

Constraints in Representation Learning

CAS AML 2024

Dr. Mykhailo Vladymyrov

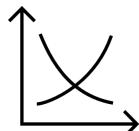
Data Science Lab
University of Bern

Representation learning

Representation Learning is an approach in machine learning where algorithms create *new data representations* from the data, that are easier to process.



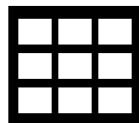
- Melting rate
- Temperature
- ...
- Poetic description



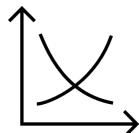
- Prediction
- Anomaly discovery
- ...
- Understanding system dynamics

Representation learning

Representation Learning is an approach in machine learning where algorithms create *new data representations* from the data, that are easier to process.



- Melting rate
- Temperature
- ...
- Poetic description



- Prediction
- Anomaly discovery
- ...
- Understanding system dynamics

These representations can be designed for interpretability, reveal hidden features, or be used for transfer learning.

The data – defines “quality” of the representation

The model learns patterns in your data.

If the main pattern of your data is variety of noise features – this will be dominant feature learned.

If your dataset is too small – the model would learn features specific to your samples, rather than to the distribution they are sampled from.

this is often an issue
in medical data (i.e. noise)
in MRI machines)

The algorithm – defines properties of the representation

- ① The representation will have correlation with my original task
Weights on the last layer can be thought of as prototypes of the classes of my original classification task.
- ② Representation is not biased by the original classes as data is unlabelled.
(somehow unsup. learning)
there is statistical constraints here.
- ③ the model could be trained / arranged to de-noise samples
and i.e. define if a new sample less likelyhood to be real or not.
example: model trained on generation of butterflies, then used to determine the likelihood of a new butterfly picture to be real or false⁵.

The algorithm – defines properties of the representation

Let's discuss what does this imply for:

- Classification model (1) \rightarrow bias toward specific task
- AE (and VAE, VQVAE) (2) \rightarrow not biased tow. specific task but statistically constrained.
- DDM (3) \rightarrow extend applicability of a trained model
For a new unrelated Task

The task – defines desired property of the representation

The different algorithms do for sure different things, but other elements play a role as well, for example the kind of loss that we use. The loss used has a strong impact on the performance of the model vs the task.

Perceptive Loss:

Called perceptive as usually it's used for image domain.

Compare representation of the model at the lower levels with the latent representation.

Practically we end up comparing activation function of different representation (of the feature maps). We don't compare the output, but the activation of the values in a different model.

Representation
produce by a low
level of a \neq model
(pretrained)

→ a low level representation will
have only high frequency components
(i.e. edges)

→ to compare with
the outcome of the AE.

} this approach force the AE
to pay more attention to the
high frequency components of
image

(butterfly
with
aws)

The task – defines desired property of the representation

e.g. consider the difference:

- AE with L2 vs L1 vs Perceptive Loss
- AE vs VAE vs VQVAE \rightarrow imply some stat prop. of the latent space.
- DDM in original vs latent space

↓

Depending on what we do in the latent space there'll be as well a difference. DDM will give a gradient of like likelihood.

There should be a difference between a property being unlikely and a property being unknown by the model. It should be possible to disentangle & encode and define which property are unknown and which are unlikely.

perceptive loss, formula example:

$$\underset{\theta}{\operatorname{argmin}} \mid X - AE(X| \theta) \mid$$

We have in principle ∞ way to represent the data and consequently infinite way to train the AE

space. Close together sample in the latent space, far samples will be forced apart.

So how can we steer it?

How can we make the representations
the way we want them?

We have:

- Constraints
- Transformations
- Conditions

Constraints

Constraints

→ Are defined by the model:

- Architecture
- Loss function
- Activation function

- Size of representation
- Statistical properties – e.g. using variational inference (VAE handle static of the sample more than the sample itself).
- Space in which representations “Live”:
 - Vector
 - Matrix
 - Vector field → diffusion (the model learn likelihood of feature in the vector space).

Spacial dimension of the latent space are continuous.

The representation can be different depending on what are we using. can be a vector of size 2 if we rotate in a 2D vector, but can be also something different. If we increase the dimensions, those new dimensions will have some constraints in the vector space that depends from the model.

