



# DD2437 – Artificial Neural Networks and Deep Architectures (annda)

## Representation learning and deep generative models

Pawel Herman

Computational Science and Technology (CST)  
KTH Royal Institute of Technology

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

# Recap on deep networks & representation learning

- What is the motivation & inspiration for deep network architectures?
  - expressive power (*expressibility*) and compactness (*efficiency*) – exponential gain
  - hierarchical brain (cortex) organisation
  - multiple levels of abstraction
  - multiple levels of representations suitable for multi-task learning

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

# Recap on deep networks & representation learning

- What is the motivation & inspiration for deep network architectures?
  - expressive power (*expressibility*) and compactness (*efficiency*) – exponential gain
  - hierarchical brain (cortex) organisation
  - multiple levels of abstraction
  - multiple levels of representations suitable for multi-task learning
- Why are deep networks hard to train with gradient descent methods?
  - Unstable (mostly *vanishing*) gradient problem

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

# Recap on deep networks & representation learning

- What is the motivation & inspiration for deep network architectures?
  - expressive power (*expressibility*) and compactness (*efficiency*) – exponential gain
  - hierarchical brain (cortex) organisation
  - multiple levels of abstraction
  - multiple levels of representations suitable for multi-task learning
- Why are deep networks hard to train with gradient descent methods?
  - Unstable (mostly *vanishing*) gradient problem
- How was this challenge originally addressed?
  - *greedy layer-wise unsupervised pre-training*: stacked autoencoders and DBNs
  - two-phase learning: *unsupervised pre-training* and *supervised tuning* with gradient descent based optimisation (the entire network or the top layers only).

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

# Recap on deep networks & representation learning

- Hypotheses about the role of unsupervised pre-training (still not well understood)
  - *regularisation*: an implicit penalisation term, minimisation of variance
  - *optimisation*: good initial condition for optimisation (areas that could otherwise be difficult to find)

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

# Recap on deep networks & representation learning

- Hypotheses about the role of unsupervised pre-training (still not well understood)
  - *regularisation*: an implicit penalisation term, minimisation of variance
  - *optimisation*: good initial condition for optimisation (areas that could otherwise be difficult to find)
- BUT: Contemporary trend to avoid pre-training
  - employ *ReLU* units (less risk for overfitting and dealing with unstable gradients)
  - new regularisation approaches: *dropout*, *batch normalisation*

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

# Recap on deep networks & representation learning

- Hypotheses about the role of unsupervised pre-training (still not well understood)
  - *regularisation*: an implicit penalisation term, minimisation of variance
  - *optimisation*: good initial condition for optimisation (areas that could otherwise be difficult to find)
- BUT: Contemporary trend to avoid pre-training
  - employ *ReLU* units (less risk for overfitting and dealing with unstable gradients)
  - new regularisation approaches: *dropout*, *batch normalisation*
- Focus on variations of CNN and LSTM architectures
  - mainly application-driven developments
  - rebirth of interest in stacked autoencoders and less focus on DBNs

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

# Recap on deep networks & representation learning

- Why do we think DL works so well?
  - “*cheap learning*” (lower dimensional nature of problems to model with inherent constraints)
  - “*no-flattening*” theorems (huge flattening costs when training: accuracy vs compute time)
  - capturing intrinsic hierarchical structure of the physical world
  - ability to *learn* rich *representations* of data



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

# Lecture overview

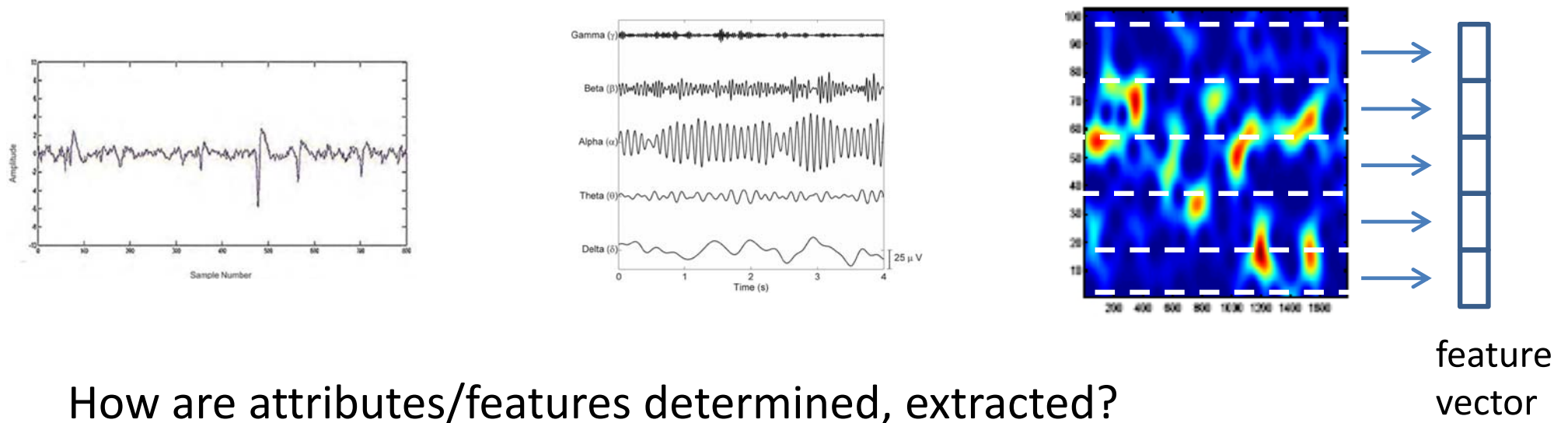
1. Data representations: desirable characteristics and key concepts.
2. The notion of hierarchical and distributed sparse representations.
3. Learning representations in deep neural networks.
4. Transfer, multi-task/modal learning.
5. Generative deep learning models

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

# Data representations

- Multiple ways of representing information – what is the difference? Why should we care?

From low-level data description to higher-order representations

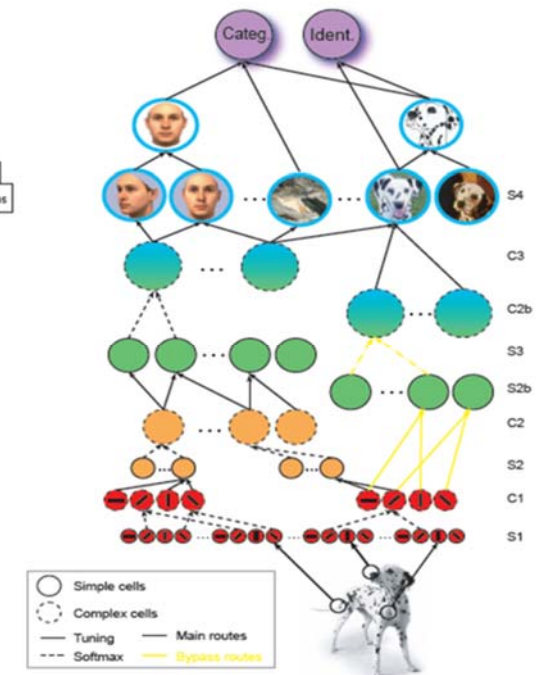


- # Data representations

- # Data parameterisation

003\_1\_1.inorm.bindata

003\_1\_1.dct.bindata

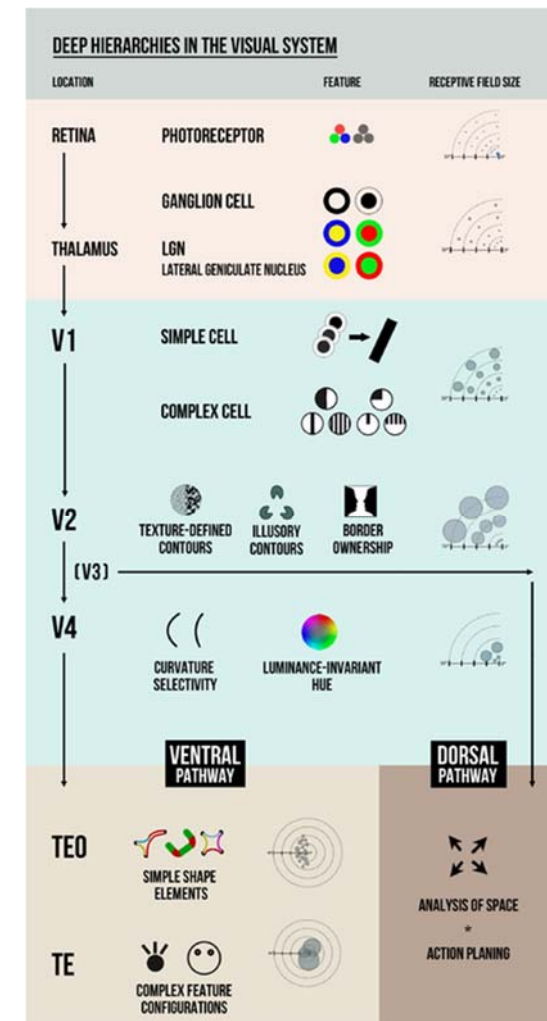


KTH

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

# Representations in the brain as an inspiration

- Sensory information represented by neural activity
  - neurons with different response properties (selectivity, tuning curves)
  - *distributed* nature of neural representations in populations vs grandmother cell concept
  - *sparseness*, redundancy
- Hierarchical representations
  - sensory pathways are organised into *hierarchies*
  - hierarchy of *abstraction* levels

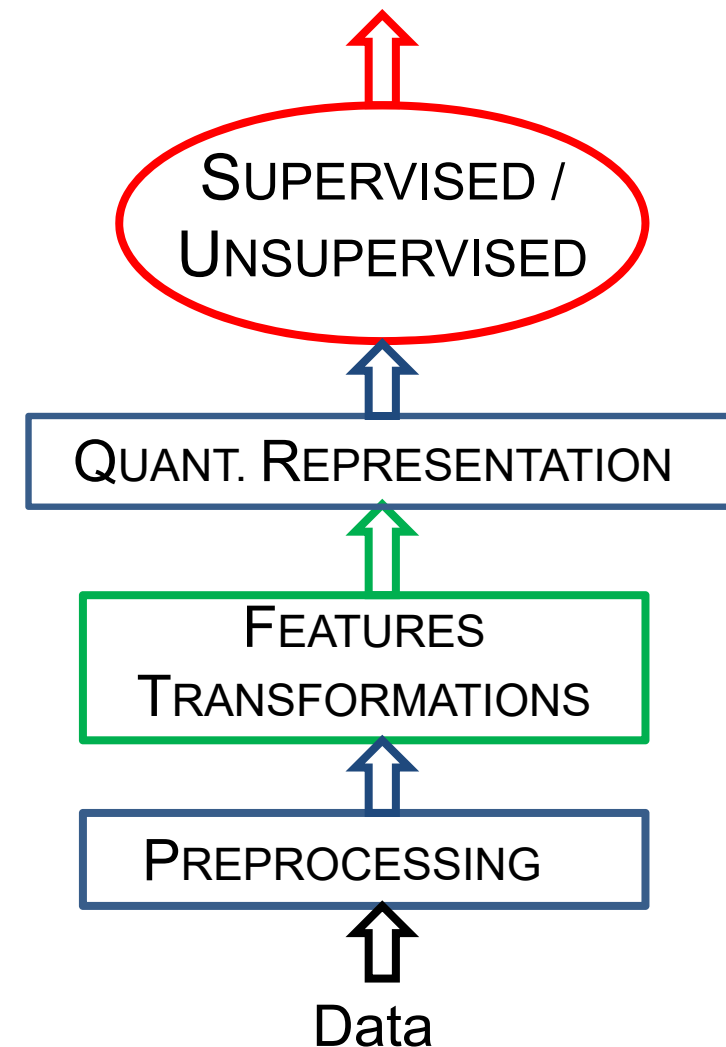


Wikibooks

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

# Representation learning problem

- Importance in the machine learning or pattern recognition context
- there is a trade-off between minimising “information” loss and obtaining “nice” properties
- What makes representation good?  
What is desirable/useful information?

