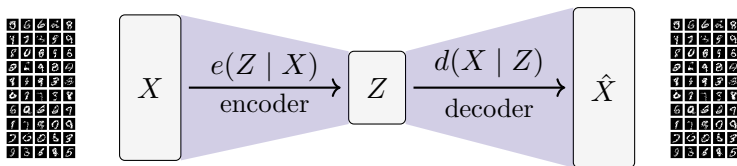
- ▶ For each  $x$ , want to minimize  $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$  “reconstruction error”
- ▶ Also have a prior  $p(Z)$  that we want encodings of  $x$  to follow.
- ▶ Add terms & negate to get

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

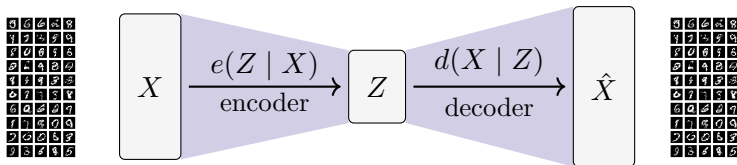
- ▶ For each  $x$ , want to minimize  $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$  (“reconstruction error”)
- ▶ Also have a prior  $p(Z)$  that we want encodings of  $x$  to follow.
- ▶ Add terms & negate to get

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



- Objective:

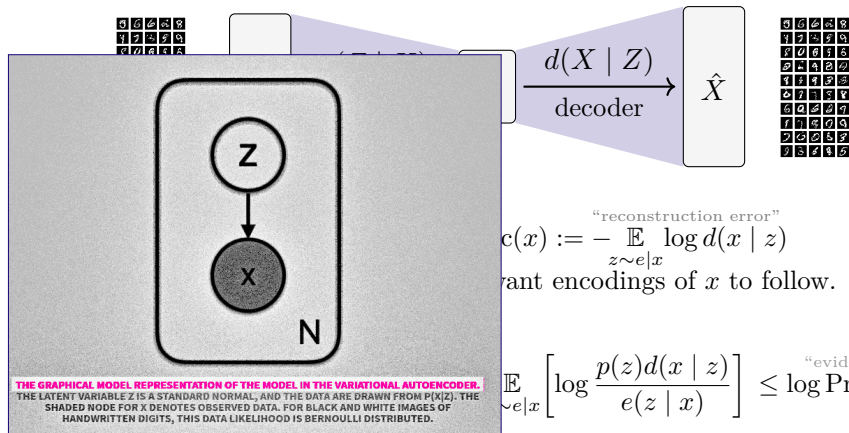
- ▶ For each  $x$ , want to minimize  $\text{Rec}(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$  “reconstruction error”
- ▶ Also have a prior  $p(Z)$  that we want encodings of  $x$  to follow.
- ▶ Add terms & negate to get

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):



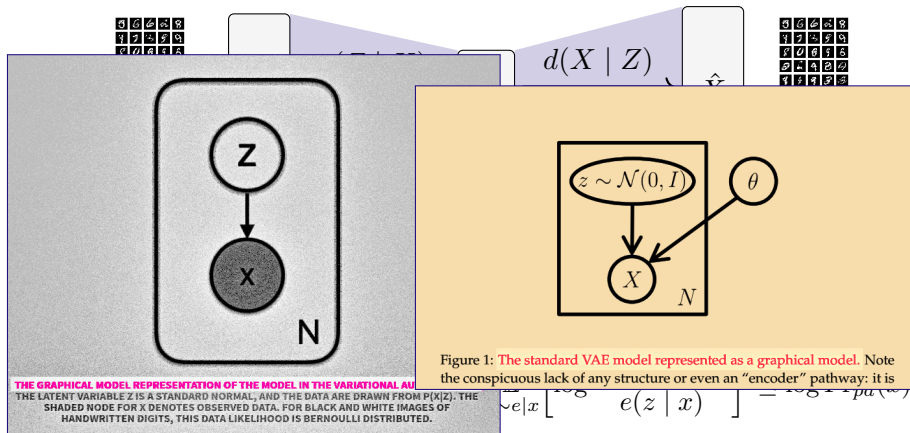
“reconstruction error”  
 $c(x) := - \mathbb{E}_{z \sim e|x} \log d(x | z)$   
 want encodings of  $x$  to follow.