Transformations

Input transformations

the very augmentation that we apply
will have an impact on the training
of the model

Augmentations:

- Translation/rotation
- Noise
- Derivate explicit representation

⇒ help the model to generalize
but at the same time tells the model
that those properties are not important.
(teach actually the model to ignore
them).

Representation transformations

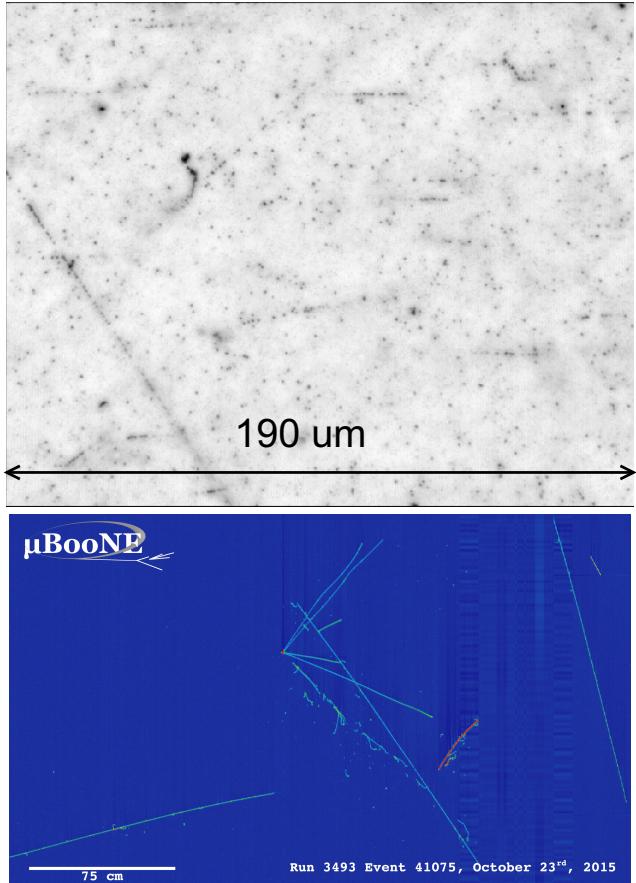
- Sampling from modelled distribution
- Equivariance
- Invariance

We somehow already saw: in the VAE we sample in a gaussian distribution to smooth out to a gaussian distib.

Some transformation in the latent space will not change the reconstruction.

EXAMPLE of APPLICATION :

Tracking: finding lines



Algorithms:

- Leverage small part of available information
- Tedious adaptation to new conditions

Supervised Machine Learning

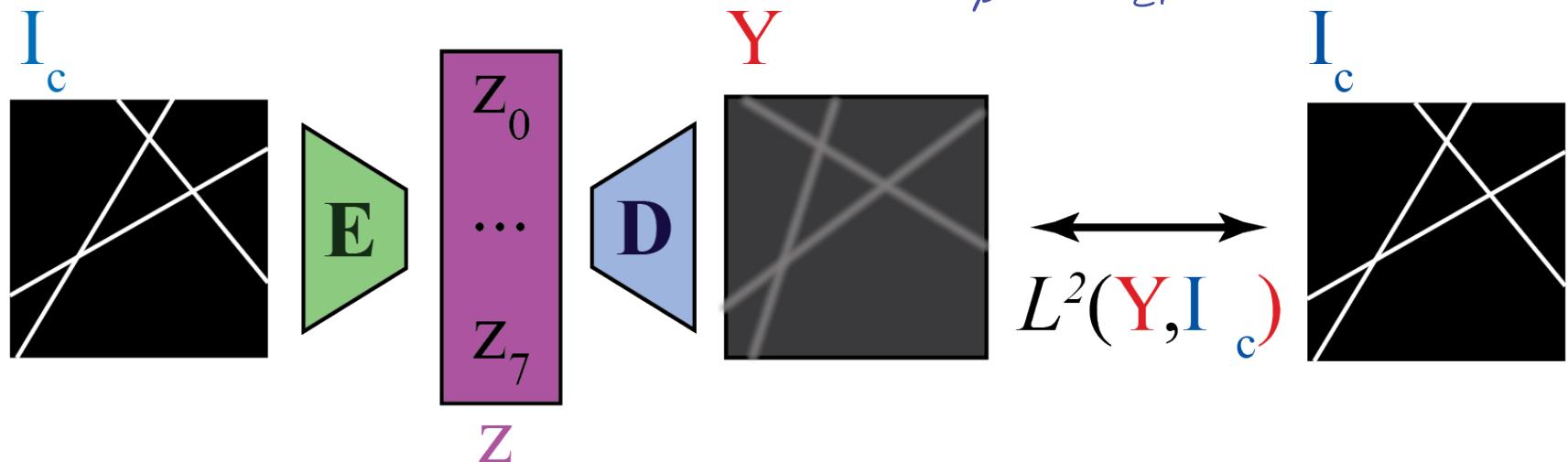
- Leverage all low-level information
- Requires labeled data / training on simulations

