# PROBABILISTIC AND DIFFERENTIABLE PROGRAMMING V1: INTRODUCTION

#### Özgür L. Özçep Universität zu Lübeck Institut für Informationssysteme



#### Differentiable Programming and Probabilistic Programming for Machine Learning<sup>1)</sup>

1) Yes, this is a footnote on a slide, believe it or not. The three lines summarizing the topic of the course is the optimal outcome w.r.t my subjective measure - using a non-gradient optimization procedure starting from the original course name: Probabilistic Differential Programming -> Probabilistic and Differentiable Programming ->Differentiable and Probabilistic Programming



#### What this lecture V<sub>1</sub> is about

#### Agenda

- 2. Differentiable Programming and
- 3. Probabilistic Programming for
- 1. Machine Learning<sup>1)</sup>

Pointers to lectures in this fancy format.

1) Yes, this is a footnote on a slide, believe it or not. The three lines summarizing the topic of the course is the optimal outcome w.r.t my subjective measure - using a non-gradient optimization procedure starting from the original course name: Probabilistic Differential Programming -> Probabilistic and Differentiable Programming ->Differentiable and Probabilistic Programming



# MACHINE LEARNING



#### What We Mean by "Learning"

Machine learning (ML) is programming algorithms for

- optimizing a performance criterion
- using example (training) data
- by constructing general(izable) models
- that are good approximations of the data

**Role of Mathematics** 

- Building mathematical model
- core task is inference from a sample

- solve the optimization problem
- represent and evaluate the model for inference



#### Differentiable Programming

Machine learning (ML) is programming algorithms for

- optimizing a performance criterion
- using example (training) data
- by constructing general(izable) models
- that are good approximations of the data

#### **Role of Mathematics**

- Programming differentiable model
- core task is inference from a sample

- solve the optimization problem
- represent and evaluate the model for inference

#### Probabilistic Programming

Machine learning (ML) is programming algorithms for

- optimizing a performance criterion
- using example (training) data
- by constructing general(izable) models
- that are good approximations of the data

#### **Role of Mathematics**

- Programming probabilistic model
- core task is inference from a sample

- solve the optimization problem
- represent and evaluate the model for inference

## Types of learning (classically)

- Supervised Learning learn to predict an output for input vector after training with labelled data
- Unsupervised Learning discover a good internal representation of the input
- Reinforcement Learning
   learn to select an action to
   maximize the expectation of
   future rewards (payoff)





#### Subtypes of unsupervised I. (in Deep Learning context)

- Self-supervised (Self-taught) Learning - learn with targets induced by a prior on the unlabelled training data
- Semi-supervised Learning learn with few labelled examples and many unlabelled ones (same distribution for labelled & non-labelled data)





#### Generative vs. Discriminative/descriptive

- Many unsupervised and self-supervised models can be classed as 'generative models'.
  - Given unlabelled data X, a unsupervised generative model learns full joint probability distribution P(X,Y).
  - These are characterised by an ability to 'sample' the model to 'create' new data
- In contrast: Discriminative models learn P(Y|X) (which can be calculated in a generative model, too, using Bayes's rule but not vice versa)

(X: observations, data, Y: categories, classes, non-observed)



#### **Example Supervised Learning: Classification**

- Class C of a "family car"
  - Prediction: Is car *x* a family car?
  - Knowledge extraction: What do people expect from a family car?
- Output:

Positive (+) and negative (–) examples

• Input representation by two features:

*x*<sub>1</sub>: price, *x*<sub>2</sub>: engine power



#### Training set X





#### Class C





#### Hypothesis class H



#### **Example Supervised Learning: Regression**



Price of a used car		
<i>x</i> :	car attribute	
<i>y</i> :	price	
$\hat{y} = g(x \mid \theta)$ :	hypothesis	
g():	linear model	
g(x)	$) = w_1 x + w_0$	
heta:	parameters	
	(here w <sub>1</sub> , w <sub>2</sub> )	



#### **Example Supervised Learning: Regression**



Price of a used car x: car attribute y: price  $\hat{y} = g(x | \theta)$ : hypothesis g(): quadratic model  $g(x) = w_2 x^2 + w_1 x + w_0$   $\theta$ : parameters (here  $w_0, w_1, w_2$ )



