

PROBABILISTIC DEPENDENCY GRAPHS

Oliver E. Richardson Joseph Y. Halpern

Cornell University
Department of Computer Science

AAAI, February 2021

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

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

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

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

Why do we need another one?

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

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

Why do we need another one?

- To resolve inconsistency, we must first model it.

YET ANOTHER PROBABILISTIC GRAPHICAL MODEL

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

Why do we need another one?

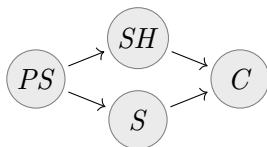
- To resolve inconsistency, we must first model it.
- In doing so, we get much more ...

TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

Qualitative BN, \mathcal{G}

an independence relation on variables

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



TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

Qualitative BN, \mathcal{G}

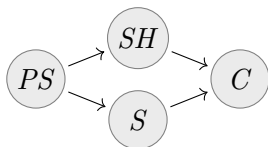
an independence relation on variables

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

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

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

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



TWO ASPECTS OF BAYESIAN NETWORKS (BNs)

Qualitative BN, \mathcal{G}

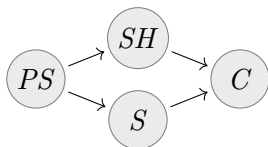
an independence relation on variables

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

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

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

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



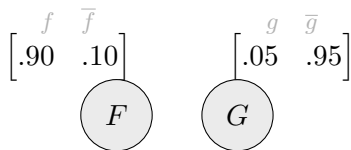
MODELING EXAMPLE: FLOOMPS AND GUNS

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

MODELING EXAMPLE: FLOOMPS AND GUNS

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

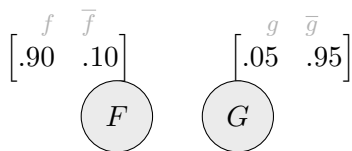
BN



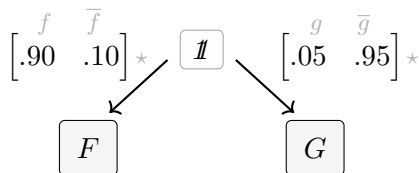
MODELING EXAMPLE: FLOOMPS AND GUNS

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

BN

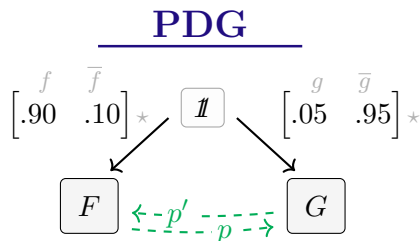
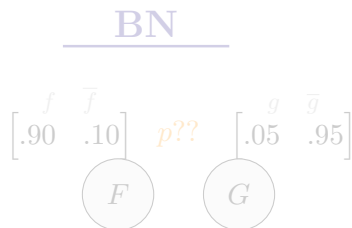


PDG



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

MODELING EXAMPLE: FLOOMPS AND GUNS

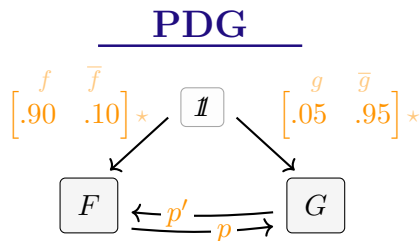
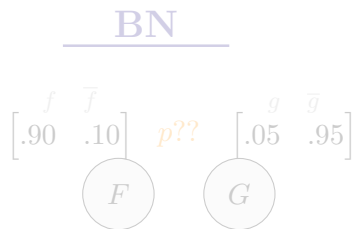


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.

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

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

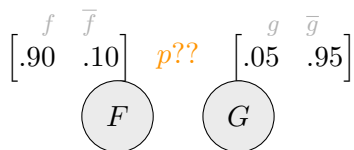
MODELING EXAMPLE: FLOOMPS AND GUNS



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

MODELING EXAMPLE: FLOOMPS AND GUNS

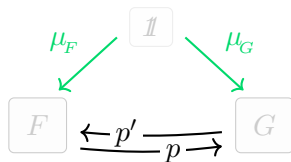
BN



Incorporated:

μ_F μ_G p p'

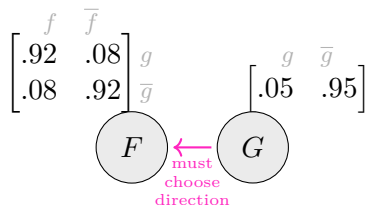
PDG



- 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.

MODELING EXAMPLE: FLOOMPS AND GUNS

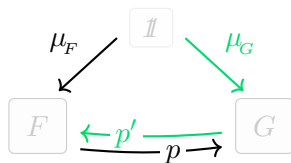
BN



Incorporated:

μ_F μ_G p p'

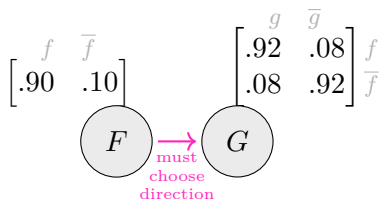
PDG



- 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 irreversibly lose information.

