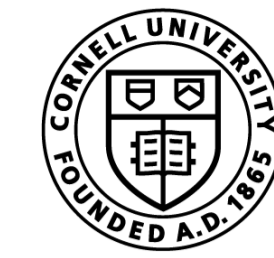


THE LOCAL INCONSISTENCY RESOLUTION ALGORITHM



Oliver E Richardson

A generic algorithm for learning and (approximate) inference, with an intuitive epistemic interpretation. Unifies many important algorithms.

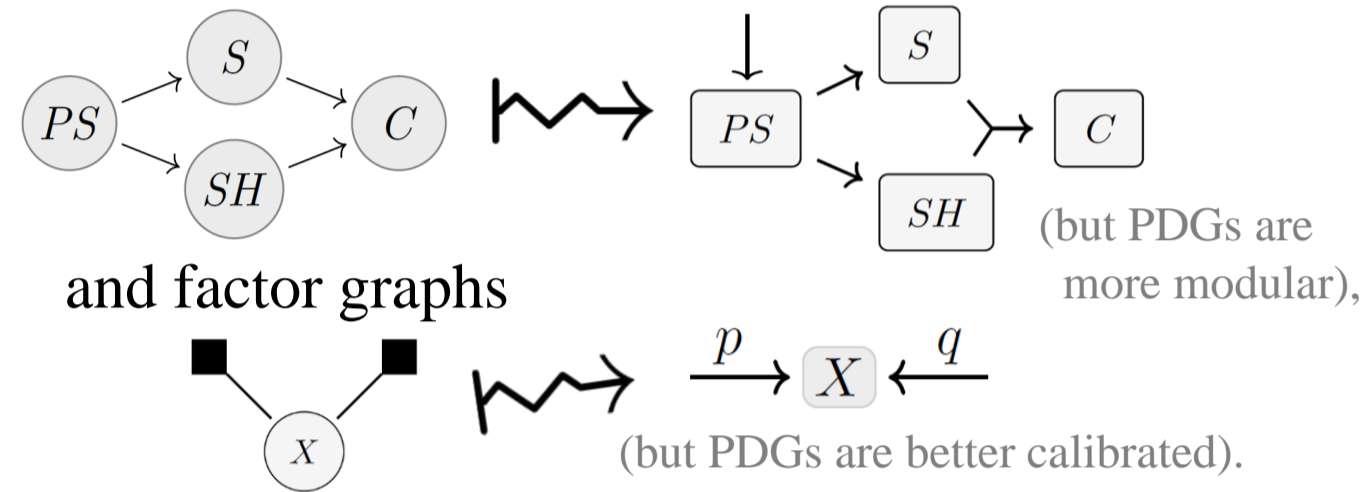
What causes changes in beliefs? Some say it is internal conflict. But identifying inconsistencies is difficult. So in practice, we resolve them *locally*: looking only at a small part of the picture, and changing only another small part at a time.

Key Representation: Probabilistic Dependency Graphs (PDGs) are directed (hyper) graphs with probabilities and confidences attached to edges.

PDGs can capture:

- ❖ inconsistent beliefs, providing a natural way to measure the degree of this inconsistency;

- ❖ **graphical models**, such as Bayesian networks



- ❖ **learning settings** and their **loss functions**, e.g.,

- variational autoencoders (VAEs)

$$\left\langle \left\langle \begin{array}{c} p \\ \rightarrow \\ Z \end{array} \begin{array}{c} d \\ \rightarrow \\ X \end{array} \leftarrow x \right\rangle \right\rangle = -\text{ELBO}_{p,e,d}(x)$$
 ... including their standard loss function, as inconsistency
- statistical divergences

