



DD2437 – Artificial Neural Networks and Deep Architectures (annda)

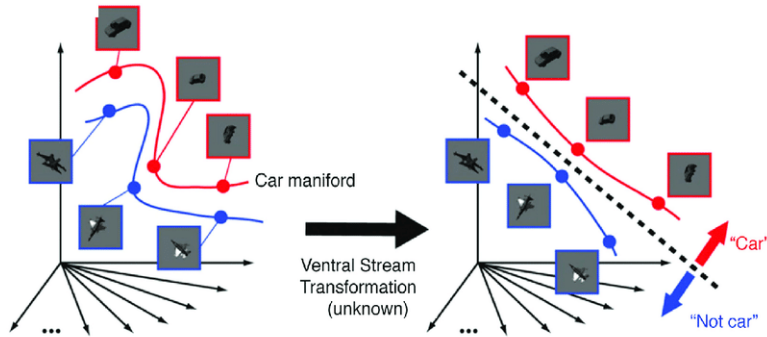
Lecture 11: **Learning deep representations, deep generative models**

Pawel Herman

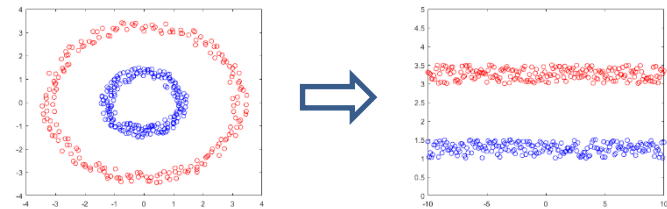
Computational Science and Technology (CST)

KTH Royal Institute of Technology

The importance of representations

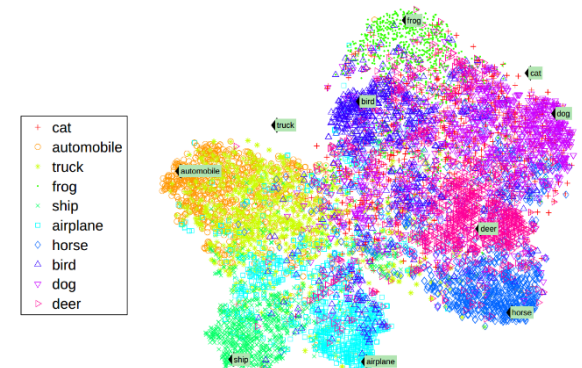
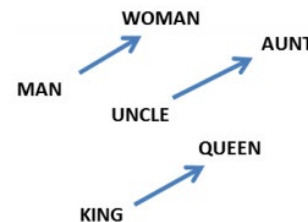


Good representations can make certain computations cheaper and more robust



“woman” $\rightarrow [-0.3, 0.2, 0.5, 0, -0.1, 0...]$

“woman” – “man” \approx “queen” – “king”



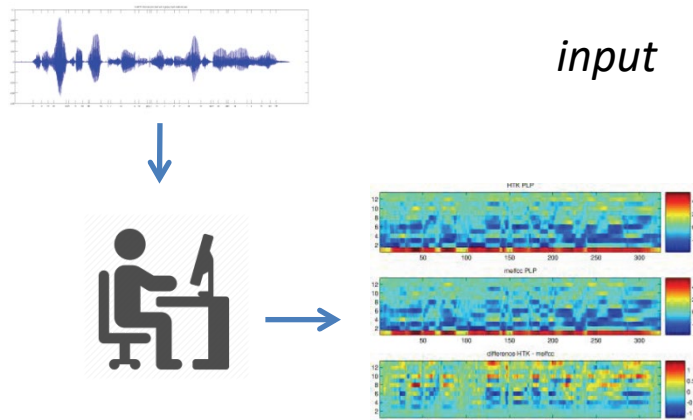
Learning representations + complex inference = future AI?

Learning representations as a hallmark of DL

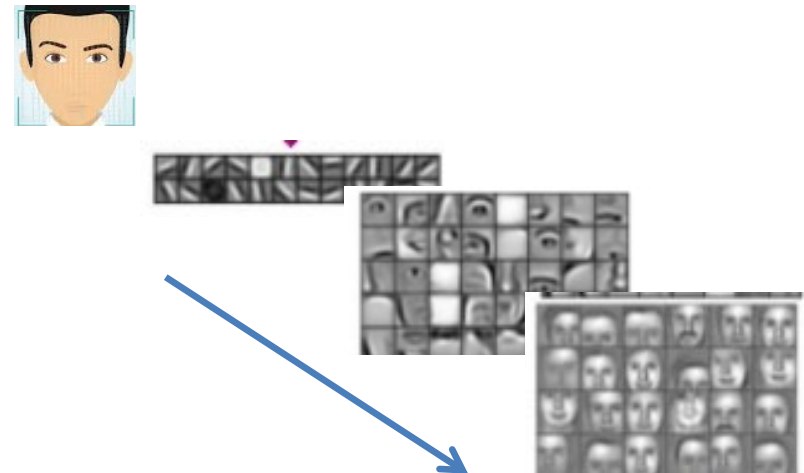
Learning features as part of DL modus operandi

- Traditional pattern recognition VS deep neural network approach

Hand-engineered features in a traditional pattern recognition approach



End-to-end networks with learned features spaces, data representations

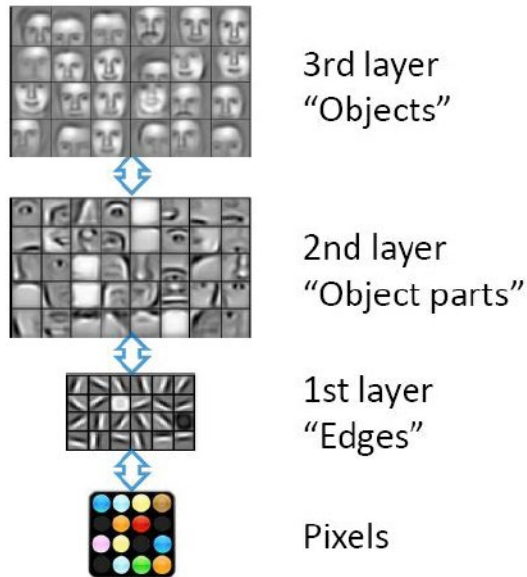


features, representations

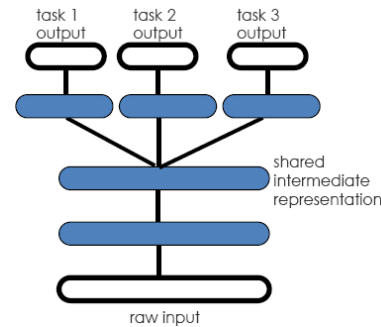
Learning representations as a hallmark of DL

Learning features as part of DL modus operandi
with many implications.....

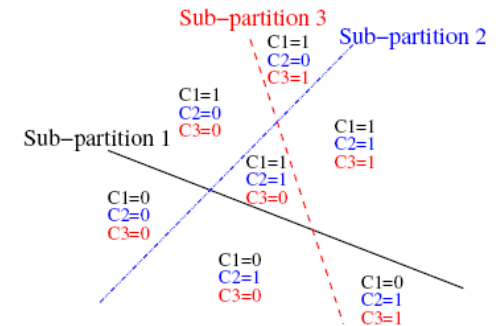
hierarchy of abstraction levels



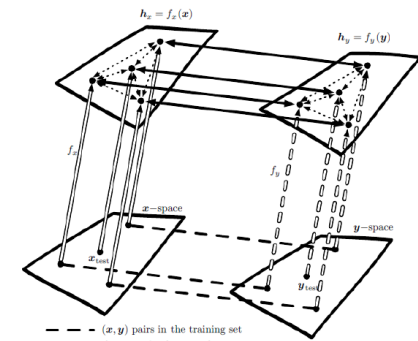
multi-tasking and transfer learning



multi-clustering

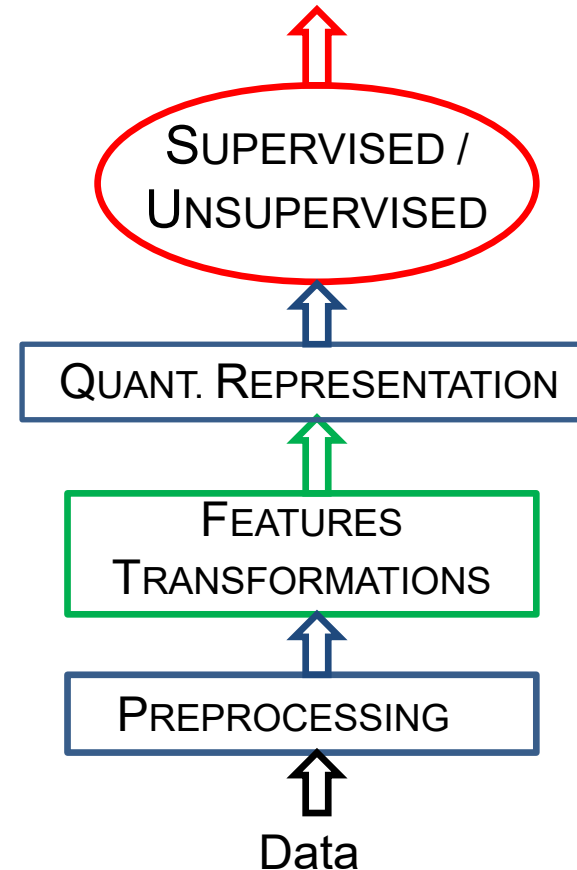
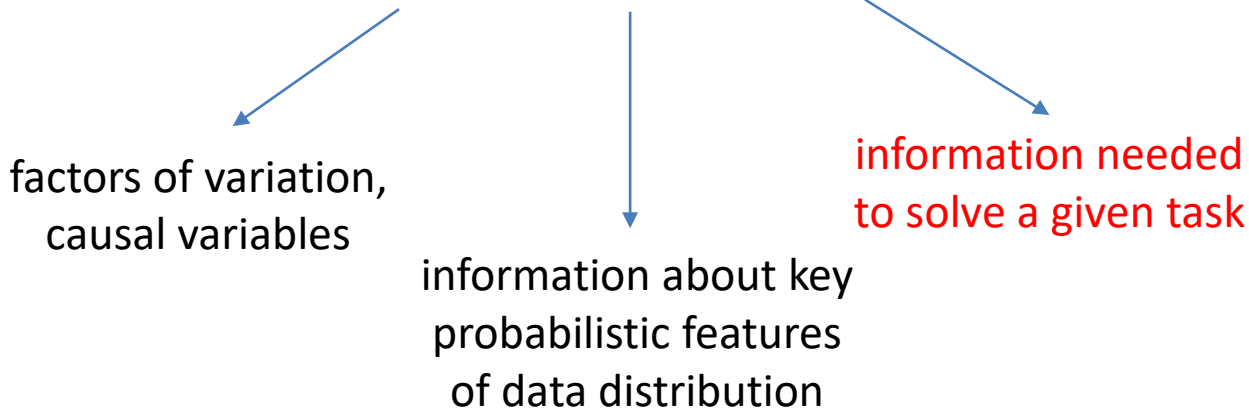


zero-shot learning



Representation learning problem

What makes representation good? What is desirable/useful information?