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

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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

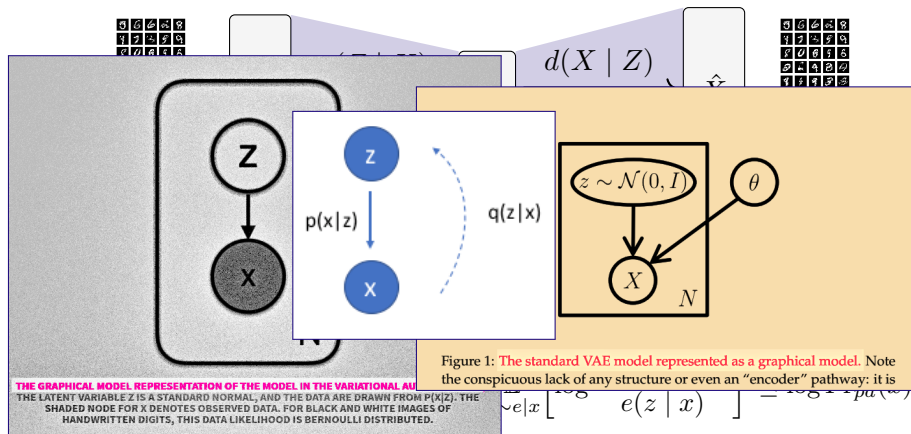
- Structure consists of two neural networks (cpds):



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

# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

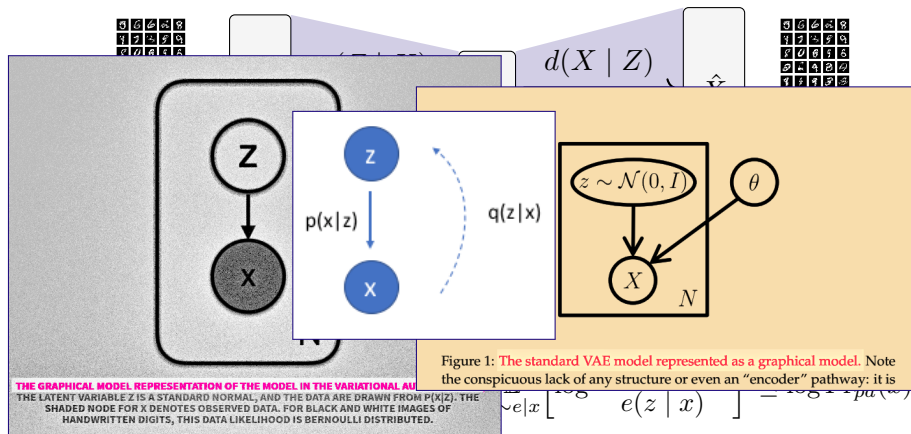


Urge to use graphical models (even if can't quite capture *entire* VaE)



# VARIATIONAL AUTO-ENCODERS, TAKE 1

- Structure consists of two neural networks (cpds):