#### **Example Supervised Learning: Regression**



Calculating the gradient  $\nabla E$ analytically NOT feasible for thousands of parameters

- > Differentiable programming



$$X = \{ (x^{t}, r^{t}) \}_{t=1}^{N} \quad r^{t} = f(x^{t}) \in \mathbb{R}$$

Mean squared error for general and linear hypothesis g  $E(g|X) = (1/N) \sum_{t=1}^{N} (g(x^{t}) - r^{t})^{2}$  $E(w_{1}, w_{0}|X) = (1/N) \sum_{t=1}^{N} (w_{1}x^{t} + w_{0}) - r^{t})^{2}$ 

#### Optimization:

$$\nabla E = \left(\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}\right) = (0, 0)$$
$$w_1 = \frac{\sum_t x^t r^t - \overline{xr}N}{\sum_t (x^t)^2 - N \ \overline{x}^2}$$

 $w_0 = \bar{r} - w_1 \, \bar{x}$ 

17

# DIFFERENTIABLE PROGRAMMING



#### What is Differentiable Programming (DP)?

"Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers. The important point is that people are now building a new kind of software by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization....It's really very much like a regular program, except it's parameterized, automatically differentiated, and trainable/optimizable....

(Part of a post of Yann Lecun, somewhere in Facebook, found at https://gist.github.com/halhenke/872708ccea42ee8cafd950c6c2069814)



",Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers. The important point is that people are now building a new kind of software by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization.....It's really very much like a regular program, except it's parameterized, automatically differentiated, and trainable/optimizable. ...

(Part of a post of Yann Lecun, somewhere in Facebook, found at https://gist.github.com/halhenke/872708ccea42ee8cafd950c6c2069814)



#### What is Deep Learning (DL)?



- Deep learning is based on function composition
  - Feedforward networks:  $\mathbf{y} = f(g(\mathbf{x}, \theta_g), \theta_f)$ Often with relatively simple functions (e.g.  $f(\mathbf{x}, \theta_f) = \sigma(\mathbf{x}^T \theta_f)$ )
  - Recurrent networks:

 $\mathbf{y}_{t} = f(\mathbf{y}_{t-1}, \mathbf{x}_{t}, \mathbf{\theta}) = f(f(\mathbf{y}_{t-2}, \mathbf{x}_{t-1}, \mathbf{\theta}), \mathbf{x}_{t}, \mathbf{\theta}) = \dots$ 

- In early days focus of DL on functions for classification
- Nowadays the functions are much more general in their inputs and outputs.



#### Network view of composed functions





V2

#### Deep networks







"Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers. The important point is that people are now building a new kind of software by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization....It's really very much like a regular program, except it's parameterized, automatically differentiated, and trainable (antimizable)

differentiated, and trainable/optimizable. ...

(Part of a post of Yann Lecun, somewhere in Facebook, found at https://gist.github.com/halhenke/872708ccea42ee8cafd950c6c2069814)



#### Gradient Descent

• Total loss

$$L = -\sum_{(x,y)\in D} l(g(x,\theta), y)$$

for some loss function I, dataset D and model g with parameters  $\theta$ 

- Define how many passes (epochs) over the data to make
- learning rate  $\eta$
- Gradient Descent: update  $\theta$  by gradient in each epoch  $\theta \leftarrow \theta - \eta \nabla_{\theta} L$









#### **Pros and Cons**

- Gradient Descent has good statistical properties (very low variance)
- But is very data inefficient (particularly when data has many similarities)
- Doesn't scale to effectively infinite data (e.g. with augmentation)



#### Stochastic Gradient Descent (SGD)

- Define loss function l, dataset D and model g with learnable parameters θ.
- Define how many passes over the data to make (each one known as an Epoch)
- Define a learning rate  $\eta$
- Stochastic Gradient Descent updates the parameters θ by moving them in the direction of the negative gradient with respect to the loss of a single item I by the learning rate η multiplied by the gradient:
- for each Epoch:
   for each (x,y)∈D:



#### Backprop: efficient implementation of gradient descent



Backpropagation idea

 Generate error signal that measures difference between predictions and target values



UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME (b) Backward pass

#### Deep networks



#### Problem: Many, many parameters, no structure

![](_page_28_Picture_3.jpeg)

#### What is Deep Learning?

V3

- *Deep learning* systems are neural network models similar to those popular in the '80s and '90s, with:
  - 1. some architectural and algorithmic innovations (e.g. many layers, ReLUs, dropout, LSTMs)
  - 2. vastly larger data sets (web-scale)
  - 3. vastly larger-scale compute resources (GPU, cloud)
  - 4. much better software tools (Theano, Torch, TensorFlow)
  - 5. vastly increased industry investment and media hype

![](_page_29_Picture_8.jpeg)

#### Deep Learning (ad 1.)

![](_page_30_Picture_1.jpeg)

![](_page_30_Figure_2.jpeg)

UNIVERSITAT ZU LÜBECK INSTITUT FÜR INFORMATI Adapted from https://www.xenonstack.com/blog/static/public/uploads/media/machine-learning-vs-deep-learning.phg FOCUS DAS LEBEN

![](_page_31_Picture_1.jpeg)

Example: Convolutional Neural Networks (CNN)

- More structure: local receptive fields
- Less parameters: weight tying, pooling

![](_page_31_Figure_5.jpeg)

http://deeplearning.stanford.edu/wiki/index.php/Feature\_extraction\_using\_convolution

![](_page_31_Picture_7.jpeg)

![](_page_32_Figure_0.jpeg)

https://www.asimovinstitute.org/neural-network-zoo/

![](_page_32_Picture_2.jpeg)

![](_page_33_Picture_1.jpeg)

"Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers. The important point is that people are now building a new kind of software by assembling networks of parameterized functional blocks and by training them from examples using some form of gradient-based optimization....It's really very much like a regular program, except it's parameterized, automatically differentiated, and trainable/optimizable....

(Part of a post of Yann Lecun, somewhere in Facebook, found at https://gist.github.com/halhenke/872708ccea42ee8cafd950c6c2069814)

![](_page_33_Picture_4.jpeg)

#### Automatic Differentiation (AD)

![](_page_34_Picture_1.jpeg)

• AD is a mix of

ERSITÄT ZU LÜBECK

- symbolic differentiation (SD) (rules s.a. chain rule, product rule)
- numerical differentiation (ND): use  $\frac{dy}{dx} \approx \frac{\Delta y}{\Delta x}$

$$\frac{d(f(x) \cdot g(x))}{dx} = \frac{d f(x)}{dx} g(x) + \frac{d g(x)}{dx} f(x) \quad \text{(Product rule)}$$
$$- h(x) := g(x) \cdot f(x)$$

- $-\frac{dh(x)}{dx}$  and *h* have two components in common
- This may also be the case for f.
- Symbollically calculating f won't profit from common parts of f and  $\frac{df(x)}{dx}$

![](_page_34_Picture_9.jpeg)

![](_page_35_Picture_0.jpeg)

![](_page_35_Figure_1.jpeg)

![](_page_35_Picture_2.jpeg)

# PROBABILITIES

![](_page_36_Picture_1.jpeg)

# *I* The third wave of differentiable programming

![](_page_37_Figure_1.jpeg)

Kristian Kersting - Sum-Product Networks: The Third Wave of Differentiable Programming

**//** 1)

1) Yes, a slide, quoting a slide

![](_page_37_Picture_5.jpeg)

#### Problems with deep (neural) networks (Ghahramani)

- Very data hungry (e.g. often millions of examples)
- Very compute-intensive to train and deploy (cloud GPU resources)
- Poor at representing uncertainty
- Easily fooled by adversarial examples
- Finicky to optimise: non-convex + choice of architecture, learning procedure, initialisation, etc, require expert knowledge and experimentation
- Uninterpretable black-boxes, lacking in trasparency, difficult to trust

![](_page_38_Picture_7.jpeg)

#### Bayes rule to rule them all ...

![](_page_39_Picture_1.jpeg)

 ...then inverse probability (i.e. Bayes rule) allows us to infer unknown quantities, adapt our models, make predictions and learn from data.