Representation learning problem

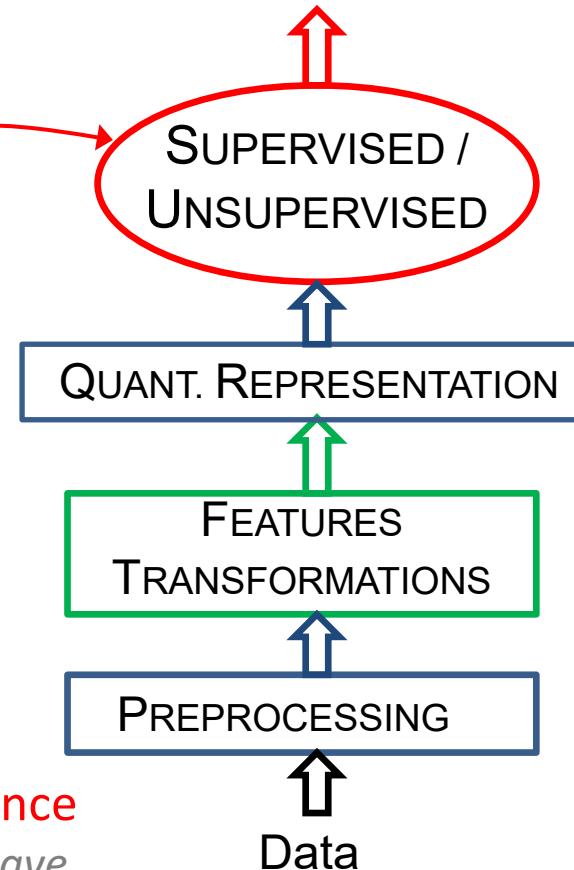
What makes representation good? What is desirable/useful information?

factors of variation,
causal variables

information about key
probabilistic features
of data distribution

information needed
to solve a given task

Facilitate the subsequent learning task, maximise performance
(easiest to define for supervised learning problems but does not have
to lead to “good” representations)



Supervised vs unsupervised approach

- Supervised – distilling information relevant to a concrete task defined by labels
 - very useful when solving particular tasks
 - strongly relying on the abundance of labelled data and prone to excessive information bottleneck

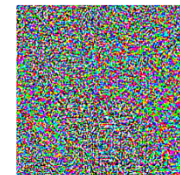
- lacking “common sense”



“panda”

57.7% confidence

+ .007 ×



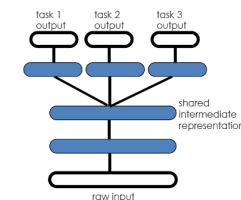
noise

=



“gibbon”

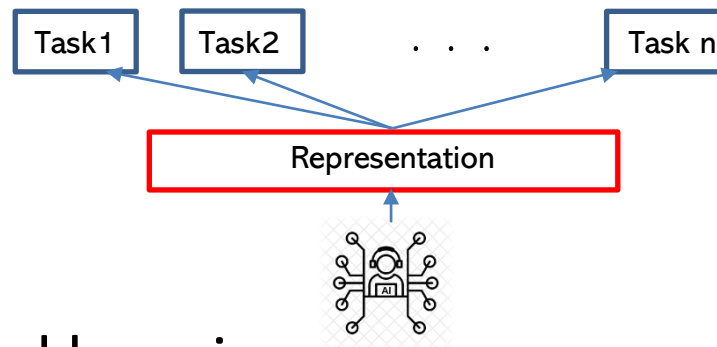
99.3% confidence



- different ways to “improve” representations

Supervised vs unsupervised approach

- Supervised – distilling information relevant to a concrete task defined by labels
- Unsupervised learning
 - “The idea of learning to represent the world before learning a task – this is what babies do” (*LeCun*)
 - It appears as a more generic approach less susceptible to inf. Loss



- Semi-supervised learning

Representation learning

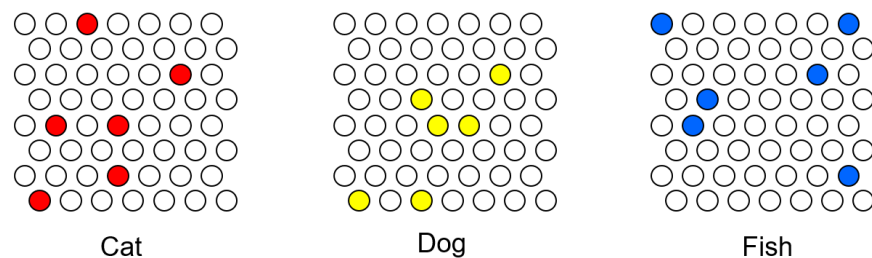
- Computational perspective: disentangling unknown factors causing relevant variation in the data (*factors of variation*)
 - *causes explain* the observed data (discriminative context, both unsupervised and supervised aspects)
 - factors in combination can be used to generate data (generative context)
- Probabilistic perspective
 - *Classical approach: density estimation* – learn probability distribution for data with the use of latent variables (PCA, ICA, GMM etc.) -> explain data
 - $P(\text{data} | \text{latent var})$ for generation and $P(\text{latent var} | \text{data})$ for recognition
 - full probabilistic model advocated by generative models, $P(\text{data}, \dots)$

Distributed vs local representations

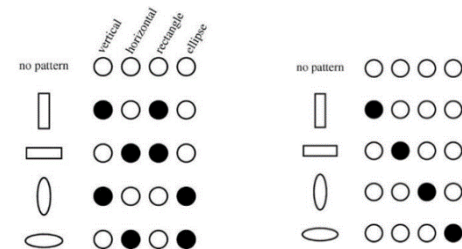
Information is distributed across many units that account for information about features that are not mutually exclusive.....

... unlike in clustering (cluster centres act as prototypes) with distinct regions where *local generalisation* is observed.

Locality in input space implies different behaviour of the learned function in different regions of data space (local or symbolic representations).



DISTRIBUTED vs **LOCAL**



Generalisation due to shared attributes and semantic proximity.

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Sparse vs dense representations

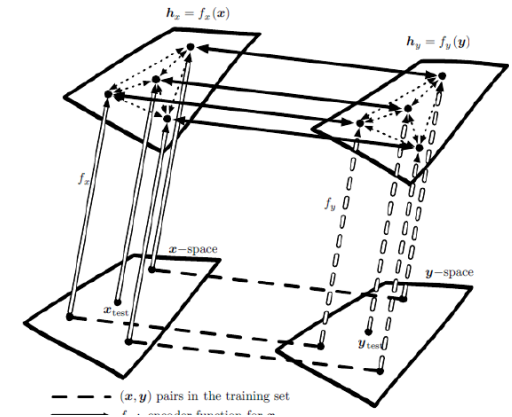
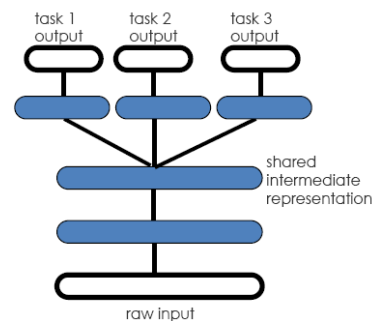
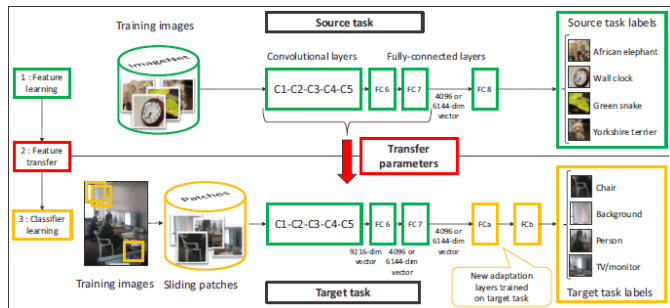
Sparse representations

- orthogonalisation/decorrelation – more separable in high-dimensional spaces
- “metabolic” efficiency
- neural selectivity (vs coarse coding with broad tuning)
- balance between sparse local representations that suffer from the *curse of dimensionality* and dense representations that *entangle* factors and can be hard to interpret

sparse not distributed	not sparse distributed	sparse distributed
0 .2 0 0 0	.1 .8 .7 .5 .7	0 .8 0 .5 0
0 0 0 0 .1	.8 .9 .6 .2 .4	0 0 .6 0 .4
0 0 0 .4 0	.3 .1 .6 .3 .3	.3 0 0 .3 0

Functional implications of distributed codes

- Transfer learning and domain adaptation
- Multi-task learning
- Zero-shot learning



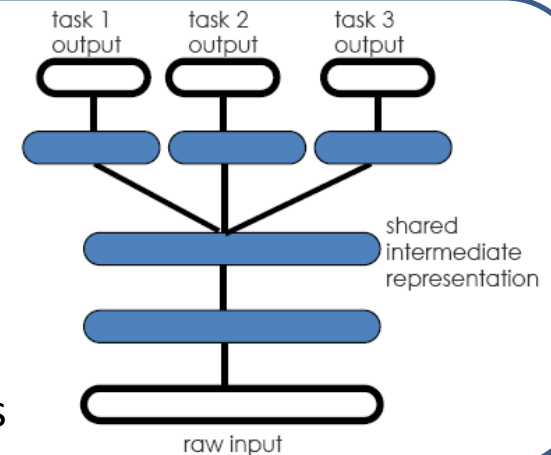
Transfer learning

A more general concept in ML: a model developed for one task is reused (re-purposed) as the starting point for a model on another task

- also referred to as *domain adaptation*
- related to both *multi-tasking* and *concept drift*

Sharing factors across tasks

- assumption that factors explaining the variations in different tasks are shared/common
- especially low-level features are expected to be the same
- supported by hierarchical and distributed representations



Transfer learning

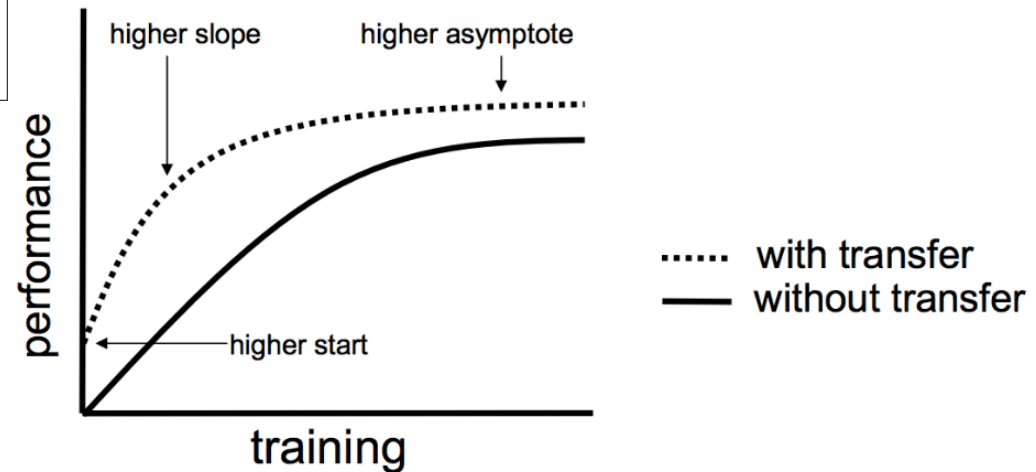
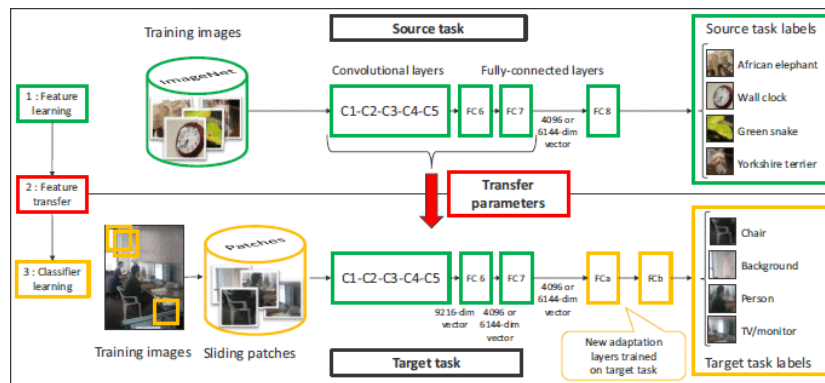
A more general concept in ML: a model developed for one task is reused (re-purposed) as the starting point for a model on another task