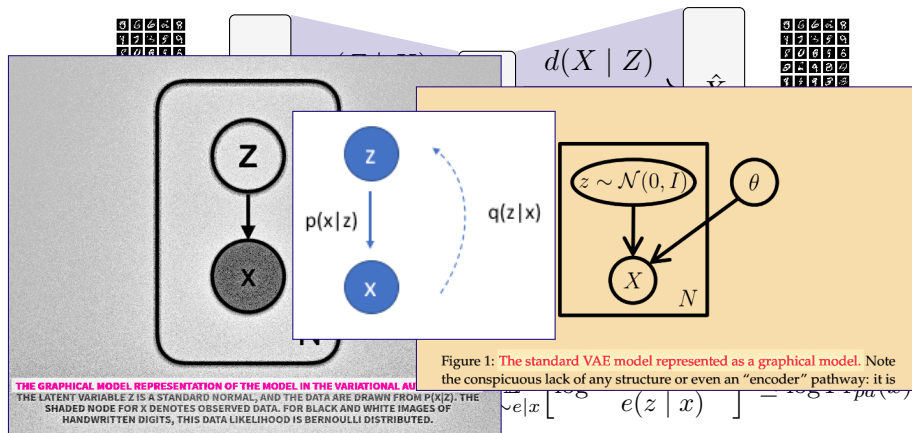


Urge to use graphical models (even if can't quite capture *entire* VaE)

- $e(Z | X)$  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 2

- Structure:

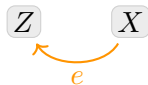
$Z$

$X$

# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$  : encoder

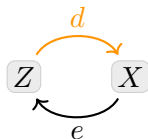


# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder



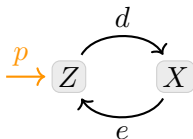
# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior



# VARIATIONAL AUTO-ENCODERS, TAKE 2

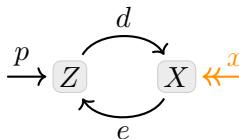
- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$



# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

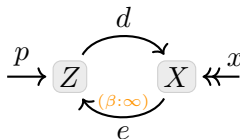
$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$

▶ and trust encoding





# VARIATIONAL AUTO-ENCODERS, TAKE 2

- Structure:

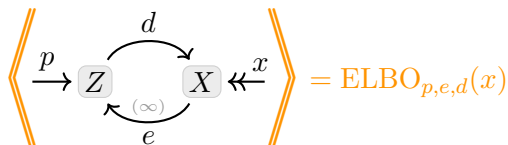
$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$ 
  - ▶ and trust encoding

Objective function is free:



# VARIATIONAL AUTO-ENCODERS, TAKE 2

Objective function is free:

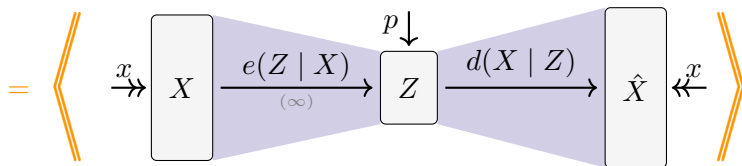
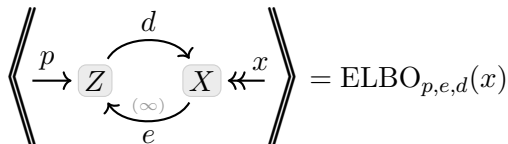
- Structure:

$e(Z | X)$  : encoder

$d(X | Z)$  : decoder

$p(Z)$  : prior

- observe a sample  $x$ 
  - ▶ and trust encoding



# A VERY USEFUL FACT

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

## Lemma (monotonicity of inconsistency)

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

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

# VISUAL PROOF: THE VARIATIONAL BOUND

# VISUAL PROOF: THE VARIATIONAL BOUND

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

# VISUAL PROOF: THE VARIATIONAL BOUND

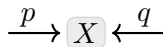
$$-\log \Pr_{p,d}(X=x) = \left\langle \begin{array}{c} p \rightarrow Z \xrightarrow{d} X \leftarrow x \end{array} \right\rangle \left\langle \begin{array}{c} p \rightarrow Z \xrightleftharpoons[e!]{d} X \leftarrow x \end{array} \right\rangle = -\text{ELBO}_{p,e,d}(x).$$

