

Autoencoders



DIG-UM Training course

"Advanced Deep learning"

28.11. - 1.12.22, Meinerzhagen

Autoencoders

Michael Krämer

(RWTH Aachen University)

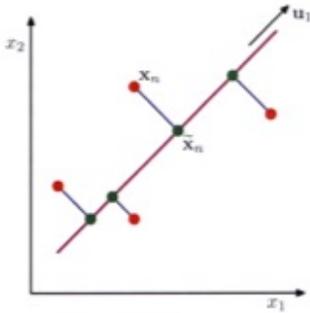
Content:

- Unsupervised learning
- Autoencoders
- Variational autoencoders
- Anomaly detection

Principle component analysis

→ data reduction

→ decorrelation of features



← minimize sum of squares of projection errors

from book by Bishop

Data: $\{\bar{x}^{(1)}, \dots, \bar{x}^{(m)}\}$ with $\bar{x}^{(i)} \in \mathbb{R}^n$

want to compress data $\bar{x}^{(i)} \rightarrow \bar{z}^{(i)}$ with $\bar{z}^{(i)} \in \mathbb{R}^l$ and $l < n$
but minimal loss of information

Try to find encoding and decoding functions:

$$f(\bar{x}) = \bar{z} \quad \text{and} \quad \bar{x} \approx g(\bar{z}) = g(f(\bar{x}))$$

Simple choice: $g(\bar{z}) = D\bar{z}$; $D \in \mathbb{R}^{n \times l}$

Constraint: $D_{:,i}$ (columns of D) are orthogonal to each other; $D_{:,i}$ have unit norm.

Proceed in 2 steps: 1) find optimal c^*
2) find optimal D^*

$$1) \text{ Find } \vec{z}^* = \operatorname{argmin}_{\vec{z}} \|\vec{x} - g(\vec{z})\|_2$$

$$\text{or equivalently } \vec{z}^* = \operatorname{argmin}_{\vec{z}} \|\vec{x} - g(\vec{z})\|_2^2$$

$$= \operatorname{argmin}_{\vec{z}} (\vec{x} - g(\vec{z}))^T (\vec{x} - g(\vec{z}))$$

$$\rightarrow (\vec{x} - g(\vec{z}))^T (\vec{x} - g(\vec{z})) = \underbrace{\vec{x}^T \vec{x}}_{\text{not relevant for finding } \vec{z}^*} - \vec{x}^T g(\vec{z}) - \underbrace{g(\vec{z})^T \vec{x}}_{= (g(\vec{z})^T \vec{x})^T = \vec{x}^T g(\vec{z})} + g(\vec{z})^T g(\vec{z})$$

$$\rightarrow \text{consider } \vec{z}^* = \operatorname{argmin}_{\vec{z}} (-2 \vec{x}^T g(\vec{z}) + g(\vec{z})^T g(\vec{z}))$$

Recall that $g(\vec{z}) = D\vec{z}$:

$$\vec{z}^* = \operatorname{argmin}_{\vec{z}} (-2 \vec{x}^T D\vec{z} + \underbrace{\vec{z}^T D^T D \vec{z}}_{= \mathbb{1}_e})$$

$$= \operatorname{argmin}_{\vec{z}} (-2 \vec{x}^T D\vec{z} + \vec{z}^T \vec{z})$$

Minimize a function $f(\vec{z})$:

$$\vec{\nabla}_{\vec{z}} (-2 \vec{x}^T D\vec{z} + \vec{z}^T \vec{z}) = 0$$

$$\left[\frac{\partial}{\partial c_i} \left(-2 \sum_{j=1}^n x_j D_{j,i} c_i + \sum_j c_j^2 \right) = -2 \sum_{j=1}^n x_j D_{j,i} \frac{\partial c_j}{\partial c_i} + 2 \sum_j c_j \frac{\partial c_j}{\partial c_i} \right]$$

$$= -2 \sum_j x_j D_{j,i} + 2 c_i$$

$$\Rightarrow \vec{\nabla}_{\vec{z}} (-2 \vec{x}^T D\vec{z} + \vec{z}^T \vec{z}) = -2 D^T \vec{x} + 2 \vec{z} \perp$$

$$\Rightarrow \vec{z} = D^T \vec{x}$$

\rightarrow complete reconstruction: $r(\vec{x}) = g(f(\vec{x})) = D\vec{z} = DD^T \vec{x}$

2) Find optimal matrix D : data point $i, i=1, \dots, m$

$$D^* = \operatorname{argmin}_D \sqrt{\sum_{i,j} (x_j^{(i)} - r(\vec{x}^{(i)})_j)^2} \quad \text{subject to } D^T D = \mathbb{1}_e$$

Frobenius norm j-component of vectors $\in \mathbb{R}^n$

Consider simplest case: $l=1$

$$\rightarrow D \in \mathbb{R}^{n \times l} \rightarrow \vec{d} \in \mathbb{R}^n$$

In the simplest case $l=1$ we consider

$$\begin{aligned} \vec{d}^* &= \arg \min_{\vec{d}} \sum_i \|\vec{x}^{(i)} - \underbrace{\vec{d} \vec{d}^T \vec{x}^{(i)}}\|_2^2 \quad \text{subject to } \|\vec{d}\|_2 = 1 \\ &= \vec{d}^T \vec{x}^{(i)} \vec{d} = (\vec{d}^T \vec{x}^{(i)})^T \vec{d} \\ &= \vec{x}^{(i)T} \vec{d} \vec{d} \end{aligned}$$

Introduce some helpful notation: design matrix $X \in \mathbb{R}^{m \times n}$:

$$X_{i,:} = \vec{x}^{(i)} \quad \text{or explicitly} \quad X = \begin{bmatrix} X_1^{(1)} & \dots & X_n^{(1)} \\ X_1^{(2)} & \dots & X_n^{(2)} \\ \vdots & & \vdots \\ X_1^{(m)} & \dots & X_n^{(m)} \end{bmatrix}$$

$$\text{Then } \vec{d}^* = \arg \min_{\vec{d}} \|X - X \vec{d} \vec{d}^T\|_F^2 \quad \text{subject to } \vec{d}^T \vec{d} = 1$$