Unsupervised Machine Learning

- Leverage all low-level information
- No / small amount of labeled data

Proposed model

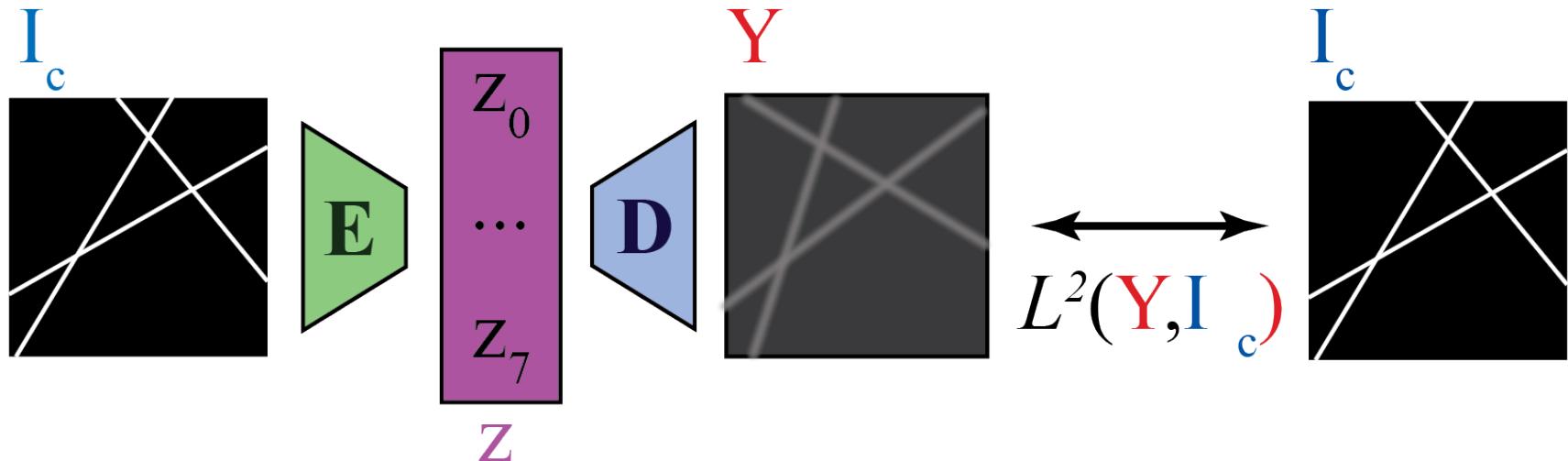
Autoencoder input image with some lines, decoder reconstruct from latent space. We want this (parcels/lines) to be a *part cluster*



- parameters that shall correspond to the i^{th} track

Proposed model

line as parameters are a constraint.