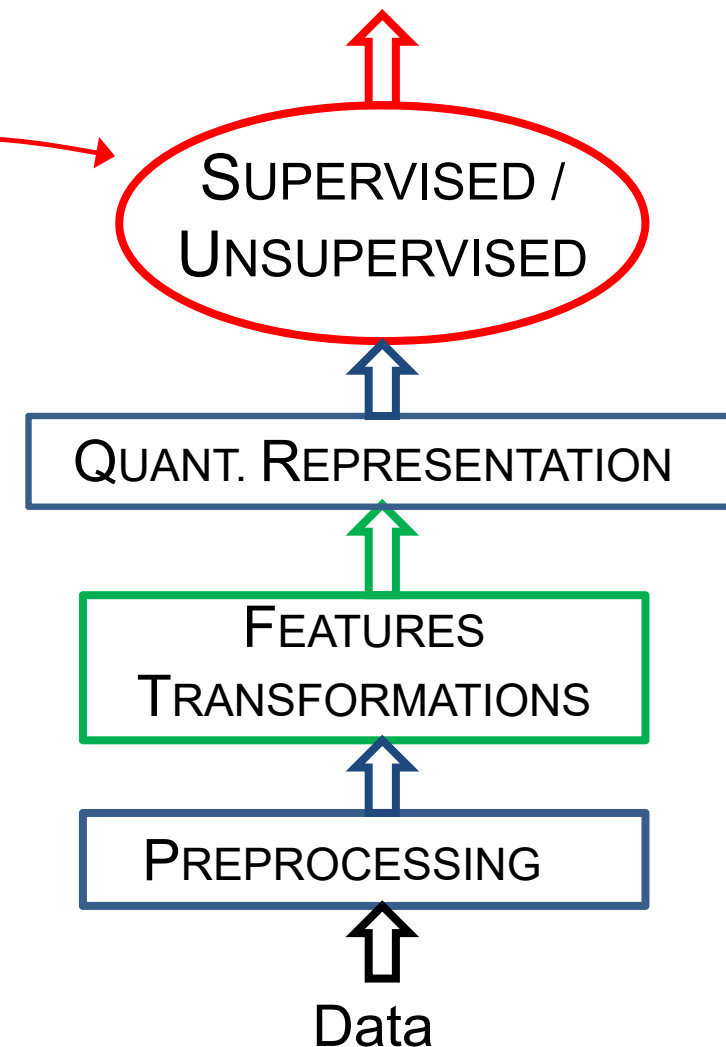


- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Representation learning problem

- Importance in the machine learning or pattern recognition context
- there is a trade-off between minimising “information” loss and obtaining “nice” properties
- What makes representation good?  
What is desirable/useful information?

Facilitate the subsequent learning task,  
maximise its performance



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

# Representation learning

- Computational perspective: disentangling unknown factors causing relevant variation in the data
  - *causes explain* the observed data (discriminative context, both unsupervised and supervised aspects)

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

# Representation learning

- Computational perspective: disentangling unknown factors causing relevant variation in the data
  - *causes explain* the observed data (discriminative context, both unsupervised and supervised aspects)

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.

$$p(\mathbf{h}|\mathbf{x}) = p(\mathbf{x}|\mathbf{h}) p(\mathbf{h})$$

$$\text{Ideally: } p(\mathbf{h}) = \prod_i p(h_i)$$



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

# Representation learning

- Computational perspective: disentangling unknown factors causing relevant variation in the data
  - *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
  - *density estimation* – learn probability distribution for data with the use of latent variables (PCA, ICA etc.) -> explain data
  - $P(\text{data} | \text{latent var})$  for recognition and  $P(\text{latent var} | \text{data})$  for generation

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

# Representation learning

- Computational perspective: disentangling unknown factors causing relevant variation in the data
  - *causes explain* the observed data (discriminative context, both unsupervised and supervised aspects)
  - factors in *discriminative context*
- Probabilistic
  - *density estimation* with the use of latent variables
  - $P(\text{data} | \text{latent var})$  for recognition and  $P(\text{latent var} | \text{data})$  for generation

Can we implicitly guide the *unsupervised* learning to discover features corresponding to underlying/causal factors?

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

# Distributed representations

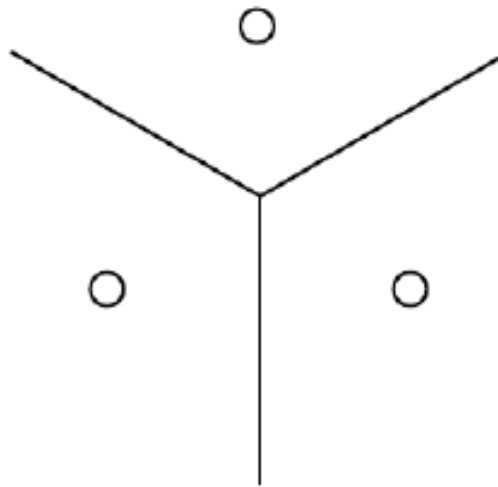
Information is distributed across many units that account for information about features that are not mutually exclusive.....

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

# Distributed representations

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

... unlike in clustering with distinct regions where *local generalisation* is observed.



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

# Distributed representations

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

... unlike in clustering 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).

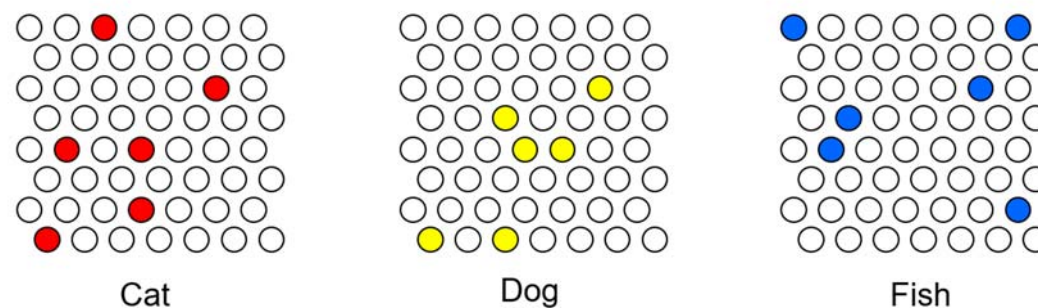
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Distributed representations

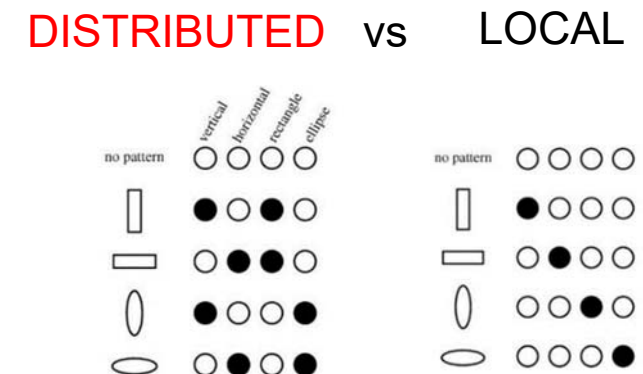
Information is distributed across many units that account for information about features that are not mutually exclusive.....

... unlike in clustering 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).



Generalisation due to shared attributes and semantic proximity.



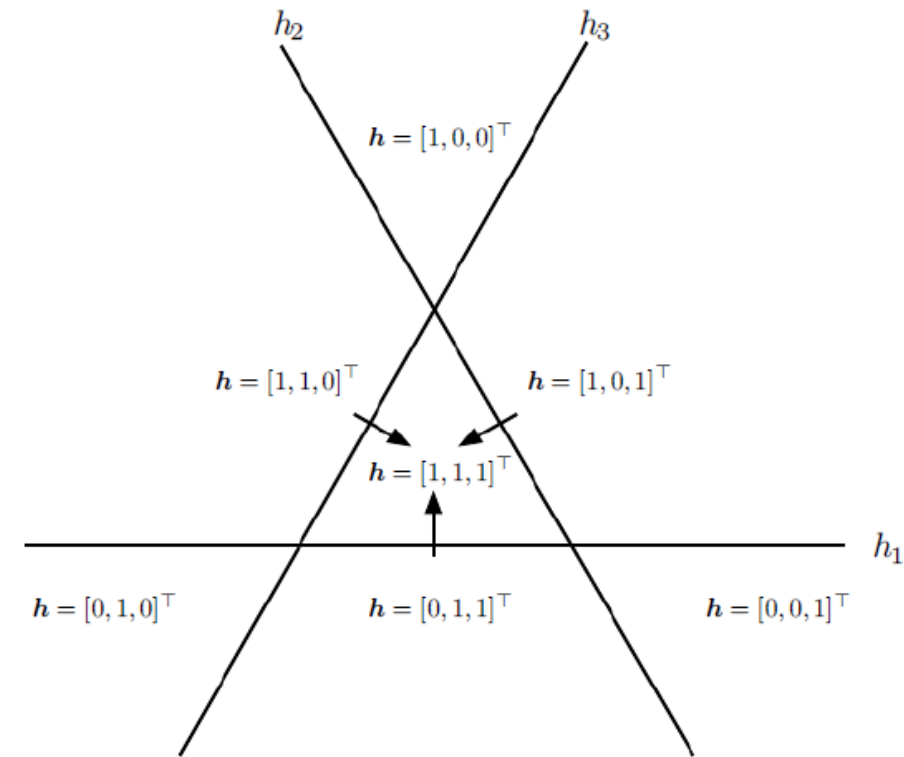
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# The power of distributed representations

In summary:

- expressiveness ( $n$  features with  $k$  values each can describe  $k^n$  concepts)
- the combination of powerful distributed representations with weak classifiers could be a strong regulariser

fault tolerance



Goodfellow et al.  
Bengio et al., 2009

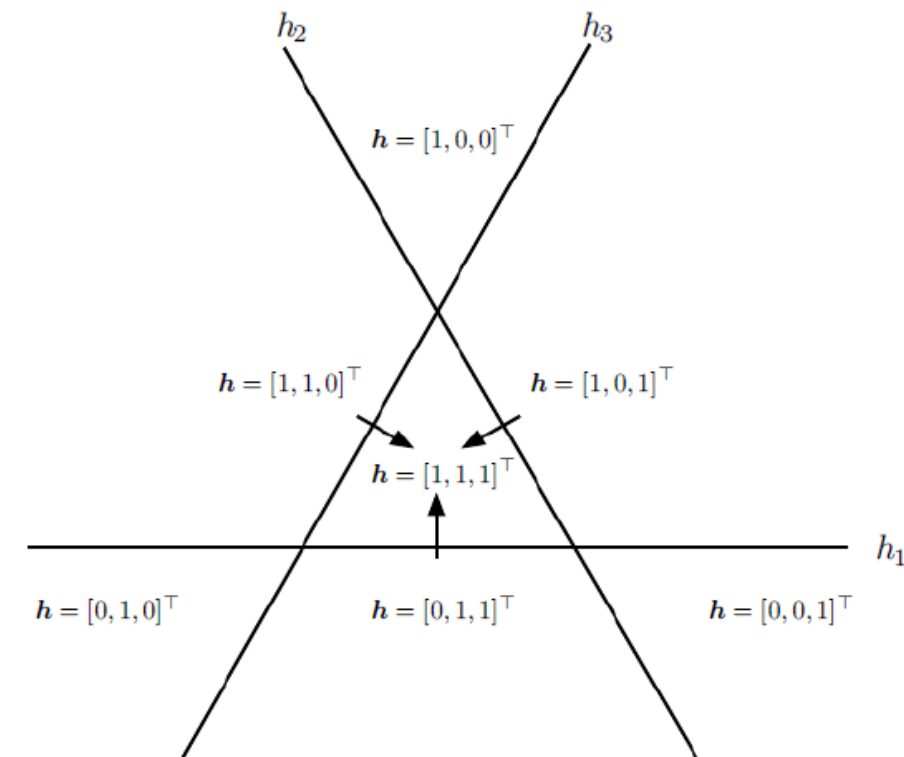
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# The power of distributed representations