MODELING EXAMPLE: FLOOMPS AND GUNS

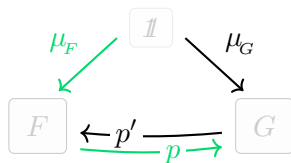
BN



Incorporated:

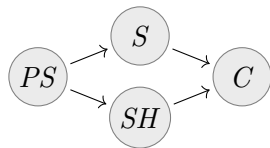
μ_F μ_G p p'

PDG

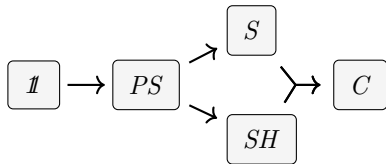
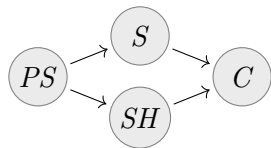


- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent
 - ▶ ... but BNs must resolve inconsistency first, which may break symmetry and irreversibly lose information.

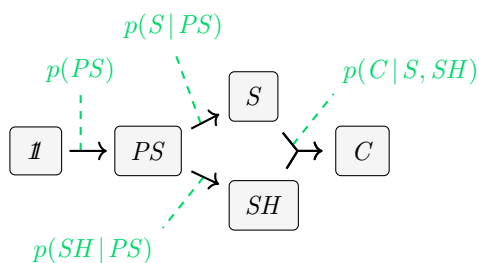
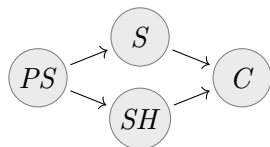
BAYESIAN NETWORKS AS PDGs



BAYESIAN NETWORKS AS PDGs



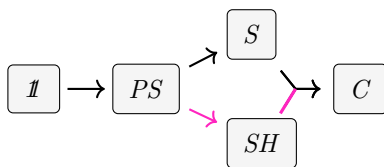
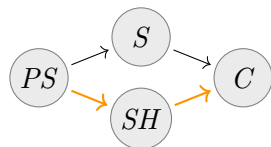
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

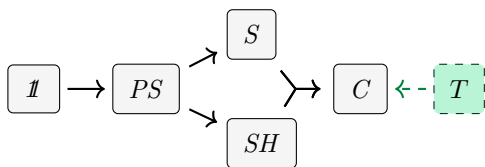
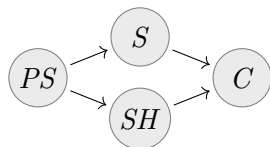
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

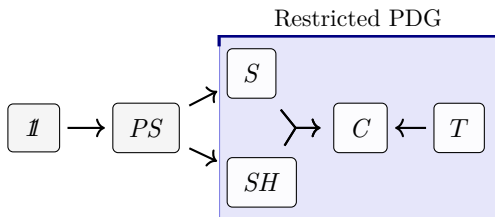
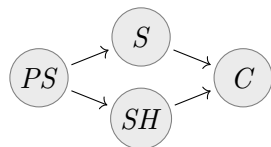
BAYESIAN NETWORKS AS PDGs



In contrast with BNs:

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

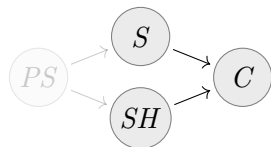
BAYESIAN NETWORKS AS PDGs



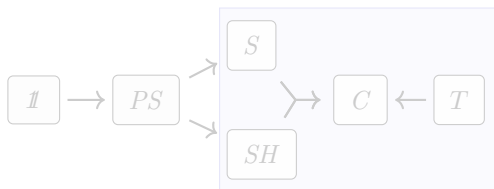
In contrast with BNs:

- 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



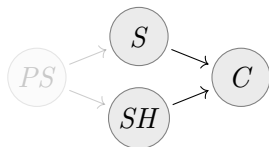
Must now give distributions on SH and S , or distinguish them as “observed” (a *conditional* BN).



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).

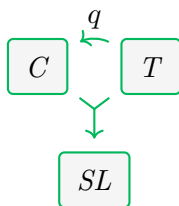
In a qualitative BN: *removing data results in new knowledge: $A \perp\!\!\!\perp C$.*



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!

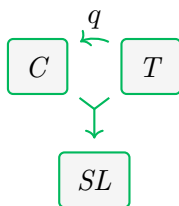
COMBINING PDGs



Grok wants to be supreme leader (SL).

- She notices that those who use tanning beds have more power,

COMBINING PDGs

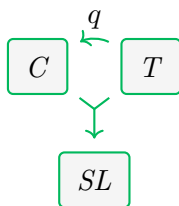


Grok wants to be supreme leader (SL).

- She notices that those who use tanning beds have more power,

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

COMBINING PDGs



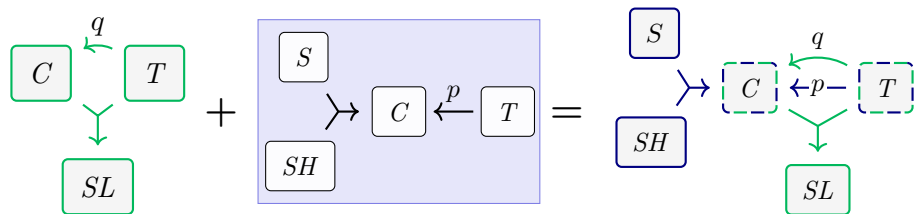
Grok wants to be supreme leader (SL).