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)} = \frac{P(D|H) \cdot P(H)}{\sum_{h} P(D|h)P(h)}$$

H = hypothesis, model D = data, observation

IVERSITÄT ZU LÜBECK

INFORMATIONSSYSTEME

**Bayes Rule** 

## Probabilistic graphical models

# Encode efficiently full joint probabilities

- Directed graphs
   (Bayesian networks, Hidden Markov models ...)
- undirected graphs
   (Markov networks...)
- Mixed models
- Factor graphs

![](_page_40_Figure_6.jpeg)

**Requires Normalization** 

$$P(B = b, E = e, A = a, j, m) = \frac{1}{z}\phi_{JA}(a, j)\phi_{MA}(a, m)\phi_{AB}(a, b), \phi_{AE}(a, e)\phi_{B}(b)$$

$$Z = \sum_{x} \prod_{j} \phi_{j} \qquad Par$$

Partition function

![](_page_40_Picture_11.jpeg)

# PROBABILISTIC PROGRAMMING

![](_page_41_Picture_1.jpeg)

## Why then not stick to probabilities

- Problem 1: Probabilistic model development and the derivation of inference algorithms is time-consuming and error-prone.
- Problem 2: Exact (and approximate inference) hard due to normalization: partition function Z)
- Solution to 1
  - Develop Probabilistic Programming Languages for expressing probabilistic models as computer programs that generate data (i.e. simulators).
  - Derive Universal Inference Engines for these languages that do inference over program traces given observed data (Bayes rule on computer programs).

![](_page_42_Picture_6.jpeg)

V8

#### Comparison

![](_page_43_Picture_1.jpeg)

#### Probabilistic Programming Example

![](_page_44_Figure_1.jpeg)

#### Hidden markov model in Julia

![](_page_44_Picture_3.jpeg)

V8

# ADEQUATE DEEP STRUCTURES

![](_page_45_Picture_1.jpeg)

#### Problem 2 of probabilistic graphical models

- Exact (and even approxiate) inference not tractable for general probabilistic models (problem: normalization function Z).
- Restricting the models in expressivity is possible (thin junction trees and so on) but not desirable
- Find a better compromise of expressivity and feasibility: sum-product networks/probabilisitc boolean circuits

![](_page_46_Picture_4.jpeg)

#### Programming with Adequate Deep Structures

Machine learning (ML) is programming algorithms for

- optimizing a performance criterion
- using example (training) data
- by constructing general(izable) models
- that are good approximations of the data

**Role of Mathematics** 

- Building mathematical model
- core task is inference from a sample

- solve the optimization problem
- represent (expressively) and evaluate (feasibily) the model for inference

![](_page_48_Picture_1.jpeg)

$$P(X_1, ..., X_n) = \frac{1}{Z} \prod_{j} \phi_j(X_1, ..., X_n)$$

- Bottleneck: Summing out variables
- E.g.: Partition function

Sum of exponentially many products

$$Z = \sum_{x} \prod_{j} \phi_{j}$$

![](_page_48_Picture_7.jpeg)

#### **Alternative Representation**

![](_page_49_Picture_1.jpeg)

![](_page_49_Figure_2.jpeg)

$$P(X) = 0.4 \cdot X_1 \cdot X_2$$
  
+ 0.2 \cdot X\_1 \cdot \overline{X}\_2  
+ 0.1 \cdot \overline{X}\_1 \cdot X\_2  
+ 0.3 \cdot \overline{X}\_1 \cdot \overline{X}\_2

![](_page_49_Picture_4.jpeg)

#### Sum Out Variables

![](_page_50_Picture_1.jpeg)

			$e: X_1 = 1$
<i>X</i> <sub>1</sub>	<i>X</i> <sub>2</sub>	P(X)	$P(e) = 0.4 \cdot X_1 \cdot X_2$
1	1	0.4	$+0.2 \cdot X_1 \cdot \overline{X}_2$
1	0	0.2	$+ 0 1 \overline{V} V$
0	1	0.1	$+ 0.1 \cdot A_1 \cdot A_2$
0	0	0.3	$+0.3 \cdot X_1 \cdot X_2$