In summary:

- expressiveness ( $n$  features with  $k$  values each can describe  $k^n$  concepts)
- the combination of powerful distributed representations with weak classifiers could be a strong regulariser
- similarity (topological) space with a distributed code – semantically close objects are close in distance
- generalisation due to shared attributes

content addressability



Goodfellow et al.  
Bengio et al., 2009

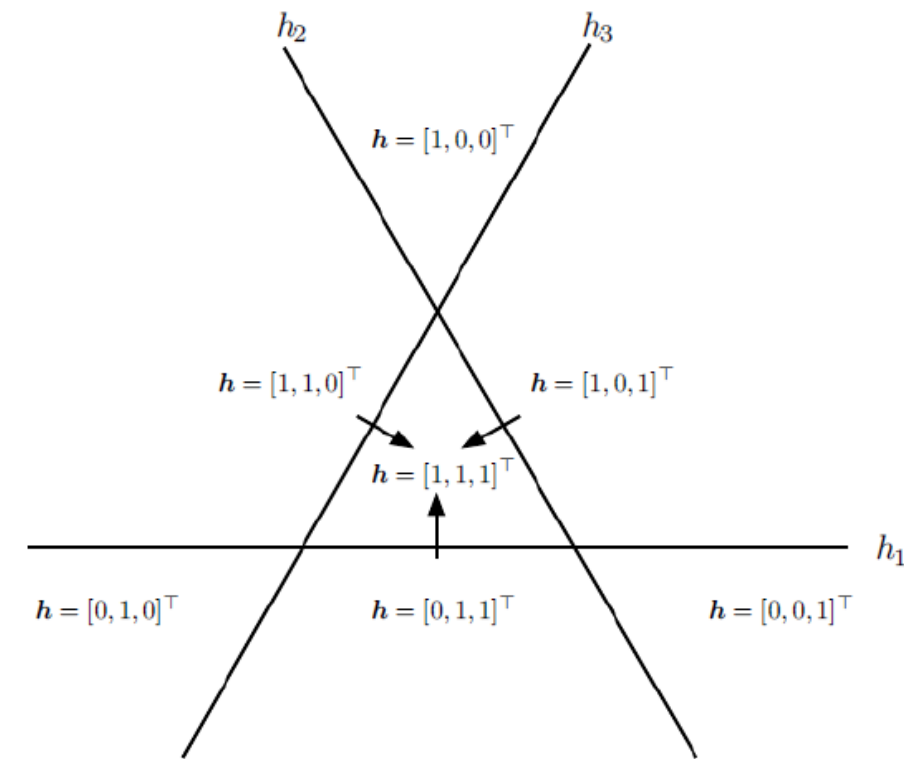


- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# The power of distributed representations

In summary:

- expressiveness ( $n$  features with  $k$  values each can describe  $k^n$  concepts)
- the combination of powerful distributed representations with weak classifiers could be a strong regulariser
- similarity (topological) space with a distributed code – semantically close objects are close in distance
- generalisation due to shared attributes
- in line with the idea that hidden units can learn to represent the underlying causal factors as different variables (here: directions in the representation space)



Goodfellow et al.  
Bengio et al., 2009

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

# Sparse vs dense representations

## Sparse representations

- promoting memory capacity
- orthogonalisation/decorrelation
- “metabolic” efficiency
- neural selectivity (vs coarse coding with broad tuning)

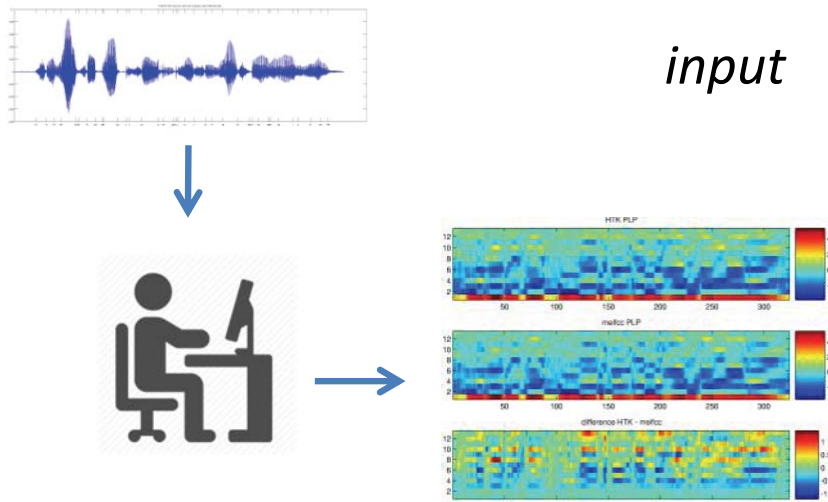
sparse not distributed	not sparse distributed	sparse distributed
0 .2 0 0 0	.1 .8 .7 .5 .7	0 .8 0 .5 0
0 0 0 0 .1	.8 .9 .6 .2 .4	0 0 .6 0 .4
0 0 0 .4 0	.3 .1 .6 .3 .3	.3 0 0 .3 0

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

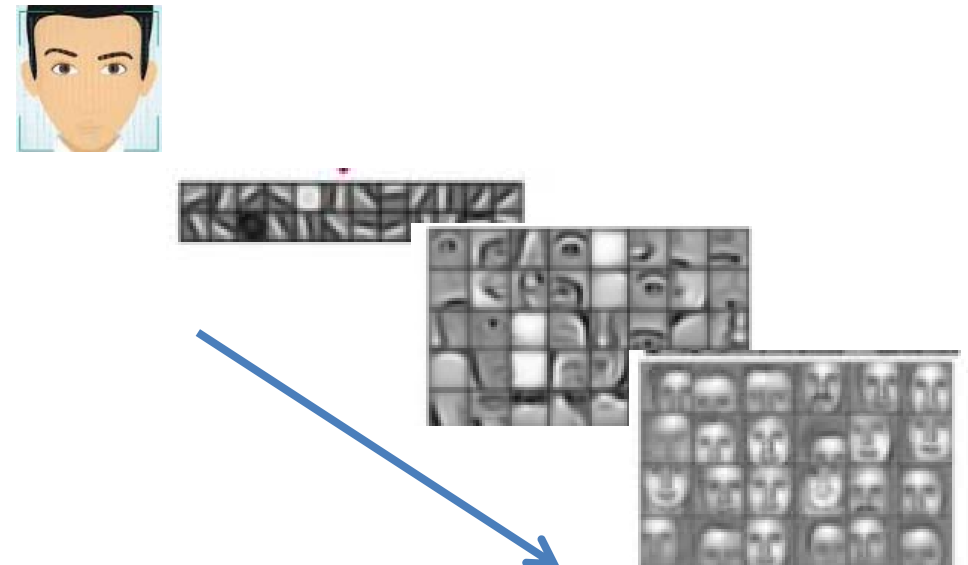
# Representation learning in deep models

- The composition of multiple non-linear transformations with the expectation for the hierarchy of abstraction levels

Hand-engineered features in a traditional pattern recognition approach



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

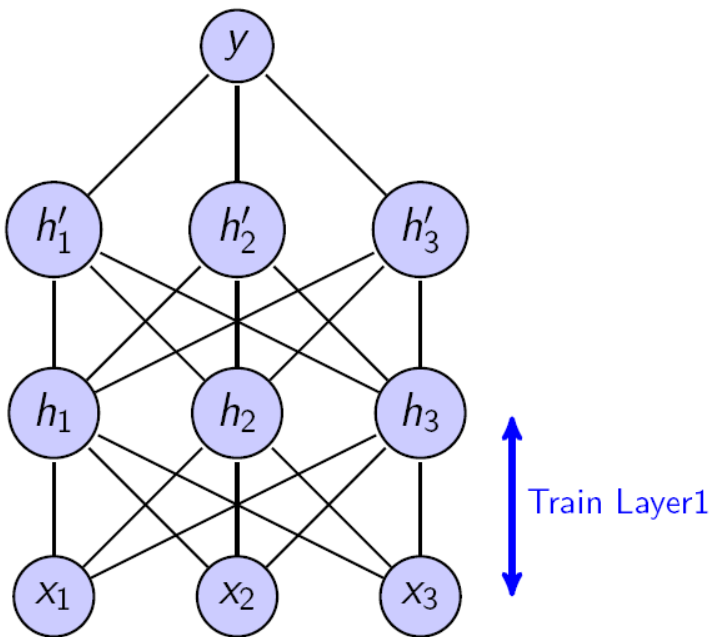


*features, representations*

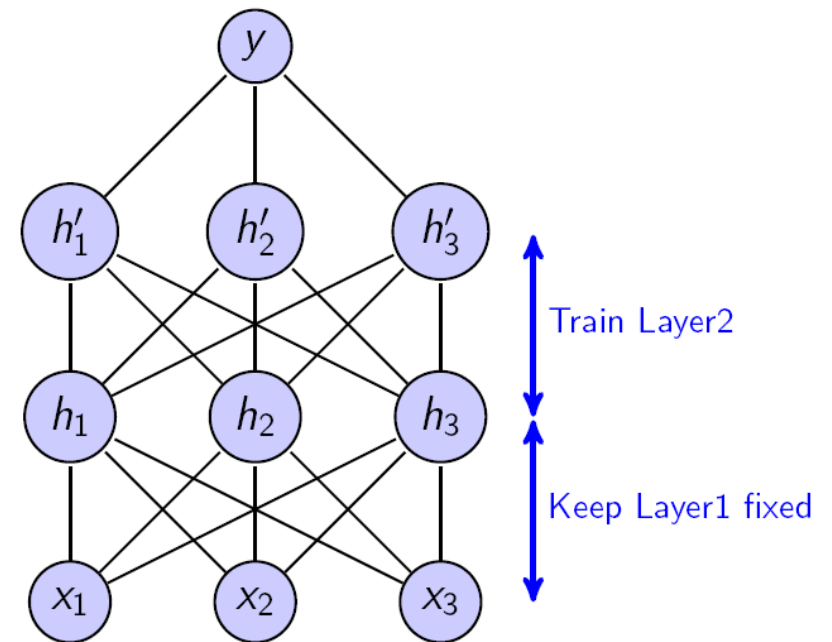
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Representation learning in deep models

