# Interpretable Machine Learning and Hyperbolic Geometry

#### **Ullrich Köthe**

Visual Learning Lab, Heidelberg University joint work with Lynton Ardizzone, Jakob Kruse, Stefan Radev, Jens Müller, Felix Draxler, Peter Sorrenson, Carsten Rother

Seminar Hyperbolic Geometry & Data Science, December 2020









#### Inductive vs. Transductive Reasoning

• Traditional science: inductive reasoning



- Problem: did not work well for machine learning until recently
  - available models were too weak to theoretically explain complex phenomena
  - ML community gave up on it around 1990



- Why trying to understand the world, when one is only required to make good predictions?
- Such black-box models are very general and work surprisingly well in practice
- Most ML successes rest on transductive approach



#### Inductive Solution of Inverted Triple Pendulum: Optimal Control based on Equations of Motion





#### Transductive Solution of Inverted Double Pendulum: Reinforcement Learning



UNIVERSITÄT PADERBORN

Swing-up and balancing of the double pendulum on a cart by reinforcement learning





#### The Brain's Solution of Complex Balancing: Transductive or Inductive?





6

#### **Generative Modelling**

- Problem: transductive models are hard to interpret, explain, validate
  - not acceptable for applications with critical consequences
  - nor acceptable if ML shall be used as a scientific research tool
  - the return to inductive modelling is now within reach!

#### ⇒ Generative modelling

- Transductive models: to predict y from x, learn the posterior distribution
- Generative models turn the problem around: learn the data generating process, i.e. likelihood p(x | y)
- Solve the original prediction task via Bayes theorem

$$p(y \mid x) = \frac{p(x \mid y) p(y)}{p(x)}$$

- This is much more difficult and requires **insight**, i.e. theoretical understanding

#### Feynman: "What I cannot create, I do not understand."

 $p(y \mid x)$ 

#### **Neural Networks**

- Standard definition of a neuron:  $\mathbb{R}^D \to \mathbb{R}$ :  $z' = \phi(w^T z + b)$ 
  - weights w, bias b, non-linear activation function  $\phi$ , e.g.  $\phi(t) = \tanh(t)$  or  $\phi(t) = \max(t, 0)$  (ReLU)
- Fully connected networks: combine neurons in parallel and in series (= layers)  $-f: \mathbf{z}_0 \in \mathbb{R}^D \to \mathbf{z}_L \in \mathbb{R}^M$ : recursive definition  $\mathbf{z}_l = \phi_l(\mathbf{W}_l^T \mathbf{z}_{l-1} + \mathbf{b}_l)$  with weight matrices  $\mathbf{W}_l$  and element-wise application of  $\phi_l$ 
  - deep learning: large number L of layers (10 ... 1000)
- Convolutional networks: weight matrices  $\boldsymbol{W}_l$  are Toeplitz matrices
  - Additionally: some layers execute "max pooling": reduce spatial resolution by  $\frac{1}{2}$  in every direction by taking the maximum of certain subsets of the input  $z_{l-1}$  (e.g. of 2x2 blocks of pixels)
  - Image classification networks: convolutional and pooling layers followed by fully connected layers
  - Image segmentation networks: only convolutional and pooling layers
- Most successful networks are transductive ("discriminative models")

#### **Deep Learning Example: Semantic Segmentation**

- Neural network classifies each pixel as "Street", "Sidewalk", "Person", "Car" etc.
  - Prerequisite for autonomous driving, but not as long as it is an unverifiable black-box





#### **Opening-up the Black-Box in Retrospect**

- Neural networks learn suitable features for the task lacksquare
  - better than traditional hand-crafted features
  - early layers detect local edges, later layers entire objects
  - feature visualization shows what each layer is looking for:



objects



Ohla et al.(2017)



#### **Opening-up the Black-Box in Retrospect**

- Emergent structure evolves in the feature spaces
- Can identify linear "concept axes"
- Extrapolation along a concept axis changes image in well-defined way
  - transform data  $x \rightarrow z$
  - $z' = z + \tau \cdot \text{concept}$
  - reconstruct  $z' \rightarrow x'$
- Active concepts in z serve as explanations



age axis

Upchurch et al.(2016)



beard axis

11



#### **Generative Modeling**

• Given features, learn to generate corresponding data



#### **Deep Generative Modelling**

- Generative Adversarial Networks (GANs):
  - two player game: generator learns to create realistic images discriminator learns to recognize synthetic images
  - converge to a state where synthetic and real images are indistinguishable



