# PROBABILISTIC AND DIFFERENTIABLE PROGRAMMING

## V13: ROUND-UP

Özgür L. Özçep

**Universität zu Lübeck**

**Institut für Informationssysteme**

# What this course was about

Differentiable Programming and

Probabilistic Programming for
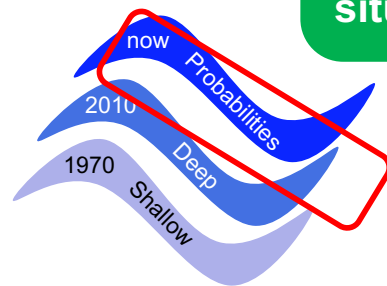
Machine Learning

# What this lecture $V_{13}$ is about

Nearly the same as V1, but even shorter

**„ The third wave of differentiable programming**

**Getting deep systems that know when they do not know and, hence, recognise new situations and adapt to them**

now
Probabilities
2010
Deep
1970
Shallow

**"** 1)

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

# Gradient Descent
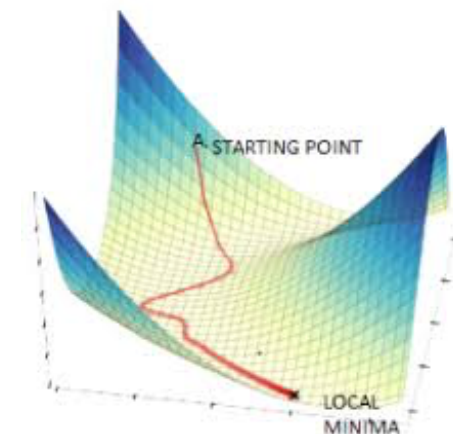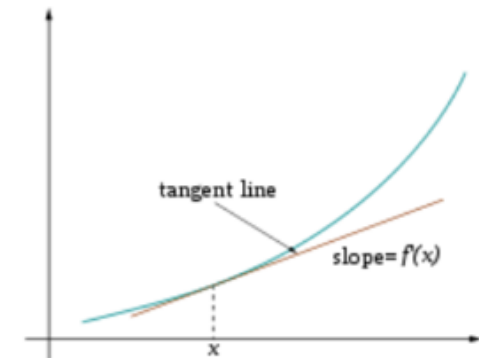
- Total loss

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

  for some loss function l, dataset D
  and model g with parameters θ



tangent line

slope=$f'(x)$

x

- Define how many passes (epochs) over the data to make

- learning rate η

- Gradient Descent: update θ by gradient in each epoch   $\theta \leftarrow \theta - \eta\nabla_\theta L$



STARTING POINT

LOCAL MINIMA

# Backprop: efficient implementation of gradient descent

| Sample labeled data (**batch**) | → | **Forward** it through the network, get predictions | → | **Back-propagate** the errors | → | **Update** the network weights |
|---|---|---|---|---|---|---|

Backpropagation idea

- Generate **error signal** that measures difference between predictions and target values

- Use error signal to change the weights and get more accurate predictions backwards
- Underlying mathematics: chain rule

**Chain rule (1-dim)**

$$\frac{dh}{dx} = \frac{df}{dg}\frac{dg}{dx}$$

( for $h(x) = f(g(x))$ )

(a) Forward pass



$w_1$
$\partial E/\partial w_1$

$x_1$     $y_1$     $w_5$
$w_2$          $\partial E/\partial w_5$
$\partial E/\partial w_2$     $\partial E/\partial y_1$

$w_3$     $y_3$     $E(y_3, t)$
$\partial E/\partial w_3$          $\partial E/\partial E$

$x_2$     $y_2$     $\partial E/\partial y_3$

$w_4$     $w_6$
$\partial E/\partial w_4$     $\partial E/\partial w_6$
$\partial E/\partial y_2$

(b) Backward pass

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

IM FOCUS DAS LEBEN

A mostly complete chart of
# Neural Networks
©2019 Fjodor van Veen & Stefan Leijnen   asimovinstitute.org