– parameters that *shall* correspond to the i^{th} track

We need to 1. disentangle different tracks; 2. disentangle track parameters

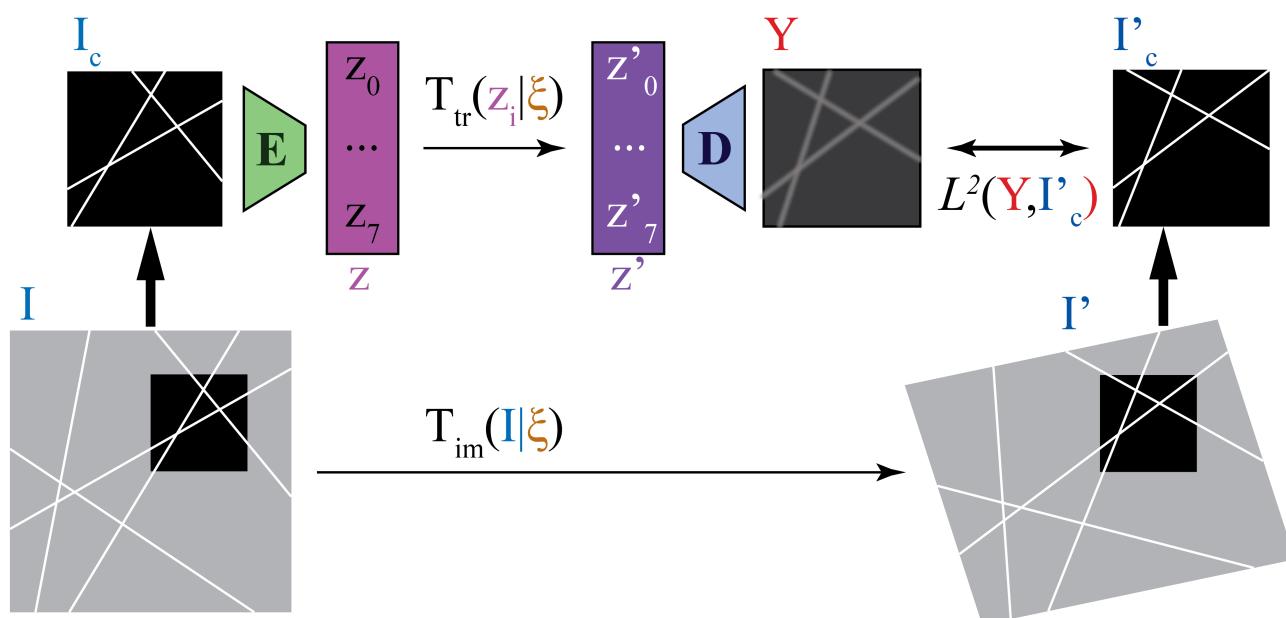
Equivariance constraint

, – geometrical transformations in the image and the latent representation space, parametrized by the same parameter set

Demand equivariance of encoder and decoder:

->

Equivariance constraint



$$E, D = \underset{E, D}{\operatorname{arg\,min}} \# L^2(Y, I'_c) = \underset{E, D}{\operatorname{arg\,min}} \underset{\xi}{\mathbb{E}} [L^2(D(T_{tr}(E(C(I)) | \xi)), C(T_{im}(I | \xi)))]$$

There is relation b/w if the image rotates, the latent space should rotate too

$I' = T_{im}(I | \xi)$, $Z_t^i = T_{tr}(z_t | \xi)$ – geometrical transformations in the image and the latent representation space, parametrized by the same parameter set

Demand equivariance of encoder and decoder:

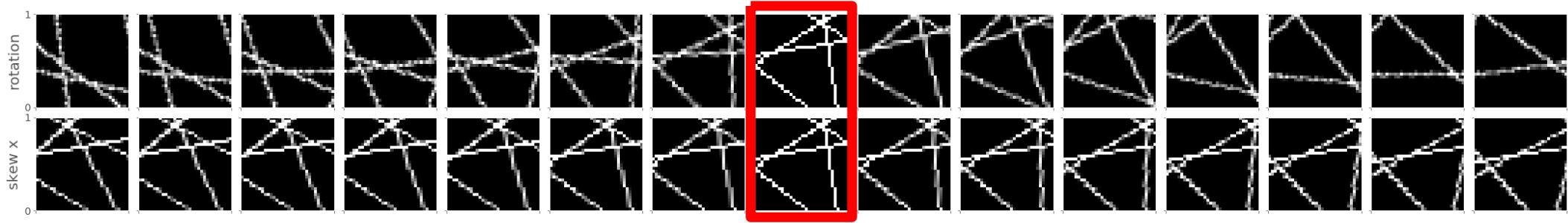
$$\begin{aligned} T_{tr}(E(I) | \xi) &= E(T_{im}(I | \xi)); \forall \xi \\ D(T_{tr}(z | \xi)) &= T_{im}(D(z) | \xi); \forall \xi \end{aligned}$$