Set 
$$X_1 = 1, \overline{X}_1 = 0, X_2 = 1, \overline{X}_2 = 1$$

Easy: Set both indicators to 1

Easy: Partition function: Set all indicators to 1

![](_page_50_Picture_6.jpeg)

#### **Graphical Representation**

![](_page_51_Picture_1.jpeg)

![](_page_51_Figure_2.jpeg)

But in general may lead to exponentially large networks (e.g. parity). Solution: Make a deep dive (reuse computations) with Sum-Product networks

![](_page_51_Picture_4.jpeg)

![](_page_52_Picture_0.jpeg)

#### Sum-Product Networks (SPNs)

- Rooted DAG
- Nodes: Sum, product, input indicator
- Weights on edges from sum to children
- More general class: Probabilistic Boolean Circuits

![](_page_52_Figure_6.jpeg)

![](_page_52_Picture_7.jpeg)

![](_page_52_Figure_8.jpeg)

# **NEARLY THE END**

![](_page_53_Picture_1.jpeg)

#### Topic progress of course in short

- It's from discriminative to generative models
- It's from pure functions to algorithms to algorithms over semi-declarative structures (and some logic)
- It's from non-probabilities to probabilities (and some logic)

![](_page_54_Picture_4.jpeg)

Uhhh, a lecture with a hoepfully useful

## **APPENDIX**

![](_page_55_Picture_2.jpeg)

## Todays lecture is based on the following slides

- Jonathon Hare: Lecture 1,2 of course "COMP6248 Differentiable Programming (and some Deep Learning") <u>http://comp6248.ecs.soton.ac.uk/</u>
- Zoubin Ghahramani: Probabilistic Machine Learning and AI, Microsoft AI Summer School Cambridge 2017 <u>http://mlss.tuebingen.mpg.de/2017/speaker\_slides/Zoubin1.pdf</u>
- Hoifung Poon: Sum-Product Networks: A New Deep Architecture <u>https://www.microsoft.com/en-us/research/wp-</u> <u>content/uploads/2017/05/spn11.pdf</u>
- E. Alpaydin: Course on machine learning, introductory slides, <u>https://www.cmpe.boun.edu.tr/~ethem/i2ml2e/2e\_v1-0/i2ml2e-chap1-v1-0.pptx</u>
- I. Lorentzou: Introduction to Deep Learnin, <u>link</u>
- F. Wood: Probabilistic Programming, PPAML Summer School, Portland 2016, <u>link</u>

![](_page_56_Picture_7.jpeg)

#### Color Convention in this course

- Formulae, when occurring inline
- Newly introduced terminology and definitions
- Important results (observations, theorems) as well as emphasizing some aspects
- Examples are given with standard orange with possibly light orange frame
- Comments and notes
- Algorithms

![](_page_57_Picture_7.jpeg)

## Books for topics covered in this lecture (1)

- Nielsen: Neural Networks and Deep Learning. http://neuralnetworksanddeeplearning.com/
- Zhang et al.: Dive into Deep Learning https://d2l.ai/
- I. Goodfellow, Y. Bengio, and A. Courville. Deep Learning. MIT Press, 2016.
- D. Koller and N. Friedman. Probabilistic Graphical Models: Principles and Techniques - Adaptive Computation and Machine Learning. The MIT Press, 2009.
- L. D. Raedt, K. Kersting, and S. Natarajan. Statistical Relational Artificial Intelligence: Logic, Probability, and Computation. Morgan & Claypool Publishers, 2016.

#### Books for topics covered in this lecture (2)

- J.-W. van de Meent, B. Paige, H. Yang, and F. Wood. An Introduction to Probabilistic Programming. arXiv eprints, arXiv:1809.10756, Sept. 2018.
- U. Naumann. The Art of Differentiating Computer Programms. Siam, 2012.
- K. Murphy. Machine Learning: A Probabilistic Perspective. Adaptive Computation and Machine Learning series. MIT Press, 2012.
- S. J. Russell and P. Norvig. Artificial Intelligence A Modern Approach. Prentice Hall, 1995.

![](_page_59_Picture_5.jpeg)