Let's work out $\|X - X \vec{d} \vec{d}^T\|_F^2$

$$\begin{aligned} \|X - X \vec{d} \vec{d}^T\|_F^2 &= \text{Tr} \left(\underbrace{(X - X \vec{d} \vec{d}^T)^T (X - X \vec{d} \vec{d}^T)} \right) \\ &= \underbrace{X^T X}_{\text{not relevant}} - X^T X \vec{d} \vec{d}^T - \vec{d} \vec{d}^T X^T X + \vec{d} \vec{d}^T X^T X \vec{d} \vec{d}^T \\ &= -\text{Tr}(X^T X \vec{d} \vec{d}^T) - \underbrace{\text{Tr}(\vec{d} \vec{d}^T X^T X)} + \underbrace{\text{Tr}(\vec{d} \vec{d}^T X^T X \vec{d} \vec{d}^T)} \\ &= -2 \text{Tr}(X^T X \vec{d} \vec{d}^T) + \underbrace{\text{Tr}(X^T X \vec{d} \vec{d}^T \vec{d} \vec{d}^T)}_{=1} \\ &= -2 \text{Tr}(X^T X \vec{d} \vec{d}^T) + \text{Tr}(X^T X \vec{d} \vec{d}^T) \end{aligned}$$

$$\rightarrow \vec{d}^* = \arg \min_{\vec{d}} (-\text{Tr}(X^T X \vec{d} \vec{d}^T)) \quad \text{subject to } \vec{d}^T \vec{d} = 1$$

$$= \arg \max_{\vec{d}} (+\text{Tr}(X^T X \vec{d} \vec{d}^T)) \quad \text{subject to } \vec{d}^T \vec{d} = 1$$

$$= \arg \max_{\vec{d}} \text{Tr}(\underbrace{\vec{d}^T X^T X \vec{d}}_{\text{this is a scalar}}) \quad \text{--- " ---}$$

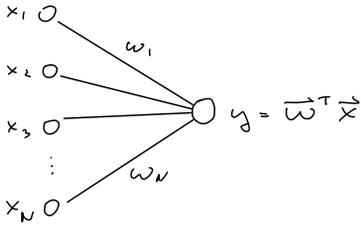
$$= \arg \max_{\vec{d}} (\vec{d}^T X^T X \vec{d}) \quad \text{subject to } \vec{d}^T \vec{d} = 1$$

\rightarrow problem of maximizing quadratic form $\vec{d}^T A \vec{d}$, where $A = X^T X$ is symmetric \rightarrow optimal \vec{d} is given by eigenvector of $X^T X$ with largest eigenvalue.

Unsupervised learning

Machine learning for unlabeled data: find patterns, detect outliers, generate data, ...

let us start with a simple example: Oja's rule



consider data distribution $p(\bar{x})$ and i.i.d. drawn samples $\bar{x} = (x_1, \dots, x_N)^T$

Weight update: $\bar{w} \rightarrow \bar{w} + \delta \bar{w}$ with $\delta \bar{w} = \epsilon y \bar{x}$
learning rate
 \rightarrow output $|y|$ becomes the larger the more often input feature occurs in data distribution.

can write this as a differential equation:

$$\tau \frac{d\bar{w}}{dt} = y \bar{x} \quad \text{since} \quad \tau \frac{\bar{w}(t+\Delta t) - \bar{w}(t)}{\Delta t} = y \bar{x}$$
$$\Rightarrow \bar{w}(t+\Delta t) = \bar{w}(t) + \frac{\Delta t}{\tau} y \bar{x}$$

Is this stable?

$$\frac{d\|\bar{w}\|^2}{dt} = 2\bar{w}^T \frac{d\bar{w}}{dt} = 2\bar{w}^T \frac{1}{\tau} y \bar{x} = \frac{2}{\tau} y^2 > 0$$

$\rightarrow \bar{w}$ grows without bound

\rightarrow add some weight decay:

$$\bar{w} \rightarrow \bar{w} + \epsilon y (\bar{x} - \alpha y \bar{w}) \quad \text{or} \quad \tau \frac{d\bar{w}}{dt} = y \bar{x} - \alpha y^2 \bar{w} \quad [\text{Oja 82}]$$

Is Oja's rule stable?

$$\begin{aligned}\frac{d\|\vec{w}\|^2}{dt} &= 2\vec{w}^T \frac{d\vec{w}}{dt} = \frac{2}{\tau} \vec{w}^T (y\vec{x} - \alpha y^2 \vec{w}) = \frac{2}{\tau} (y^2 - \alpha y^2 \vec{w}^T \vec{w}) \\ &= \frac{2}{\tau} y^2 (1 - \alpha \|\vec{w}\|^2)\end{aligned}$$

→ converges to $\|\vec{w}\|^2 = 1/\alpha$

What does Hebbian learning do?

Consider average over input data:

$$\tau \frac{d\vec{w}}{dt} = \langle y\vec{x} \rangle_{\vec{x}} = \langle \vec{x} \vec{x}^T \vec{w} \rangle_{\vec{x}} = \langle \vec{x} \vec{x}^T \rangle_{\vec{x}} \vec{w} = \underbrace{C}_{\text{data covariance matrix}}$$

for input data with zero mean

Let us write \vec{w} in terms of the eigenvectors of C :

$$\vec{w}(t) = \sum_i c_i(t) \vec{u}_i \quad [C \text{ is real \& symmetric} \rightarrow \vec{u}_i \text{ form basis}]$$

$$\Rightarrow \tau \sum_i \frac{dc_i(t)}{dt} \vec{u}_i = \sum_i c_i(t) \lambda_i \vec{u}_i$$

$$\Rightarrow \tau \frac{dc_i(t)}{dt} = c_i(t) \lambda_i \Rightarrow c_i(t) = c_i(0) \exp(\lambda_i t / \tau)$$

$$\Rightarrow \vec{w}(t) = \sum_i c_i(t) \vec{u}_i = \sum_i c_i(0) \exp(\lambda_i t / \tau) \vec{u}_i$$

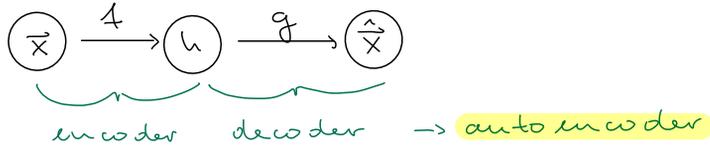
For $t \gg \tau$, the term corresponding to the largest eigenvalue (say λ_1) dominates:

$$\vec{w}(t) \propto \vec{u}_1(t) \quad [\text{or } \vec{w}(t) \propto \vec{u}_1 / \sqrt{\alpha} \text{ for Oja's rule}]$$

⇒ Hebbian learning implements the principle component analysis!

let us now try to implement **feature learning on unlabeled data with non-linear neural networks.**

Start with data i.i.d. drawn from $p(\mathcal{X})$, and set up neural network with a hidden layer to map data to itself:

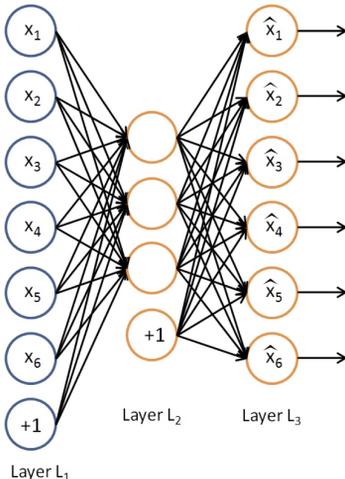


learning the identity mapping is not very useful.
instead construct

undercomplete autoencoders:

- h has lower dimension than \mathcal{X}
- f or g has low capacity (e.g. linear g)
- discard information in h

or autoencoders with regularization.



undercomplete AE

($\dim h < \dim \mathcal{X}$)

want $\hat{\mathcal{X}} \approx \mathcal{X}$

→ AE needs to learn
compressed representation
of data

→ cf. PCA

→ extract salient features
of data!