# VISUAL PROOF: THE VARIATIONAL BOUND

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

# DIVERGENCES AND INCONSISTENCY

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

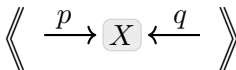




# DIVERGENCES AND INCONSISTENCY

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

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



# DIVERGENCES AND INCONSISTENCY

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

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

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

# DIVERGENCES AND INCONSISTENCY

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

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

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

# DIVERGENCES AND INCONSISTENCY

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

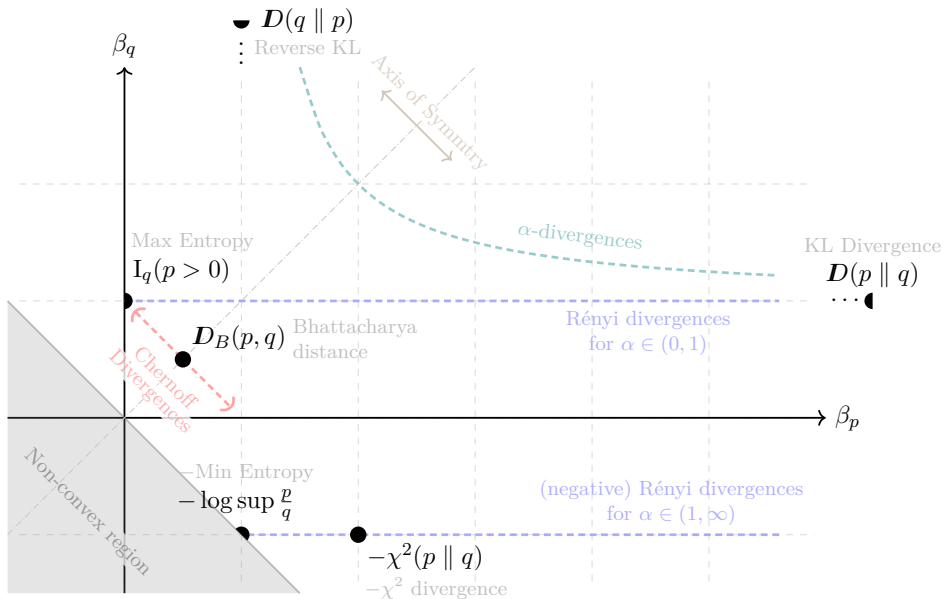
Your inconsistency: a divergence between  $p$  and  $q$ ?

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

## Lemma

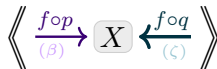
$$D_{(r,s)}^{\text{PDG}}(p, q) = -(r+s) \log \sum_x \left( p(x)^r q(x)^s \right)^{\frac{1}{r+s}}.$$

# DIVERGENCES AS INCONSISTENCIES



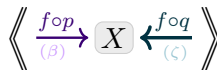
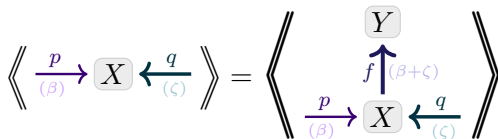
# VISUAL PROOF: DATA PROCESSING INEQUALITY

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



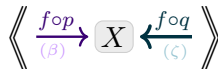
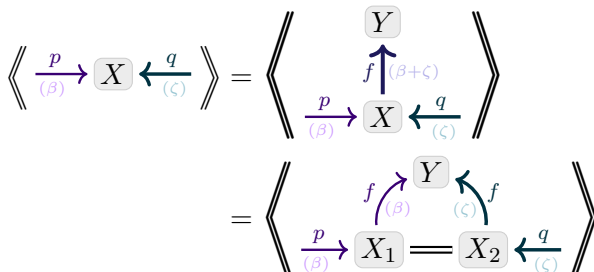
# VISUAL PROOF: DATA PROCESSING INEQUALITY