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



Single layer at a time

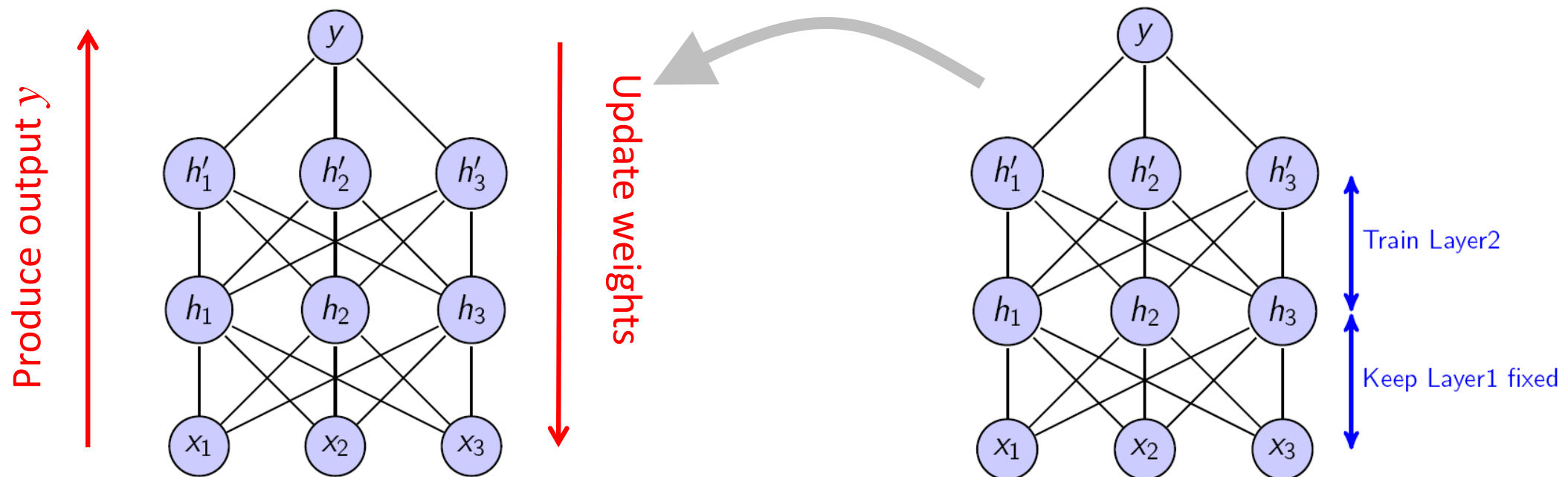


Train another layer while keeping the lower layer fixed

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

# Representation learning in deep models

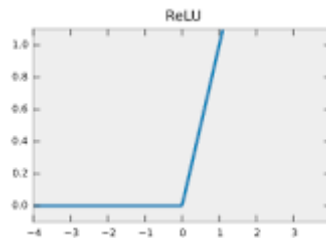
- The concept of layer-by-layer pretraining
  - greedy layer-wise unsupervised representation learning
    - intuitively, learning about the input distribution should help in learning the *mapping* between the input and output space
    - BUT having two separate phases has *disadvantages*



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

# Representation learning in deep models

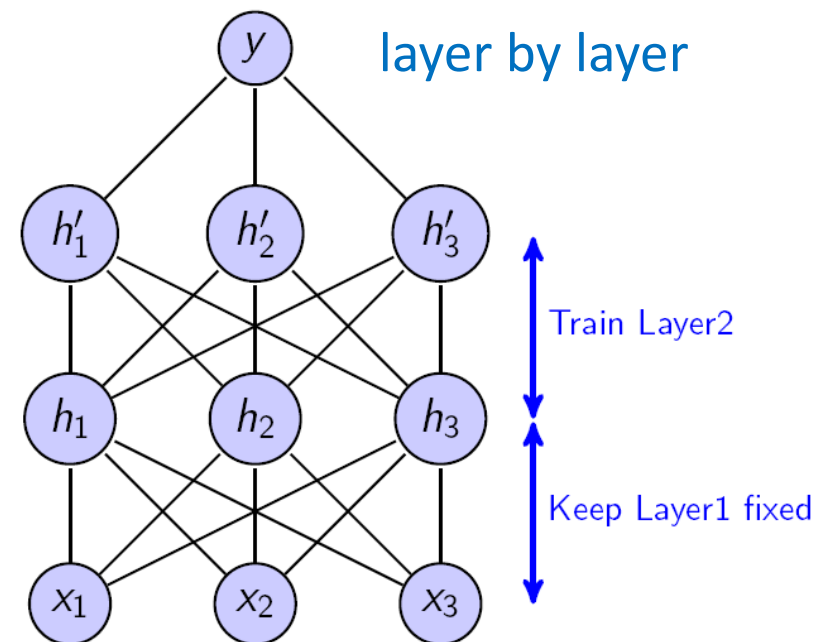
- The concept of layer-by-layer pretraining
  - greedy layer-wise unsupervised representation learning
    - intuitively, learning about the input distribution should help in learning the *mapping* between the input and output space
    - BUT having two separate phases has *disadvantages*
  - ULTIMATELY, the approach with unsupervised pretraining is largely **abandoned** (except word embeddings in NLP)
    - new regularisation techniques: dropout, batch normalisation
    - smaller datasets -> Bayesian methods
    - units with ReLU activation



- 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
    - restricted BM (RBM), 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*)



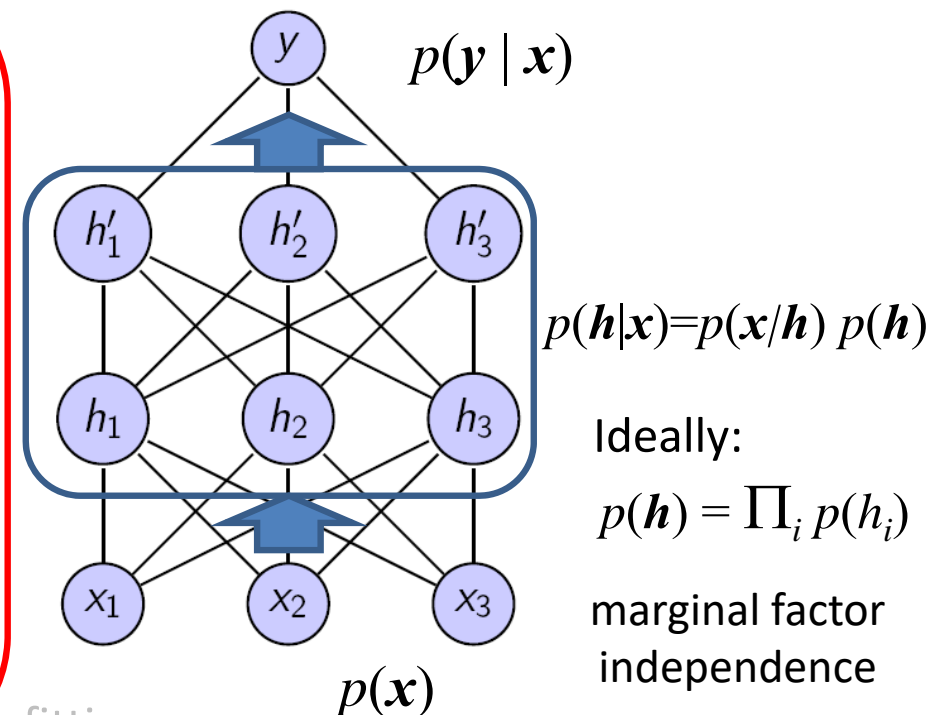
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# 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.



- 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

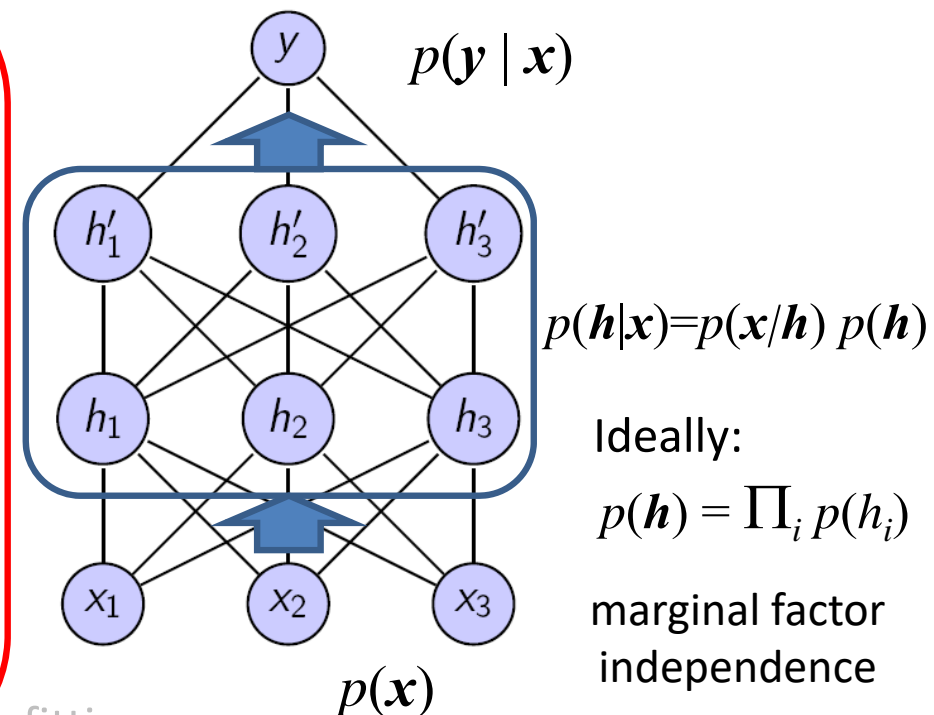
# Representation learning in deep models

- 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*)

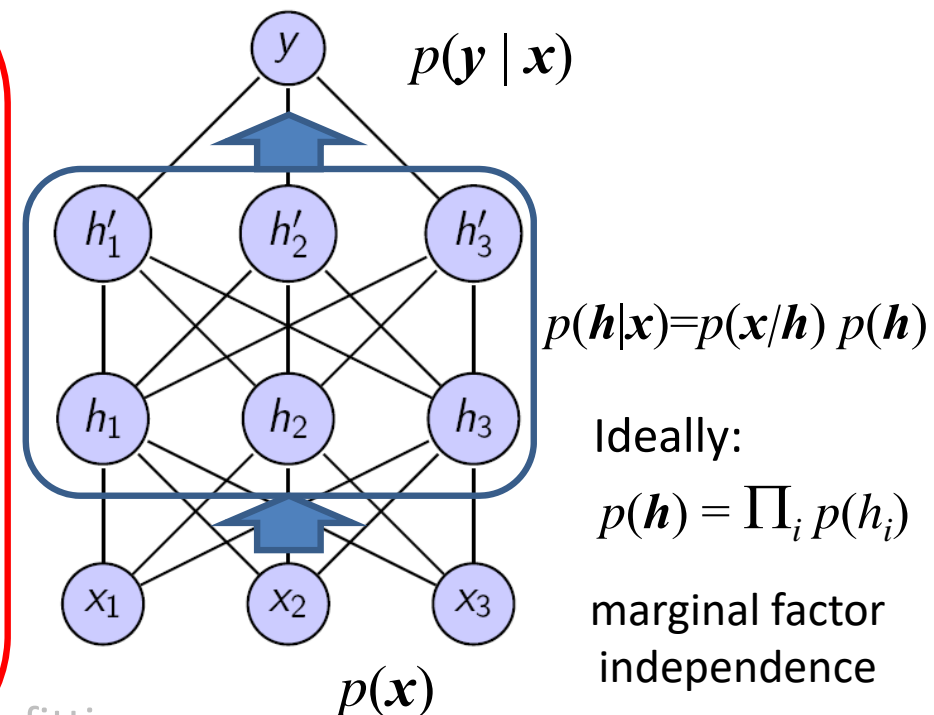
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# 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.