Regularised autoencoders, e.g. sparse autoencoders or denoising autoencoders:

Sparse autoencoders:

$f(x) = h$ and $g(h) = \hat{x}$; minimize loss function:

$$\underset{f, g}{\operatorname{arg\,min}} \sum_{x \in \text{data}} \sum_{m=1}^M \mathcal{L}(x_m, g(f(x)_m)) + \underbrace{\Omega(h)}_{\text{sparsity penalty on hidden layer}}$$

can derive $\Omega(h)$ through modeling joint distribution

$$p_{\text{model}}(\hat{x}, t) = p(\hat{x}|t) p(t) \quad \curvearrowright \text{prior over latent variables}$$

Maximize log-likelihood:

$$\ln p_{\text{model}}(\hat{x}, t) = \ln p_{\text{model}}(\hat{x}|t) + \ln p_{\text{model}}(t)$$

A Laplace prior on latent variables:

$$p_{\text{model}}(h_i) = \frac{\lambda}{2} e^{-\lambda |h_i|}$$

leads to penalty: $\Omega(t) = \lambda \sum_i |h_i|$

Alternative way to enforce sparsity:

keep average activation in each hidden layer node

small: $\hat{h}_j = \frac{1}{M} \sum_{m=1}^M h_j(x^{(m)}) = h_0$ (e.g. = 0.05) (Andrew Ng)

\leadsto add penalty term:

$$\Omega(h) = \sum_{j=1}^S \left(h_0 \ln(h_0 / \hat{h}_j) + (1 - h_0) \ln\left(\frac{1 - h_0}{1 - \hat{h}_j}\right) \right)$$

\curvearrowright sum over hidden units

$$= \sum_{j=1}^S \text{KL}(h_0 \| \hat{h}_j)$$

Denoising autoencoders: add noise to input and

$$\text{minimize } \mathcal{L}(x, g(f(\tilde{x})))$$

copy of \tilde{x} corrupted by noise

→ autoencoder must learn to undo corruption of data.

Contractive autoencoder: regularize $h = f(x)$ by penalizing

derivatives of f : $\mathcal{R}(h) = \lambda \left\| \frac{\partial f}{\partial x} \right\|_F^2$ $\left\| \cdot \right\|_F$ is Frobenius norm

→ forces autoencoder to learn function that does not change much when \tilde{x} changes slightly.

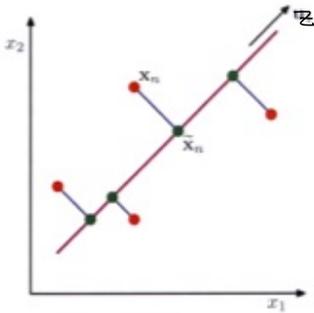
Autoencoders belong to the general class of

latent variable models

LVMs map between observation space $\vec{x} \in \mathbb{R}^D$ and latent space $\vec{z} \in \mathbb{R}^Q$: $f_\theta: \vec{x} \rightarrow \vec{z}$; $g_\phi: \vec{z} \rightarrow \hat{\vec{x}}$

- one latent variable gets associated with each data point in training set: $\vec{x}^{(m)} \rightarrow \vec{z}^{(m)}$
- latent vectors are smaller than observations, $Q < D \Rightarrow$ compression
- models can be linear or non-linear, deterministic or stochastic

Example of a linear, deterministic model:
principal component analysis:

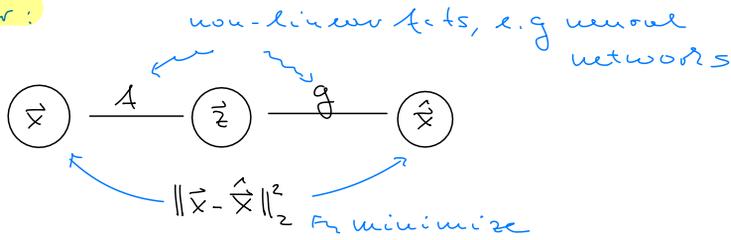


Here: find linear mapping $f: \vec{x} \rightarrow \vec{z}$ with $\vec{x} \in \mathbb{R}^2$ and $\vec{z} \in \mathbb{R}$ and $g: \vec{z} \rightarrow \hat{\vec{x}}$ such as to minimize $\sum_{i=1}^M \|\vec{x}^{(i)} - \hat{\vec{x}}^{(i)}\|_2^2$
 $= \sum_{i=1}^M \|\vec{x}^{(i)} - g(f(\vec{x}^{(i)}))\|_2^2$

Some classification:

	deterministic	non-deterministic
linear	PCA	probabilistic PCA
non-linear	Autoencoder	Variational AE
- -, no encoder		Gen. adv. networks

Autoencoder:



Both the encoder f and the decoder g are deterministic.

In this lecture we want to combine the idea of an autoencoder with the concept of generative modeling.

- want to determine models of probability distributions $p(\vec{x})$ over data points \vec{x} ;
- need to capture structural regularities in the data, e.g. correlations between pixels in images;
- generative latent variable models capture structure of data in distribution of latent variables;
- today we will discuss variational autoencoders, which only approximate $p(\vec{x})$, but allow to draw samples from $p(\vec{x})$.

Bayesian view of generative latent variable models:

$$p(\vec{x}) = \int \underbrace{p(\vec{z}) p(\vec{x}|\vec{z})}_{\text{generative process}} d\vec{z} = \int p(\vec{x}, \vec{z}) d\vec{z} = \mathbb{E}_{\vec{z} \sim p(\vec{z})} p(\vec{x}|\vec{z})$$

- $p(\vec{z})$: prior over latent variable $\vec{z} \in \mathbb{R}^d$
- $p(\vec{x}|\vec{z})$: likelihood (decoder)
- $p(\vec{x})$: marginal likelihood, or evidence

Goal: maximize $p(\vec{x}) = \bar{p}_\theta(\vec{x})$ by learning $p_\theta(\vec{z})$ and $\bar{p}_\theta(\vec{x}|\vec{z})$.

