

# Safe Semi-Supervised Learning of Sum-Product Networks

Martin Trapp<sup>1,2</sup>, Tamas Madl<sup>2</sup>, Robert Peharz<sup>3</sup>, Franz Pernkopf<sup>1</sup>, and Robert Trapp<sup>2</sup>

<sup>1</sup>Signal Processing and Speech Communication Lab., Graz University of Technology, Austria

<sup>2</sup>Austrian Research Institute for Artificial Intelligence, Austria

<sup>3</sup>Computational and Biological Learning Lab., University of Cambridge, UK

## Motivation

- In several domains, obtaining class labels is expensive
- Most semi-supervised approaches enforce restrictive assumptions on the data distribution
- **Goal:** Non-restrictive semi-supervised learning

## Sum-Product Networks (SPNs)

- Deep probabilistic model capturing expressive variable interactions, while guaranteeing exact and efficient inference
- Generative parameter learning using efficient Expectation Maximisation (EM) [1, 2]
- Discriminative parameter learning using backpropagation for cond. llh maximisation [3]

## Contrastive Pessimistic Likelihood Estimation [4]

- Semi-supervised learning for generative linear models
- Maintains soft-labels for each unlabelled observation
- Alternates between optimistic parameter updates and pessimistic soft-label updates

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmax}} \underset{q}{\operatorname{argmin}} L(\theta | \mathcal{X}, \mathcal{U}, q) - L(\theta^+ | \mathcal{X}, \mathcal{U}, q)$$

- Results in a safe solution, since the supervised solution can allways be used in the worst case

## Our Approach (MCP-SPN)

- First semi-supervised learning approach for SPNs
- Allows generative and discriminative learning of non-linear decision boundaries
- Guarantees that adding unlabelled data can increase, but not degrade, the training performance (safe)

$$L(\theta^* | \mathcal{X}, \mathcal{U}, q) \geq L(\theta^+ | \mathcal{X}, \mathcal{U}, q)$$

- Is computationally efficient (scales linearly)
- Does not enforce restrictive assumptions on the data distribution

## Generative Training

- SPNs with soft-labels

$$\mathcal{S}[\mathbf{u}, \mathbf{q} | \theta] = \sum_{k=1}^K q_k w_k S_k[\mathbf{u} | \mathbf{q}, \theta]$$

- Optimisation given labelled and unlabelled data

$$L(\theta | \mathcal{X}, \mathcal{U}, q) = \sum_{n=1}^N \log \mathcal{S}[\mathbf{x}_n, \mathbf{y}_n | \theta] + \sum_{m=1}^M \log \mathcal{S}[\mathbf{u}_m, \mathbf{q}_m | \theta]$$

- Pessimistic update of soft-labels

$$\nabla q_{mk} = \frac{\partial L(\theta^* | \mathcal{X}, \mathcal{U}, q)}{\partial q_{mk}} - \frac{\partial L(\theta^+ | \mathcal{X}, \mathcal{U}, q)}{\partial q_{mk}}$$