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



# VISUAL PROOF: DATA PROCESSING INEQUALITY

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





# VISUAL PROOF: DATA PROCESSING INEQUALITY

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

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

# OUTLINE FOR SECTION 8

## 1 INTRODUCTION

## 2 MODELING EXAMPLES

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

## 3 SYNTAX

- Formal Definitions of PDGs

## 4 SEMANTICS

## 5 CAPTURING OTHER GRAPHICAL MODELS

- Bayesian Networks
- Factor Graphs

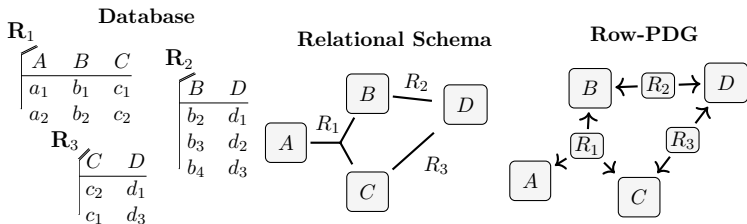
## 6 INFERENCE

## 7 INCONSISTENCY AS LOSS

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

## 8 OTHER ASPECTS OF PDGs

- Databases
- Open Problems + Future Work



## 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.

# OPEN PROBLEMS AND FUTURE WORK

- Technical tool: PDGs with incomplete cpds
- Encoding preferences, and understanding preference changes
- Trace Semantics: and the probabilistic automaton generated by a PDG
- \* \* Do PDGs capture Dependency Networks? \* \*
- Multi-agent systems

# TECHNICAL RECAP

⟨ update with second half ⟩

PDGs...

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

# TECHNICAL RECAP

⟨ update with second half ⟩

PDGs...

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

# TECHNICAL RECAP

⟨ update with second half ⟩

PDGs...

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

# TECHNICAL RECAP

⟨ update with second half ⟩

PDGs...

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



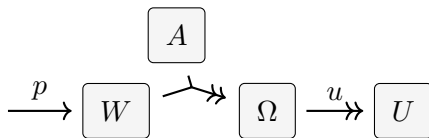
# TECHNICAL RECAP

⟨ update with second half ⟩

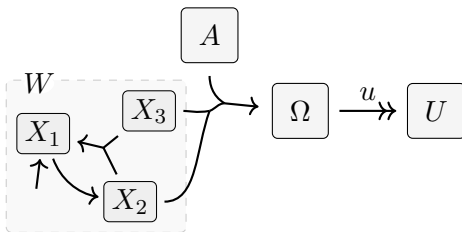
PDGs...

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

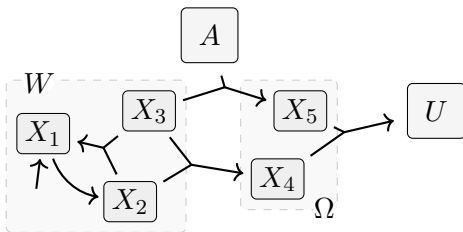
*But there is much more to be done!*



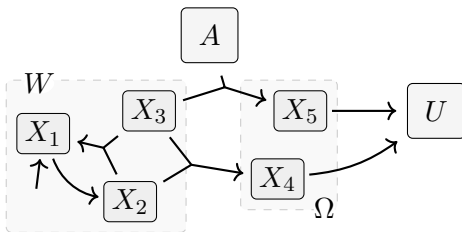
- decompose states and beliefs (like PGMs)
- driven by pursuit of coherent identity, not “favorite number go up”



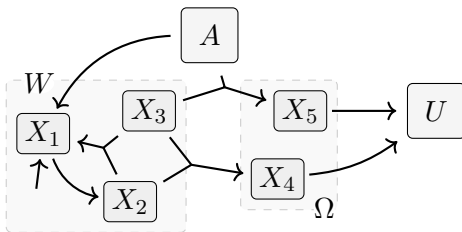
- decompose states and beliefs (like PGMs)
- driven by pursuit of coherent identity, not “favorite number go up”