- also referred to as *domain adaptation*
- related to both *multi-tasking* and *concept drift*
- popular in DL as it offers an opportunity to save computational costs
- two main approaches: develop model and **pre-trained model**
 - Google Inception model, Oxford VGG model (ImageNet)
 - Stanford's GloVe Model, Google's word2vec Model
 - check out *Caffe Model Zoo*

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Transfer learning

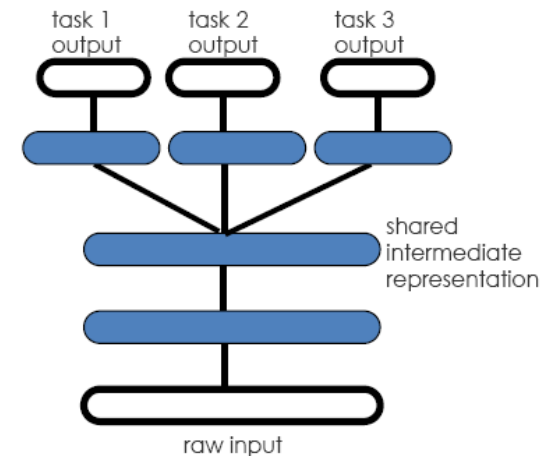
“Transfer learning and domain adaptation refer to the situation where what has been learned in one setting ... is exploited to improve generalization in another setting”



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Multi-task learning

- Excessive focus on a single task may ignore information that can help us perform better on the metric we care about
- Effectively, simultaneously optimizing different cost functions implements multi-task learning
- Form of inductive transfer learning with inductive bias/regularisation mediated by another task
- “Multi-task learning improves generalization by leveraging the domain-specific information contained in the training signals of related tasks” (*Caruana, 1998*)

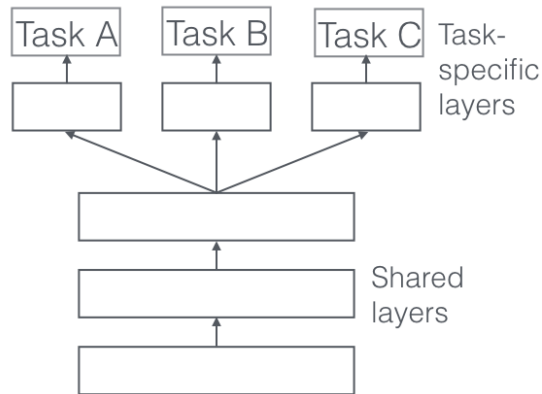


S. Ruder (2017). An Overview of Multi-Task Learning in Deep Neural Networks

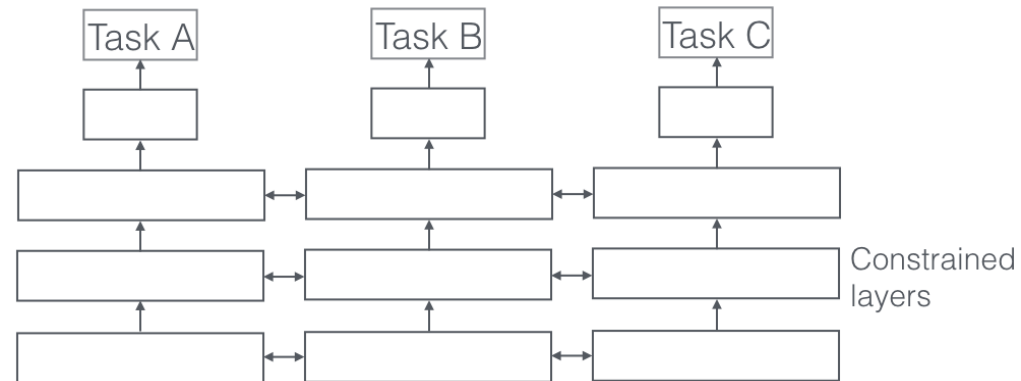
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Multi-task learning

Hard parameter sharing



Soft parameter sharing



Underlying mechanisms: implicit data augmentation, attention focusing, representation bias, regularisation

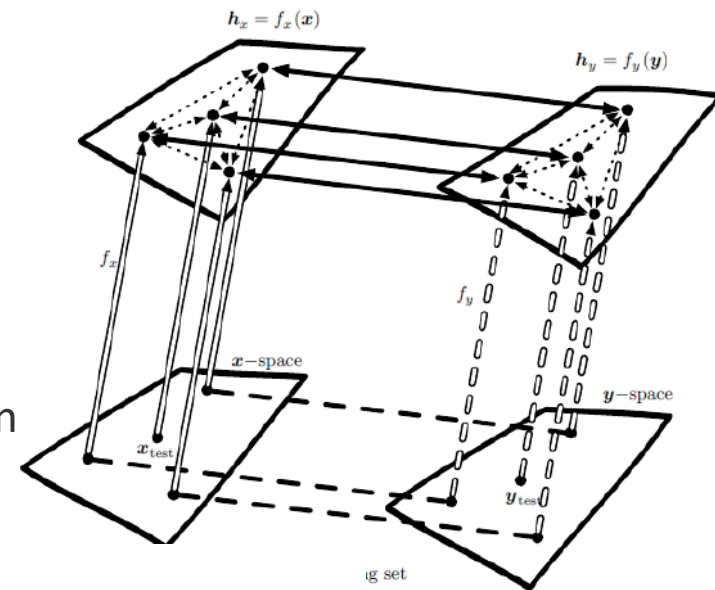
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Zero-shot learning

- Zero-shot learning as a specific form of *multi-modal learning* (capturing the relationship between representations in different modalities)
- In zero-shot classification we want to recognize objects from classes that the model has not seen during training.

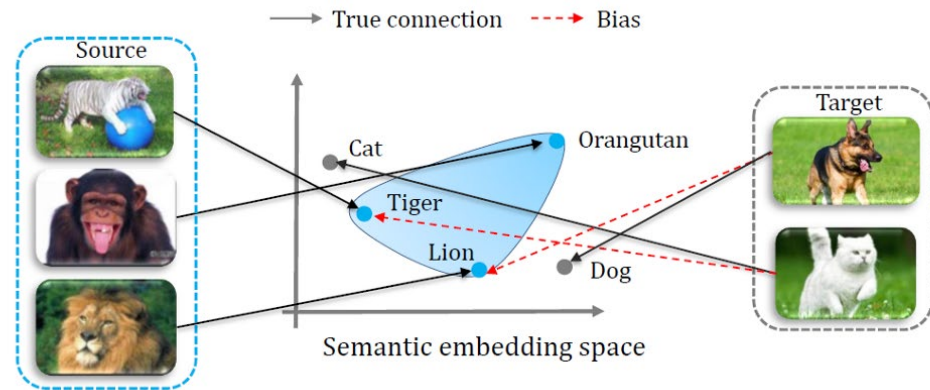
Available data:

- Seen classes** (with labels)
- Unseen classes:** (no labels available during training)
- Auxiliary information:** descriptions/semantic attributes/word embeddings for both seen and unseen classes at train time - a bridge between seen and unseen classes.



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Zero-shot learning



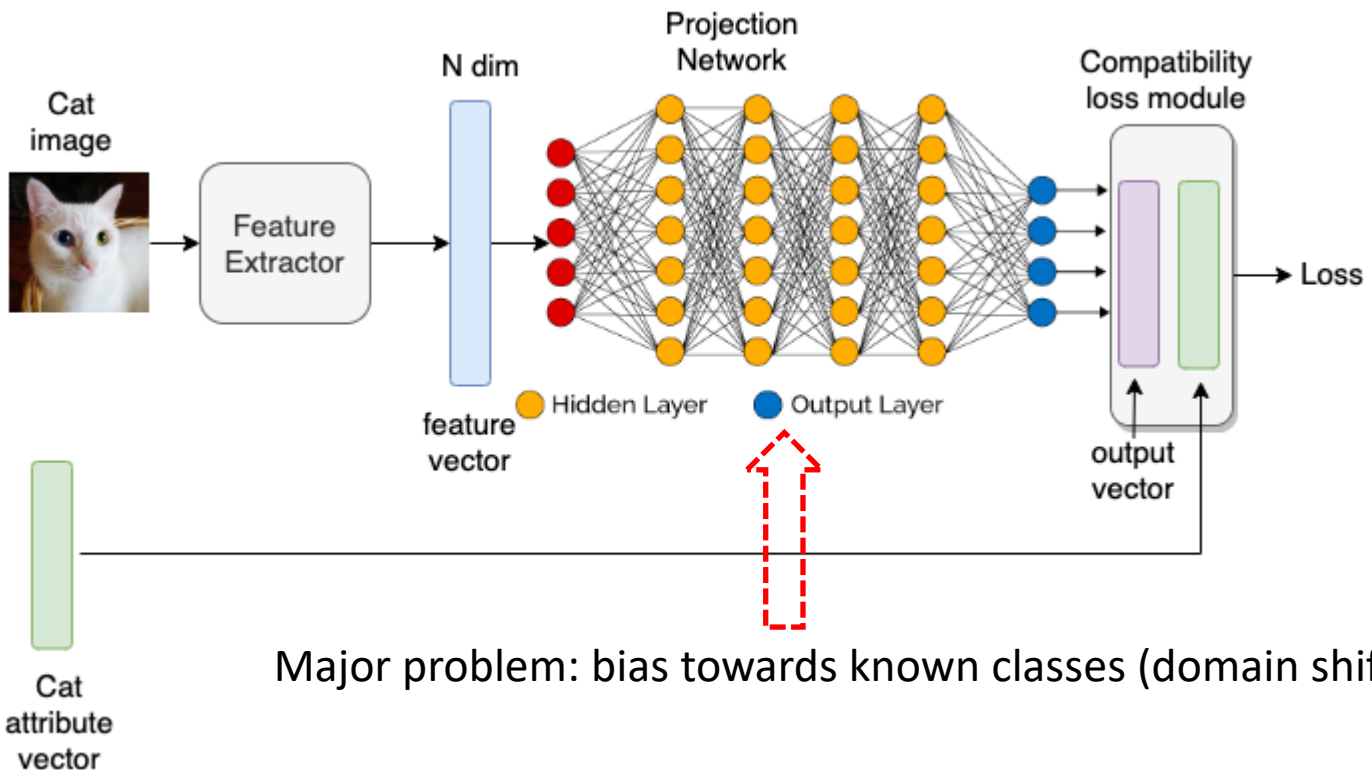
Available data:

- Seen classes** (with labels)
- Unseen classes:** (no labels available during training)
- Auxiliary information:** descriptions/semantic attributes/word embeddings for both seen and unseen classes at train time - **a bridge between seen and unseen classes.**