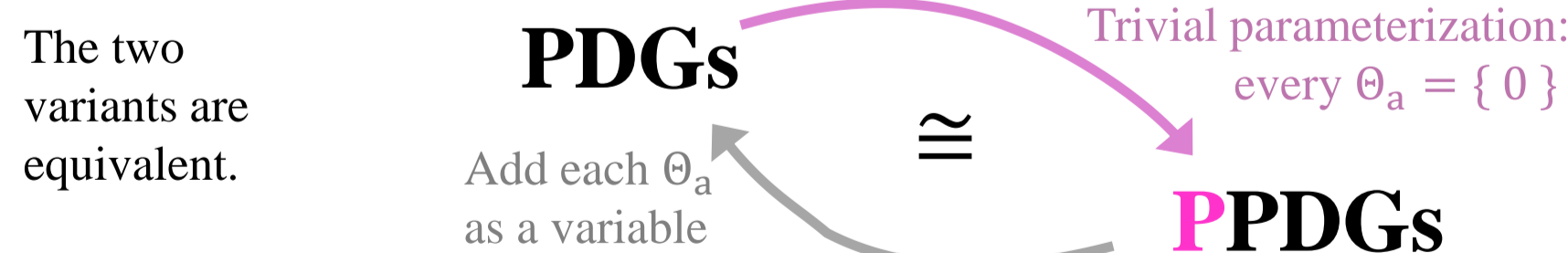
$$\left\langle \left\langle \begin{array}{c} p \\ \rightarrow \\ X \end{array} \leftarrow \begin{array}{c} q \\ \rightarrow \\ X \end{array} \right\rangle \right\rangle$$
 Generates Rényi divergences, reverse KL, conditional divergences.
- regularizers as priors, accuracy, MSE, ...

FORMALISM: PARAMETERIZED PDGs

$\mathcal{M}(\Theta)$ = variables \mathcal{X} connected by arcs \mathcal{A} ;
each $(S \xrightarrow{a} T) \in \mathcal{A}$ is associated with:

- a convex parameter space $\Theta_a \subseteq \mathbb{R}^n$
- a conditional probability $\mathbb{P}_a(T|S, \Theta_a)$,
- two confidences: β_a (observational) and α_a (structural).

Fix a parameter setting $\theta \in \prod_{a \in \mathcal{A}} \Theta_a$, to get an (ordinary) PDG $\mathcal{M}(\theta)$.



Inconsistency semantics.

A joint probability $\mu(\mathcal{X})$ can be incompatible with a PDG in two ways:

Observational Incompatibility with (\mathbb{P}, β)

$$\sum_{S \xrightarrow{a} T \in \mathcal{A}} \beta_a \mathcal{D}(\mu(T, S) \parallel \mathbb{P}_a(T|S)\mu(S))$$

Structural Deficiency with (\mathcal{A}, α)

$$\mathbb{E}_{\mu} \left[\log \frac{\mu(\mathcal{X})}{\lambda(\mathcal{X})} \prod_{S \xrightarrow{a} T} \left(\frac{\lambda(T|S)}{\mu(T|S)} \right)^{\alpha_a} \right]$$

Degree of inconsistency $\llbracket \mathcal{M} \rrbracket_{\gamma} := \inf_{\mu} \left(\text{OIncm}(\mu) + \gamma \text{SDefm}(\mu) \right)$ placing weight $\gamma \geq 0$ on the structural information is the smallest possible incompatibility with any $\mu(\mathcal{X})$.

Algorithm: Local Inconsistency Resolution (LIR)

Input: context PDG Ctx ,
mutable memory $\mathcal{M}(\Theta)$.

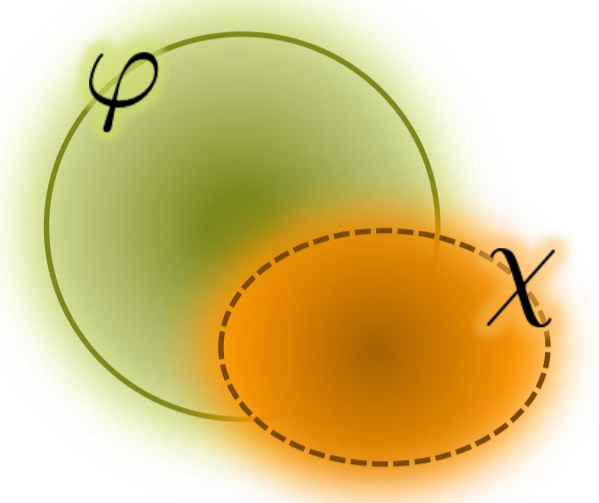
Initialize $\theta^{(0)}$;
for $t = 0, 1, 2, \dots$ **do**

1. $Ctx \leftarrow \text{REFRESH}(Ctx)$;
 $\varphi, \chi, \gamma \leftarrow \text{REFOCUS}()$;
 Write $\exp_{\theta}(t X)$ for the path following vector field X for time t , starting at θ .
 Calculate the inconsistency of the combined context and memory, weighted by attention.
Gradient Flow of $f: \Theta \rightarrow \mathbb{R}$ starting at θ :
 $t \mapsto \exp_{\theta}(t \nabla_{\Theta} f(\Theta))$
2. $\theta^{(t+1)} \leftarrow \exp_{\theta^{(t)}} \left\{ -\chi \odot \nabla_{\Theta} \left\langle \left\langle \varphi \odot (Ctx + \mathcal{M}(\Theta)) \right\rangle \right\rangle_{\gamma} \right\}$;

Reduce this inconsistency by (an approximation to) gradient flow, starting at previous state $\theta^{(t)}$, changing each parameter in proportion to our control of it.