Find model parameters Θ by minimizing negative log-likelihood:

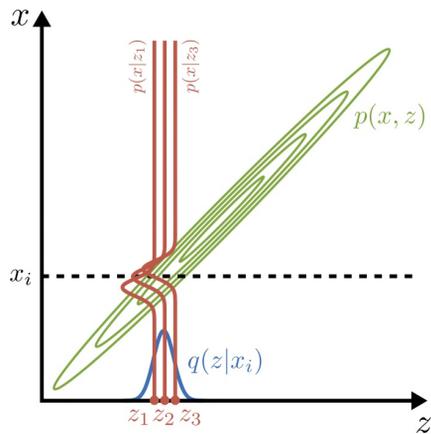
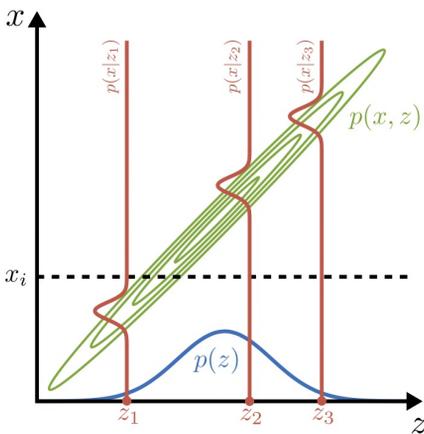
$$\begin{aligned}\Theta^* &= \operatorname{argmin}_{\Theta} \mathbb{E}_{\vec{x} \sim p_{\text{data}}} (-\ln p_{\Theta}(\vec{x})) \\ &= \operatorname{argmin}_{\Theta} \mathbb{E}_{\vec{x} \sim p_{\text{data}}} [-\ln (\mathbb{E}_{\vec{z} \sim p_{\Theta}(\vec{z})} p_{\Theta}(\vec{x}|\vec{z}))] \\ &\propto \operatorname{argmin}_{\Theta} \sum_{i=1}^N -\ln (\mathbb{E}_{\vec{z} \sim p_{\Theta}(\vec{z})} p_{\Theta}(\vec{x}^{(i)}|\vec{z}))\end{aligned}$$

Unfortunately, calculating $p_{\Theta}(\vec{x}^{(i)})$ is in general intractable

→ introduce a so-called recognition model $q_{\Theta}(\vec{z}|\vec{x})$ to approximate true posterior $p_{\Theta}(\vec{z}|\vec{x})$

In AE-terminology: $p_{\Theta}(\vec{x}|\vec{z}) \rightarrow$ decoder
 $p_{\Theta}(\vec{z}|\vec{x}) \rightarrow$ encoder

Example why computing $p_{\Theta}(\vec{x}) = \mathbb{E}_{\vec{z} \sim p_{\Theta}(\vec{z})} p_{\Theta}(\vec{x}|\vec{z})$ is hard.



from DL course Tübingen (Prof. Griger)

How can we use the recognition model $q(\tilde{z}|\vec{x})$ to maximize likelihood?

$$\begin{aligned}\ln p(\vec{x}) &= \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \ln \left(p(\vec{x}) \frac{p(\tilde{z}|\vec{x})}{p(\tilde{z}|\vec{x})} \right) \\ &= \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \left(\ln \left(\frac{p(\tilde{z}, \vec{x})}{q(\tilde{z}|\vec{x})} \right) + \ln \left(\frac{q(\tilde{z}|\vec{x})}{p(\tilde{z}|\vec{x})} \right) \right) \\ &= \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \ln \left(\frac{p(\tilde{z}, \vec{x})}{q(\tilde{z}|\vec{x})} \right) + \underbrace{\text{KL} \left(q(\tilde{z}|\vec{x}) \parallel p(\tilde{z}|\vec{x}) \right)}_{\geq 0} \\ &\geq \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \ln \left(\frac{p(\tilde{z}, \vec{x})}{q(\tilde{z}|\vec{x})} \right) \\ &\quad \sim \text{evidence lower bound (ELBO)}\end{aligned}$$

Note: $q(\tilde{z}|\vec{x})$ is variational approximation to true posterior $p(\tilde{z}|\vec{x})$

KL-div measures approximation error

Rather than maximizing log-likelihood, minimize negative log-likelihood:

$$\begin{aligned}-\ln p(\vec{x}) &\leq \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \ln \left(\frac{q(\tilde{z}|\vec{x})}{p(\tilde{z}, \vec{x})} \right) \\ &= \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \ln \left(\frac{q(\tilde{z}|\vec{x})}{p(\tilde{z}) p(\vec{x}|\tilde{z})} \right) \\ &= \mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} \left(\ln \left(\frac{q(\tilde{z}|\vec{x})}{p(\tilde{z})} \right) - \ln p(\vec{x}|\tilde{z}) \right) \\ &= \text{KL} \left(q(\tilde{z}|\vec{x}) \parallel p(\tilde{z}) \right) - \underbrace{\mathbb{E}_{\tilde{z} \sim q(\tilde{z}|\vec{x})} (-\ln p(\vec{x}|\tilde{z}))}_{\sim \text{reconstruction error}}\end{aligned}$$

\sim autoencoder structure: $\vec{x} \xrightarrow{q} \tilde{z} \xrightarrow{p} \vec{x}$

Variational autoencoder:

Minimize bound to negative log likelihood:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{\vec{x} \sim p_{\text{data}}} (-\ln p_{\theta}(\vec{x}))$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N (-\ln p_{\theta}(\vec{x}^{(i)}))$$

$$\approx \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \left[\text{KL}(q_{\theta}(\vec{z}|\vec{x}^{(i)}) \| p(\vec{z})) + \mathbb{E}_{\vec{z} \sim q_{\theta}(\vec{z}|\vec{x}^{(i)})} (-\ln p_{\theta}(\vec{x}^{(i)}|\vec{z})) \right]$$

In VAE, $q_{\theta}(\vec{z}|\vec{x})$ is a multivariate Gaussian, parametrized by neural network:

$$q_{\theta}(\vec{z}|\vec{x}) = \frac{1}{(2\pi)^{D_z}} \frac{1}{|\Sigma_{\theta}(\vec{x})|^{D_z/2}} \exp\left(-\frac{1}{2} (\vec{z} - \vec{\mu}_{\theta}(\vec{x}))^T \Sigma_{\theta}^{-1}(\vec{x}) (\vec{z} - \vec{\mu}_{\theta}(\vec{x}))\right)$$

→ mean $\vec{\mu}$ and covariance Σ are functions of the data and determined by a neural network.