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Zero-shot learning

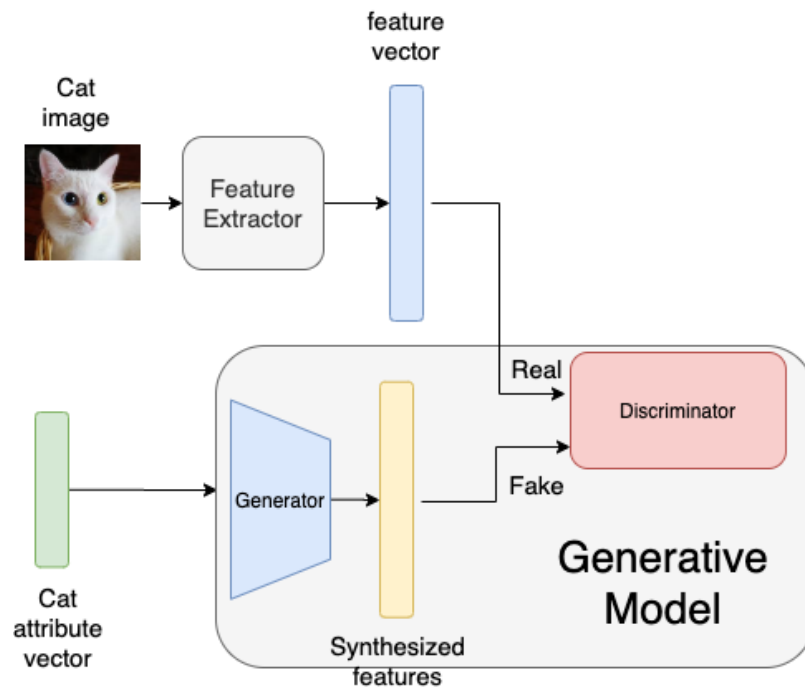
Embedding based methods



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Zero-shot learning

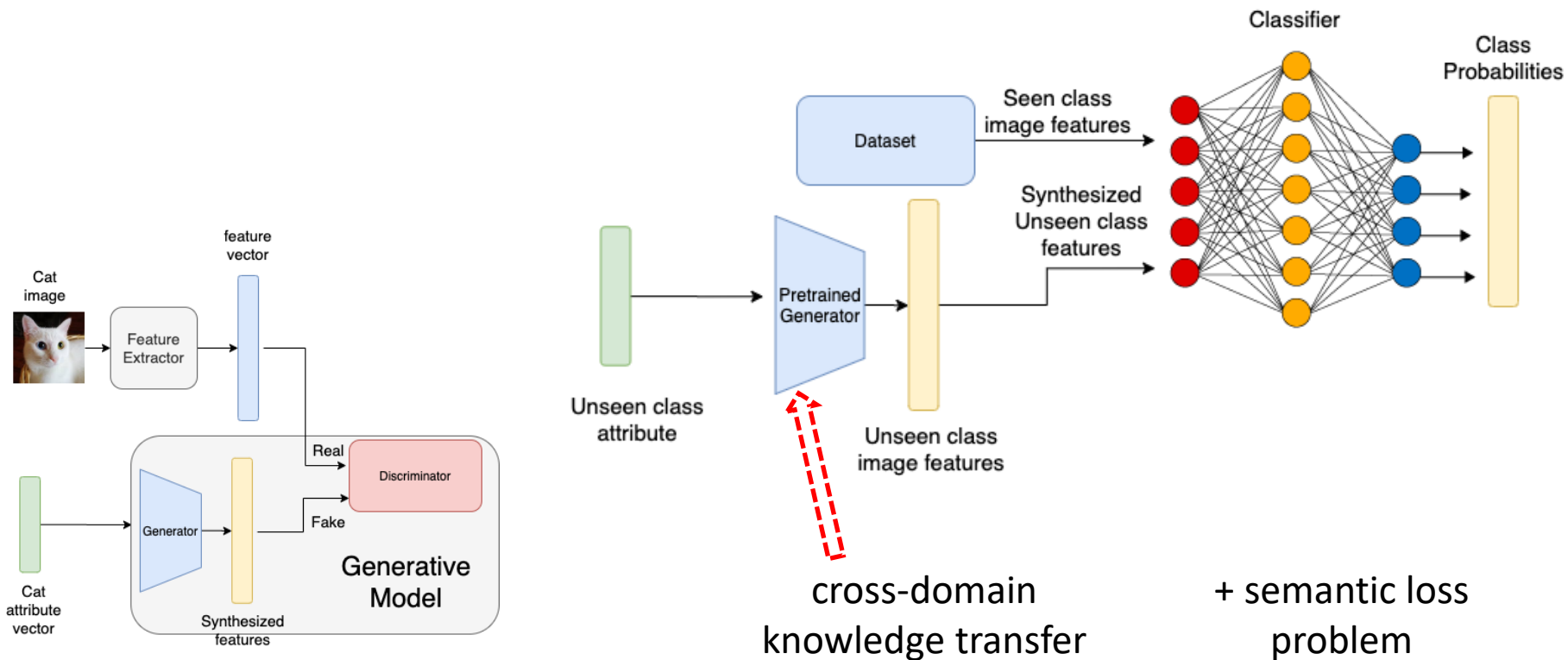
Generative model based methods



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Zero-shot learning

Generative model based methods



Representation learning in deep models

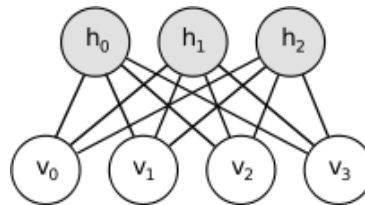
General philosophy and approach

- supervised vs unsupervised
- the concept of unsupervised *pre-training*
- *probabilistic models* (latent variables that describe the distribution) vs *direct encoding* (learning a parametric map from input to representation)
- historically important *manifold learning* and *generative* models
- growing interest in *unsupervised* representation learning with *generative*, *contrastive losses* and *self-supervised learning* approaches

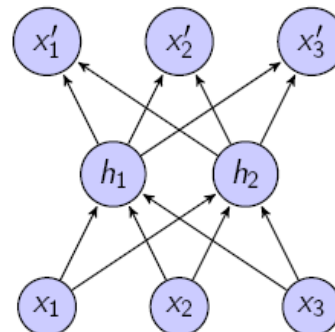
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Representation learning in deep models

- Pre-training for fundamental computational blocks
 - regularized Boltzmann machines (RBMs)



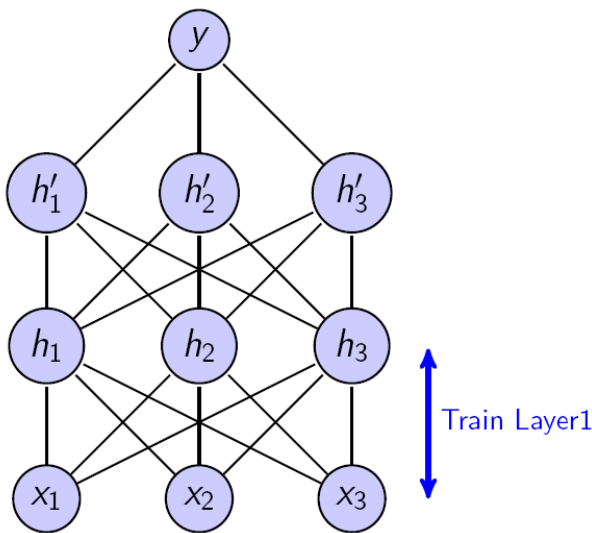
- autoencoders



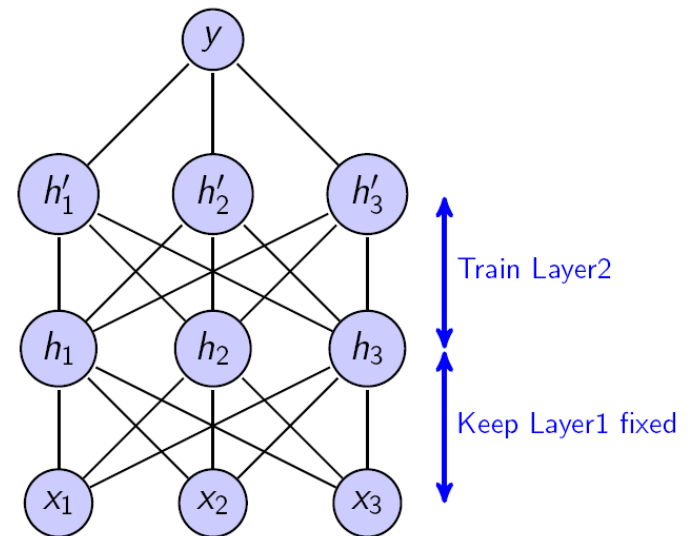
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Representation learning in deep models

- The concept of layer-by-layer pretraining
 - greedy layer-wise unsupervised representation learning



Single layer at a time

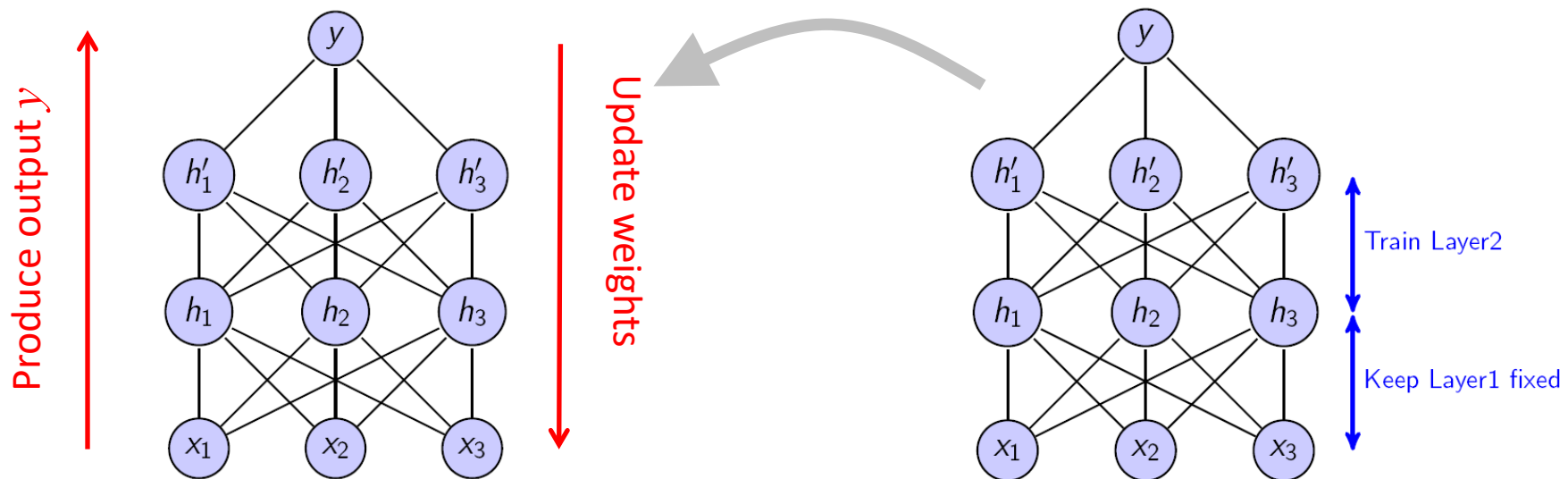


Train another layer while keeping the lower layer fixed

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Representation learning in deep models

- The concept of layer-by-layer pretraining
 - greedy layer-wise unsupervised representation learning
 - intuitively, learning about the input distribution should help in learning the *mapping* between the input and output space
 - BUT having two separate phases has *disadvantages*

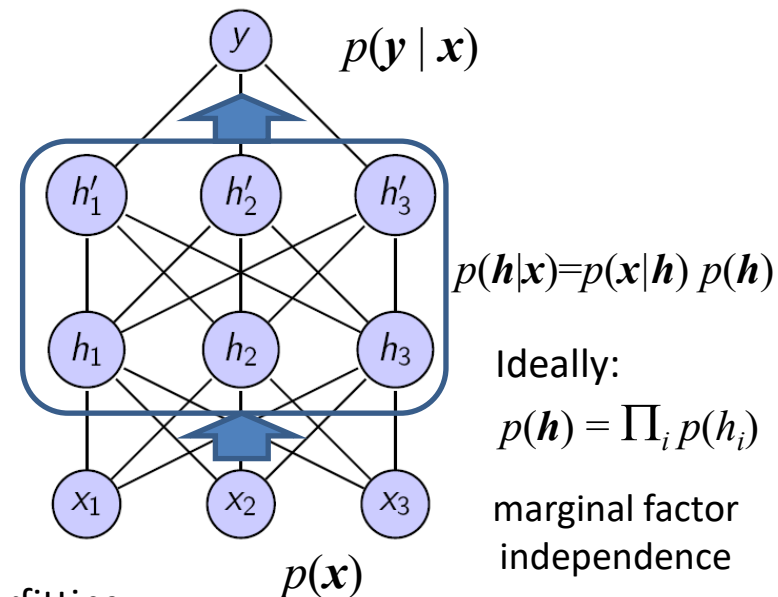


- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Representation learning in deep models