- She notices that those who use tanning beds have more power,

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

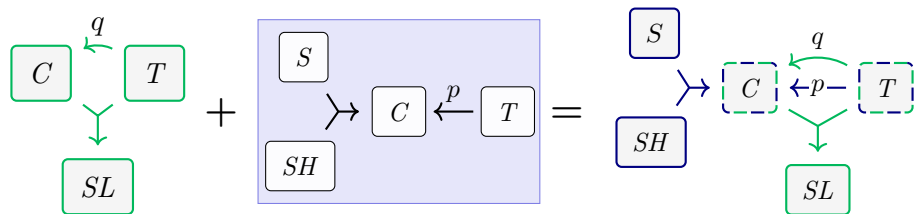
- Grok worries getting cancer from a tanning bed will make SL impossible.

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 (e.g., p, q), which directly conflict.

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 $\mathbf{m} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

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

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

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

Definition (Probabilistic Dependency Graph)

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

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

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

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

(hyper-edges)

Definition (Probabilistic Dependency Graph)

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

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

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

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

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

(hyper-edges)

Definition (Probabilistic Dependency Graph)

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

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

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

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

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

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

(hyper-edges)

Definition (Probabilistic Dependency Graph)

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

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

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

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

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

(hyper-edges)

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

$\beta_L \in (0, \infty)$ 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 $\{X \xrightarrow{L} Y\}$,

(hyper-edges)

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

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

$\alpha_L \in [0, \infty)$ a confidence in the functional dependence $X \rightarrow Y$

$\beta_L \in (0, \infty)$ a confidence in the reliability of \mathbf{p}_L .

PDG SEMANTICS

- $\{m\}$ The set of joint distributions consistent with m ;
- $[[m]]_\gamma$ A function, scoring distributions by compatibility with m ;
- $[[m]]^*$ The “best” joint distribution.

PDG SEMANTICS

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

$$\left\{ \mu \in \Delta[\mathcal{V}(\mathcal{m})] : \text{for all } X \xrightarrow{L} Y \in \mathcal{E}. \mu(Y|X) = \mathbf{p}_L(Y|X) \right\}$$

$[[\mathcal{m}]_\gamma]$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;

$[[\mathcal{m}]^*]$ The “best” joint distribution.

PDG SEMANTICS

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

$$\{\mu \in \Delta[\mathcal{V}(\mathcal{m})] : \text{for all } X \xrightarrow{L} Y \in \mathcal{E}. \mu(Y|X) = \mathbf{p}_L(Y|X)\}$$

$[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;

tradeoff parameter $\gamma \geq 0$

$$[[\mathcal{m}]]_\gamma(\mu) := \underbrace{Inc_{\mathcal{m}}(\mu)}_{\substack{\text{(quantitative)} \\ \text{term}}} + \gamma \underbrace{IDef_{\mathcal{m}}(\mu)}_{\substack{\text{(qualitative)} \\ \text{term}}}$$

$[[\mathcal{m}]]^*$ The “best” joint distribution.

PDG SEMANTICS

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

$$\left\{ \mu \in \Delta[\mathcal{V}(\mathcal{m})] : \text{for all } X \xrightarrow{L} Y \in \mathcal{E}. \mu(Y|X) = \mathbf{p}_L(Y|X) \right\}$$

$[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;

$$[[\mathcal{m}]]_\gamma(\mu) := \underbrace{Inc_{\mathcal{m}}(\mu)}_{\substack{\text{(quantitative)} \\ \text{term}}} + \gamma \underbrace{IDef_{\mathcal{m}}(\mu)}_{\substack{\text{(qualitative)} \\ \text{term}}}$$

$[[\mathcal{m}]]_\gamma^*$ The “best” joint distribution.

$$[[\mathcal{m}]]_\gamma^* := \arg \min_{\mu} [[\mathcal{m}]]_\gamma(\mu)$$

THE SCORING FUNCTION

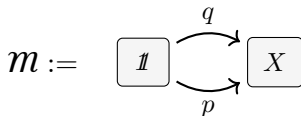
$$\llbracket \mathbf{m} \rrbracket_{\gamma}(\mu) := \text{Inc}_{\mathbf{m}}(\mu) + \gamma \text{IDef}_{\mathbf{m}}(\mu)$$

THE SCORING FUNCTION

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

Intuition: beyond stating whether or not μ is consistent with \mathcal{M} , we score μ 's compatibility with \mathcal{M} .

MOTIVATING EXAMPLES.

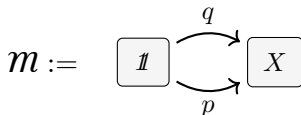


THE SCORING FUNCTION

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

Intuition: beyond stating whether or not μ is consistent with \mathcal{M} , we score μ 's compatibility with \mathcal{M} .

MOTIVATING EXAMPLES.



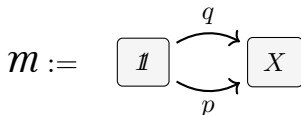
- If $p = \begin{bmatrix} x_1 & x_2 \\ .4 & .6 \end{bmatrix}^* = q$, then \mathcal{M} is consistent, and compatible with the joint distribution $\mu(X) = p$.

THE SCORING FUNCTION

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

Intuition: beyond stating whether or not μ is consistent with \mathcal{M} , we score μ 's compatibility with \mathcal{M} .

MOTIVATING EXAMPLES.