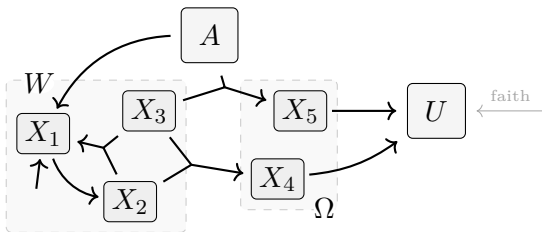
- decompose states and beliefs (like PGMs)
- driven by pursuit of coherent identity, not “favorite number go up”



- decompose states and beliefs (like PGMs)
- driven by pursuit of coherent identity, not “favorite number go up”



- decompose states and beliefs (like PGMs)
- driven by pursuit of coherent identity, not “favorite number go up”



- decompose states and beliefs (like PGMs)
- driven by pursuit of coherent identity, not “favorite number go up”

# OUTLINE FOR SECTION 9

## 9 HYPER-GRAPHS

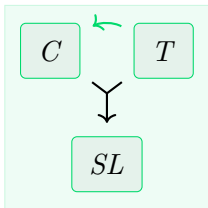
## 10 THE INFORMATION DEFICIENCY

## 11 CATEGORY THEORY

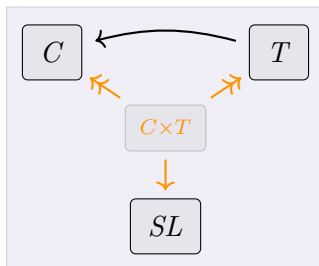
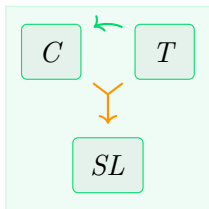
- PDGs as diagrams of the Markov Category



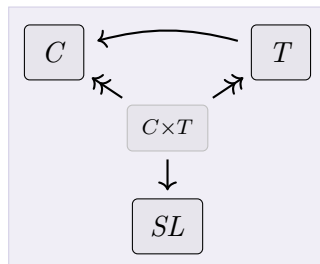
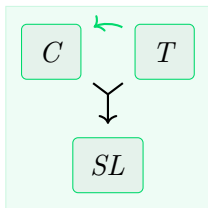
# HYPER-GRAPHS? OR MERELY GRAPHS?



# HYPER-GRAPHS? OR MERELY GRAPHS?

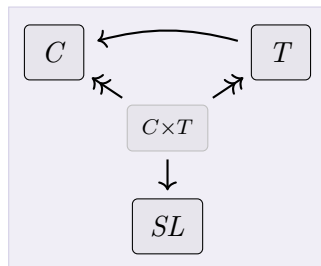
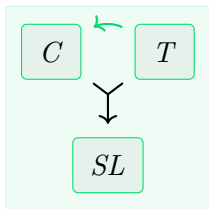


# HYPER-GRAPHS? OR MERELY GRAPHS?



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

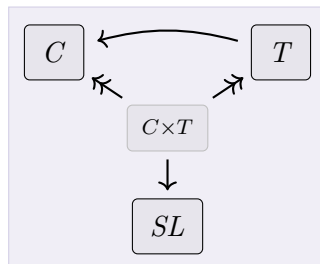
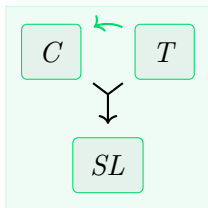
# HYPER-GRAPHS? OR MERELY GRAPHS?