- The concept of layer-by-layer pretraining
 - greedy layer-wise unsupervised representation learning

- RBMs, autoencoders
- leads to lower test classification error
- pretraining as an initialisation scheme
 - prior to supervised fine-tuning
 - initialisation for other unsupervised algorithms such as DBM, DBN etc.
- *optimisation vs regularisation hypothesis*
 - lower variance in learning, less risk for overfitting
 - as a regulariser, it urges the learning algorithm to discover **features that explain underlying causes that generate the data** (also, causal factors often remain *invariant*)



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

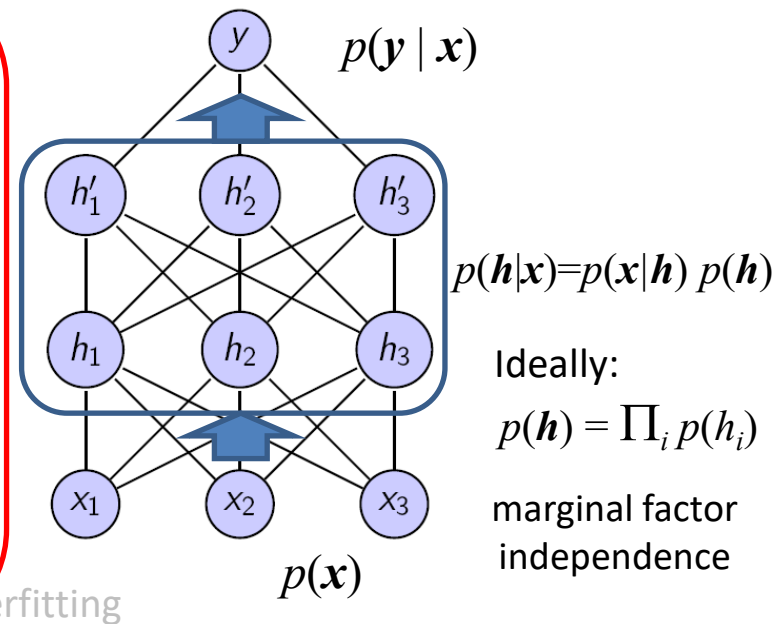
Representation learning in deep models

- The concept of layer-by-layer pretraining
 - greedy layer-wise unsupervised representation learning

Representation learning should strive towards uncovering latent factors, \mathbf{h} , which capture underlying causes in \mathbf{x} .

Then, if \mathbf{y} is one of them, i.e. $\mathbf{y} = \mathbf{h}_i$, it should be easy to learn to predict \mathbf{y} from this representation.

So, how to make representation encode relevant/salient factors?



- as a regulariser, it urges the learning algorithm to discover **features that explain underlying causes** that generate the data (also, causal factors often remain *invariant*)

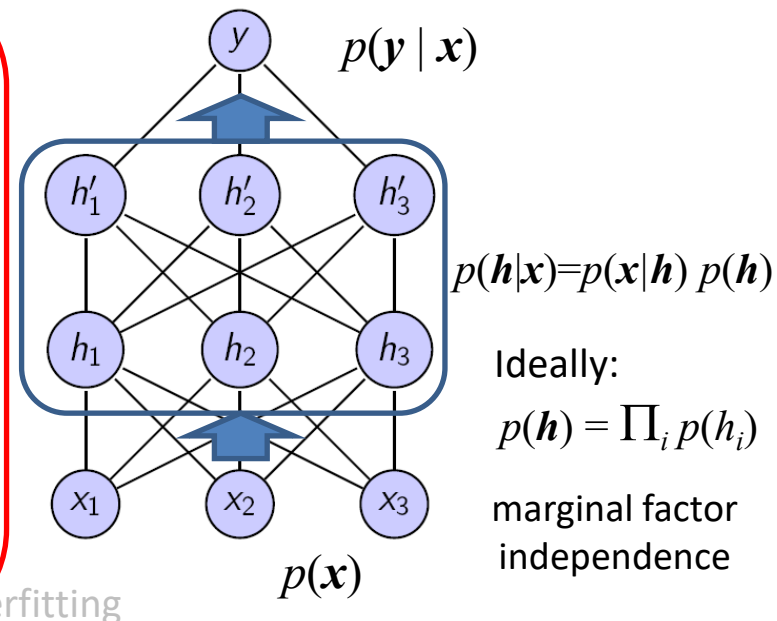
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Representation learning in deep models

- The concept of layer-by-layer pretraining
 - greedy layer-wise unsupervised representation learning

So, how to make representation encode relevant/salient factors?

- 1) Guide unsupervised pretraining with a supervised learning signal (e.g. autoencoders).
- 2) Rely on massive representations with purely unsupervised learning (e.g. RBMs).
- 3) Redefine the meaning of salience.

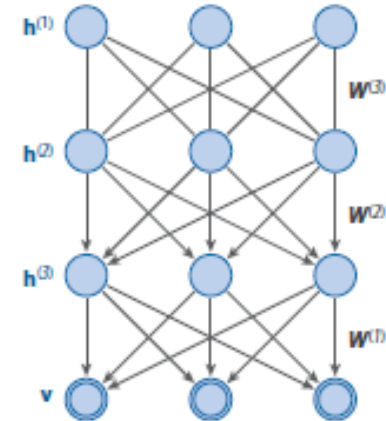
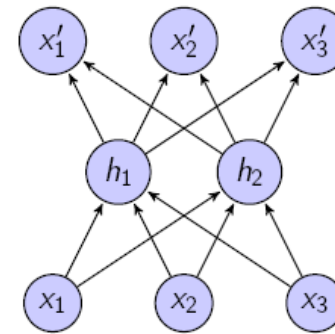
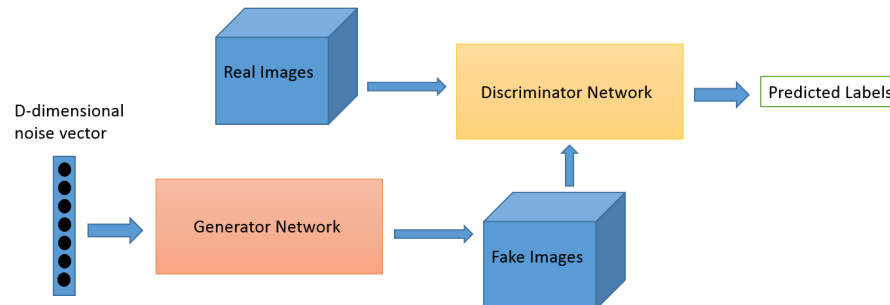
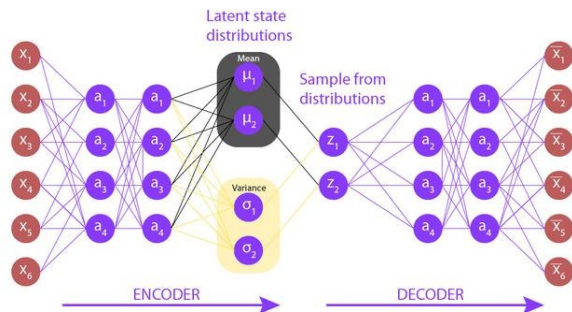


- as a regulariser, it urges the learning algorithm to discover **features that explain underlying causes that generate the data** (also, causal factors often remain *invariant*)

Representation learning in deep models

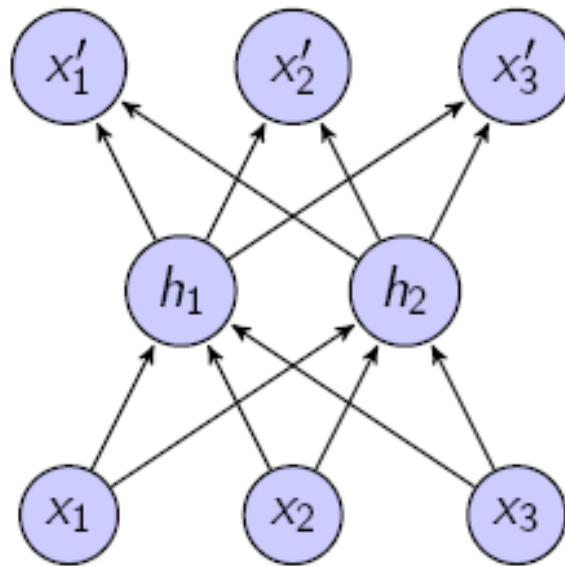
DNN architectures of special interest

- stacked autoencoders (manifold learning)
- generative models: DBNs, DBMs or Generative Adversarial Networks (GANs), Variational Autoencoders (VAE)



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Autoencoders



Decoder: $x' = f(x)$

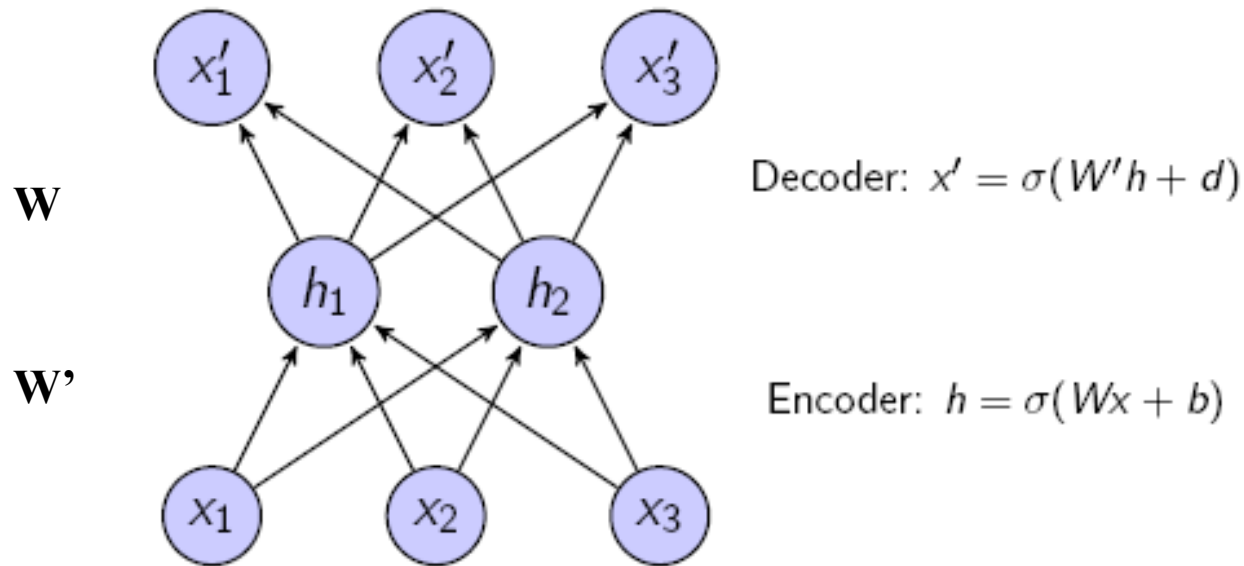
Encoder: $h = g(f(x))$

The idea is to minimise the loss function, L :

$$L(x, g(f(x)))$$

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Autoencoders



Encourage h to give small reconstruction error:

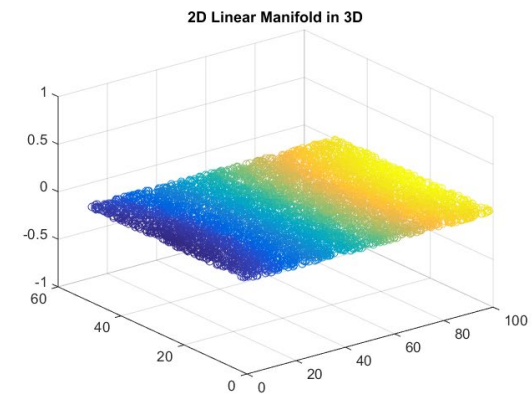
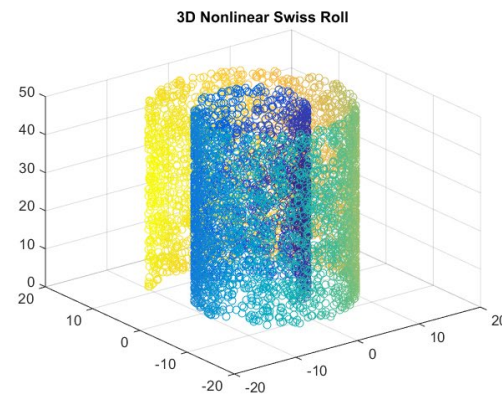
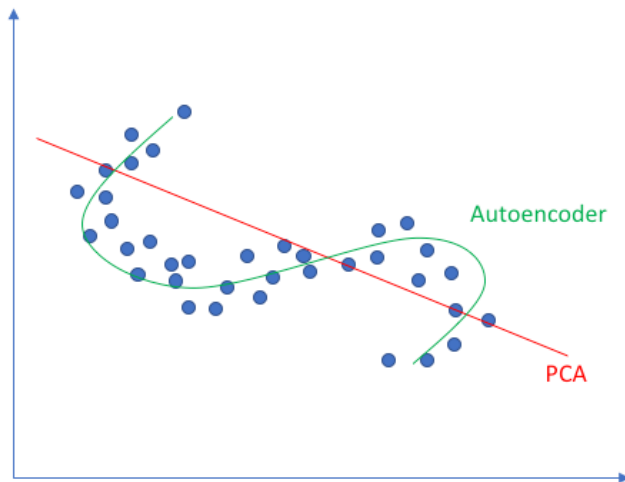
- e.g. $Loss = \sum_m \|x^{(m)} - DECODER(ENCODER(x^{(m)}))\|^2$
- Reconstruction: $x' = \sigma(W'\sigma(Wx + b) + d)$

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Autoencoders

Learning nonlinear manifolds

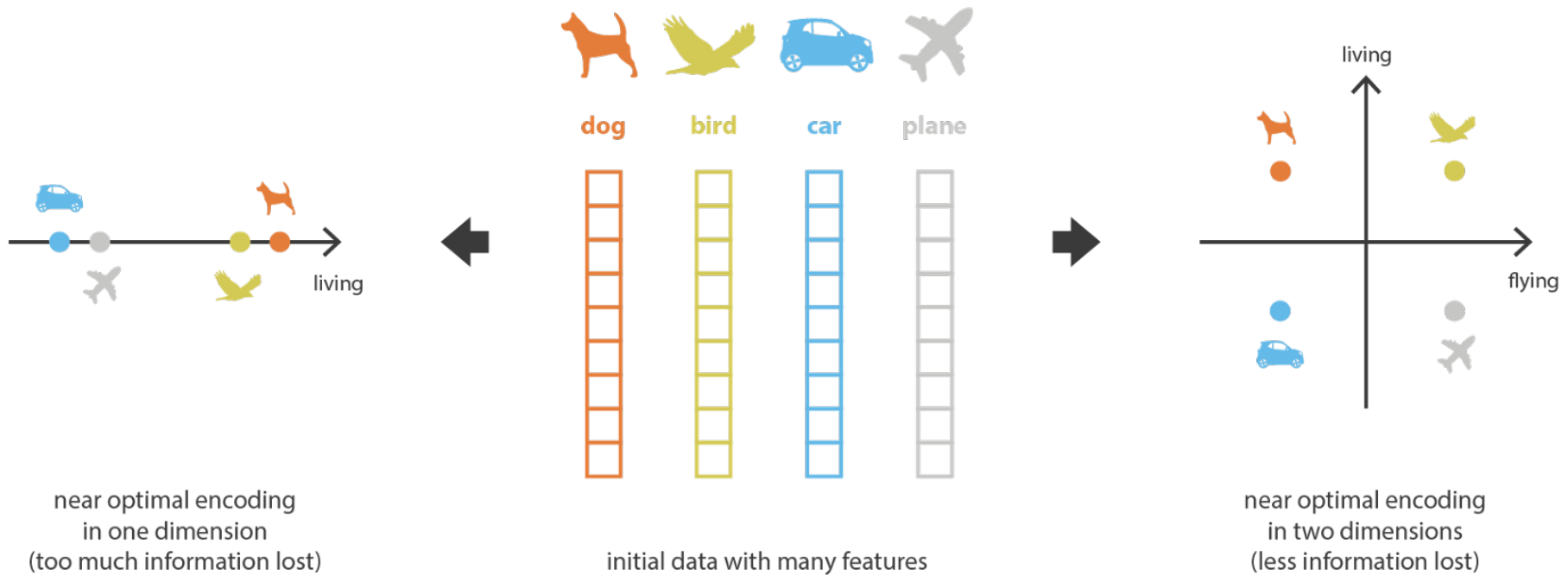
Linear vs nonlinear dimensionality reduction



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Autoencoders

Learning nonlinear manifolds



Different types of autoencoders

- Undercomplete autoencoders
 - hidden layer is smaller than the input dimensionality
- Overcomplete regularised autoencoders
 - Larger hidden layer size with the regularisation (to avoid overfitting and copying input to the output)

$$L(x, g(f(x))) + \Omega(h), \quad h = f(x)$$

- Sparse autoencoders, denosing autoencoders
- Deep autoencoders

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Sparse autoencoders

- Penalizing non-sparse solutions can be seen as adding latent variables with a prior and maximising likelihood

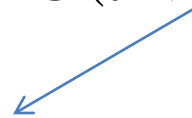
$$\log p_{\text{model}}(\mathbf{h}, \mathbf{x}) = \log p_{\text{model}}(\mathbf{h}) + \log p_{\text{model}}(\mathbf{x} | \mathbf{h})$$

for example: $p_{\text{model}}(\mathbf{h}) = \prod_i \frac{\lambda}{2} e^{-\lambda |h_i|} \Rightarrow \boxed{\Omega(\mathbf{h}) = \lambda \sum |h_i|}$

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

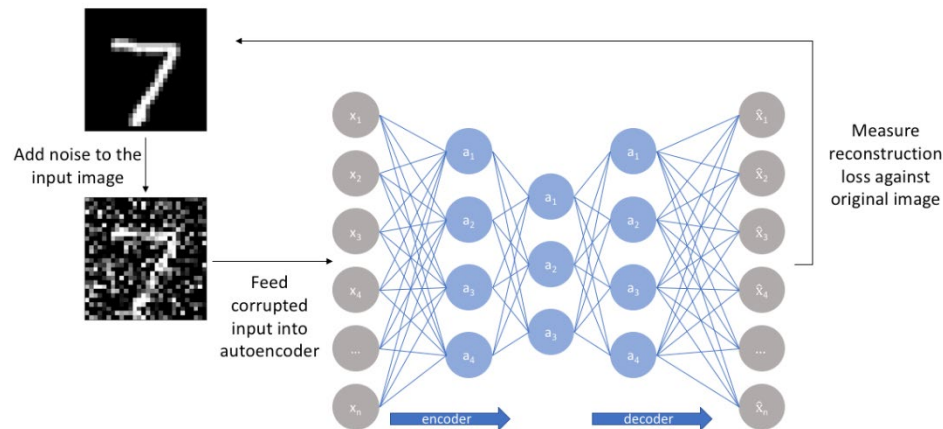
Denoising autoencoders

$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}}))), \quad h = f(\tilde{\mathbf{x}})$$



Corrupted copy of x

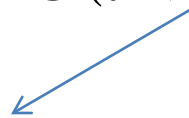
Autoencoders have to undo this corruption beyond simply coping the input.



- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Denoising autoencoders

$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}}))), \quad h = f(\tilde{\mathbf{x}})$$



Corrupted copy of x

Autoencoders have to undo this corruption beyond simply coping the input.

1. A training sample is sampled from the training data.
2. A corrupted version of the sample \mathbf{x} is drawn from some corruption process

$$C(\tilde{\mathbf{x}} \mid \mathbf{x} = \mathbf{s})$$

3. $(\mathbf{x}, \tilde{\mathbf{x}})$ is used as a training sample to estimate the autoencoder's reconstruction distribution $p_{\text{reconstruction}}(\tilde{\mathbf{x}} \mid \mathbf{x}) = p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h}), \quad \mathbf{h} = f(\tilde{\mathbf{x}})$

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Denoising autoencoders

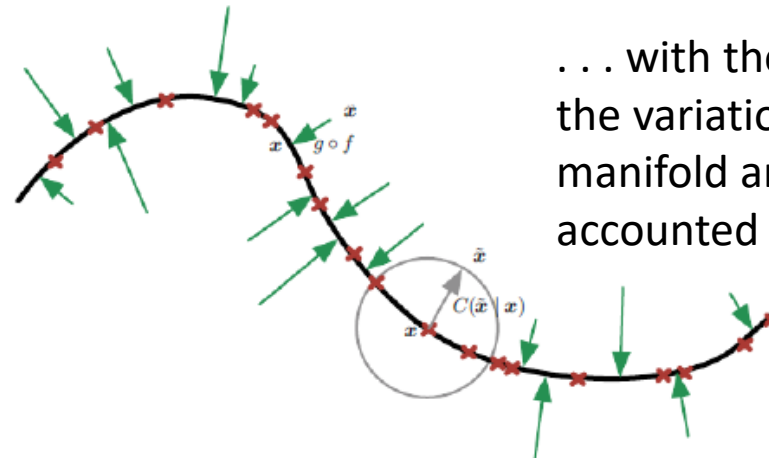
$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}}))), \quad h = f(\tilde{\mathbf{x}})$$



Corrupted copy of x

Autoencoders have to undo this corruption beyond simply coping the input.

Learning a *vector field* around
a *low-dimensional manifold* . . .



. . . with the principle that only
the variations tangent to the
manifold around \mathbf{x} should be
accounted for by changes in \mathbf{h}

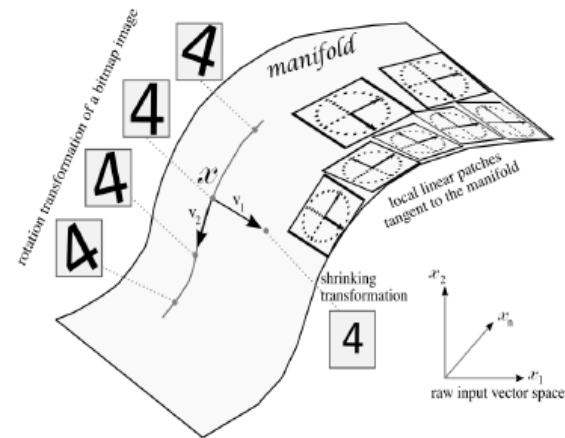
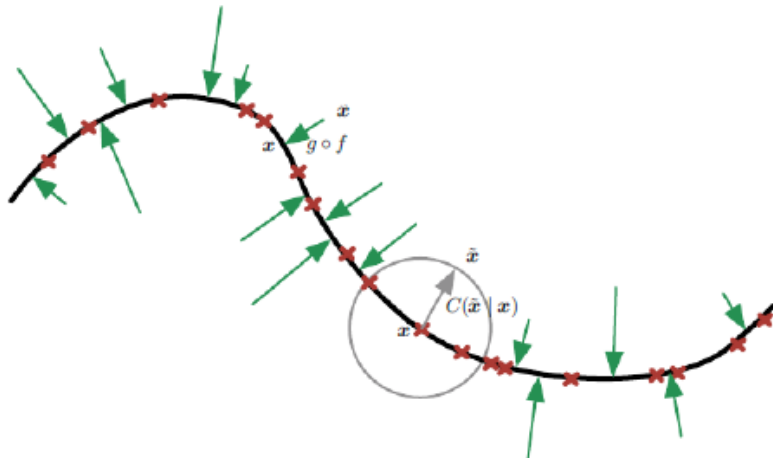
Goodfellow et al.