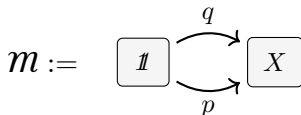
- If $p = \begin{bmatrix} x_1 & x_2 \\ .4 & .6 \end{bmatrix}^* = q$, then \mathcal{M} is consistent, and compatible with the joint distribution $\mu(X) = p$.
- If $p = \begin{bmatrix} x_1 & x_2 \\ .4 & .6 \end{bmatrix}^*$ and $q = \begin{bmatrix} x_1 & x_2 \\ .5 & .5 \end{bmatrix}^*$, then \mathcal{M} is not consistent, but $\mu = \begin{bmatrix} .45 & .55 \end{bmatrix}$ matches better than $\mu = \begin{bmatrix} .9 & .1 \end{bmatrix}$.

THE SCORING FUNCTION

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

Intuition: beyond stating whether or not μ is consistent with \mathcal{M} , we score μ 's compatibility with \mathcal{M} .

MOTIVATING EXAMPLES.



- If $p = \begin{bmatrix} x_1 & x_2 \\ .4 & .6 \end{bmatrix}^* = q$, then \mathcal{M} is consistent, and compatible with the joint distribution $\mu(X) = p$.
- If $p = \begin{bmatrix} x_1 & x_2 \\ .4 & .6 \end{bmatrix}^*$ and $q = \begin{bmatrix} x_1 & x_2 \\ .5 & .5 \end{bmatrix}^*$, then \mathcal{M} is not consistent, but $\mu = \begin{bmatrix} .45 & .55 \end{bmatrix}$ matches better than $\mu = \begin{bmatrix} .9 & .1 \end{bmatrix}$.
- If $p = \begin{bmatrix} .4 & .6 \end{bmatrix}$ and $q = \begin{bmatrix} 0 & 1 \end{bmatrix}$, then \mathcal{M} is much more inconsistent than before, even though $\llbracket \mathcal{M} \rrbracket = \emptyset$ in both cases.

THE SCORING FUNCTION

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

Definition (*Inc*)

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

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

THE SCORING FUNCTION

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

Definition (*Inc*)

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

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

THE SCORING FUNCTION

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

Definition (*Inc*)

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

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

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

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 μ with \mathbf{m} is given by

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

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

THE SCORING FUNCTION

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

Definition (*Inc*)

The *incompatibility* of a joint distribution μ with \mathcal{M} is given by

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

The *inconsistency* of \mathcal{M} is the smallest possible incompatibility,

$$\text{Inc}(\mathcal{M}) := \inf_{\mu \in \Delta \mathcal{V}(\mathcal{M})} \text{Inc}_{\mathcal{M}}(\mu).$$

THE SCORING FUNCTION

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

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

THE SCORING FUNCTION

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

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

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

THE SCORING FUNCTION

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

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

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

$H(\mu)$

THE SCORING FUNCTION

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

Definition (*IDef*)

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

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

THE SCORING FUNCTION

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

Definition (*IDef*)

The *information deficit* of a distribution μ with respect to \mathcal{M} is

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

(a) # bits needed to determine all variables

THE SCORING FUNCTION

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

Definition (*IDef*)

The *information deficit* of a distribution μ with respect to \mathcal{M} is

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

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

(a) # bits needed to determine all variables

THE SCORING FUNCTION

$$[[\mathcal{m}]]_{\gamma}(\mu) := Incm(\mu) + \gamma IDef_{\mathcal{m}}(\mu)$$

EXAMPLES

• $m_0 =$ X Y

$IDef_{m_0}(\mu) = -H_{\mu}(X, Y)$
 (optimal μ maximizes entropy of X, Y)

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

THE SCORING FUNCTION

$$[[\mathcal{m}]]_{\gamma}(\mu) := Incm(\mu) + \gamma IDef_{\mathcal{m}}(\mu)$$

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

EXAMPLES

• $\mathcal{m}_0 =$ X Y

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

(optimal μ maximizes entropy of X, Y)

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

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

(optimal μ maximizes entropy of X)

THE SCORING FUNCTION

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

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

EXAMPLES

- $\mathcal{M}_0 = \boxed{X} \quad \boxed{Y}$

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

(optimal μ maximizes entropy of X, Y)

- $\mathcal{M}_1 = \boxed{X} \longrightarrow \boxed{Y}$

$$\text{IDef}_{\mathcal{M}_1}(\mu) = -H_{\mu}(X)$$

(optimal μ maximizes entropy of X)

- $\mathcal{M}_2 = \boxed{X} \rightleftarrows \boxed{Y}$

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

(optimal μ maximizes entropy for X , and makes Y a function of X)

THE SCORING FUNCTION

$$[[\mathcal{M}]]_{\gamma}(\mu) := Incm(\mu) + \gamma IDef_{\mathcal{M}}(\mu)$$

Definition ($IDef$)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

EXAMPLES

- $\mathcal{M}_0 = \boxed{X} \quad \boxed{Y}$
 $IDef_{\mathcal{M}_0}(\mu) = -H_{\mu}(X, Y)$
 (optimal μ maximizes entropy of X, Y)

- $\mathcal{M}_1 = \boxed{X} \rightarrow \boxed{Y}$
 $IDef_{\mathcal{M}_1}(\mu) = -H_{\mu}(X)$
 (optimal μ maximizes entropy of X)

- $\mathcal{M}_2 = \boxed{X} \rightleftarrows \boxed{Y}$
 $IDef_{\mathcal{M}_2}(\mu) = -H_{\mu}(X) + H_{\mu}(Y | X)$
 (optimal μ maximizes entropy for X , and makes Y a function of X)