predicted angle of the lines

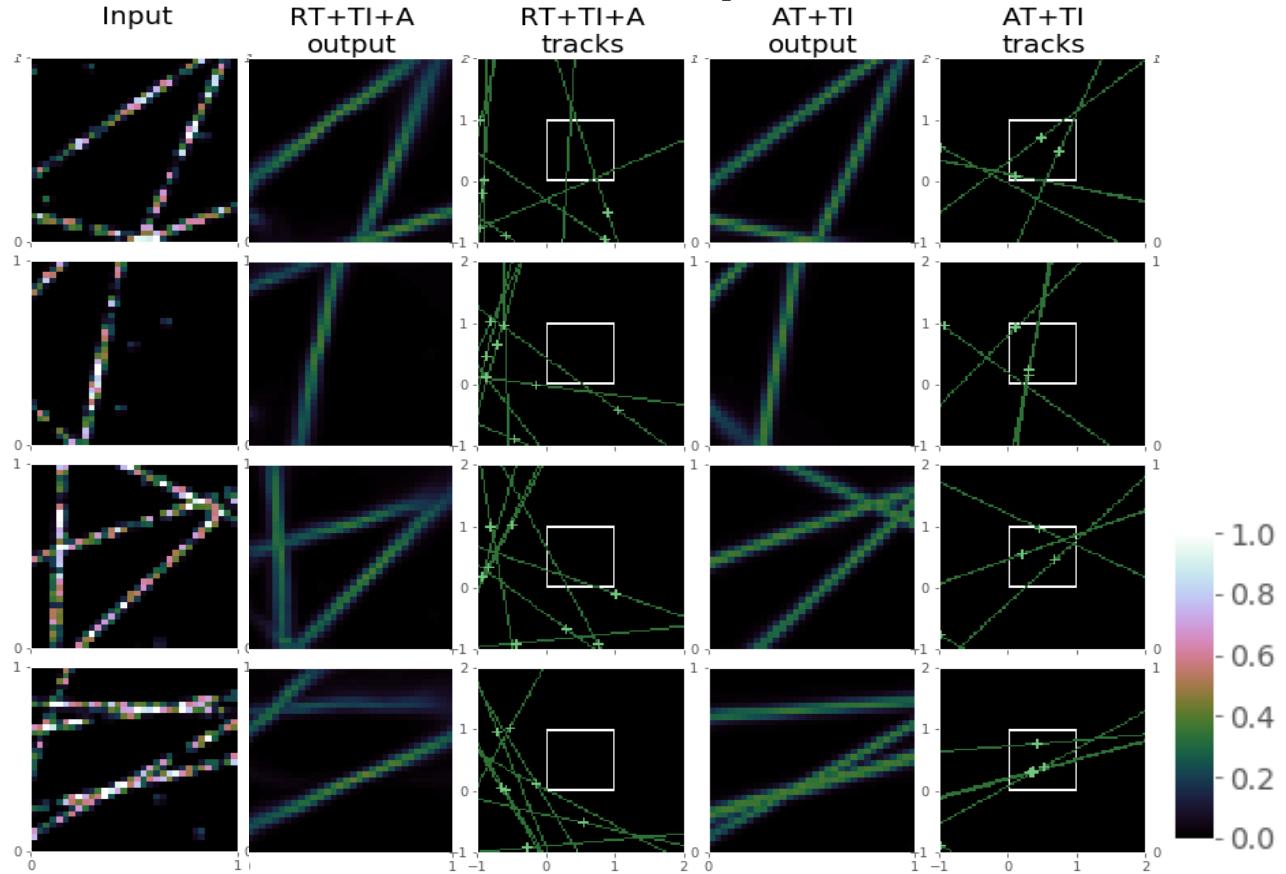
Representation transformations

By applying some coherent transformation we can obtain meaningful representation also without labelling data

Affine Transformations: rotation, translation, skew, scale



Autoencoder performance



Imposing equivariance constraint on the autoencoder under geometrical transformations in the image and latent representation domains enables the model to “discover” the existence of multiple lines in the presented images in a fully unsupervised manner.