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Denoising autoencoders

When training autoencoders there is a **compromise**

- I. Need to approximately recover \mathbf{x} – *reconstruction* force
- II. Need to satisfy the regularization term – *regularisation* force.



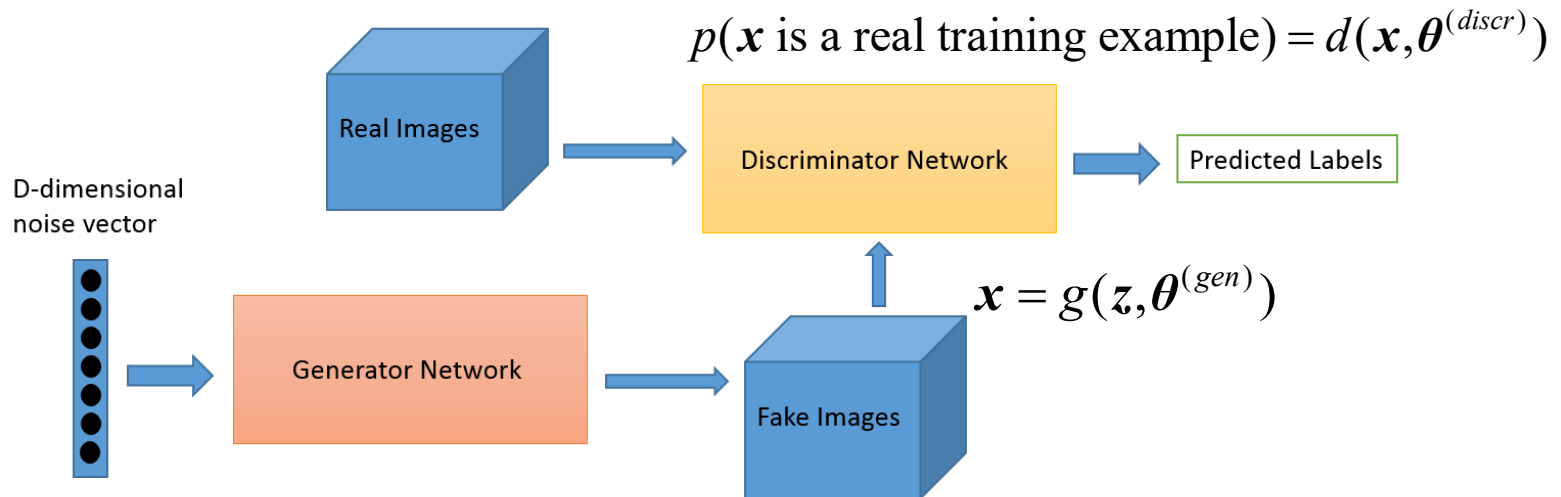
Goodfellow et al.

Generative modelling

- Unsupervised learning
 - no labels
 - the objective is to capture the hidden structure in the data
 - density estimation, clustering, feature learning, dimensionality reduction
- Generative approach allows even for generating samples
 - probabilistic in nature: learning **the joint $P(\mathbf{x}, \mathbf{y})$** -> ambitious
 - training data $\sim P_{data}(x)$ -> generated samples $\sim P_{model}(x)$
 - capable for uncovering underlying latent variables
 - could be used for many purposes, e.g. debiasing or outlier detection
 - some flagship examples of deep latent variable models: DBN, GAN, VAE

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Generative adversarial networks (GANs)



Discriminator network received some payoff v and the generator receives $-v$, so it is a zero-sum game. Both attempt to maximise their own payoff, so at the convergence:

$$g^* = \arg \min_g \max_d v(g, d)$$

$$v(\boldsymbol{\theta}^{(g)}, \boldsymbol{\theta}^{(d)}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log d(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_{\text{model}}} \log (1 - d(\mathbf{x}))$$

- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

Generative adversarial networks (GANs)

- What can GANs be used for?
 - data augmentation
 - creating art
 - image-to-image translation
 - creating images with higher resolution than the original



2014



2015



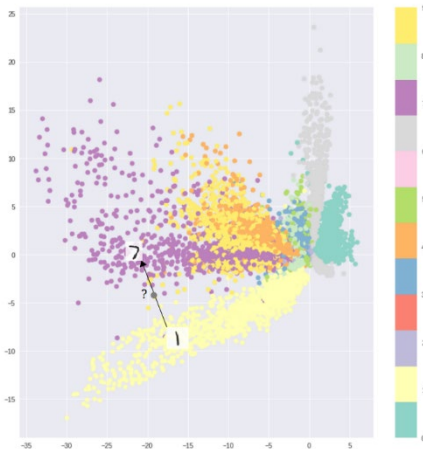
2016



2017

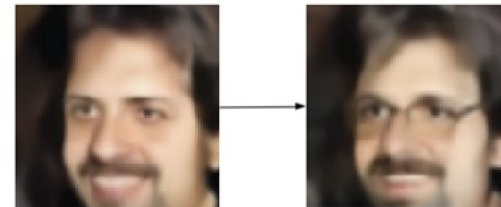
Variational autoencoders (VAE)

- The expectation from the generative model
 - capability to generate new samples from the learned distribution

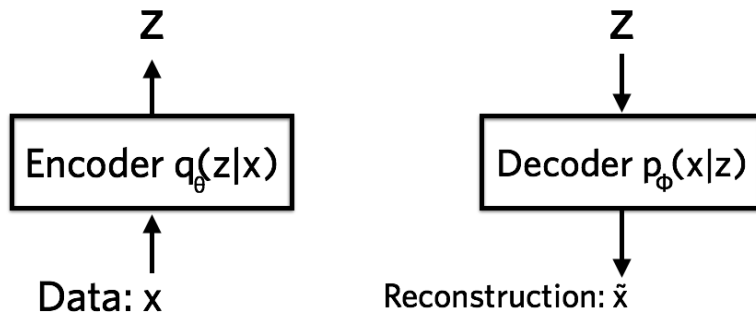


Classical autoencoders (that are not generative models) offer latent space that is often discontinuous and does not allow easy interpolation.

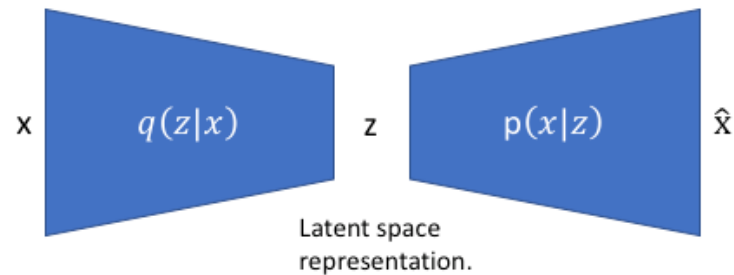
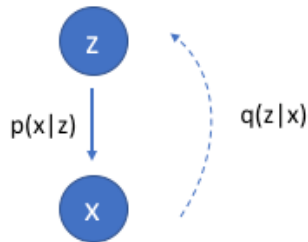
- alter and explore variations on the existing data



VAE mechanics

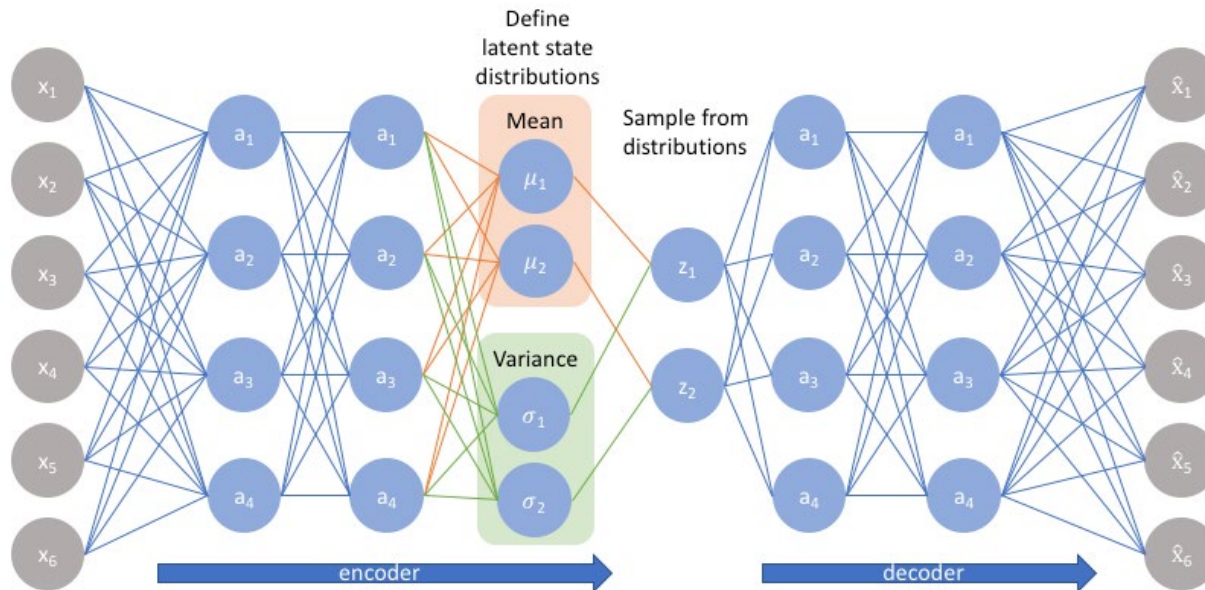


$$L(x, \hat{x}) + \sum_j KL(q_j(z|x) \| p(z))$$



VAE as a network implementation

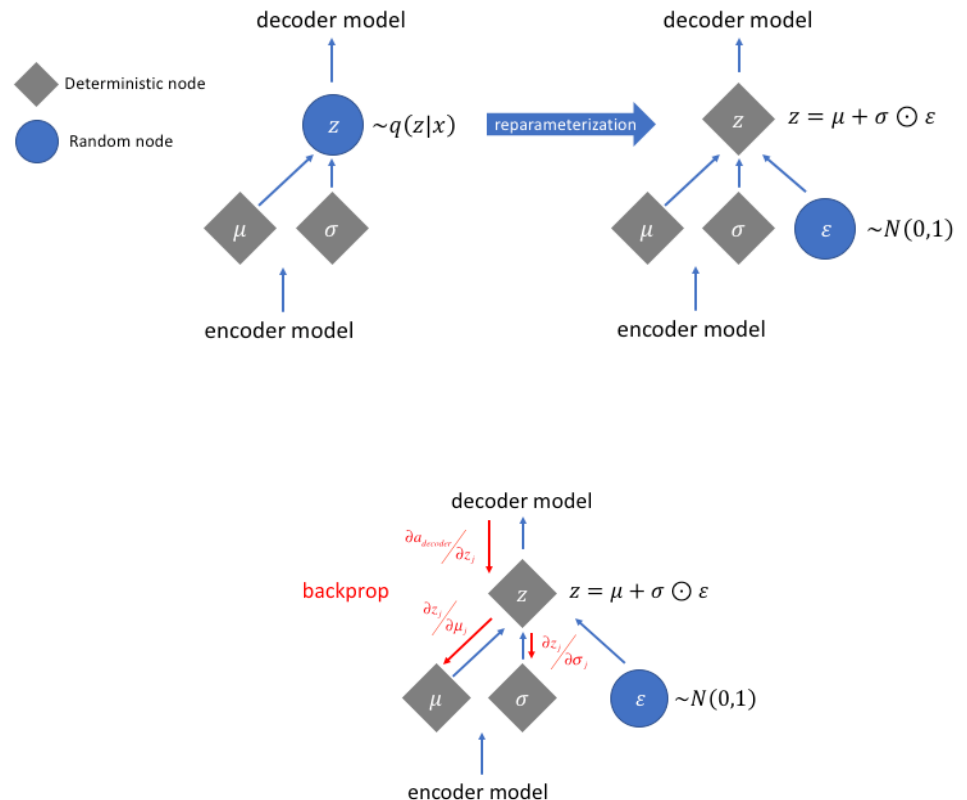
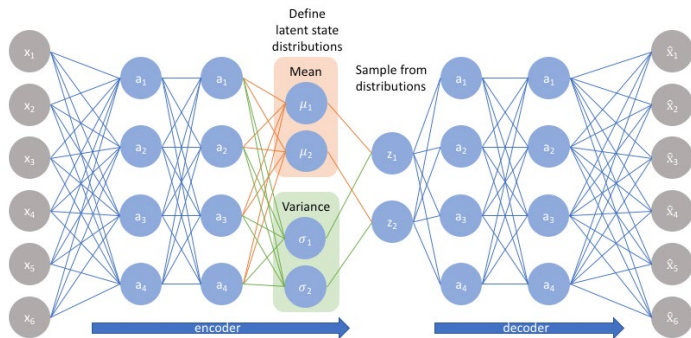
- How do we go about this in practice?