- $\mathcal{M}_3 = \boxed{X} \rightleftarrows \boxed{Y}$
 $IDef_{\mathcal{M}_3}(\mu) = -I_{\mu}(X; Y)$
 (opt. μ makes X, Y share information)

THE SCORING FUNCTION

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

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

EXAMPLES

- $\mathbf{m}_0 = \boxed{X} \quad \boxed{Y}$
 $\text{IDef}_{\mathbf{m}_0}(\mu) = -H_\mu(X, Y)$
 (optimal μ maximizes entropy of X, Y)

- $\mathbf{m}_1 = \boxed{X} \longrightarrow \boxed{Y}$
 $\text{IDef}_{\mathbf{m}_1}(\mu) = -H_\mu(X)$
 (optimal μ maximizes entropy of X)

- $\mathbf{m}_2 = \boxed{X} \rightleftarrows \boxed{Y}$
 $\text{IDef}_{\mathbf{m}_2}(\mu) = -H_\mu(X) + H_\mu(Y | X)$
 (optimal μ maximizes entropy for X , and makes Y a function of X)

- $\mathbf{m}_3 = \boxed{X} \rightleftarrows \boxed{Y}$
 $\text{IDef}_{\mathbf{m}_3}(\mu) = -I_\mu(X; Y)$
 (opt. μ makes X, Y share information)

Information Diagrams

THE SCORING FUNCTION

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

tradeoff parameter $\gamma \geq 0$

Definition (*Inc*)

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

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

The *inconsistency* of \mathcal{M} is

$$\text{Inc}(\mathcal{M}) := \inf_{\mu \in \Delta \mathcal{V}(\mathcal{M})} \text{Inc}_{\mathcal{M}}(\mu).$$

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

THE SCORING FUNCTION

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

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

tradeoff parameter $\gamma \geq 0$

Definition (*Inc*)

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

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

The *inconsistency* of \mathcal{M} is

$$\text{Inc}(\mathcal{M}) := \inf_{\mu \in \Delta \mathcal{V}(\mathcal{M})} \text{Inc}_{\mathcal{M}}(\mu).$$

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

THE SCORING FUNCTION

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

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

tradeoff parameter $\gamma \geq 0$

Definition (*Inc*)

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

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

The *inconsistency* of \mathcal{M} is

$$\text{Inc}(\mathcal{M}) := \inf_{\mu \in \Delta \mathcal{V}(\mathcal{M})} \text{Inc}_{\mathcal{M}}(\mu).$$

Definition (*IDef*)

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

bits to separately determine each target, knowing the source

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

bits to determine all vars

PDG SEMANTICS

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

$$\left\{ \mu \in \Delta[\mathcal{V}(\mathcal{m})] : \text{for all } X \xrightarrow{L} Y \in \mathcal{E}. \mu(Y|X) = \mathbf{p}_L(Y|X) \right\}$$

$[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;

$$[[\mathcal{m}]]_\gamma(\mu) := \underbrace{Inc_{\mathcal{m}}(\mu)}_{\substack{\text{(quantitative)} \\ \text{term}}} + \gamma \underbrace{IDef_{\mathcal{m}}(\mu)}_{\substack{\text{(qualitative)} \\ \text{term}}}$$

$[[\mathcal{m}]]_\gamma^*$ The “best” joint distribution.

$$[[\mathcal{m}]]_\gamma^* := \arg \min_{\mu} [[\mathcal{m}]]_\gamma(\mu)$$

PDG SEMANTICS

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

$$\left\{ \mu \in \Delta[\mathcal{V}(\mathcal{m})] : \text{for all } X \xrightarrow{L} Y \in \mathcal{E}. \mu(Y|X) = \mathbf{p}_L(Y|X) \right\}$$

$[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;

Proposition (uniqueness for small γ)

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

$[[\mathcal{m}]]_\gamma^*$ The “best” joint distribution.

$$[[\mathcal{m}]]_\gamma^* := \arg \min_{\mu} [[\mathcal{m}]]_\gamma(\mu)$$

PDG SEMANTICS

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

$$\left\{ \mu \in \Delta[\mathcal{V}(\mathcal{m})] : \text{for all } X \xrightarrow{L} Y \in \mathcal{E}. \mu(Y|X) = \mathbf{p}_L(Y|X) \right\}$$

$[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;

Proposition (uniqueness for small γ)

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

$[[\mathcal{m}]]^*$ The (unique) “best” joint distribution (in the quantitative limit).

$$[[\mathcal{m}]]^* := \lim_{\gamma \rightarrow 0} \arg \min_{\mu} [[\mathcal{m}]]_\gamma(\mu)$$

PDG SEMANTICS

1. $\{\mathcal{m}\}$ The set of joint distributions consistent with \mathcal{m} ;
2. $[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;
3. $[[\mathcal{m}]]^*$ The “best” joint distribution.

Proposition (*the second semantics extends the first*)

$$\{\mathcal{m}\} = \{\mu : [[\mathcal{m}]]_0(\mu) = 0\}.$$

PDG SEMANTICS

1. $\{\mathcal{m}\}$ The set of joint distributions consistent with \mathcal{m} ;
2. $[[\mathcal{m}]]_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathcal{m} ;
3. $[[\mathcal{m}]]^*$ The “best” joint distribution.