Standard set-up: choose $q(\vec{z}|\vec{x}) = \mathcal{N}(\vec{z}|\vec{\mu}, \Sigma)$ and prior $p(\vec{z}) = \mathcal{N}(0, \mathbb{1})$ as Gaussian, and $\Sigma = \text{diag}(\sigma^2)$

$$\begin{aligned} \text{Then: } \text{KL}(q(\vec{z}) \| p(\vec{z})) &= \int q(z) (\ln q(z) - \ln p(z)) dz \\ &= \frac{1}{2} \sum_{j=1}^D (\bar{\mu}_j^2 + \sigma_j^2 - 1 - \ln \sigma_j^2) \end{aligned}$$

Can we train the VAE using backpropagation?

Need to calculate gradient of an expectation value.

Simple example:

$$\begin{aligned} \vec{\nabla}_{\theta} \mathbb{E}_{\vec{z} \sim p(\vec{z})} (f_{\theta}(\vec{z})) &= \vec{\nabla}_{\theta} \int p(\vec{z}) f_{\theta}(\vec{z}) d\vec{z} = \int p(\vec{z}) \vec{\nabla}_{\theta} f_{\theta}(\vec{z}) d\vec{z} \\ &= \mathbb{E}_{\vec{z} \sim p(\vec{z})} [\vec{\nabla}_{\theta} f_{\theta}(\vec{z})] \end{aligned}$$

But what happens if $p(\vec{z}) = p_0(\vec{z})$?

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{\vec{z} \sim p(\vec{z})} [A_{\theta}(\vec{z})] &= \nabla_{\theta} \left[\int p_0(\vec{z}) A_{\theta}(\vec{z}) d\vec{z} \right] \\ &= \mathbb{E}_{\vec{z} \sim p_0(\vec{z})} [\nabla_{\theta} A_{\theta}(\vec{z})] + \underbrace{\int A_{\theta}(\vec{z}) \nabla_{\theta} p_0(\vec{z}) d\vec{z}}_{= 0} \end{aligned}$$

Reparametrisation trick:

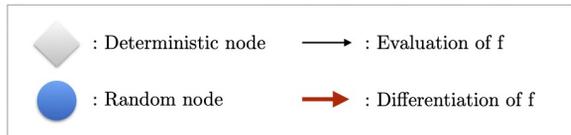
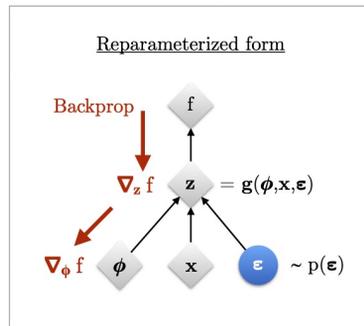
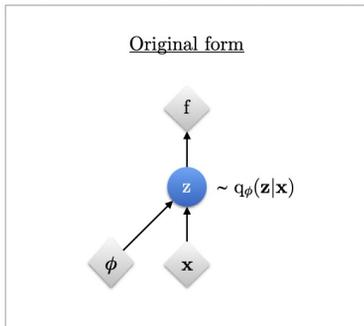
Introduce random variables $\vec{\epsilon} \sim p(\vec{\epsilon})$ and write $\vec{z} = g_{\theta}(\vec{\epsilon}, \vec{x})$ (e.g. $\vec{\epsilon} \sim \mathcal{N}(\vec{0}, \mathbb{I})$ and $\vec{z} = \vec{\mu}(\vec{x}) + \vec{\Sigma}(\vec{x}) \vec{\epsilon}$)
element-wise product

$$\text{Then } \mathbb{E}_{\vec{z} \sim p_0(\vec{z})} [A(\vec{z})] = \mathbb{E}_{\vec{\epsilon} \sim p(\vec{\epsilon})} [A(g_{\theta}(\vec{\epsilon}, \vec{x}))]$$

$$\Rightarrow \nabla_{\theta} \mathbb{E}_{\vec{z} \sim p_0(\vec{z})} [A(\vec{z})] = \nabla_{\theta} \mathbb{E}_{\vec{\epsilon} \sim p(\vec{\epsilon})} [\dots]$$

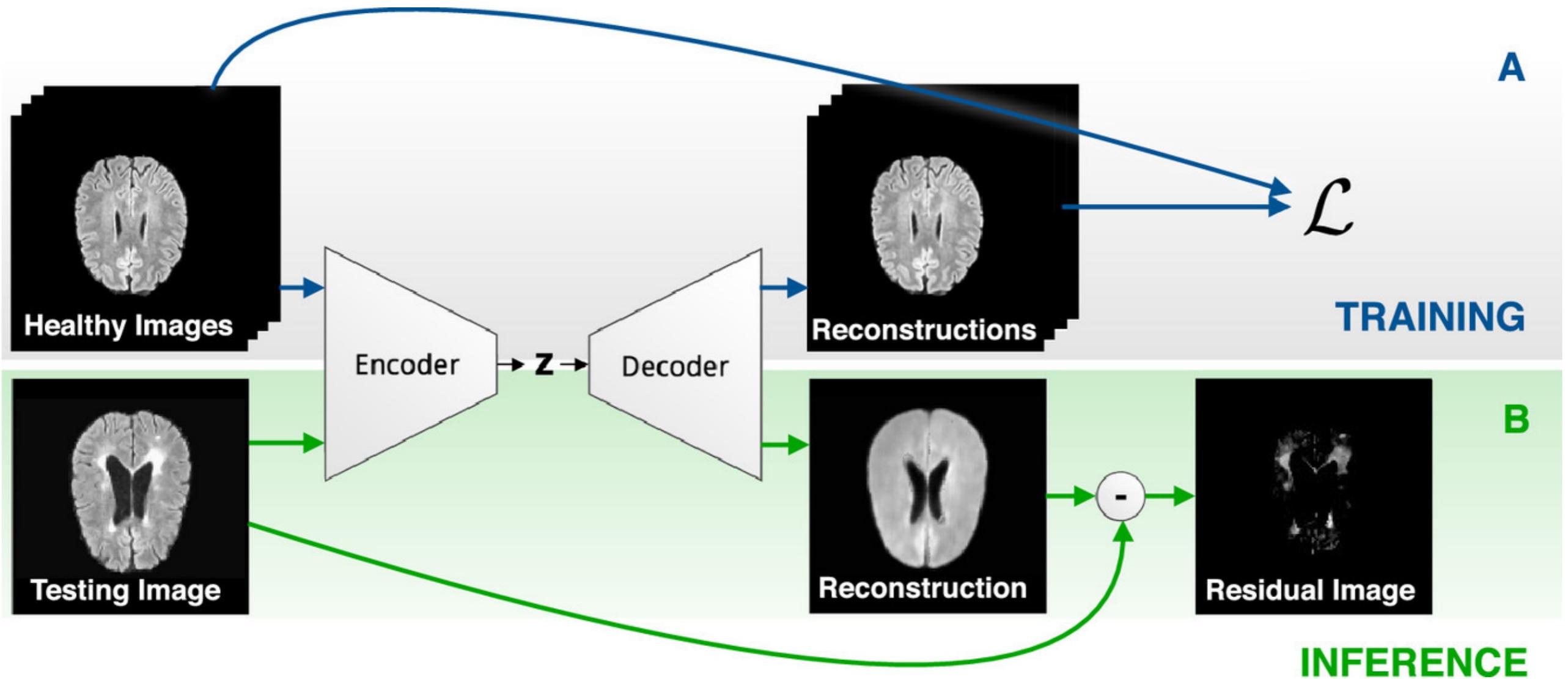
$$= \mathbb{E}_{\vec{\epsilon} \sim p(\vec{\epsilon})} [\underbrace{\nabla_{\theta} A(g_{\theta}(\vec{\epsilon}, \vec{x}))}_{\text{evaluate with MC methods}}]$$

evaluate with MC methods: $\frac{1}{L} \sum_{l=1}^L \nabla_{\theta} A(g_{\theta}(\vec{\epsilon}^{(l)}, \vec{x}))$