**Legend:**
- Input Cell
- Backfed Input Cell
- Noisy Input Cell
- Hidden Cell
- Probablistic Hidden Cell
- Spiking Hidden Cell
- Capsule Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- Gated Memory Cell
- Kernel
- Convolution or Pool

Perceptron (P)
Feed Forward (FF)
Radial Basis Network (RBF)
Deep Feed Forward (DFF)
Recurrent Neural Network (RNN)
Long / Short Term Memory (LSTM)
Gated Recurrent Unit (GRU)
Auto Encoder (AE)
Variational AE (VAE)
Denoising AE (DAE)
Sparse AE (SAE)
Markov Chain (MC)
Hopfield Network (HN)
Boltzmann Machine (BM)
Restricted BM (RBM)
Deep Belief Network (DBN)
Deep Convolutional Network (DCN)
Deconvolutional Network (DN)
Deep Convolutional Inverse Graphics Network (DCIGN)
Generative Adversarial Network (GAN)
Liquid State Machine (LSM)
Extreme Learning Machine (ELM)
Echo State Network (ESN)
Deep Residual Network (DRN)
Differentiable Neural Computer (DNC)
Neural Turing Machine (NTM)
Capsule Network (CN)
Kohonen Network (KN)
Attention Network (AN)

V3

V4

V6

(V8)

V3

V6

V4

(V3)

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

$$W_{V \times N}^{T} \times x_{cat} = v_{cat}$$

| 0.1 | 2.4 | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | 2.6 | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

Input layer

$x_{cat}$

V-dim

$W_{V \times N}^{T} \times x_{cat} = v_{cat}$

$+$

$W_{V \times N}^{T} \times x_{on} = v_{on}$

$x_{on}$

V-dim

Word embedding a la word2vec

Output layer

$\hat{v} = \dfrac{v_{cat} + v_{on}}{2}$

sat

V-dim

Hidden layer

N-dim



Knowledge graph embedding

Woman

$r$

Man

Aunt

$r$

Uncle

Queen

$r$

King

# Automatic Differentiation (AD)

- AD is a mix of

  - symbolic differentiation (SD) (rules s.a. chain rule, product rule)