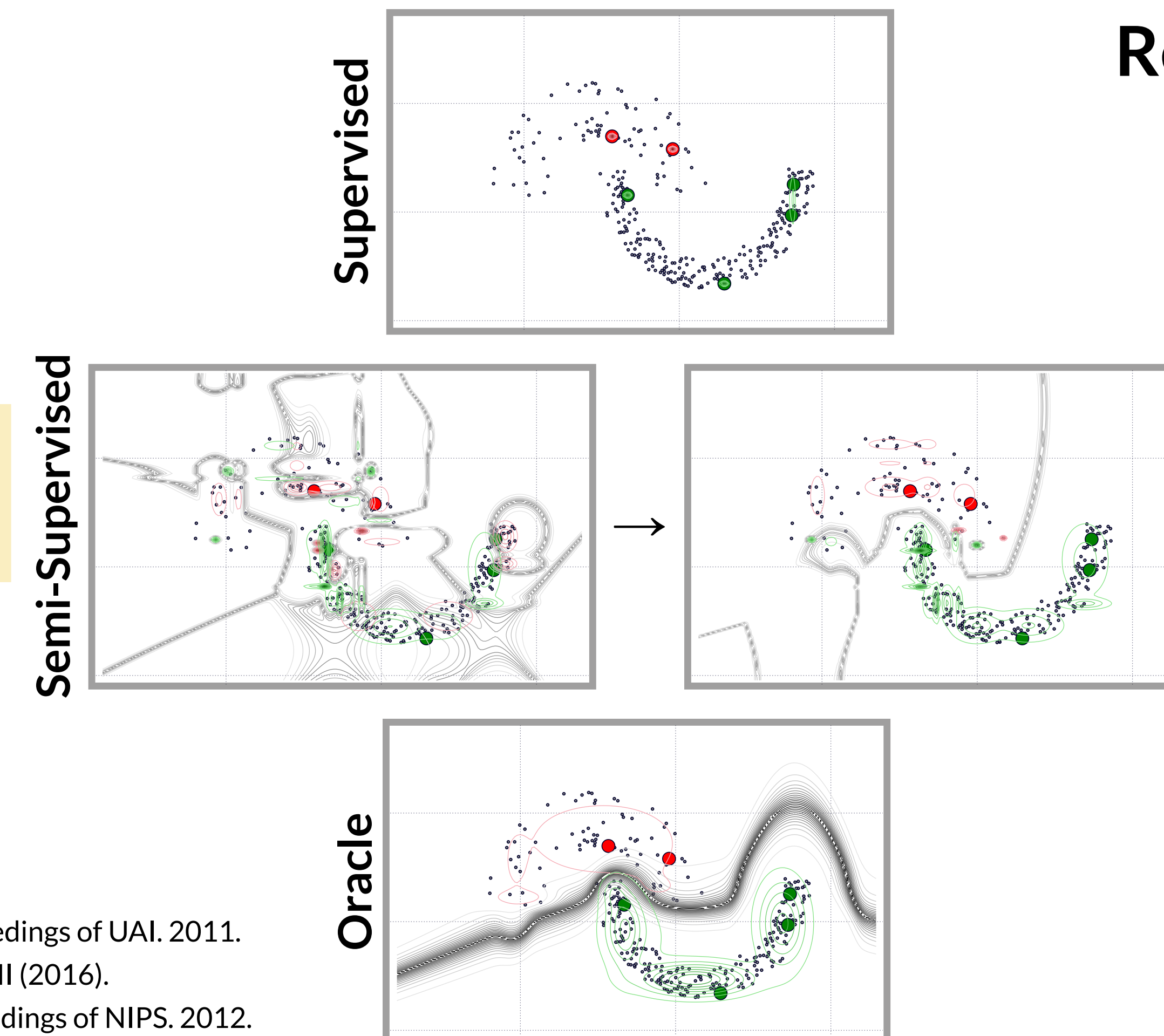
## Discriminative Training

- Optimisation of conditional log likelihood given labelled and unlabelled data using backpropagation

$$CL(\theta | \mathcal{X}, \mathcal{U}, q) = \sum_{n=1}^N \log \mathcal{S}[\mathbf{y}_n | \mathbf{x}_n, \theta] + \sum_{m=1}^M \log \mathcal{S}[\mathbf{q}_m | \mathbf{u}_m, \theta]$$

- Pessimistic update of soft-labels as in the generative setting
- Initialisation of soft-labels using optimistic labelling for discriminative training and random draws from a Dirichlet distribution for generative training

## Qualitative Results



## Quantitative Results

Data Set	Supervised	SSL	Oracle	MCPLDA
BUPA	$-438.75 \pm 7 \cdot 10^0$	$-7.31 \pm 6 \cdot 10^{-2}$	$-8.80 \pm 2 \cdot 10^{-1}$	$-9.07 \pm 3 \cdot 10^{-2}$
Fertility	$-3.31 \pm 3 \cdot 10^{-2}$	$-3.06 \pm 7 \cdot 10^{-3}$	$-3.00 \pm 6 \cdot 10^{-3}$	$-12.68 \pm 5 \cdot 10^{-2}$
Haberman	$-138.63 \pm 4 \cdot 10^0$	$-5.05 \pm 6 \cdot 10^{-2}$	$-5.14 \pm 6 \cdot 10^{-2}$	$-7.83 \pm 1 \cdot 10^{-1}$
ILPD	$-5.62 \pm 3 \cdot 10^0$	$-1.15 \pm 2 \cdot 10^{-2}$	$-1.00 \pm 1 \cdot 10^{-2}$	$-37.54 \pm 1 \cdot 10^{-1}$
Ionosphere	$-2.83 \pm 5 \cdot 10^{-2}$	$-1.61 \pm 1 \cdot 10^{-2}$	$-1.52 \pm 9 \cdot 10^{-3}$	$-46.12 \pm 5 \cdot 10^{-2}$
Iris	$-20.65 \pm 9 \cdot 10^{-1}$	$-3.78 \pm 3 \cdot 10^{-2}$	$-2.17 \pm 1 \cdot 10^{-2}$	$-2.65 \pm 5 \cdot 10^{-2}$
Parkinsons	$-1.32 \pm 4 \cdot 10^{-3}$	$-1.34 \pm 4 \cdot 10^{-3}$	$-1.30 \pm 2 \cdot 10^{-3}$	$-2.27 \pm 5 \cdot 10^{-2}$
WDBC	$-1.90 \pm 1 \cdot 10^{-3}$	$-1.93 \pm 2 \cdot 10^{-3}$	$-1.88 \pm 3 \cdot 10^{-4}$	$-10.75 \pm 1 \cdot 10^{-2}$
Wine	$-2.47 \pm 4 \cdot 10^{-3}$	$-2.47 \pm 2 \cdot 10^{-3}$	$-2.44 \pm 9 \cdot 10^{-4}$	$-15.28 \pm 2 \cdot 10^{-2}$