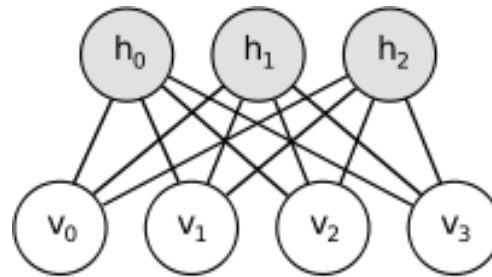


- as a regulariser, it urges the learning algorithm to discover **features that explain underlying causes** that generate the data (also, causal factors often remain *invariant*)

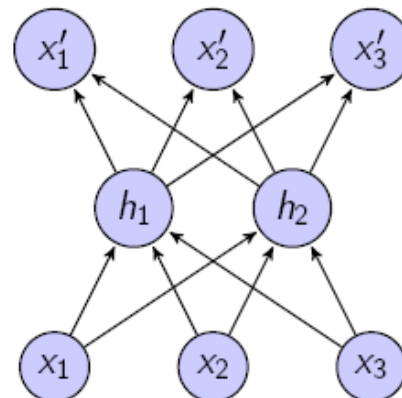
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Computational blocks for learning representations

- Two key approaches to greedy layer-wise pretraining
  - regularized Boltzmann machines (RBMs)



- autoencoders

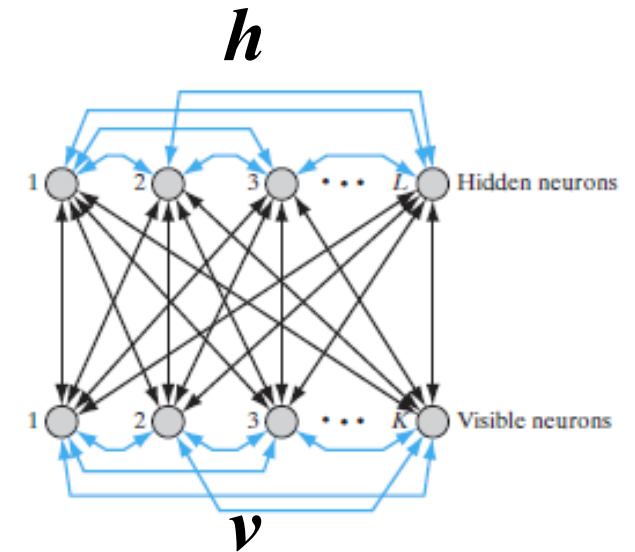


- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Recap on Boltzmann machine

$$E = -\frac{1}{2} \vec{x}^T \mathbf{W} \vec{x} = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} x_i x_j$$

$$P(\vec{x} | \mathbf{W}) = \frac{e^{-E}}{Z} = \frac{1}{Z(\mathbf{W})} \exp\left(\frac{1}{2} \vec{x}^T \mathbf{W} \vec{x}\right)$$



$$\begin{aligned} \mathbf{v}^{(p)} &= \mathbf{x}^{(p)} \\ \Downarrow \\ \mathbf{y}^{(p)} &= [\mathbf{x}^{(p)}, \mathbf{h}] \end{aligned}$$

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

# Recap on Boltzmann machine

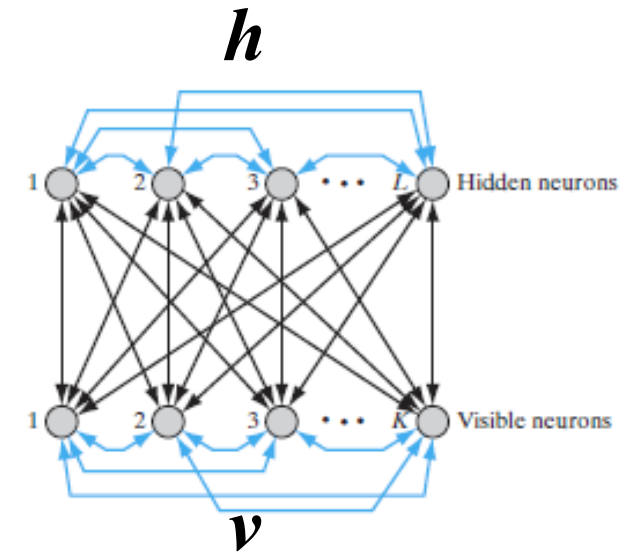
$$E = -\frac{1}{2} \vec{x}^T \mathbf{W} \vec{x} = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} x_i x_j$$

$$P(\vec{x} | \mathbf{W}) = \frac{e^{-E}}{Z} = \frac{1}{Z(\mathbf{W})} \exp\left(\frac{1}{2} \vec{x}^T \mathbf{W} \vec{x}\right)$$

The idea is to maximise log-likelihood,

$$L(\mathbf{W}) = \log(P(\mathbf{X}) | \mathbf{W})$$

$$\Delta w_{ji} = \varepsilon \frac{\partial L(\mathbf{W})}{\partial w_{ji}} \propto \langle y_i y_j \rangle_{\text{data}} - \langle y_i y_j \rangle_{\text{model}}$$



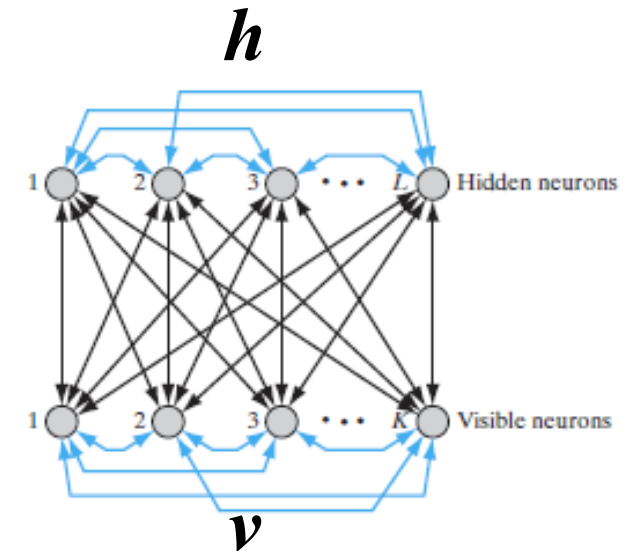
$$\begin{aligned} \mathbf{v}^{(p)} &= \mathbf{x}^{(p)} \\ \Downarrow \\ \mathbf{y}^{(p)} &= [\mathbf{x}^{(p)}, \mathbf{h}] \end{aligned}$$

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

# Recap on Boltzmann machine

$$E = -\frac{1}{2} \vec{x}^T \mathbf{W} \vec{x} = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{i,j} x_i x_j$$

$$P(\vec{x} | \mathbf{W}) = \frac{e^{-E}}{Z} = \frac{1}{Z(\mathbf{W})} \exp\left(\frac{1}{2} \vec{x}^T \mathbf{W} \vec{x}\right)$$



The idea is to maximise log-likelihood,

$$L(\mathbf{W}) = \log (P(\mathbf{X}) | \mathbf{W})$$

$$\Delta w_{ji} = \varepsilon \frac{\partial L(\mathbf{W})}{\partial w_{ji}} \propto \langle y_i y_j \rangle_{\text{data}} - \langle y_i y_j \rangle_{\text{model}}$$

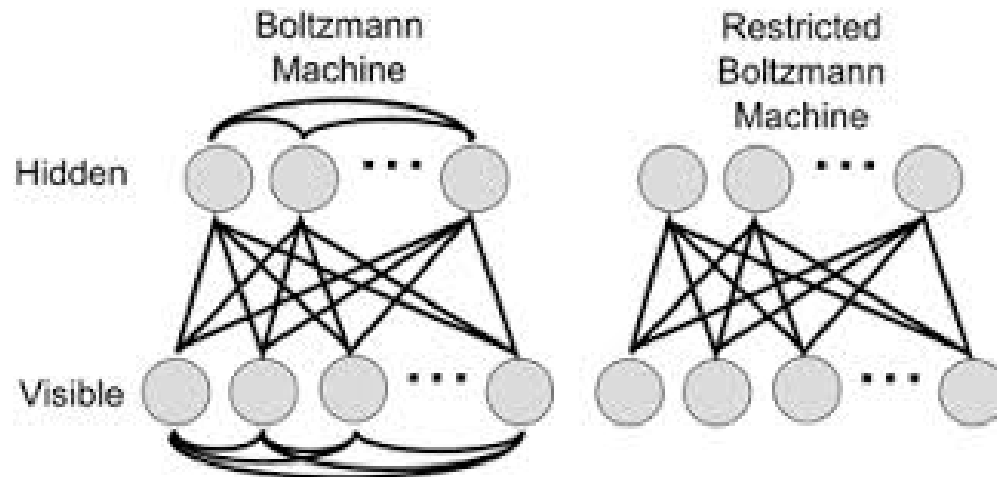
**positive:** “Hebbian learning”

**negative:** “Hebbian forgetting”

$$\begin{aligned} \mathbf{v}^{(p)} &= \mathbf{x}^{(p)} \\ \Downarrow \\ \mathbf{y}^{(p)} &= [\mathbf{x}^{(p)}, \mathbf{h}] \end{aligned}$$

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

# Restricted Boltzmann machine (RBM)



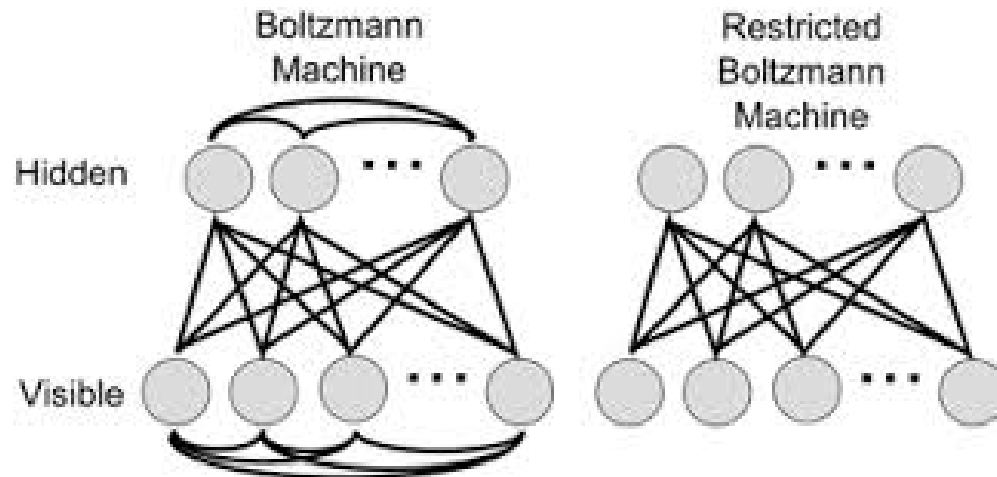
Visible and hidden units are conditionally independent given one another

$$p(\mathbf{h} | \mathbf{v}) = \prod_i p(h_i | \mathbf{v})$$

$$p(\mathbf{v} | \mathbf{h}) = \prod_j p(v_j | \mathbf{h})$$

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

# Restricted Boltzmann machine (RBM)



Visible and hidden units are conditionally independent given one another

$$p(\mathbf{h} | \mathbf{v}) = \prod_i p(h_i | \mathbf{v})$$

$$p(\mathbf{v} | \mathbf{h}) = \prod_j p(v_j | \mathbf{h})$$

Following the same principle of maximising log likelihood by means of gradient ascent, one obtains:

$$\Delta w_{ji} = \varepsilon \frac{\partial L(\mathbf{W})}{\partial w_{ji}} = \varepsilon \left( \langle v_j h_i \rangle_{\text{data}} - \langle v_j h_i \rangle_{\text{model}} \right)$$