  - numerical differentiation (ND): use $\frac{dy}{dx} \approx \frac{\Delta y}{\Delta x}$
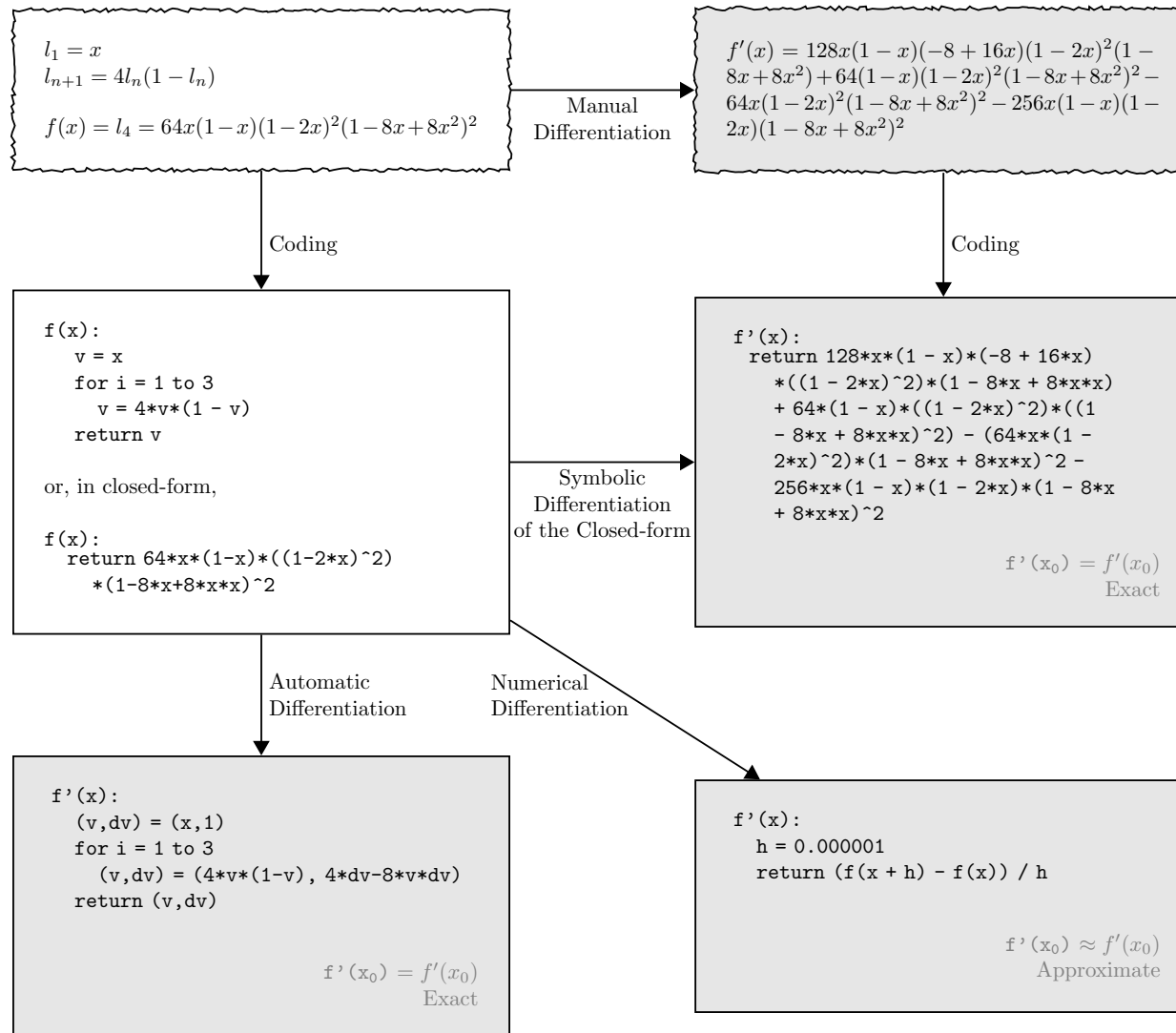
$$\frac{d(f(x) \cdot g(x))}{dx} = \frac{d\, f(x)}{dx} g(x) + \frac{d\, g(x)}{dx} f(x) \quad \text{(Product rule)}$$