Data Set	Supervised	SSL	Oracle	TSVM	ICLSC	MER
BUPA	$0.41 \pm 1 \cdot 10^{-2}$	$0.40 \pm 1 \cdot 10^{-2}$	$0.48 \pm 5 \cdot 10^{-3}$	$0.36 \pm 2 \cdot 10^{-2}$	$0.47 \pm 7 \cdot 10^{-3}$	$0.42 \pm 1 \cdot 10^{-2}$
Fertility	$0.07 \pm 2 \cdot 10^{-2}$	$0.03 \pm 1 \cdot 10^{-2}$	$0.06 \pm 2 \cdot 10^{-2}$	$0.07 \pm 2 \cdot 10^{-2}$	$0.07 \pm 2 \cdot 10^{-2}$	$0.12 \pm 2 \cdot 10^{-2}$
Haber.	$0.23 \pm 2 \cdot 10^{-2}$	$0.28 \pm 2 \cdot 10^{-2}$	$0.25 \pm 0.$	$0.20 \pm 2 \cdot 10^{-2}$	$0.33 \pm 1 \cdot 10^{-2}$	$0.27 \pm 1 \cdot 10^{-2}$
ILPD	$0.17 \pm 2 \cdot 10^{-2}$	$0.20 \pm 1 \cdot 10^{-2}$	$0.24 \pm 4 \cdot 10^{-3}$	$0.23 \pm 2 \cdot 10^{-2}$	$0.29 \pm 1 \cdot 10^{-2}$	$0.33 \pm 2 \cdot 10^{-2}$
Ionos.	$0.79 \pm 4 \cdot 10^{-3}$	$0.82 \pm 4 \cdot 10^{-3}$	$0.87 \pm 0.$	$0.66 \pm 9 \cdot 10^{-3}$	$0.61 \pm 9 \cdot 10^{-3}$	$0.70 \pm 7 \cdot 10^{-3}$
Iris	$0.73 \pm 1 \cdot 10^{-2}$	$0.88 \pm 1 \cdot 10^{-2}$	$0.93 \pm 0.$	$0.72 \pm 1 \cdot 10^{-2}$	$0.74 \pm 2 \cdot 10^{-2}$	$0.80 \pm 6 \cdot 10^{-3}$
Parkins.	$0.72 \pm 1 \cdot 10^{-2}$	$0.77 \pm 4 \cdot 10^{-3}$	$0.82 \pm 4 \cdot 10^{-3}$	$0.74 \pm 1 \cdot 10^{-2}$	$0.66 \pm 2 \cdot 10^{-2}$	$0.68 \pm 1 \cdot 10^{-2}$
PID	$0.38 \pm 1 \cdot 10^{-2}$	$0.45 \pm 1 \cdot 10^{-2}$	$0.64 \pm 8 \cdot 10^{-4}$	$0.45 \pm 1 \cdot 10^{-2}$	$0.54 \pm 7 \cdot 10^{-3}$	$0.57 \pm 9 \cdot 10^{-3}$
WDBC	$0.85 \pm 3 \cdot 10^{-3}$	$0.90 \pm 2 \cdot 10^{-3}$	$0.92 \pm 3 \cdot 10^{-4}$	$0.91 \pm 4 \cdot 10^{-3}$	$0.88 \pm 4 \cdot 10^{-3}$	$0.92 \pm 3 \cdot 10^{-3}$
Wine	$0.82 \pm 7 \cdot 10^{-3}$	$0.97 \pm 2 \cdot 10^{-3}$	$0.97 \pm 4 \cdot 10^{-3}$	$0.96 \pm 2 \cdot 10^{-3}$	$0.95 \pm 7 \cdot 10^{-3}$	$0.95 \pm 9 \cdot 10^{-3}$

[1] Poon, Hoifung and Domingos, Pedro. "Sum-product networks: a new deep architecture." In proceedings of UAI. 2011.

[2] Peharz, Robert, et al. "On the latent variable interpretation in sum-product networks." IEEE TPAMI (2016).

[3] Gens, Robert, and Pedro Domingos. "Discriminative learning of sum-product networks." In proceedings of NIPS. 2012.

[4] Loog, Marco. "Contrastive pessimistic likelihood estimation for semi-supervised classification." IEEE TPAMI (2016).