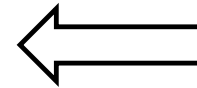
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Restricted Boltzmann machine (RBM)

Visible and hidden units are conditionally independent given one another

$$P(h_i = 1 | \mathbf{v}) = \frac{1}{1 + \exp(-bias_{h_i} - \mathbf{v}^T \mathbf{W}_{:,i})}$$

$$P(v_j = 1 | \mathbf{h}) = \frac{1}{1 + \exp(-bias_{v_j} - \mathbf{W}_{j,:} \mathbf{h})}$$



$$p(\mathbf{h} | \mathbf{v}) = \prod_i p(h_i | \mathbf{v})$$

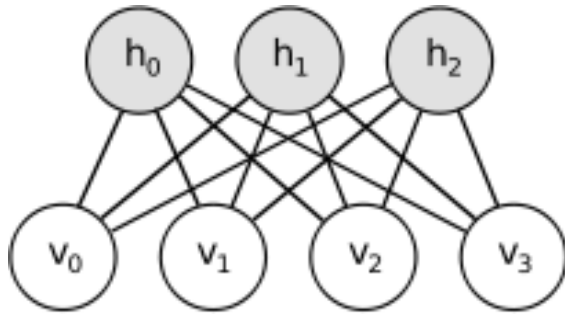
$$p(\mathbf{v} | \mathbf{h}) = \prod_j p(v_j | \mathbf{h})$$

Following the same principle of maximising log likelihood by means of gradient ascent, one obtains:

$$\Delta w_{ji} = \varepsilon \frac{\partial L(\mathbf{W})}{\partial w_{ji}} = \varepsilon \left( \langle v_j h_i \rangle_{\text{data}} - \langle v_j h_i \rangle_{\text{model}} \right)$$

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

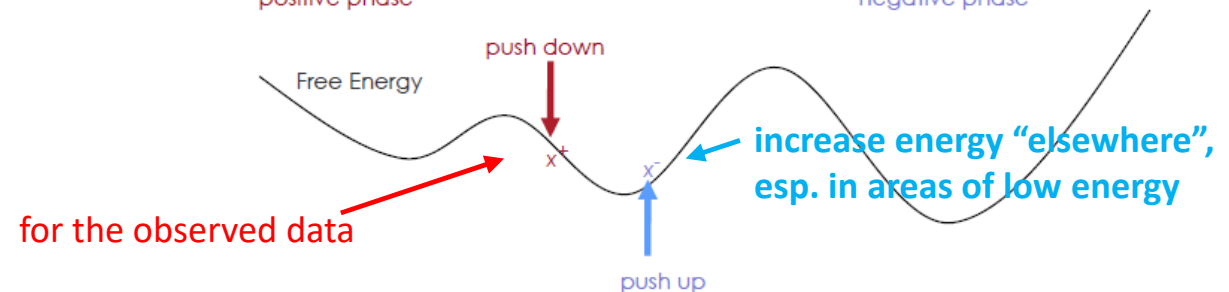
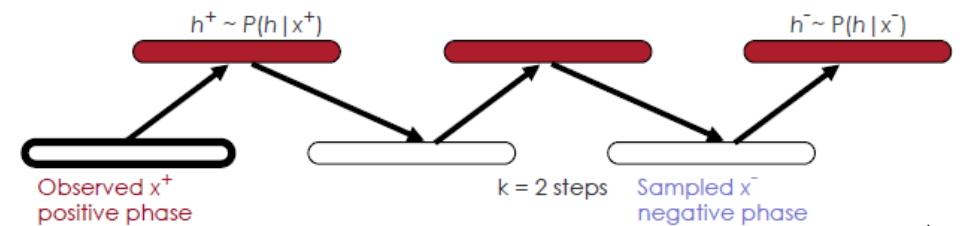
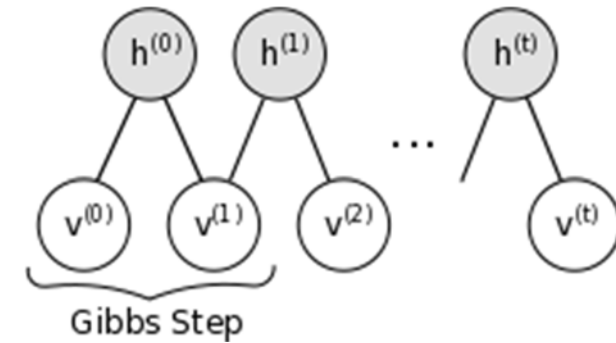
# RBM learning with Contrastive Divergence (CD)



$$P(h_i = 1 | \mathbf{v}) = \frac{1}{1 + \exp(-bias_{h_i} - \mathbf{v}^T \mathbf{W}_{:,i})}$$

$$P(v_j = 1 | \mathbf{h}) = \frac{1}{1 + \exp(-bias_{v_j} - \mathbf{W}_{j,:} \mathbf{h})}$$

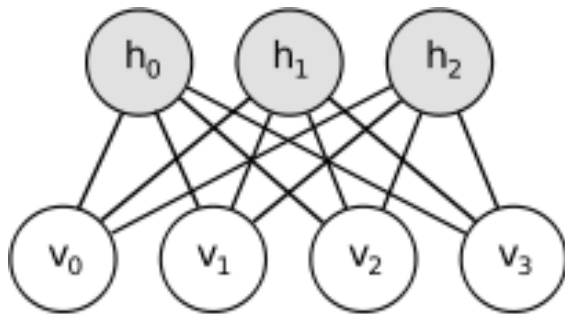
## Gibbs sampling



Hinton, 2003

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

# RBM learning with Contrastive Divergence (CD)



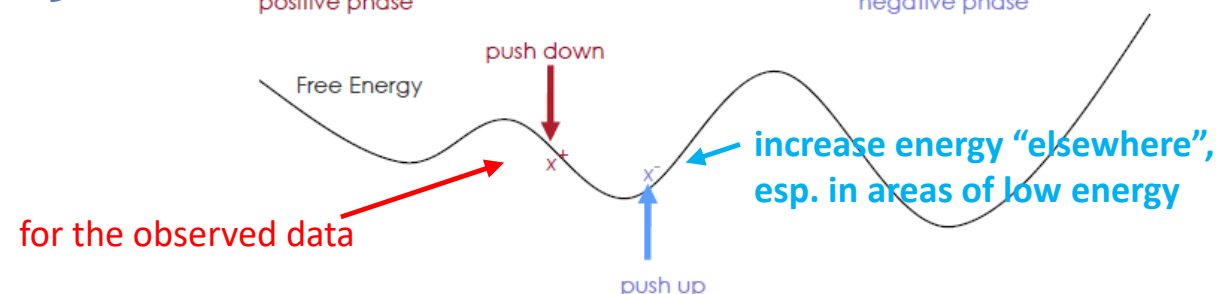
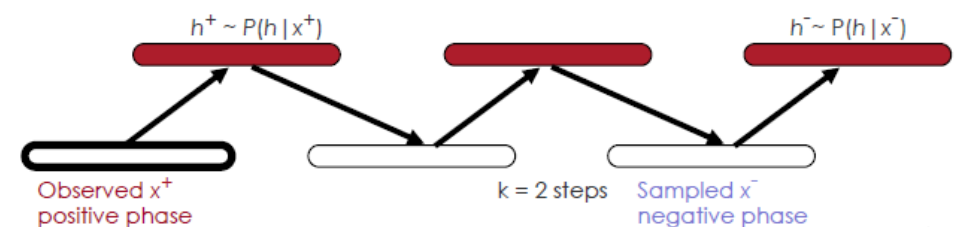
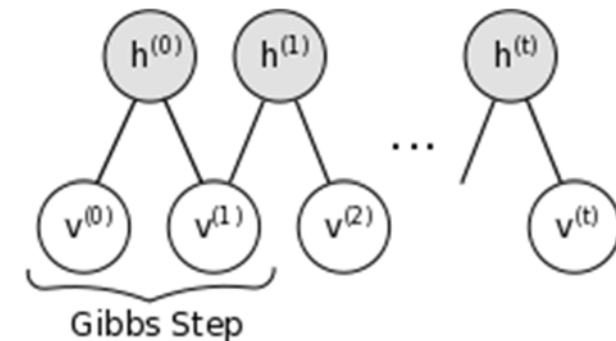
$$P(h_i = 1 | \mathbf{v}) = \frac{1}{1 + \exp(-bias_{h_i} - \mathbf{v}^T \mathbf{W}_{:,i})}$$

$$P(v_j = 1 | \mathbf{h}) = \frac{1}{1 + \exp(-bias_{v_j} - \mathbf{W}_{j,:} \mathbf{h})}$$

## GOOD TO KNOW:

Contrastive Divergence does not optimise the likelihood but it works effectively!

## Gibbs sampling

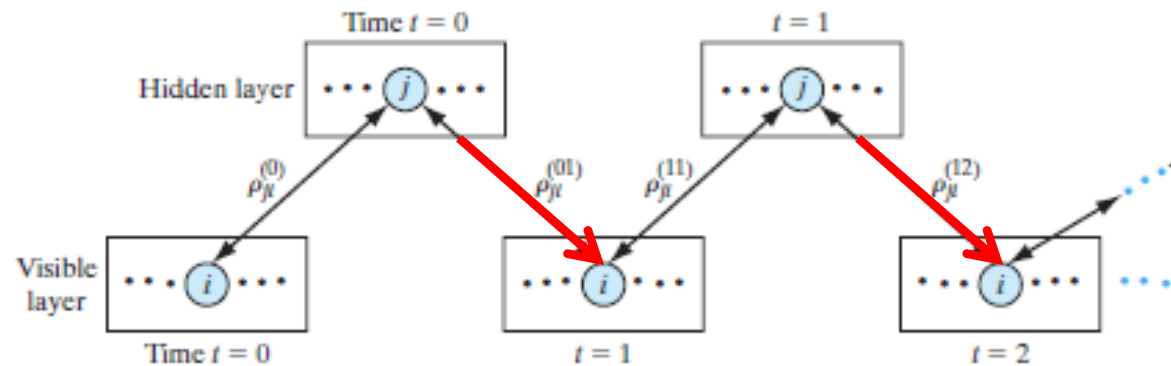


Hinton, 2003

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

# CD<sub>k</sub> recipe for training RBM

## Gibbs sampling



- 1) Clamp the visible units with an input vector and update hidden units.

$$P(h_i = 1 | \mathbf{v}) = \left( 1 + \exp \left( -bias_{h_i} - \mathbf{v}^T \mathbf{W}_{:,i} \right) \right)^{-1}$$

- 2) Update all the visible units in parallel to get a **reconstruction**.

$$P(v_j = 1 | \mathbf{h}) = \left( 1 + \exp \left( -bias_{v_j} - \mathbf{W}_{j,:} \mathbf{h} \right) \right)^{-1}$$

- 3) Collect the statistics for correlations after  $k$  steps using mini-batches and update weights:

$$\Delta w_{j,i} = \frac{1}{N} \sum_{n=1}^N \left( v_j^{(n)} h_i^{(n)} - \hat{v}_j^{(n)} \hat{h}_i^{(n)} \right)$$

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

# From RBM to Gaussian-Bernoulli RBM

Bernoulli-Bernoulli (binary-binary)

$$p(v_i = 1|\mathbf{h}) = g\left(\sum_j W_{ij}b_j + b_i\right)$$

$$p(b_j = 1|\mathbf{v}) = g\left(\sum_i W_{ij}v_i + a_j\right)$$



Gaussian-Bernoulli (real/cont.-binary)

$$p(v_i = x|\mathbf{h}) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{\left(x - b_i - \sigma_i \sum_j b_j W_{ij}\right)^2}{2\sigma_i^2}\right),$$

$$p(b_j = 1|\mathbf{v}) = g\left(b_j + \sum_i W_{ij} \frac{v_i}{\sigma_i}\right),$$



Visible units are real-valued whereas hidden units remain binary.