Proposition (*the second semantics extends the first*)

$$\{\mathcal{m}\} = \{\mu : [[\mathcal{m}]]_0(\mu) = 0\}.$$

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

$$[[\mathcal{m}]]^* \in [[\mathcal{m}]]_0^*, \text{ so if } \mathcal{m} \text{ is consistent, then } [[\mathcal{m}]]^* \in \{\mathcal{m}\}.$$

CAPTURING BAYESIAN NETWORKS

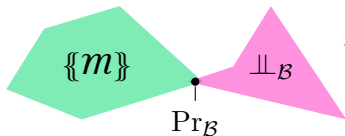
Let $\mathbf{m}_{\mathcal{B},\beta}$ be the PDG corresponding to the BN \mathcal{B} , with weights β .

Theorem (*BNs are PDGs*)

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

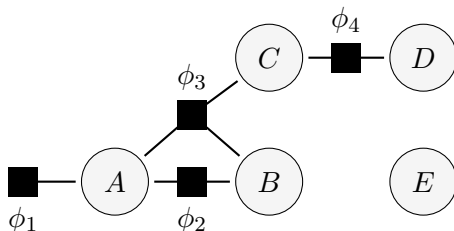
$$[[\mathbf{m}_{\mathcal{B},\beta}]_{\gamma}^* = \{\text{Pr}_{\mathcal{B}}\}, \quad \text{and thus} \quad [[\mathbf{m}_{\mathcal{B},\beta}]^* = \text{Pr}_{\mathcal{B}}.$$

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



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

FACTOR GRAPHS



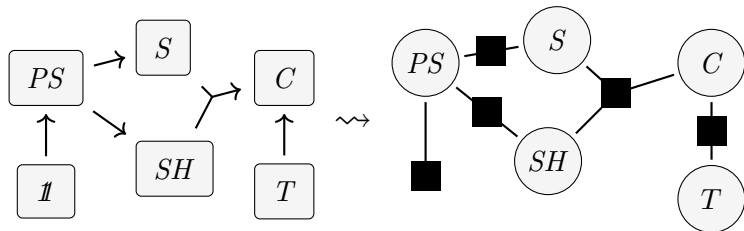
Definition

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

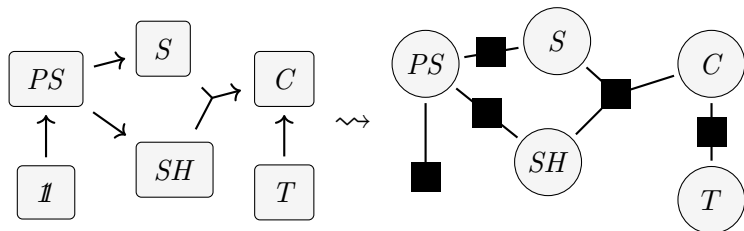
$$\Pr_{\Phi}(\vec{x}) = \frac{1}{Z_{\Phi}} \prod_{J \in \mathcal{J}} \phi_J(\vec{x}_J),$$

where Z_{Φ} is the normalization constant.

PDGs AS FACTOR GRAPHS

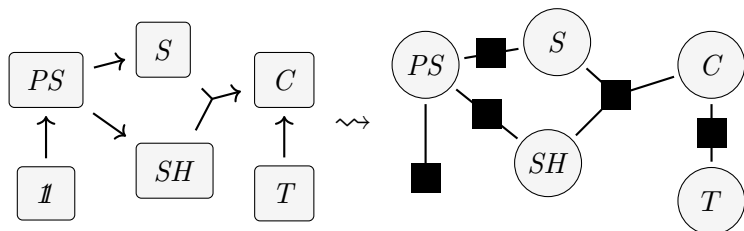


PDGs AS FACTOR GRAPHS



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

PDGs AS FACTOR GRAPHS

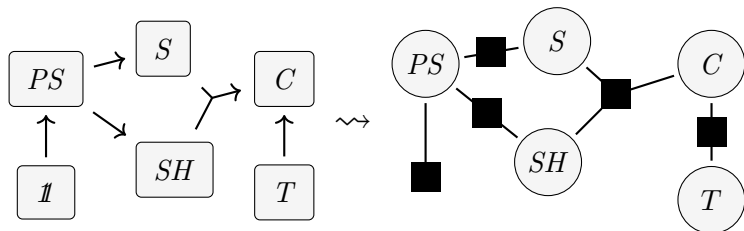


The cpds of a PDG are essentially factors. Are the semantics different?
Not for $\gamma = 1$.

Theorem

$[[\mathcal{n}]]_1^* = \Pr_{\Phi_{\mathcal{n}}}$ for all unweighted PDGs \mathcal{n} .

PDGs AS FACTOR GRAPHS



The cpds of a PDG are essentially factors. Are the semantics different?
Not for $\gamma = 1$.

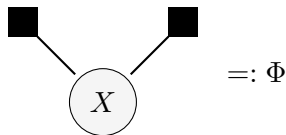
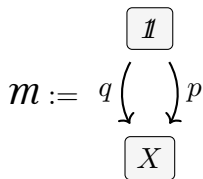
Theorem

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

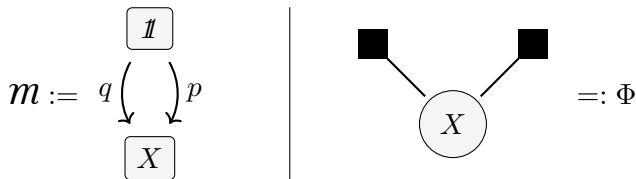
Theorem

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

WHY NOT USE FACTOR GRAPHS, THEN?