Training Data



#### **Deep Generative Modelling**

- Autoencoder
  - learn a lossy compression through a latent bottleneck
  - converge to maximum reconstruction quality by keeping only the relevant information in the bottleneck







#### **Deep Generative Modelling**

- Normalizing flows
  - learn an invertible transport map between target distribution in data space and a simple latent base distribution
  - converge to maximum reconstruction quality by keeping only the relevant information in the bottleneck



## Generative Modelling for Interpretable Machine Learning

- Idea: control training such that the latent space evolves interpretable structure
  - Happens emergently in classical training, but can we do better?
- Example 1: disentanglement
  - each latent variable has well-defined meaning: encodes a single isolated data property
  - complex transformations of the real data can be expressed by simple latent superpositions
  - analogy in physics:
     finding a good coordinate transformation is key to the analytic solution of many problems



# Disentanglement of MNIST using Invertible Neural Networks

 First 8 latent variables control global properties





Height



Width of top



Width of bottom

Sorrenson et al. "Disentanglement by Nonlinear ICA with General Incompressible-flow Networks (GIN)", ICLR 2020. 17



variability of handwriting in MNIST digits dataset



### Disentanglement of MNIST using Invertible Neural Networks

- First 8 latent variables control global properties
- Following 14 control local shape



variability of handwriting in MNIST digits dataset



Openness of lower loop



Extension to top right

Tail of 2

Sorrenson et al. "Disentanglement by Nonlinear ICA with General Incompressible-flow Networks (GIN)", ICLR 2020. 10

### Generative Modelling for Interpretable Machine Learning

True: bow tie

- Idea: control training such that the latent space evolves interpretable structure
  - Happens emergently in classical training, but can we do better?
- Example 2: Class Activation Mapping
  - Map latent interpretations back to data space
  - Heat maps: which input pixels are relevant to particular latent interpretations?
  - Visualize plausibility of interpretations and detect overfitting



 $Q_{\rm class}({\rm suit})$ 

 $Q_{\text{class}}(\text{sunglass})$ 

 $Q_{\text{class}}(\text{bow tie})$ 

### Generative Modelling for Interpretable Machine Learning

- Idea: control training such that the latent space evolves interpretable structure
  - Happens emergently in classical training, but can we do better?
- Example 3: Metric learning
  - Simple metric in latent space encodes complicated similarity relations in data space







#### **Two Perspectives on Network Interpretability**

- 1. External perspective: What black-boxes do we want for interpretable ML? Assume that networks can represent any desired mapping  $z = f_{\theta}(x)$ 
  - What types of mappings are especially suitable for interpretability?
  - Which properties should these mappings have (concept disentanglement, group equivariance, ...)
  - How can we quantify the uncertainty of such mappings (Bayesian posteriors, outlier detection, ...)
  - How can network mappings serve as a tool for scientific discovery?
- Internal perspective: How to achieve the desired mappings? Open-up the black-box
  - What network architectures realize interesting interpretable function families?
  - How can one guarantee convergence of the training?
  - Why do neural networks generalize so well?
  - How much trainings data is enough?



### Geometric Limitations of Current Generative Networks

- Most generative networks use Euclidean geometry in latent space
  - Euclidean metric for metric learning
  - Gaussian latent distributions
  - Ordinary gradients for training
- Since networks are continuous mappings, this translates to the data space
  - Cannot exactly model
    - Manifolds with arbitrary topology
    - Data supported on multiple disconnected compact regions
    - Non-Euclidean group symmetries
  - If latent dimension is high enough, accurate approximate embeddings are still possible, but
    - Representation does not have correct intrinsic dimension
    - No disentanglement into meaningful concepts
    - Representation artifacts (e.g. incorrect similarities, outliers not recognized, ...)

#### Interpretability requires that latent structure is compatible with intrinsic structure of the data

# Example: Embedding of the WordNet Hierarchy

- Famous paper by Nickel and Kiela (2017)
  - Hierarchies have intrinsically hyperbolic structure
  - ⇒ Embeddings in hyperbolic space are much more accurate