adapted from Jeremy Jordan

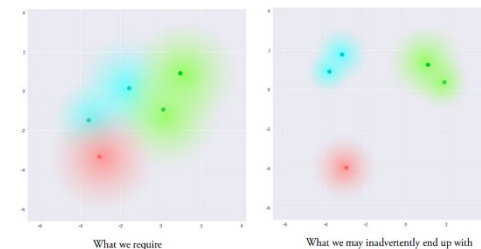
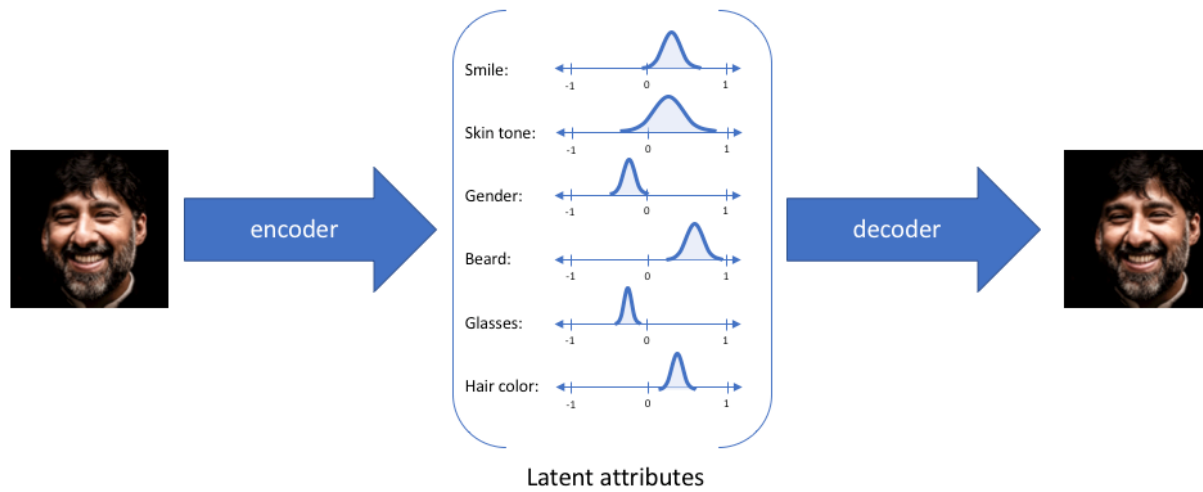
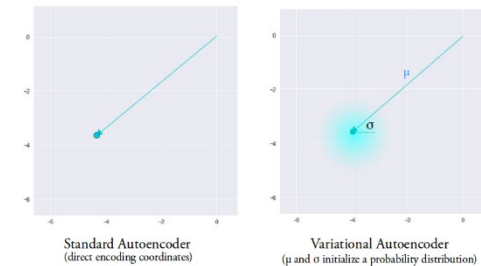
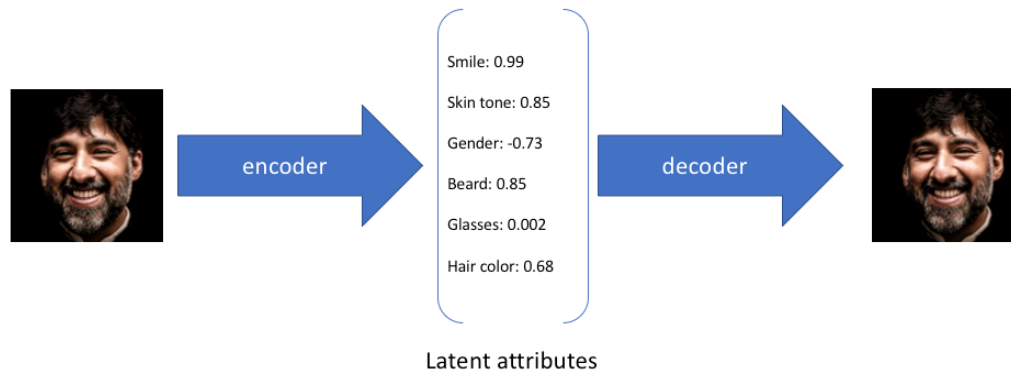
VAE as a network implementation

- How do we go about this in practice?



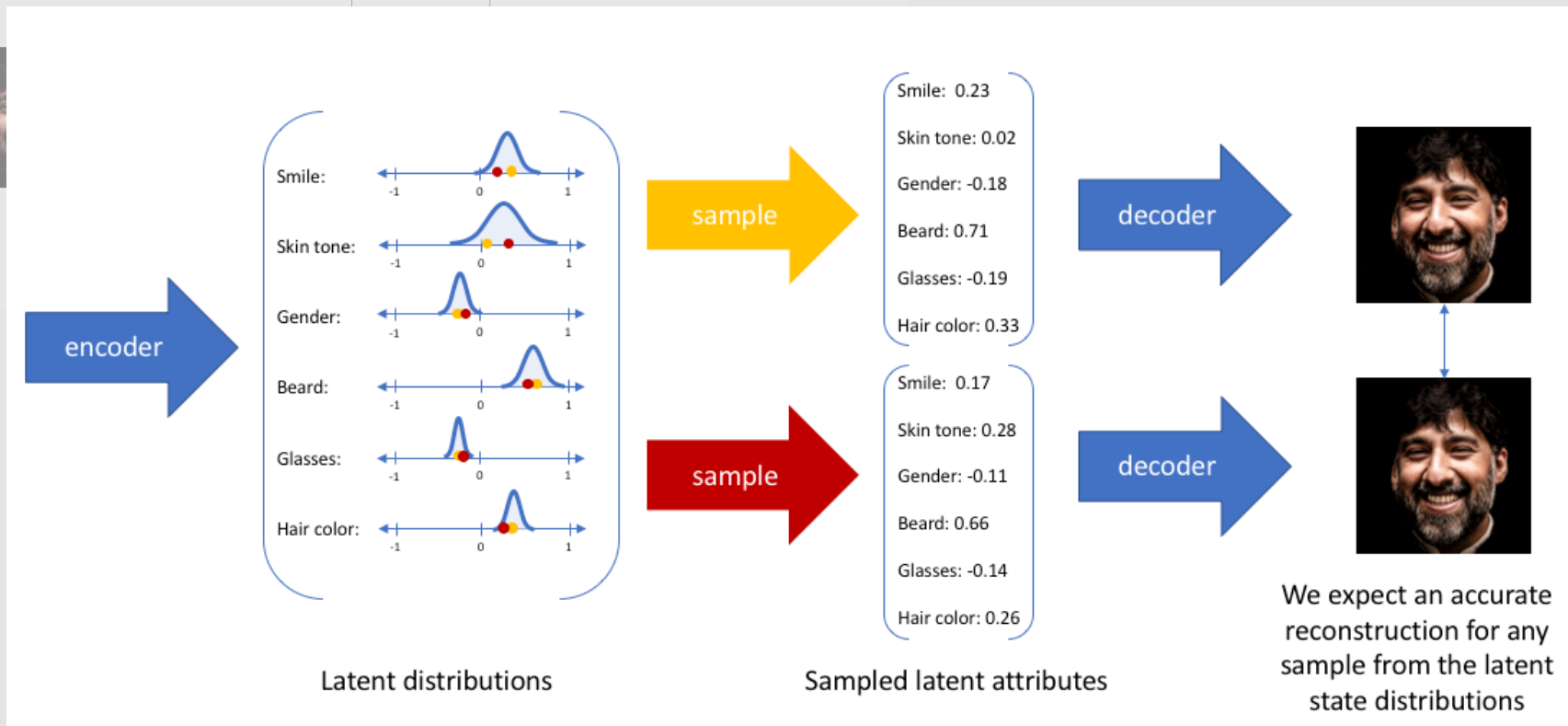
adapted from Jeremy Jordan

The intuition behind VAEs



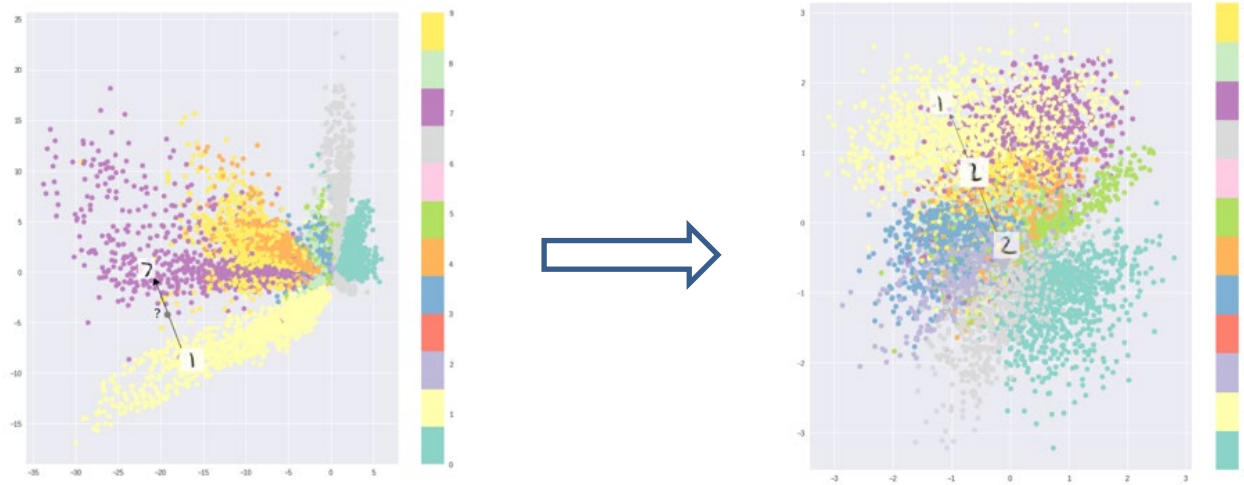
adapted from Irum Shafkat's (towards data science)

The intuition behind VAEs

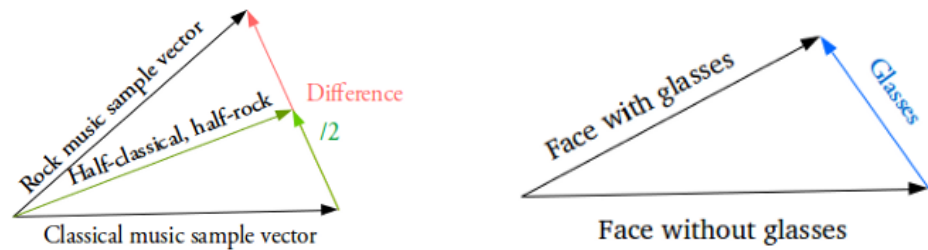


adapted from Irum Shafkat's (towards data science)

VAE implications



Vector arithmetic in the latent space



adapted from Irum Shafkat's (towards data science)

Impressive VAE applications



[Google Brain's Magenta's MusicVAE](#)
(Roberts et al., 2017)

Deep Feature Consistent Variational Autoencoder