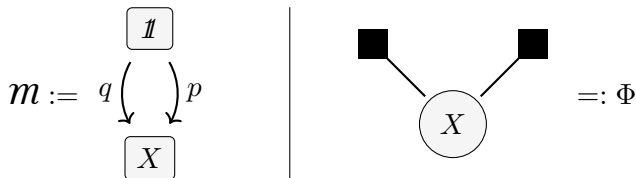


WHY NOT USE FACTOR GRAPHS, THEN?



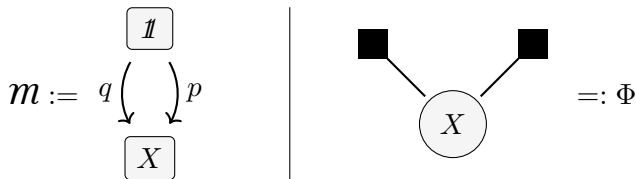
- If $p = q$, then $[[m]]^* = p = q \dots$

WHY NOT USE FACTOR GRAPHS, THEN?



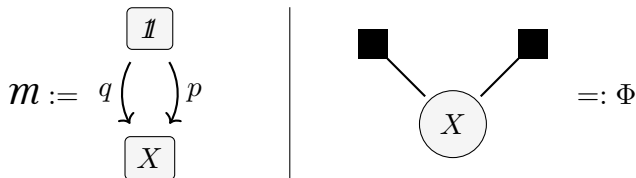
- If $p = q$, then $[[m]]^* = p = q \dots$
- ... but $\Pr_{\Phi} \propto p^2$

WHY NOT USE FACTOR GRAPHS, THEN?



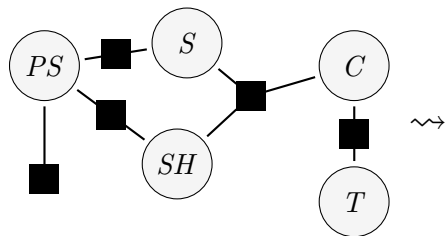
- If $p = q$, then $[[m]]^* = p = q \dots$
- \dots but $\Pr_{\Phi} \propto p^2$
- More generally, (positive) factors individually have *no meaning*,

WHY NOT USE FACTOR GRAPHS, THEN?

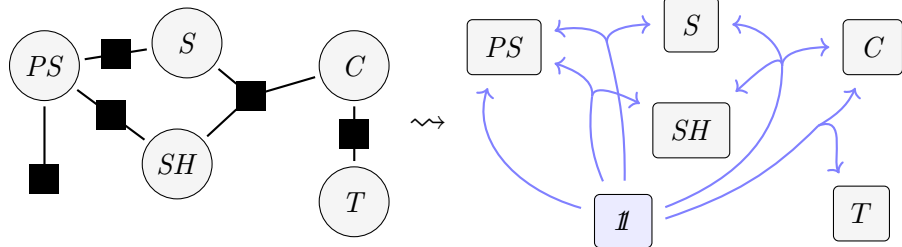


- If $p = q$, then $[[\mathcal{M}]]^* = p = q \dots$
- \dots but $\Pr_{\Phi} \propto p^2$
- More generally, (positive) factors individually have *no meaning*,
- a factor graph can fail to normalize, in which case it has no global semantics either.

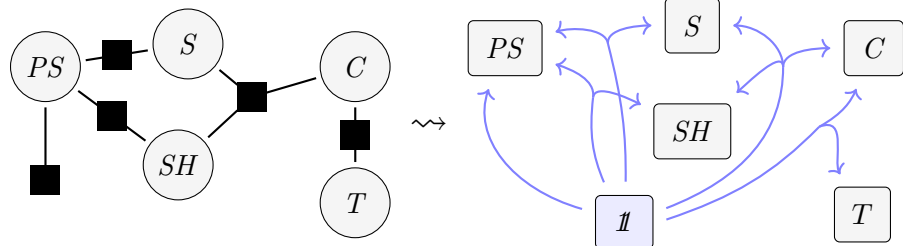
FACTOR GRAPHS AS PDGs



FACTOR GRAPHS AS PDGs



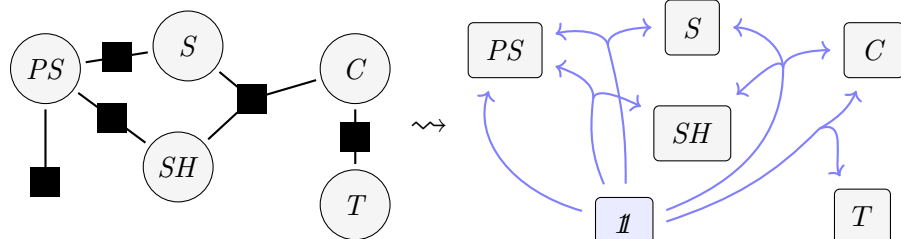
FACTOR GRAPHS AS PDGs



Theorem

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

FACTOR GRAPHS AS PDGs



Theorem

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

Theorem

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

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

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

INFERENCE AND INCONSISTENCY: A GLIMPSE.