Salakhutdinov, 2015

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

# From RBM to Gaussian-Bernoulli RBM

Bernoulli-Bernoulli (binary-binary)

$$p(v_i = 1|\mathbf{h}) = g\left(\sum_j W_{ij}b_j + b_i\right)$$

$$p(b_j = 1|\mathbf{v}) = g\left(\sum_i W_{ij}v_i + a_j\right)$$



Gaussian-Bernoulli (real/cont.-binary)

$$p(v_i = x|\mathbf{h}) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{\left(x - b_i - \sigma_i \sum_j b_j W_{ij}\right)^2}{2\sigma_i^2}\right),$$

$$p(b_j = 1|\mathbf{v}) = g\left(b_j + \sum_i W_{ij} \frac{v_i}{\sigma_i}\right),$$

Visible units are real-valued whereas hidden units remain binary.

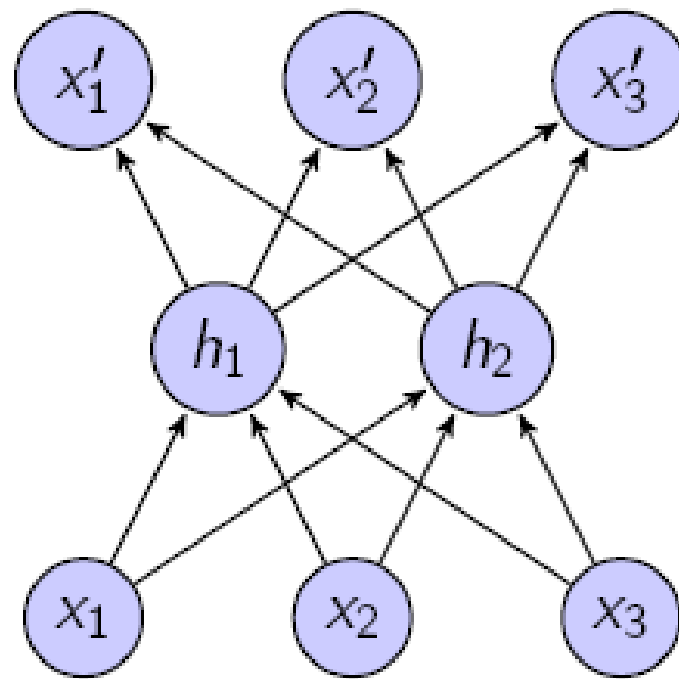
The derivative of the log-likelihood:

$$\frac{\partial \log P(\mathbf{v}; \theta)}{\partial W_{ij}} = \mathbb{E}_{P_{\text{data}}}\left[\frac{1}{\sigma_i} v_i b_j\right] - \mathbb{E}_{P_{\text{model}}}\left[\frac{1}{\sigma_i} v_i b_j\right]$$

Salakhutdinov, 2015

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

# Autoencoders – principles



Decoder:  $x' = f(x)$

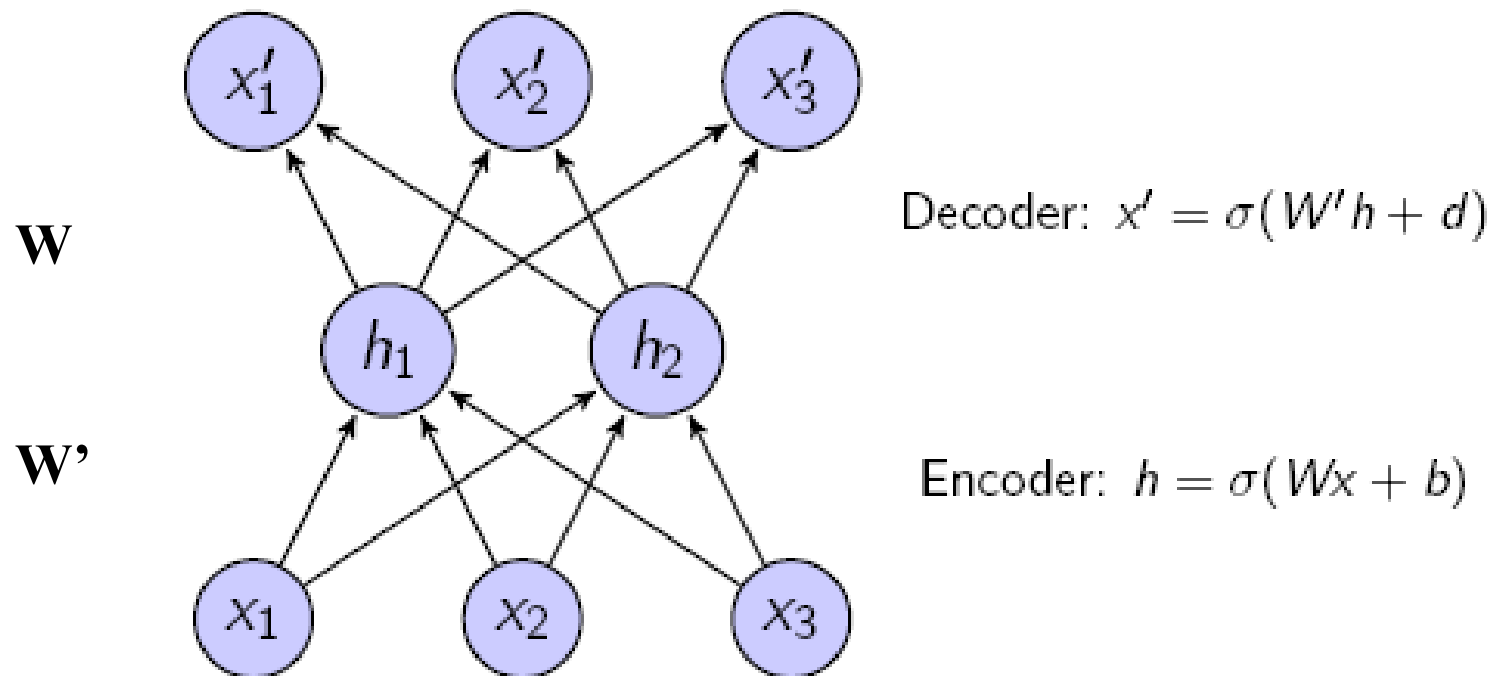
Encoder:  $h = g(f(x))$

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

# 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

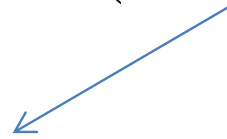
$$\log p_{\text{model}}(\mathbf{h}, \mathbf{x}) = \log p_{\text{model}}(\mathbf{h}) + \log p_{\text{model}}(\mathbf{x} | \mathbf{h})$$

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

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

# Denoising autoencoders

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



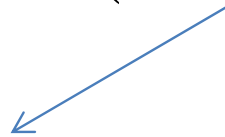
Corrupted copy of  $x$

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  $p_{reconstruction}(\tilde{\mathbf{x}} \mid \mathbf{x}) = p_{decoder}(\mathbf{x} \mid \mathbf{h}), \quad \mathbf{h} = f(\tilde{\mathbf{x}})$

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

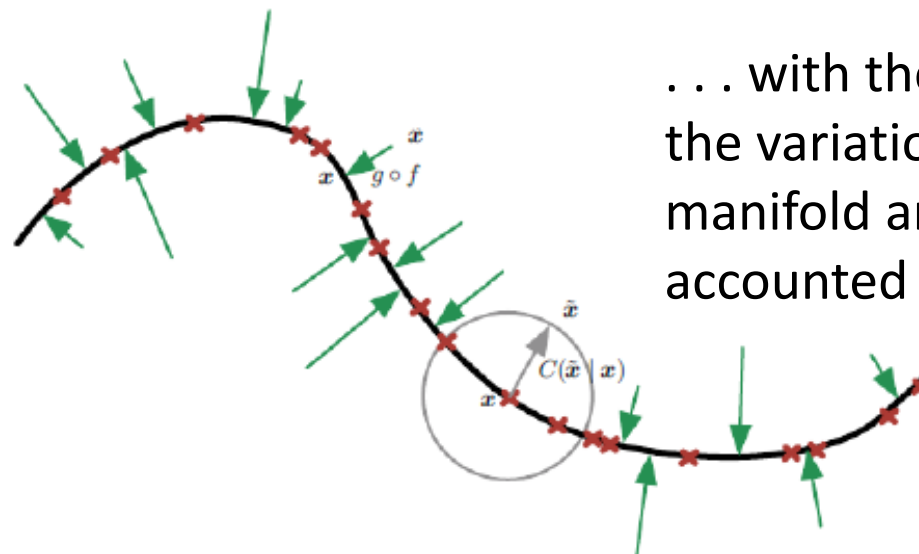
# Denoising autoencoders

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

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  $\mathbf{h}$

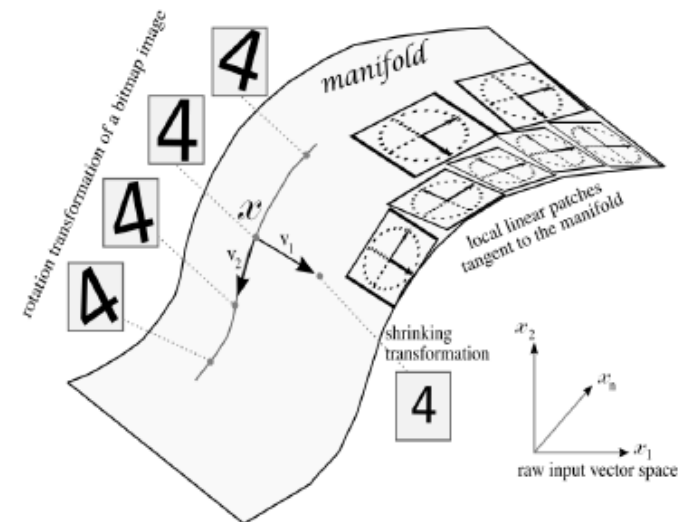
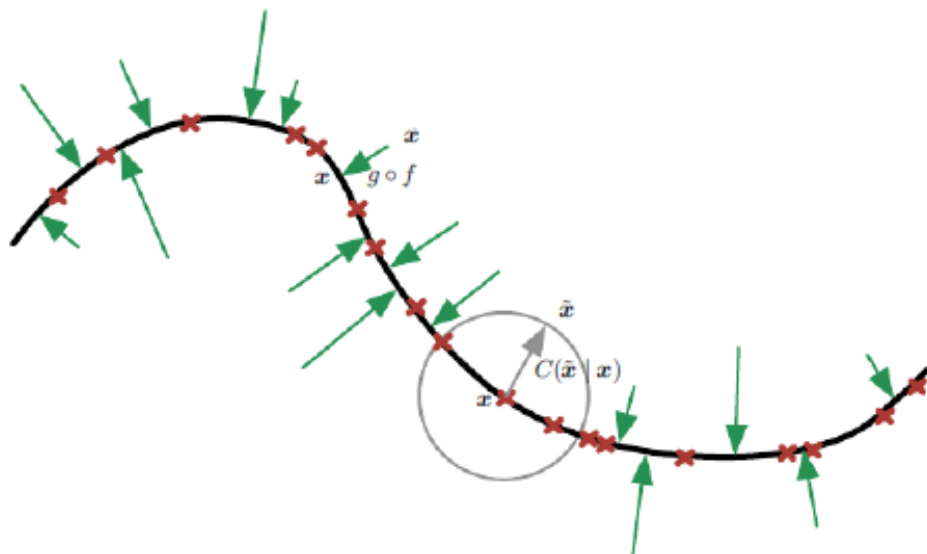
Goodfellow et al.