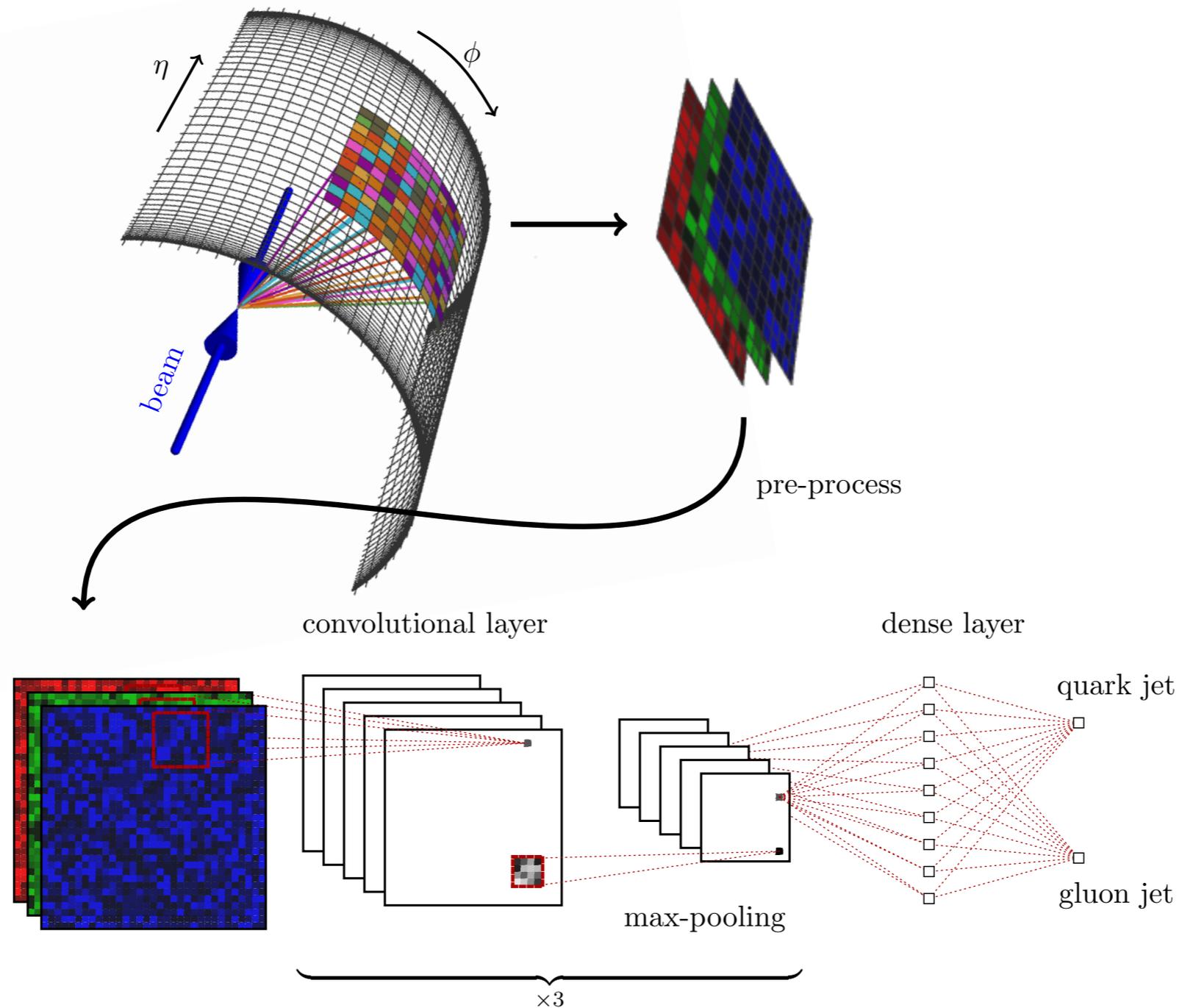
See Kingma and Welling 2014

Anomaly detection with autoencoders



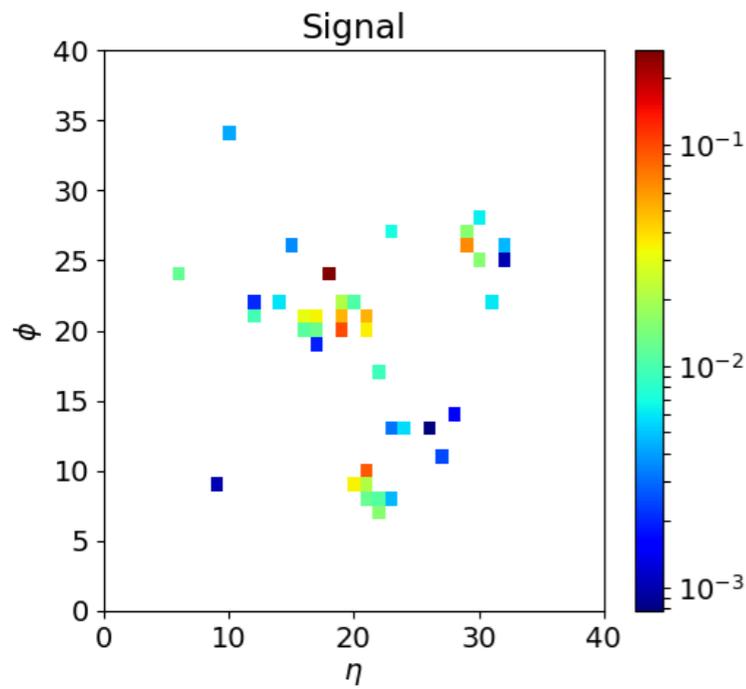
Baur et al., 2020

Convolutional neural networks for jet images