  - $h(x) := g(x) \cdot f(x)$

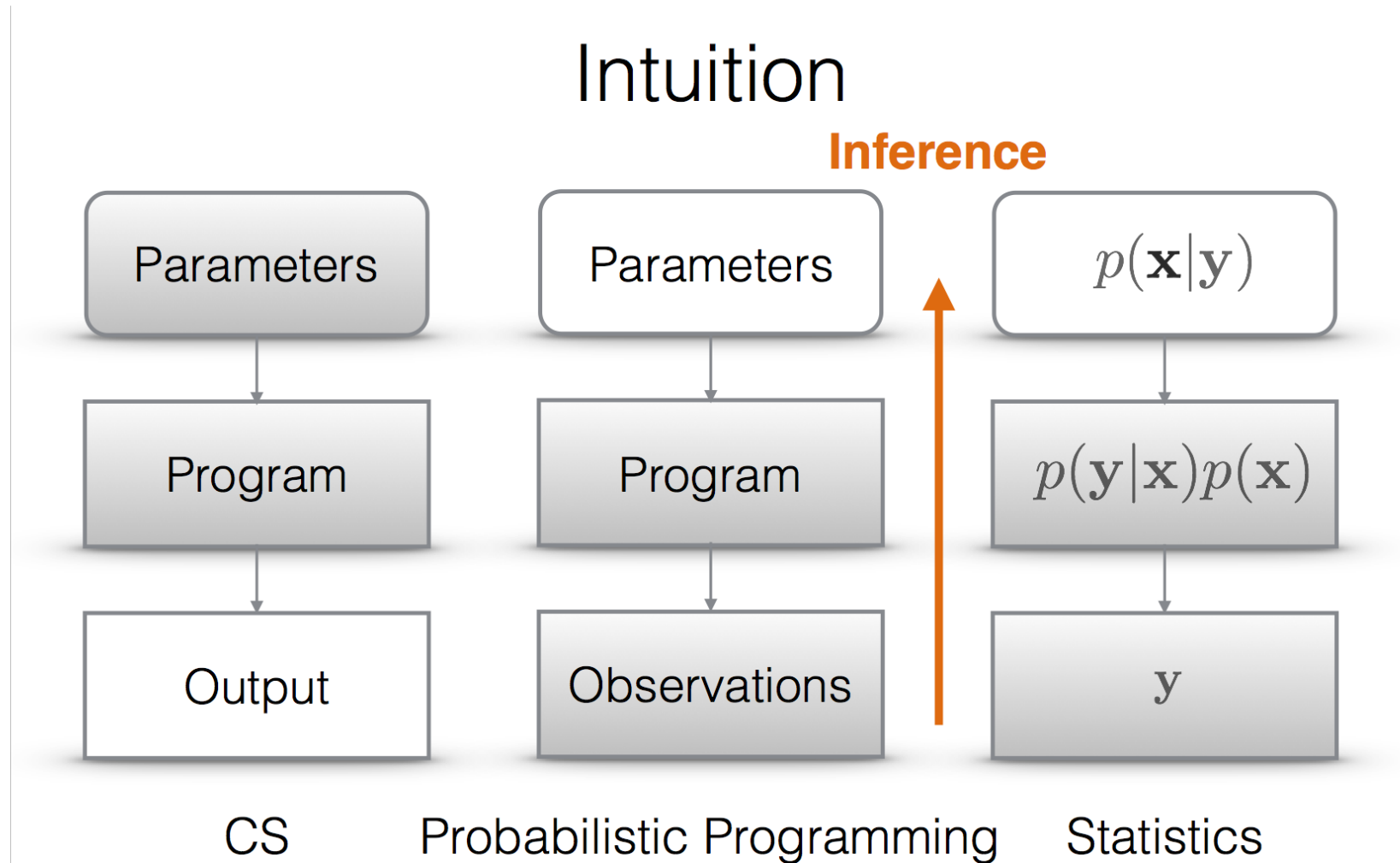  - $\frac{dh(x)}{dx}$ and $h$ have two components in common

  - This may also be the case for $f$.

  - Symbollicaly calculating $f$ won't profit from common parts of $f$ and $\frac{df(x)}{dx}$

$$l_1 = x$$
$$l_{n+1} = 4l_n(1 - l_n)$$

$$f(x) = l_4 = 64x(1-x)(1-2x)^2(1-8x+8x^2)^2$$

**Manual Differentiation** →

$$f'(x) = 128x(1-x)(-8+16x)(1-2x)^2(1-8x+8x^2)+64(1-x)(1-2x)^2(1-8x+8x^2)^2-64x(1-2x)^2(1-8x+8x^2)^2-256x(1-x)(1-2x)(1-8x+8x^2)^2$$

**Coding** ↓

```
f(x):
    v = x
    for i = 1 to 3
        v = 4*v*(1 - v)
    return v
```

or, in closed-form,

```
f(x):
    return 64*x*(1-x)*((1-2*x)^2)
        *(1-8*x+8*x*x)^2
```

**Coding** ↓

```
f'(x):
    return 128*x*(1 - x)*(-8 + 16*x)
        *((1 - 2*x)^2)*(1 - 8*x + 8*x*x)
        + 64*(1 - x)*((1 - 2*x)^2)*((1
        - 8*x + 8*x*x)^2) - (64*x*(1 -
        2*x)^2)*(1 - 8*x + 8*x*x)^2 -
        256*x*(1 - x)*(1 - 2*x)*(1 - 8*x
        + 8*x*x)^2
```