FOCUS: ATTENTION AND CONTROL

attend only to probabilities of a subset of arcs $A \subseteq \mathcal{A}$ (or attn mask φ)
control only parameters of a subset of arcs $C \subseteq \mathcal{A}$ (or ctrl mask χ)

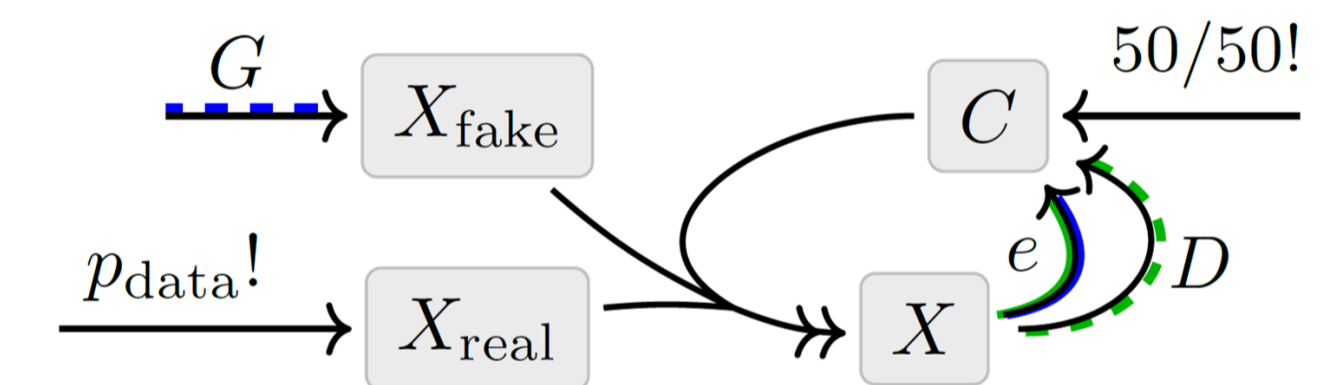


Typical use case: select focus (φ, χ) from a fixed set of foci $\mathbf{F} = \{ \blacksquare, \blacksquare, \dots \}$.

MORE EXAMPLES

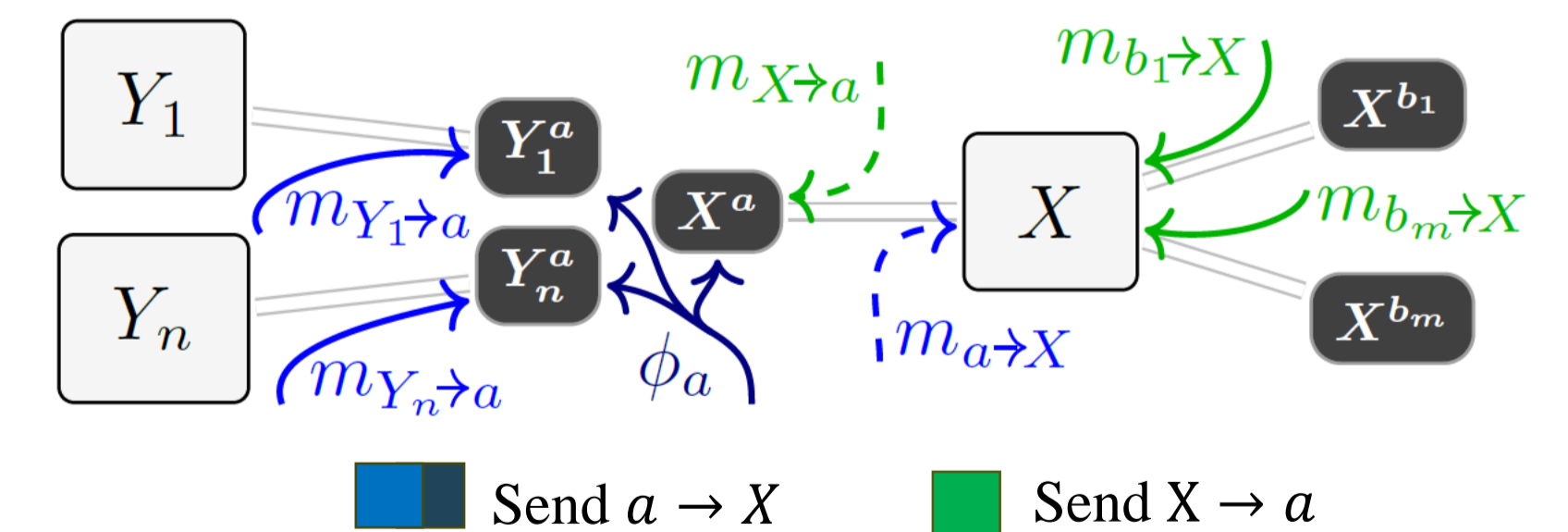
- ❖ Training Generative Adversarial Networks (GANs).