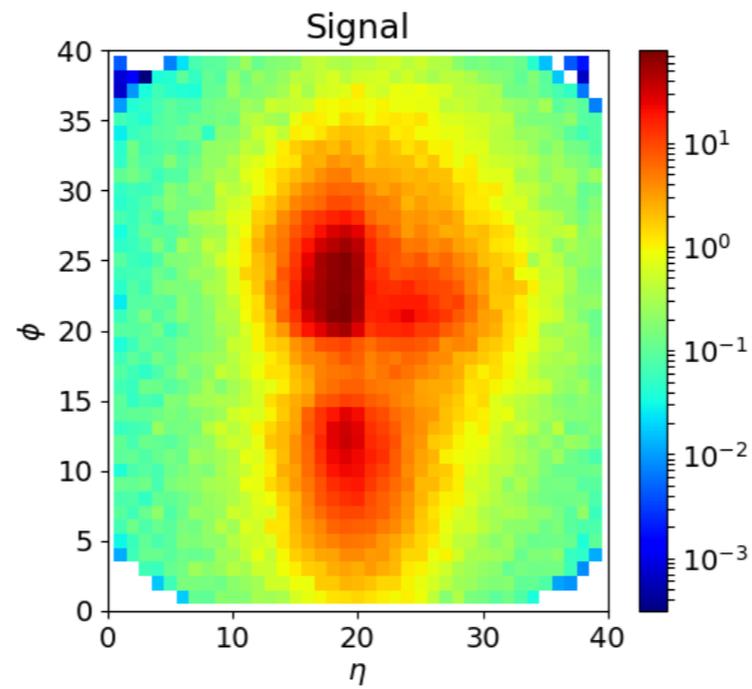


Jet images

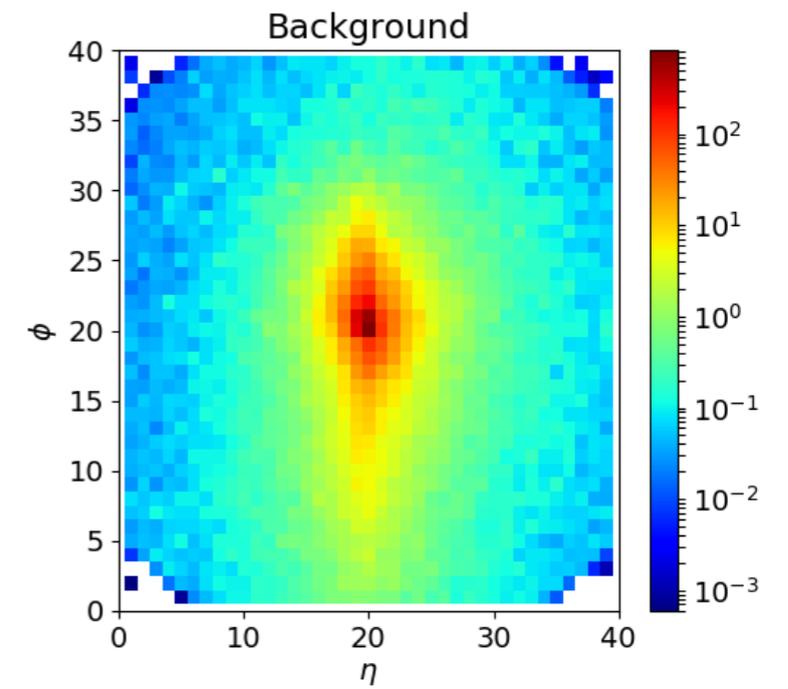
Typical single top-jet



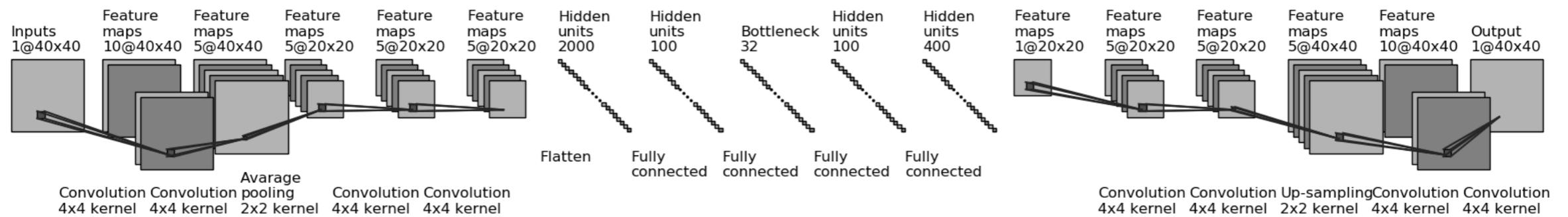
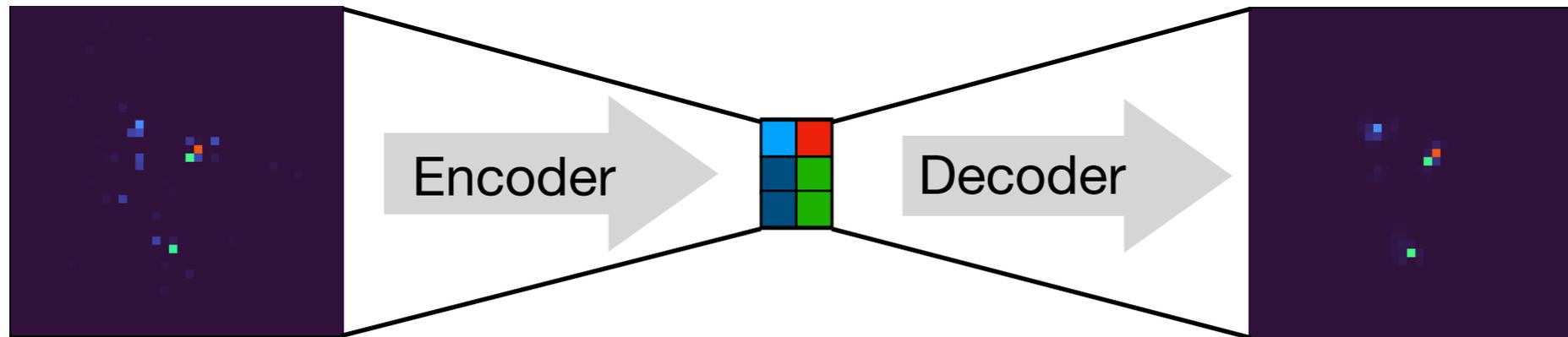
Average top-jet