Conditioning as inconsistency resolution.

To condition on $Y = y$, in \mathcal{m} , simply add the edge $\mathbb{1} \xrightarrow{\delta_y} Y$ to get $\mathcal{m}_{Y=y}$.
Then $[[\mathcal{m}_{Y=y}]]^* = [[\mathcal{m}]]^* \mid (Y = y)$.

INFERENCE AND INCONSISTENCY: A GLIMPSE.

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{1} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $[[\mathcal{M}_{Y=y}]]^* = [[\mathcal{M}]]^* \mid (Y=y)$.

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

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

INFERENCE AND INCONSISTENCY: A GLIMPSE.

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{m} , simply add the edge $\mathbb{1} \xrightarrow{\delta_y} Y$ to get $\mathcal{m}_{Y=y}$.
Then $[[\mathcal{m}_{Y=y}]]^* = [[\mathcal{m}]]^* \mid (Y=y)$.

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

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

INFERENCE AND INCONSISTENCY: A GLIMPSE.

Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{1} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $[[\mathcal{M}_{Y=y}]]^* = [[\mathcal{M}]]^* \mid (Y=y)$.

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

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

INFERENCE AND INCONSISTENCY: A GLIMPSE.

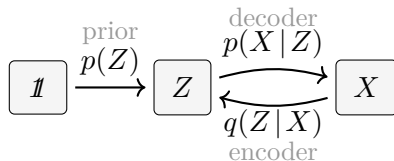
Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{1} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $[[\mathcal{M}_{Y=y}]]^* = [[\mathcal{M}]]^* \mid (Y=y)$.

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

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

This is closely related to standard variational techniques!



INFERENCE AND INCONSISTENCY: A GLIMPSE.

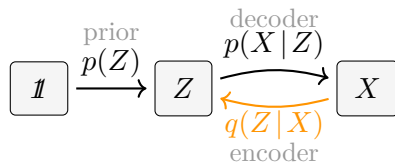
Conditioning as inconsistency resolution.

To condition on $Y=y$, in \mathcal{M} , simply add the edge $\mathbb{1} \xrightarrow{\delta_y} Y$ to get $\mathcal{M}_{Y=y}$.
Then $[[\mathcal{M}_{Y=y}]]^* = [[\mathcal{M}]]^* \mid (Y=y)$.

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

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

This is closely related to standard variational techniques!



SUMMARY

PDGs...

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

SUMMARY

PDGs...

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

SUMMARY

PDGs...

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

SUMMARY

PDGs...

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

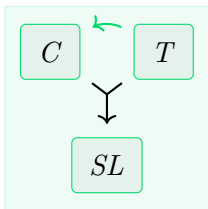
SUMMARY

PDGs...

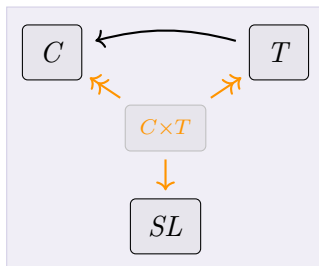
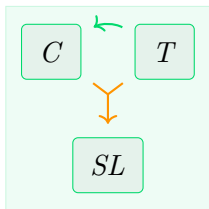
- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular; to combine info from two sources, simply take a PDG union. This incorporates new data (edge cpds) and concepts (nodes) without affecting previous information.
- cleanly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

But there is much more to be done!

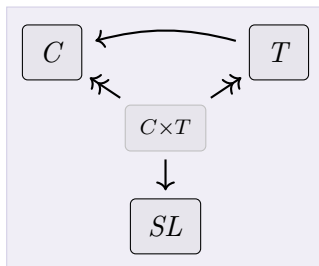
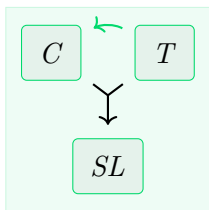
HYPER-GRAPHS? OR MERELY GRAPHS?



HYPER-GRAPHS? OR MERELY GRAPHS?

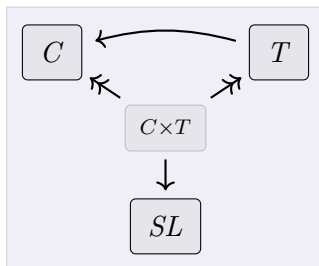
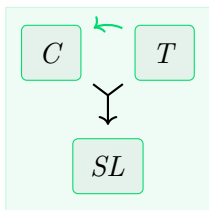


HYPER-GRAPHS? OR MERELY GRAPHS?



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

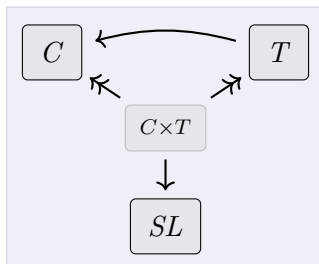
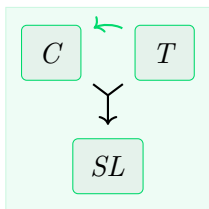
HYPER-GRAPHS? OR MERELY GRAPHS?



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

joint distributions \Leftrightarrow expanded joint distributions
satisfying coherence constraints

HYPER-GRAPHS? OR MERELY GRAPHS?



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

joint distributions \Leftrightarrow expanded joint distributions
satisfying coherence constraints

(working directly with hypergraphs is also possible)

ILLUSTRATIONS OF $IDef$