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

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

# HYPER-GRAPHS? OR MERELY GRAPHS?



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

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

(working directly with hypergraphs is also possible)

# OUTLINE FOR SECTION 10

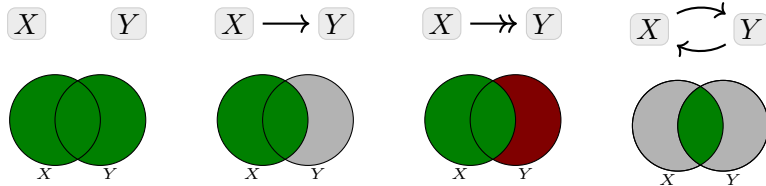
9 HYPER-GRAPHS

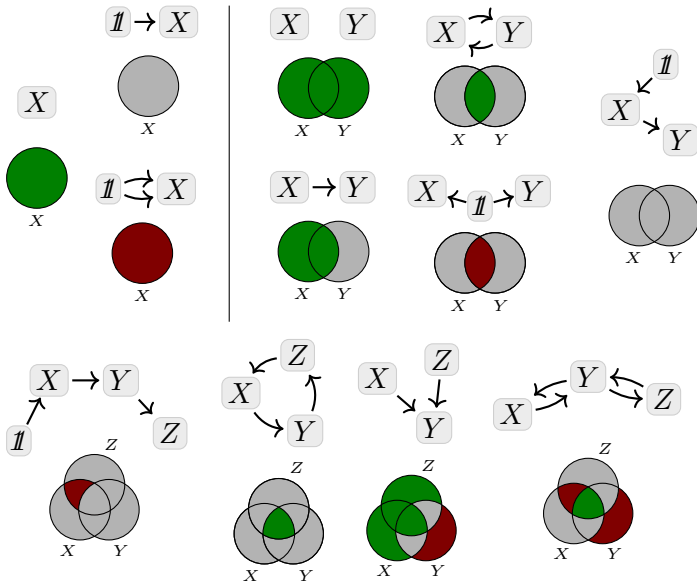
10 THE INFORMATION  
DEFICIENCY

11 CATEGORY THEORY

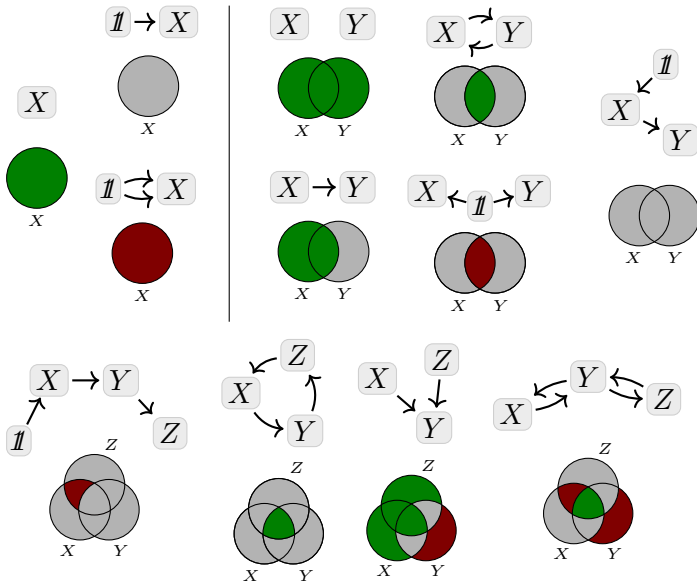
- PDGs as diagrams of the Markov Category