*M. Vladymrov, A. Ariga,
Novel tracking approach based on fully-
unsupervised disentanglement of the geometrical
factors of variation*

Conditions

Conditions

- Distribution properties
e.g. for VAE we had:

Conditions

- Sparseness – like Lasso regression
Impose by augmenting the loss function:

We can apply weight
less wangs

$$L = L_{\text{task}} + \lambda \cdot \sum_i (z_i - 1)$$

Project on a unit sphere

$$z_i \rightarrow \frac{z_i}{|z_i|}$$

We can apply some
normalization

Conditions

- Integrity:

e.g. terms sum up to 100%:

project on a unit sphere:

Conditions

- Physical properties

Any property in fact can be added similar way (known as Physics informed Neural Networks). Given equation

$$f(z) = 0$$
$$L = L_0 + \lambda \cdot f(z)$$

Coefficient balances scale of the function value range

We can, the same way, apply also any physical property.

Coeff. balances balance the order of my equation.

Conditions

- Physical properties

Any property in fact can be added similar way (known as Physics informed Neural Networks). Given equation

$$f(z) = 0$$
$$L = L_0 + \lambda \cdot f(z)$$

Coefficient balances scale of the function value range

- Static – features of the output.
- Dynamic – gradient of features

PINNs

Consider searching for a PDE solution:

$$u_t + uu_x - (0.01/\pi)u_{xx} = 0, \quad x \in [-1, 1], \quad t \in [0, 1], \\ u(0, x) = -\sin(\pi x), \\ u(t, -1) = u(t, 1) = 0.$$

$$f := u_t + uu_x - (0.01/\pi)u_{xx},$$

Train should minimize the total loss.

total loss will have the loss on the boundary condition + the loss of making $\text{eq}=0$

differential equation

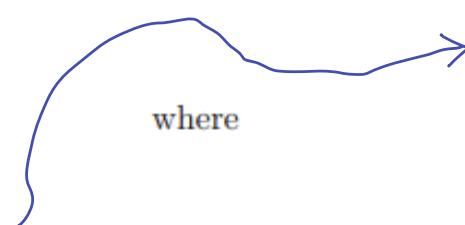
```
def u(t, x):
    u = neural_net(tf.concat([t,x],1), weights, biases)
    return u
```

\mathcal{L} with t and x the model predict u

Correspondingly, the *physics informed neural network* $f(t, x)$ takes the form

```
def f(t, x):
    u = u(t, x)
    u_t = tf.gradients(u, t)[0]
    u_x = tf.gradients(u, x)[0]
    u_xx = tf.gradients(u_x, x)[0]
    f = u_t + u*u_x - (0.01/tf.pi)*u_xx
    return f
```

The shared parameters between the neural networks $u(t, x)$ and $f(t, x)$ can be learned by minimizing the mean squared error loss



$$MSE = MSE_u + MSE_f, \quad \text{Total loss} \quad (4)$$

where

$$MSE_u = \frac{1}{N_u} \sum_{i=1}^{N_u} |u(t_u^i, x_u^i) - u^i|^2, \quad \text{boundary condition}$$

$$MSE_f = \frac{1}{N_f} \sum_{i=1}^{N_f} |f(t_f^i, x_f^i)|^2. \quad \text{use of the function } f.$$

PINNs

Instead of the NN to solve the eq, authors try to predict the hamiltonian of the system.

$$\mathcal{L}_{HNN} = \left\| \frac{\partial \mathcal{H}_\theta}{\partial p} - \frac{\partial q}{\partial t} \right\|_2 + \left\| \frac{\partial \mathcal{H}_\theta}{\partial q} + \frac{\partial p}{\partial t} \right\|_2$$

The way you shape your model will determine how your model learns.