$$\texttt{f'(x}_0\texttt{)} = f'(x_0)$$
Exact

**Symbolic Differentiation of the Closed-form** →

**Automatic Differentiation** ↓

```
f'(x):
    (v,dv) = (x,1)
    for i = 1 to 3
        (v,dv) = (4*v*(1-v), 4*dv-8*v*dv)
    return (v,dv)
```

$$\texttt{f'(x}_0\texttt{)} = f'(x_0)$$
Exact

**Numerical Differentiation** ↘

```
f'(x):
    h = 0.000001
    return (f(x + h) - f(x)) / h
```

$$\texttt{f'(x}_0\texttt{)} \approx f'(x_0)$$
Approximate

# Comparison

## Intuition

**Inference**

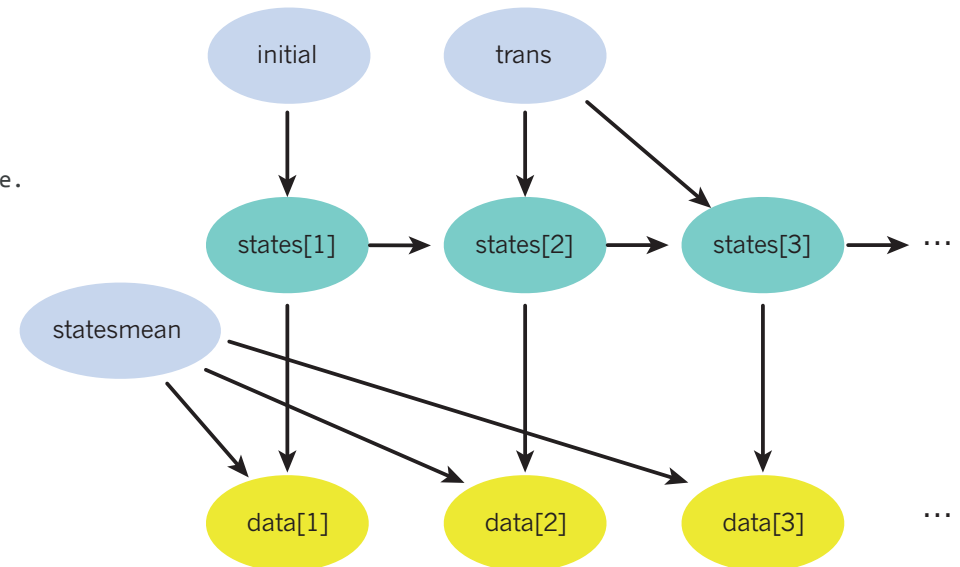| CS | Probabilistic Programming | Statistics |
|---|---|---|
| Parameters | Parameters | $p(\mathbf{x}\vert\mathbf{y})$ |
| Program | Program | $p(\mathbf{y}\vert\mathbf{x})p(\mathbf{x})$ |
| Output | Observations | $\mathbf{y}$ |