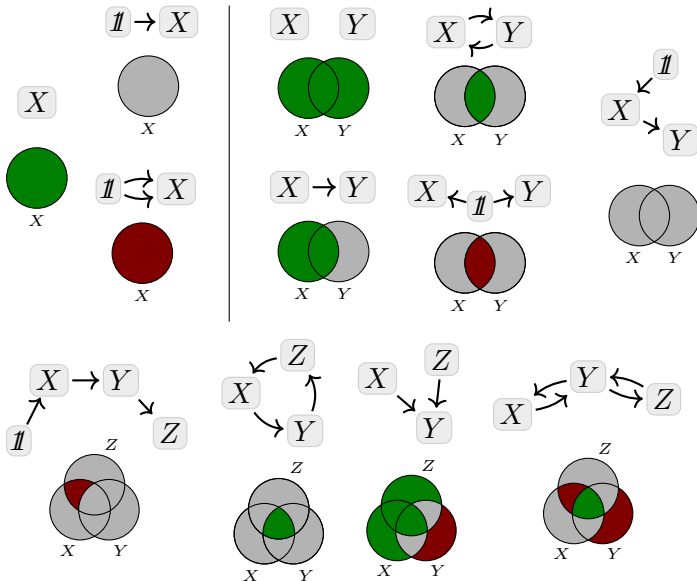
# ILLUSTRATIONS OF $IDef$

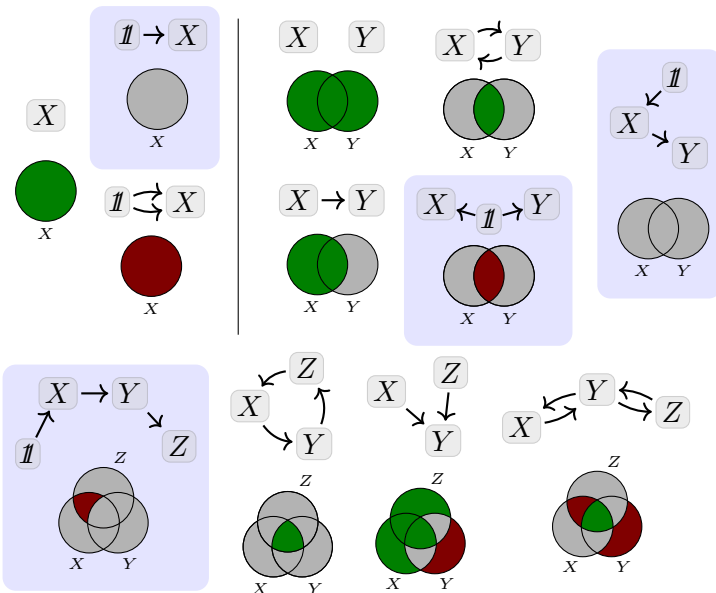












# OUTLINE FOR SECTION 11

9 HYPER-GRAPHS

10 THE INFORMATION  
DEFICIENCY

11 CATEGORY THEORY

- PDGs as diagrams of the Markov Category

## Definition (PDG)

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

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

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

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

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

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

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

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

## Definition (PDG)

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

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

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

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

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

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

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

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

## Definition (PDG)

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

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

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

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

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

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

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

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

## Definition (PDG)

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

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

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

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

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

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

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

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



## Definition (PDG)

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

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

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

## Definition (PDG)

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

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

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

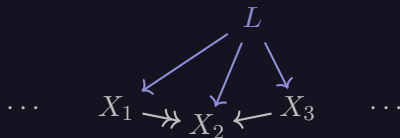
## Equivalent Categorical Definition

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

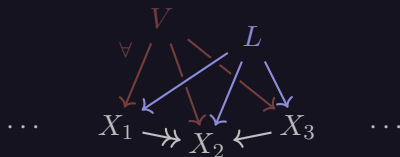
What do you do with diagrams? Take **limits** / **colimits**.

$$\cdots \quad X_1 \rightrightarrows X_2 \longleftarrow X_3 \quad \cdots$$

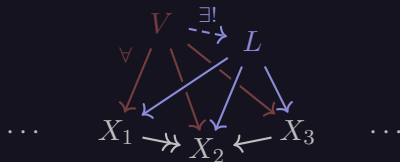
What do you do with diagrams? Take limits / colimits.



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



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



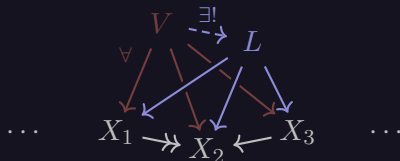
What do you do with diagrams? Take limits / colimits.



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

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

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



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



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

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

In general:

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

For a BN  $\mathcal{B}$ :

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