Average QCD-jet

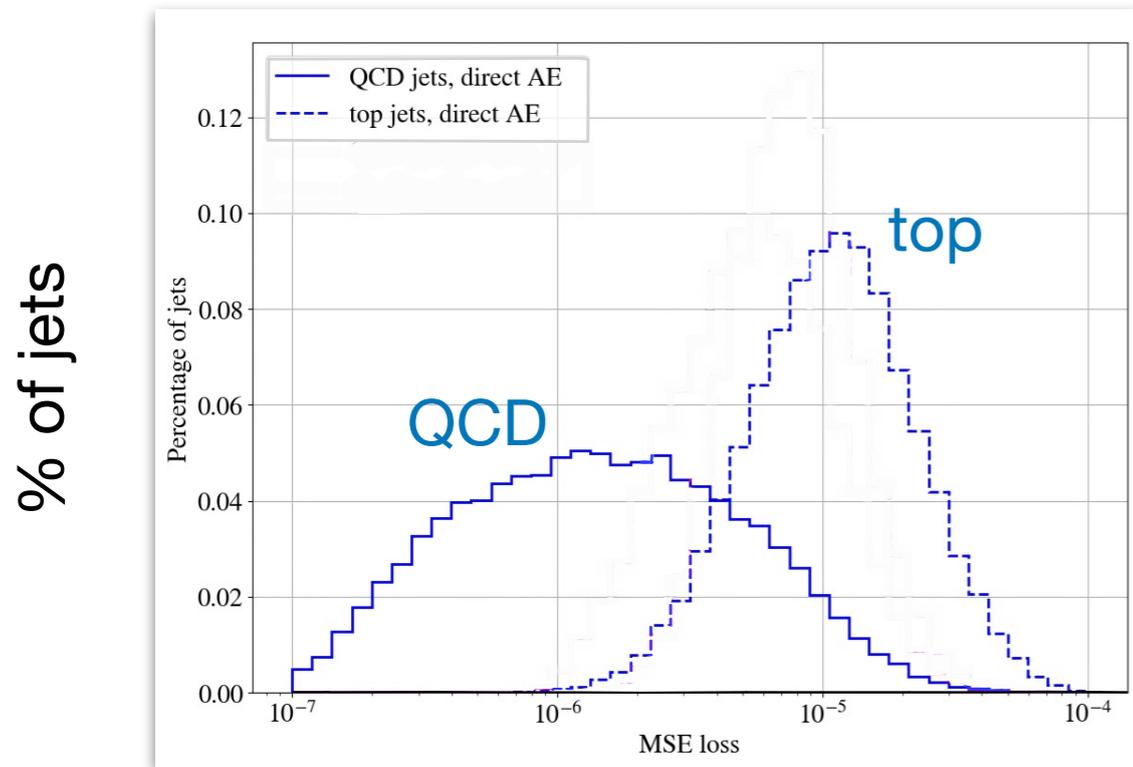
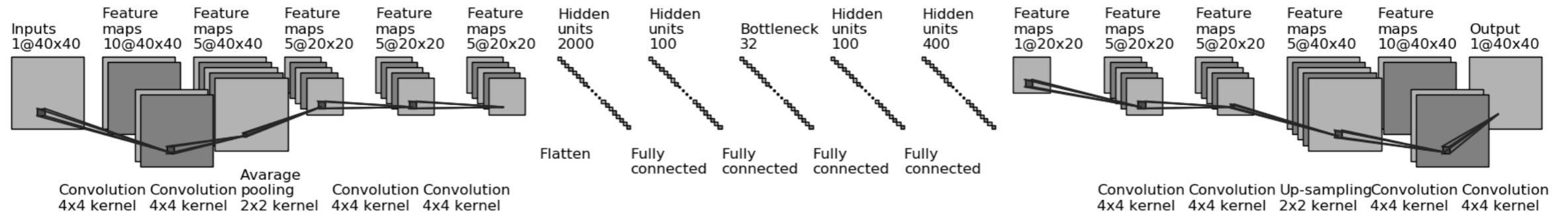


Anomaly detection with autoencoders

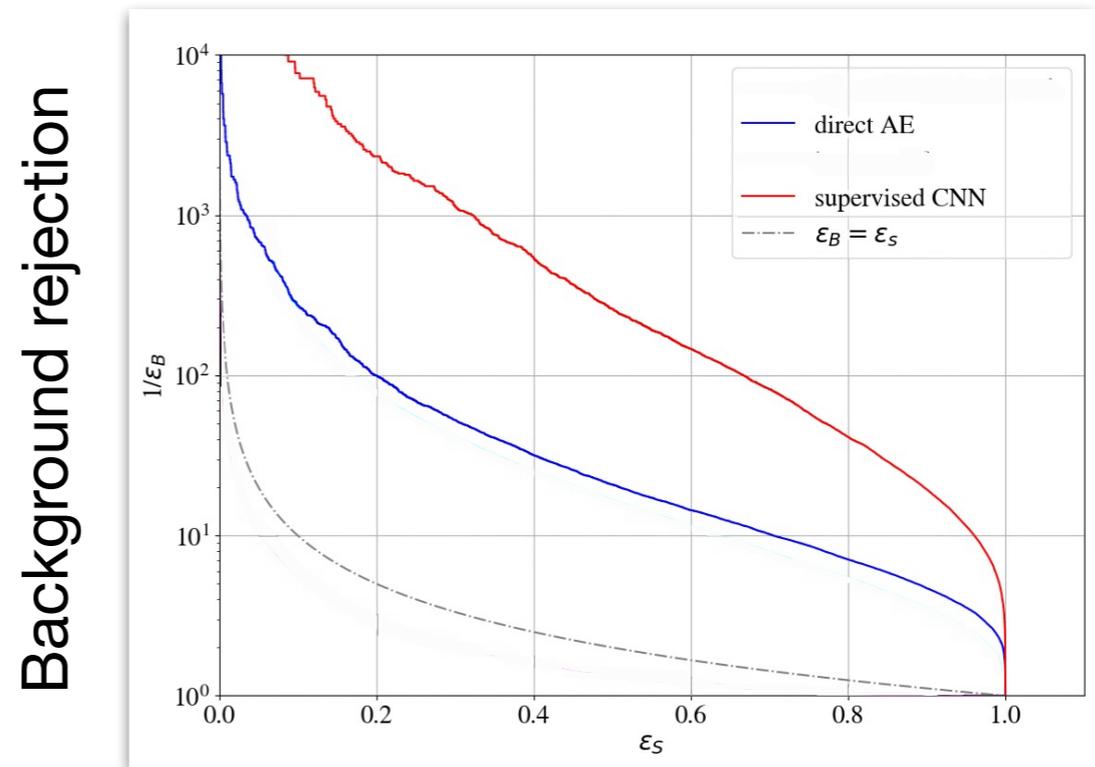


cf Heimgel, Kasieczka, Plehn, Thompson, SciPost Phys. 6, 030 (2019); Farina, Nakai, Shih, PRD 101 (2020)

Anomaly detection with autoencoders



MSE loss



Top tagging efficiency