Typically trained with minimax game: $\min_G \max_D \mathcal{L}^{GAN}$



- **Generator's focus**
inconsistency = $\text{JSD}(G, p_{\text{data}})$
or $+\mathcal{L}^{GAN}$ if disbelieves D
- **Discriminator's focus**
inconsistency = $\text{KL}(D, D^{\text{opt}})$
or $-\mathcal{L}^{GAN}$ if disbelieves e

- ❖ Message Passing: Sum-Product Belief Propagation



Observation: the message passing equations are sums of products of factors, i.e., do inference in *local* factor graphs.

Here, $\mathcal{M}(\theta)$: collection of messages (BP data structure)
 Ctx : original factor graph, as a PDG

- ❖ Variational Inference, EM algorithm, e.g., VAE training.

LIR IN THE CLASSIFICATION SETTING

Consider a discriminator $p_{\theta}(Y|X)$ and sample (x, y) .
Together, they have inconsistency

$$\left\langle \left\langle \begin{array}{c} x \\ \rightarrow \\ X \end{array} \begin{array}{c} p_{\theta} \\ \rightarrow \\ Y \end{array} \leftarrow y \right\rangle \right\rangle = \log \frac{1}{p_{\theta}(y|x)}.$$

Can resolve by modifying:

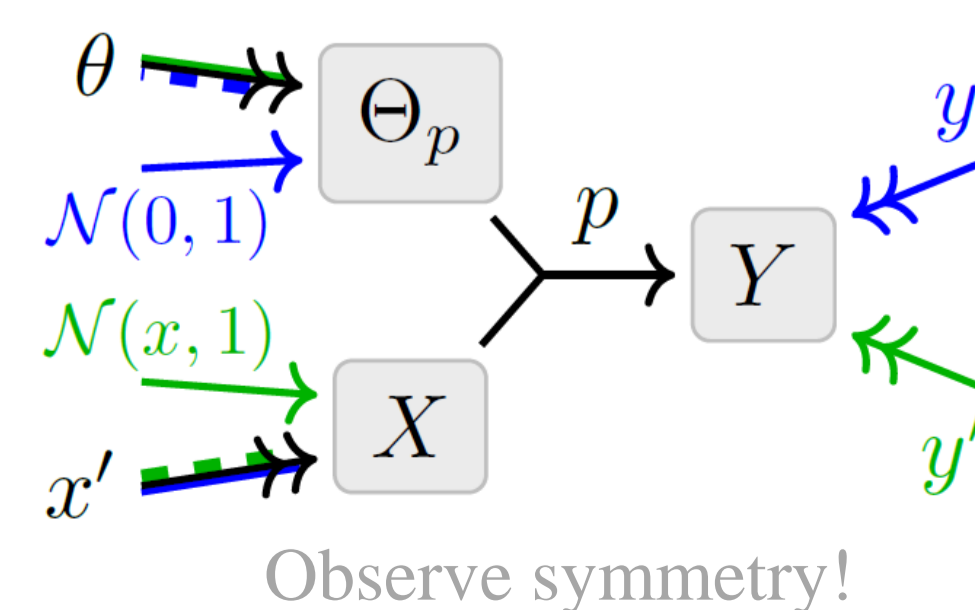
- θ , to train the discriminator
- y , resulting in a forward pass
- x , to form an adversarial example

SGD. Control over p_{θ} . Replace (x, y) with empirical distribution over a batch, and suppose $\text{REFRESH}(Ctx)$ gets a new batch. This performs SGD with learning rate $\chi(p) \cdot \varphi(p)$.

Adversarial Training.

Add discriminator params as a variable Θ_p with Gaussian prior.

- Construct attack $x' \approx x$ that p misclassifies as y'
- Patch p to classify x' as y



Observe symmetry!