Is you can in principle do everything

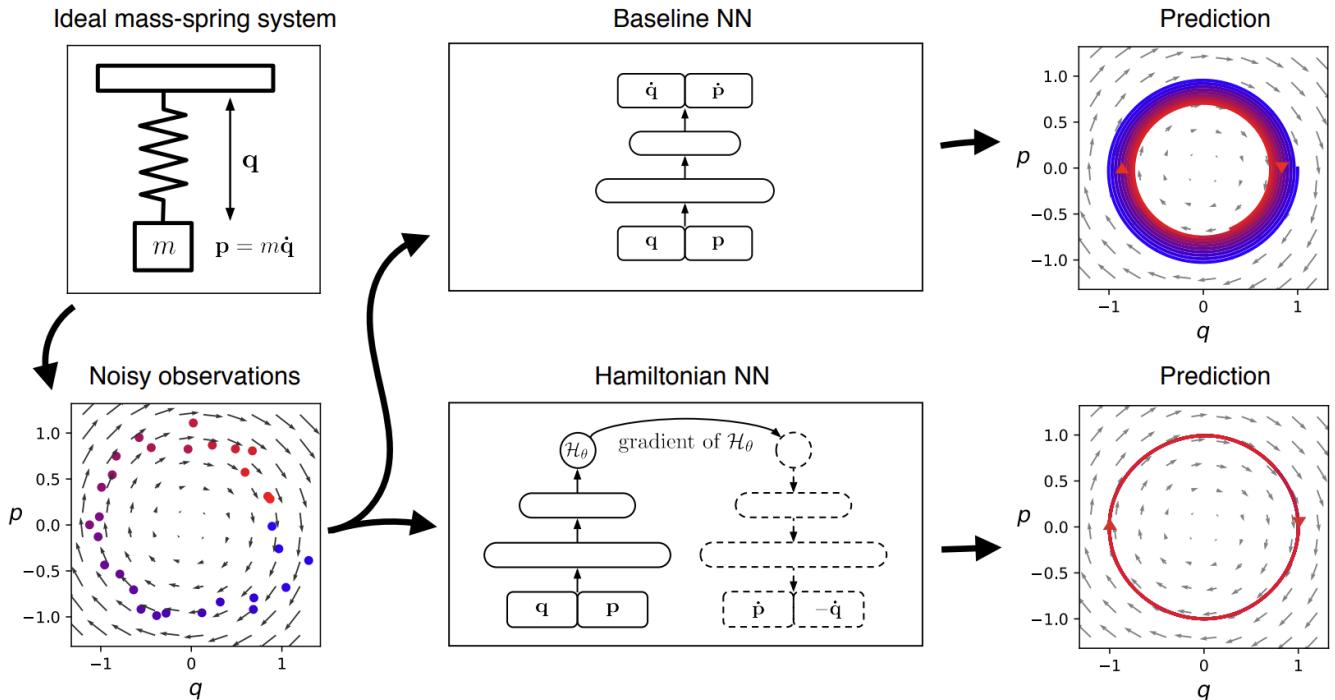


Figure 1: Learning the Hamiltonian of a mass-spring system. The variables q and p correspond to position and momentum coordinates. As there is no friction, the baseline's inner spiral is due to model errors. By comparison, the Hamiltonian Neural Network learns to *exactly* conserve a quantity that is analogous to total energy.

Further readings

- PINNs:

- <https://arxiv.org/pdf/1906.01563>
- <https://towardsdatascience.com/physics-informed-neural-networks-pinns-an-intuitive-guide-fff138069563>
- <https://arxiv.org/pdf/1711.10561>

links

- NN important tools:

- Normalizations: <https://www.pinecone.io/learn/batch-layer-normalization/>
- Gradient propagation: <https://mrinalwalia.medium.com/understanding-the-vanishing-gradient-problem-in-deep-learning-c648a4f16b05>
- Stop gradient

Tutorials, topics
to look into

- Interpreting learned features in CNNs:

- <https://distill.pub/2017/feature-visualization/>
- <https://distill.pub/2018/building-blocks/>
- <https://distill.pub/2018/differentiable-parameterizations/>

some more reading on interpretation of
the latent space & meaningful
representation

Workshop

1h + 15min

Plan the project

- What task are you attempting?
- Which representations would be useful?
- Which tools would you use?
- Which constraints are imposed by the model?
- Which conditions would you apply?

PROJECT :

- do smth about represent. learning
- ideally adding some constraints

- which could be the properties of my repres?
- which simmetry?
- how would my latent space look like?