			Dimensionality					
			5	10	20	50	100	200
WORDNET Reconstruction	Euclidean	Rank MAP	3542.3 0.024	2286.9 0.059	1685.9 0.087	1281.7 0.140	1187.3 0.162	1157.3 0.168
	Translational	Rank MAP	205.9 0.517	179.4 0.503	95.3 0.563	92.8 0.566	92.7 0.562	91.0 0.565
	Poincaré	Rank MAP	4.9 0.823	4.02 0.851	3.84 0.855	3.98 0.86	3.9 0.857	3.83 0.87
WORDNET Link Pred.	Euclidean	Rank MAP	3311.1 0.024	2199.5 0.059	952.3 0.176	351.4 0.286	190.7 0.428	81.5 0.490
	Translational	Rank MAP	65.7 0.545	56.6 0.554	52.1 0.554	47.2 0.56	43.2 0.562	40.4 0.559
	Poincaré	Rank MAP	5.7 0.825	<b>4.3</b> 0.852	4.9 0.861	4.6 <b>0.863</b>	4.6 0.856	4.6 0.855

• 5-dimensional hyperbolic embedding more accurate than 200-dimensional Euclidean

# Example: Embedding of the WordNet Hierarchy

- Famous paper by Nickel and Kiela (2017)
  - Hierarchies have intrinsically hyperbolic structure
  - ⇒ Embeddings in hyperbolic space are much more accurate



- 5-dimensional hyperbolic embedding more accurate than 200-dimensional Euclidean
- 3-dimensional hyperbolic space suffices with better initialization (Clemens Fruböse)

#### **Example: Image Pyramid**

- Classical representation in image analysis ullet
  - Alternate blurring and down-sampling
  - Get gradually coarser and more abstract images
- Two kinds of relations:  $\bullet$ 
  - Spatial: context with in single level (e.g. grid neighbors) ⇒ Euclidean
  - Scale: objects with their constituents (e.g. pyramid neighbors) ⇒ Hyperbolic (?)
  - No satisfactory "scale embedding" yet allowing, for example latent interpolation of size
  - Scale-space theory failed Wavelets are not interpretable



### Summary

- Generative models with structured latent spaces ulletare a promising approach to interpretable ML
- Most networks use Euclidean latent geometry •
- Results in sub-optimal representations when ulletthe geometry/topology of the data differs
- ⇒ Networks with other latent geometries, group equivariance, Riemannian gradient descent are hot research topics
- $\Rightarrow$  I'm especially interested in hyperbolic geometry for "learning to abstract"
- $\Rightarrow$  What types of network mappings  $\mathbf{z} = f_{\theta}(\mathbf{x})$ do we want and how can we get them?





#### Invertible Neural Networks (INNs) with Coupling Layers

Powerful generative models: RealNVP ("non-volume preserving") [Dinh et al. 2017]

- Network is a sequence of *affine coupling layers*
- Each coupling layer splits its input  $x \in \mathbb{R}^{D}$  into two halves  $x_1, x_2 \in \mathbb{R}^{D/2}$
- Upper half is subjected to an affine transformation  $\Rightarrow$  outputs  $z_1, z_2 \in \mathbb{R}^{D/2}$
- Affine coefficients are computed by standard fully connected or convolutional networks  $s_2 \in \mathbb{R}^{D/2}_+$  and  $t_2 \in \mathbb{R}^{D/2}_+$  from the lower half's data

Forward computation:  $z_1 = x_1 \odot s_2(x_2) + t_2(x_2)$ ,  $z_2 = x_2$ Inverse computation:  $x_1 = (z_1 - t_2(z_2)) \oslash s_2(z_2)$ ,  $x_2 = z_2$  nested networks  $s_2$  and  $t_2$  are always executed in the same direction

Coupling layer  $x_1$   $z_1$   $z_1$   $x_1$   $z_1$   $z_1$   $z_1$   $z_1$   $z_2$   $z_2$  z

### Deep INNs

- Concatenate many coupling layers
- Alternate with orthogonal layers Q
  - ⇒ Active (upper lane) and passive (lower lane) dimensions change in each layer
  - Random permutations or projections are good enough, learning Q is not necessary
- Surprisingly powerful despite its simplicity
- Similar to autoencoder: forward mode = encoder, backward mode = decoder
  - Encoder and decoder are merged into a single network
  - Lossless encoding due to invertibility (no bottleneck)



# FAR

### **Conditional INNs (cINNs)**

#### Model conditional posterior $p(\mathbf{x} \mid \mathbf{y})$

• Receive observation y as additional conditioning input



Observation

### **cINN** Training

Estimated conditional density  $\hat{p}(x \mid y)$  expressed via 'reparameterization trick'

- Let the cINN represent the function  $z = f_{\theta}(x; y)$
- Train cINN such that  $p_z(z) \approx \mathcal{N}(0, \mathbb{I})$
- Then  $\hat{p}(x | y) \approx p^*(x | y)$  is defined by the change-of-variables formula

$$\hat{p}(x \mid y) = p_z(z = f_\theta(x; y)) \left| \det\left(\frac{\partial f_\theta(x; y)}{\partial x}\right) \right|$$