Ex: F. Wood: Probabilistic Programming, PPAML Summer School, Portland 2016
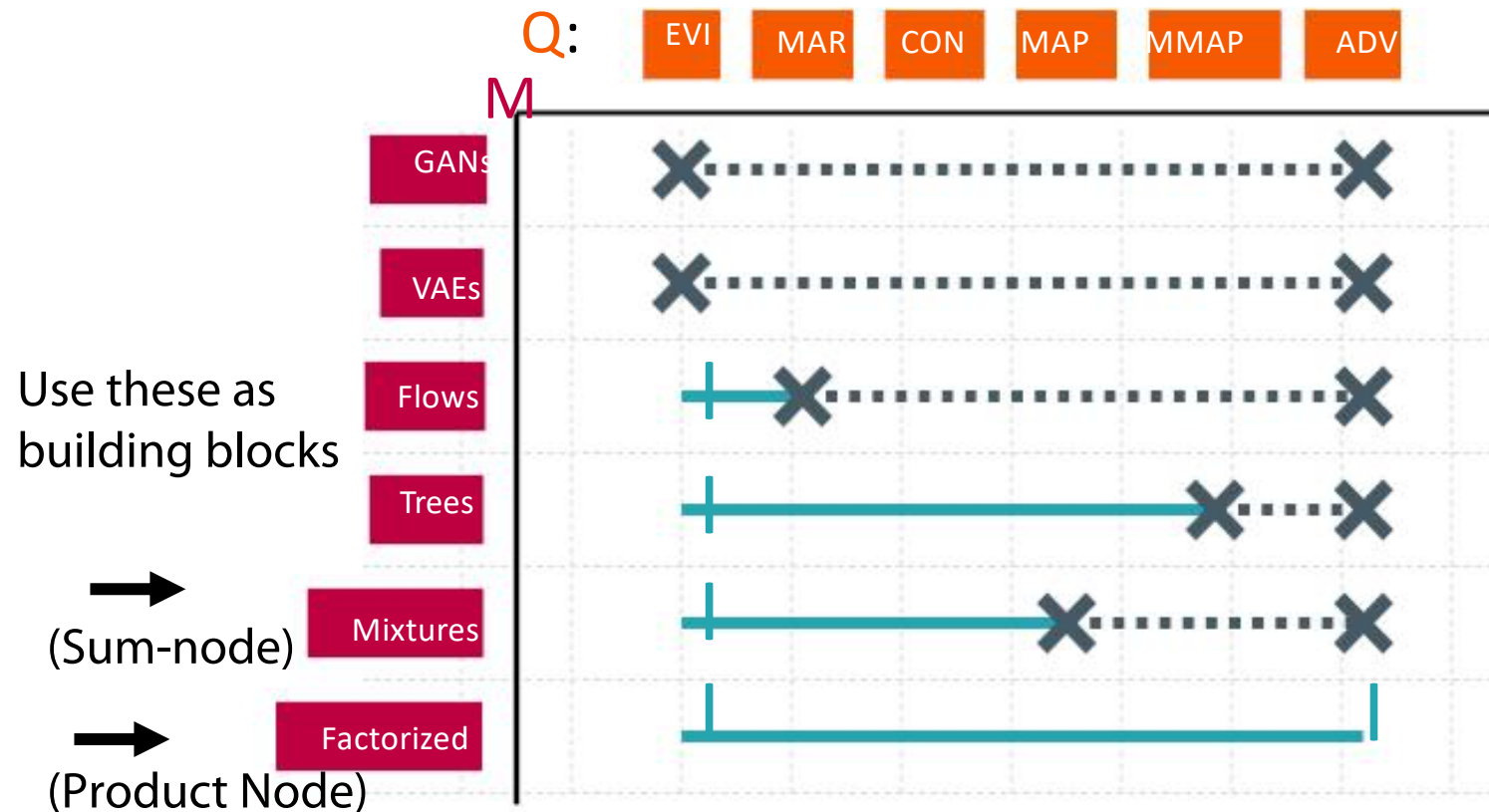
# Probabilistic Programming Example

```julia
statesmean = [-1, 1, 0]  # Emission parameters.
initial    = Categorical([1.0/3, 1.0/3, 1.0/3]) # Prob distr of state[1].
trans      = [Categorical([0.1, 0.5, 0.4]), Categorical([0.2, 0.2, 0.6]),
              Categorical([0.15, 0.15, 0.7])]   # Trans distr for each state.
data       = [Nil, 0.9, 0.8, 0.7, 0, -0.025, -5, -2, -0.1, 0, 0.13]

@model hmm begin # Define a model hmm.
 states = Array(Int, length(data))
 @assume(states[1] ~ initial)
 for i = 2:length(data)
   @assume(states[i] ~ trans[states[i-1]])
   @observe(data[i]  ~ Normal(statesmean[states[i]], 0.4))
 end
 @predict states
end
