# Differentiation and Weighted Model Integration

Pedro Zuidberg Dos Martires

**KU LEUVEN**

DTAI

Weighted Model Integration:

- Probabilistic inference in the discrete-continuous domain.
- Generalization of weighted model counting.
- Performing counting and integration in this hybrid domain.

Weighted Model Integration:

- ▶ Probabilistic inference in the discrete-continuous domain.
- ▶ Generalization of weighted model counting.
- ▶ Performing counting and integration in this hybrid domain.

Why do we need differentiation?

- ▶ Probabilistic inference is hard, many approximation schemes rely on differentiation, e.g. variational inference, Hamilton Monte Carlo
- ▶ Taking derivatives allows for gradient based optimization!

**Weighted Model Integration:**

- ▸ Probabilistic inference in the discrete-continuous domain.
- ▸ Generalization of weighted model counting.
- ▸ Performing counting and integration in this hybrid domain.

**Why do we need differentiation?**

- ▸ Probabilistic inference is hard, many approximation schemes rely on differentiation, e.g. variational inference, Hamilton Monte Carlo
- ▸ Taking derivatives allows for gradient based optimization!

**This talk:**

- ▸ Show how differentiation can be done in the weighted model integration context.
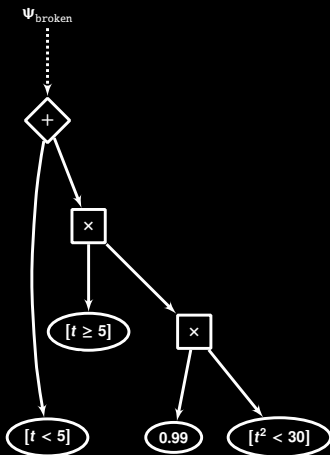- ▸ Show difficulties that lie ahead!

$$\texttt{working} \leftrightarrow (\texttt{cooling} \wedge (t^2 < 30)) \vee (t < 5)$$

$$\texttt{working} \leftrightarrow (\texttt{cooling} \land (\texttt{t}^2 < 30)) \lor (\texttt{t} < 5)$$
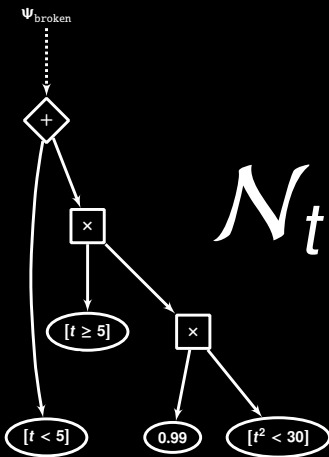
$$p(\texttt{cooling}) = 0.99$$

$$\texttt{t} \sim \mathcal{N}_\texttt{t}(20, 5)$$

Compile logic structure into arithmetic circuit.
Arithmetic circuits are related to sum-product-networks.

$$\mathrm{WMI}(\phi, w | \mathbf{x}, \mathbf{b}) = \int \sum_{\mathbf{b}_\mathcal{I} \in \mathcal{I}_{\mathbf{b},\mathbf{b_a}}(\phi_a)} \prod_{b_i \in \mathbf{b}_\mathcal{I}} \alpha_{b_i}(\mathbf{x}) w_x(\mathbf{x}) d\mathbf{x} \quad (1)$$

$$= \int \Psi(\mathbf{x}) w_x(\mathbf{x}) d\mathbf{x} \quad (2)$$

$$= \mathbb{E}_{w_x(\mathbf{x})}[\Psi(\mathbf{x})] \quad (3)$$

# Cross-Entropy *H*

Tells you how different two distributions *p* and *q* are.

$$H(p, q) = \mathbb{E}_p[-\log q] \tag{4}$$

*p* is the *true* distribution (observation of the world).
*q* is our model of the true distribution.
q depends on the parameters $\theta$.

Minimize $H(p, q)$ by learning the parameters $\theta$.