Can be trained with maximum likelihood loss:

$$\hat{\theta} = \arg \max_{\theta} \sum_{i \in \mathcal{TS}} \hat{p}(x_i \mid y_i)$$
  
=  $\arg \min_{\theta} \sum_{i \in \mathcal{TS}} \left( \frac{\|f_{\theta}(x_i; y_i)\|_2^2}{2} - \sum_{l=1}^L \log(\|s_l(x_i; y_i\|_1)) \right)$ 

Since  $f_{\theta}$  is invertible (given  $\hat{y}$ ), we get a generative model for free:

$$x \sim \hat{p}(x \mid \hat{y}) \iff z \sim \mathcal{N}(0, \mathbb{I}) \text{ and } x = f_{\hat{\theta}}^{-1}(z; \hat{y})$$

#### **Simulation-Based Inference**

- Classical simplification: reduce posterior to point estimate
  - Regularization: disambiguate inverse when forward process is surjective (not information preserving)
  - Maximum a-posteriori (MAP) inference: find only mode of posterior
  - $\Rightarrow$  No idea of solution diversity and uncertainty
- Standard solution: approximate Bayesian computation (ABC)
  - Define a distance  $d(y_{obs}, y_{sim})$  between observed and simulated data
    - for m=1,..., M:

Sample hidden parameters  $\mathbf{x}^{(m)} \sim p(\mathbf{x})$  from prior

Run the simulation  $y_{sim}^{(m)} = g(x^{(m)}; \xi)$ 

Keep  $\boldsymbol{x}^{(m)}$  if  $d\left(\boldsymbol{y}_{\text{obs}}, \boldsymbol{y}_{\text{sim}}^{(m)}\right) < \epsilon$ , reject otherwise

- Return the set of "surviving"  $x^{(m)}$  as an approximate sample from  $p(x | y_{obs})$  $\Rightarrow$  Very slow: high rejection rate, because small  $\epsilon$  are needed for accurate results
- Case-based inference
  - Efficient sampling methods (MCMC, SMC, ...) with good proposal distribution have low rejection rates
  - $\Rightarrow$  Learn a good proposal distribution for each observed data set  $y_{
    m obs}$ , i.e. on a per case basis
  - $\Rightarrow$  Still expensive if many different  $y_{
    m obs}$  must be evaluated



#### **Simulation-Based Inference**

- Classical simplification: reduce posterior to point estimate
  - Regularization: disambiguate inverse when forward process is surjective (not information preserving)
  - Maximum a-posteriori (MAP) inference: find only mode of posterior
  - $\Rightarrow$  No idea of solution diversity and uncertainty
- Standard solution: approximate Bayesian computation (ABC)
  - Define a distance  $d(y_{obs}, y_{sim})$  between observed and simulated data

#### We can do much better using invertible neural networks!

- Return the set of "surviving"  $x^{(m)}$  as an approximate sample from  $p(x | y_{obs})$  $\Rightarrow$  Very slow: high rejection rate, because small  $\epsilon$  are needed for accurate results
- Case-based inference
  - Efficient sampling methods (MCMC, SMC, ...) with good proposal distribution have low rejection rates
  - $\Rightarrow$  Learn a good proposal distribution for each observed data set  $y_{obs}$ , i.e. on a per case basis
  - $\Rightarrow$  Still expensive if many different  $y_{
    m obs}$  must be evaluated



#### **Multiple Possibilities for INNs**

Autoregressive Models

Chain rule decomposition:

 $p(x_1, ..., x_D) = \prod_i p_i(x_i \mid x_{< i})$ triangular reparameterization:  $\forall i: x_i = f_i(z_i, x_{< i})$  monoton.



inverse direction inefficient ⇒ use two complementary nets

example: parallel WaveNet

iResNets (invertible residual networks)





z = x + f(x)

is invertible when  $\|f(x)\|_{1}$ 

 $||f(x)||_{\text{Lipshitz}} < 1$ 

inverse direction is reasonably efficient (fixpoint or Newton iterations)

example: Residual Flow Net

RealNVP



$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} x_1 \cdot s_2(x_2) + t_2(x_2) \\ x_2 \end{bmatrix}$$

inverse is equally efficient:

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} (z_1 - t_2(z_2))/s(z_2) \\ z_2 \end{bmatrix}$$

example: GLOW