```
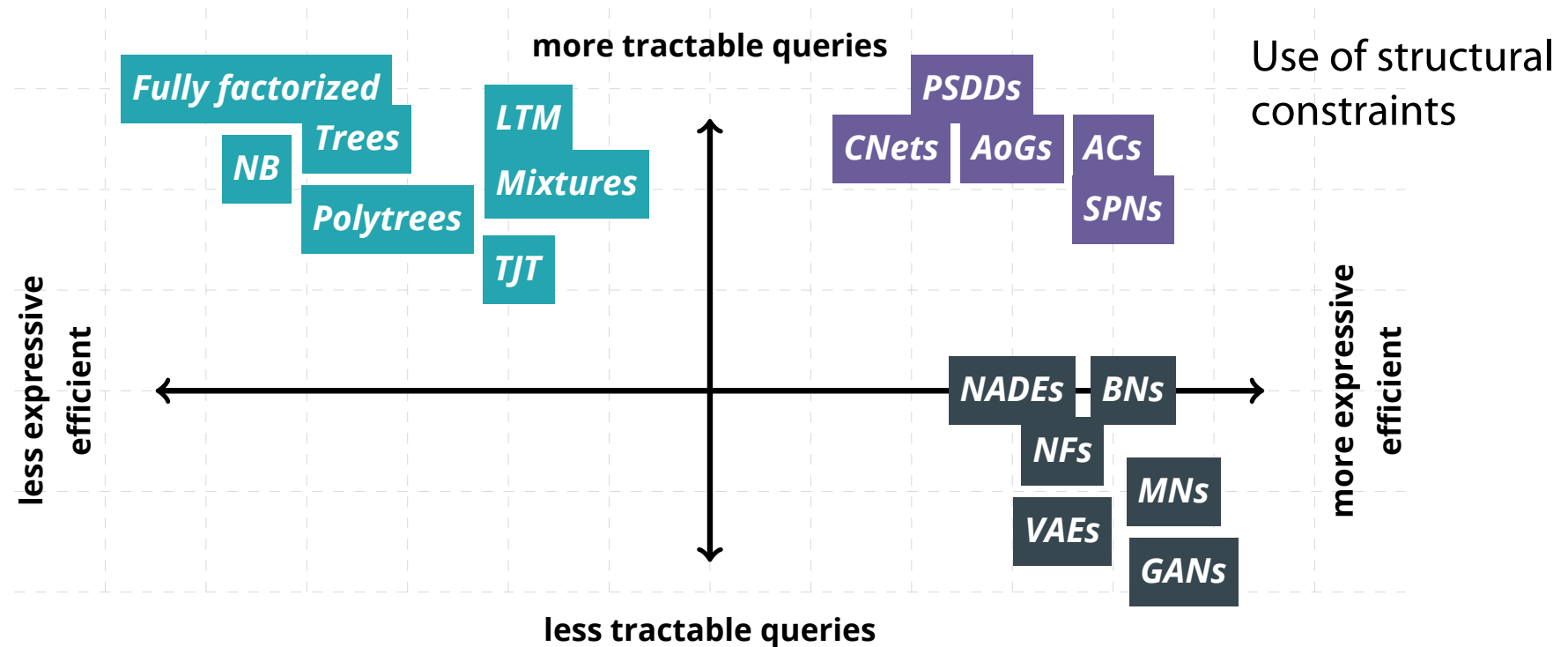
Hidden markov model  in Julia

UNIVERSITÄT ZU LÜBECK
INSTITUT FÜR INFORMATIONSSYSTEME

Q:  EVI  MAR  CON  MAP  MMAP  ADV

M

GANs

VAEs

Use these as
building blocks

Flows

Trees

→

(Sum-node)      Mixtures

→

(Product Node)      Factorized

„Eat the cake and have it"

*tractable bands*

more tractable queries

**Fully factorized**
**Trees**
**NB**
**Polytrees**
**LTM**
**Mixtures**
**TJT**

**PSDDs**
**CNets** | **AoGs** | **ACs**
**SPNs**

Use of structural constraints

less expressive efficient

more expressive efficient

**NADEs** | **BNs**
**NFs**
**MNs**
**VAEs**
**GANs**

less tractable queries

# *tractability vs expressive efficiency*

Can use efficiency also in learning

Thanks for your interest!