# Gradient Descent

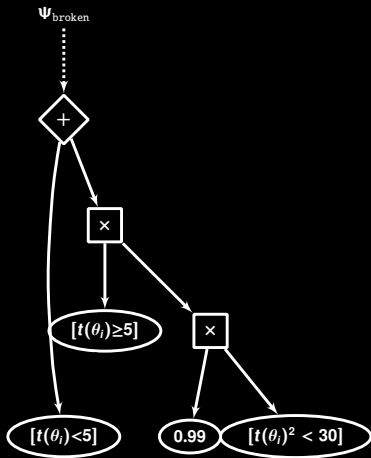$$\theta \leftarrow \theta + \eta \nabla_\theta \mathbb{E}_p[-\log q] \tag{5}$$

# Gradient of the Cross-Entropy

$$\nabla_\theta \mathbb{E}_p[-\log q] \tag{6}$$

$$= \mathbb{E}_p\Big[-\frac{1}{q(\theta)}\nabla_\theta \sum_{\mathbf{b}_\mathcal{I}\in\mathcal{I}_{\mathbf{b}.\mathbf{b}_a}(\phi_a)} \prod_{b_i\in\mathbf{b}_\mathcal{I}} \alpha_{b_i}(\theta)\Big] \tag{7}$$

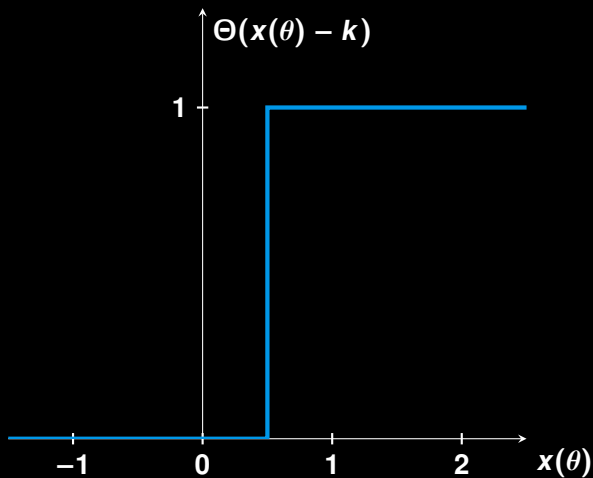$$= \mathbb{E}_p\Big[-\frac{1}{\Psi(\theta)}\nabla_\theta \Psi(\theta)\Big] \tag{8}$$

$$\frac{\partial}{\partial \theta_i}$$

Ψ_broken

+

×

[t(θ_i)≥5]

×

[t(θ_i)<5]

0.99    [t(θ_i)² < 30]

# Applying the Product Rule

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_p[-\log q] \tag{9}$$

$$= \mathbb{E}_p[-\frac{1}{q(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} \sum_{\mathbf{b}_I \in \mathcal{I}_{\mathbf{b},\mathbf{b}_a}(\phi_a)} \prod_{b_i \in \mathbf{b}_I} \alpha_{b_i}(\boldsymbol{\theta})] \tag{10}$$

$$= \mathbb{E}_p[-\frac{1}{\Psi(\boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta})] \tag{11}$$

$$= \mathbb{E}_p[-\frac{1}{\Psi(\boldsymbol{\theta})} \sum_{\mathbf{b}_I \in \mathcal{I}_{\mathbf{b},\mathbf{b}_a}(\phi_a)} \sum_{b_i \in \mathbf{b}_I} \nabla_{\boldsymbol{\theta}}(\alpha_{b_i}(\boldsymbol{\theta})) \prod_{b_j \in \mathbf{b}_I \setminus \{b_i\}} \alpha_{b_j}(\boldsymbol{\theta})] \tag{12}$$

# A Simple One Dimensional Case

$$\nabla_{\boldsymbol{\theta}} \alpha(\boldsymbol{\theta}) = \frac{\partial \alpha(\theta)}{\partial \theta} \tag{13}$$

$$= \frac{\partial \Theta(x(\theta) - k)}{\partial \theta} \tag{14}$$

# Heaviside Step Function

# A Simple One Dimensional Case

$$\nabla_{\boldsymbol{\theta}}\alpha(\boldsymbol{\theta}) = \frac{\partial\alpha(\theta)}{\partial\theta} \tag{15}$$

$$= \frac{\partial\Theta(x(\theta) - k)}{\partial\theta} \tag{16}$$

$$= \delta(x(\theta) - k)\frac{\partial x(\theta)}{\partial\theta} \tag{17}$$

# In Higher Dimensions

- ▶ Gradient is generalization of inward normal derivative.
- ▶ Leads to surface integral (boundary of indicator function)

Open questions:

- ▶ What is the computational hardness of the surface integral?
- ▶ Is it equivalent to the optimization of the 0-1 loss (NP-complete)?
- ▶ Can we use convex relaxation for a practical algorithm?
- ▶ Is it beneficial to restrict ourselves to a subclass of constraints? (very probably yes)

Where might this be useful?

- ▸ Parameter learning in hybrid probabilistic programs.
- ▸ Probabilistic inference through (stochastic) variational inference.
- ▸ Probabilistic inference through Hamilton Monte Carlo.