- 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.

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

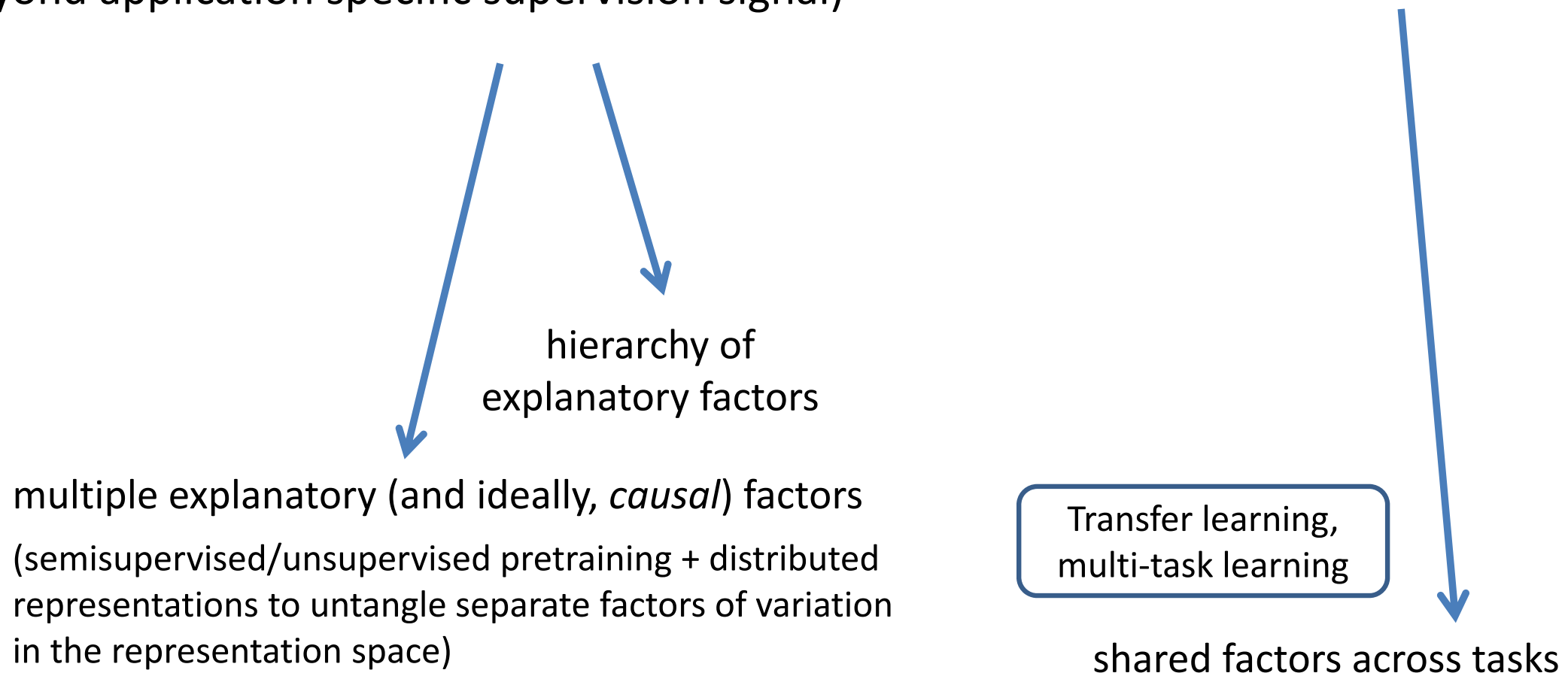
# Strategies to guide discovery of salient/causal factors

How can the discovery/identification of the underlying causal factors of variation that generate the data be further supported?  
(beyond application specific supervision signal)

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

# Generic regularisation strategies (Bengio et al., 2013)

How can the discovery/identification of the underlying causal factors of variation that generate the data be further supported? (beyond application specific supervision signal)

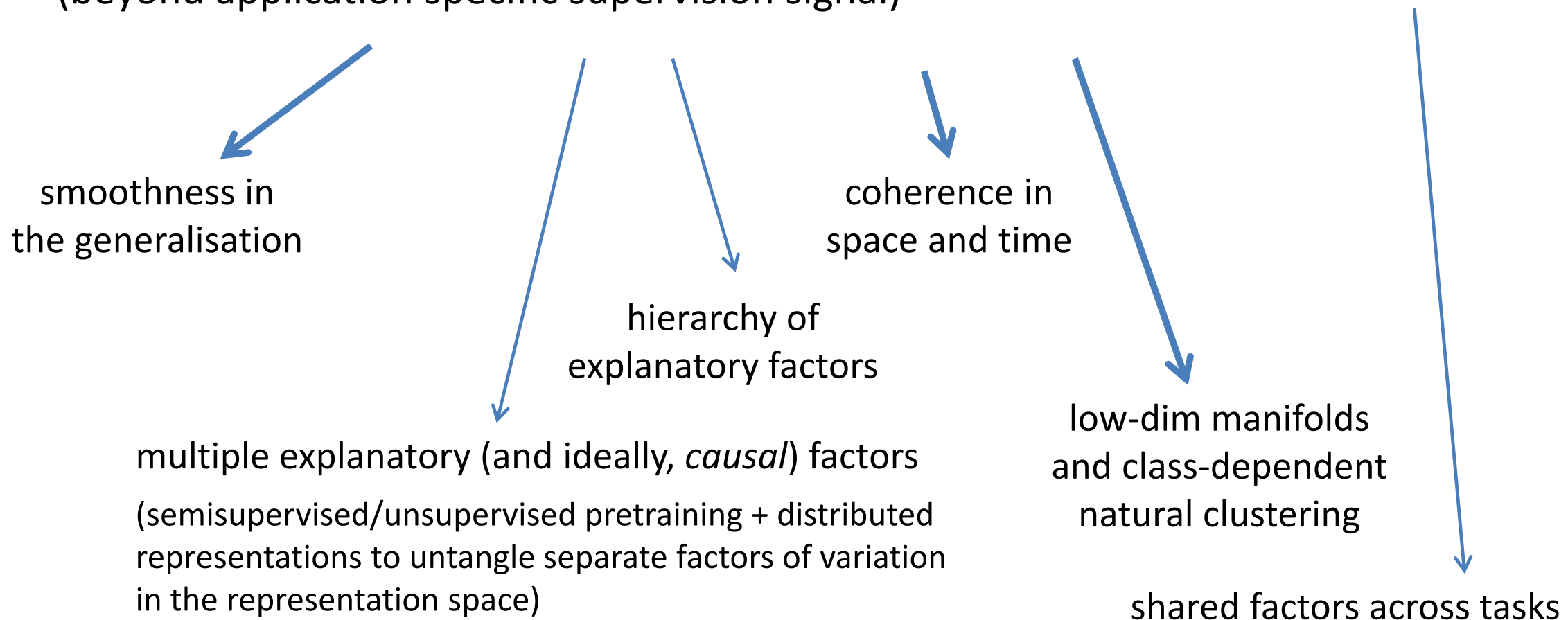




- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Generic regularisation strategies (Bengio et al., 2013)

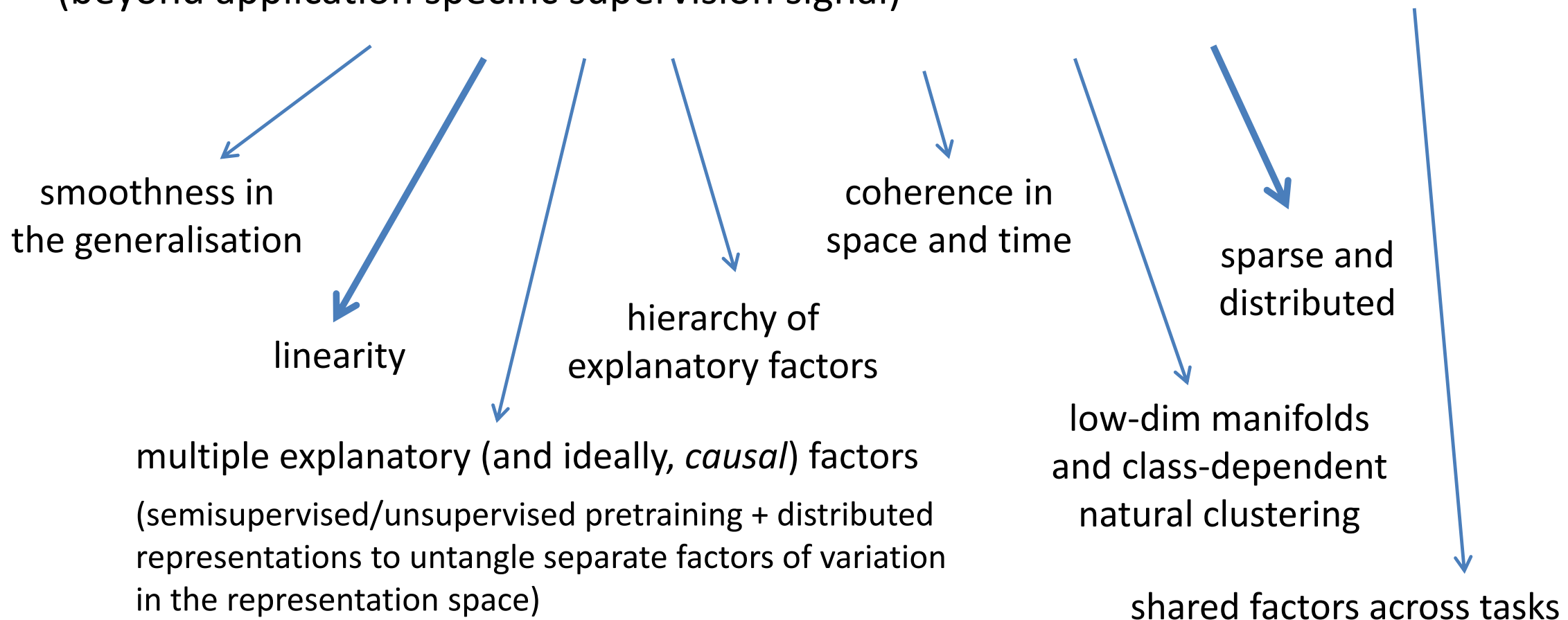
How can the discovery/identification of the underlying causal factors of variation that generate the data be further supported? (beyond application specific supervision signal)



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

# Generic regularisation strategies (Bengio et al., 2013)

How can the discovery/identification of the underlying causal factors of variation that generate the data be further supported? (beyond application specific supervision signal)

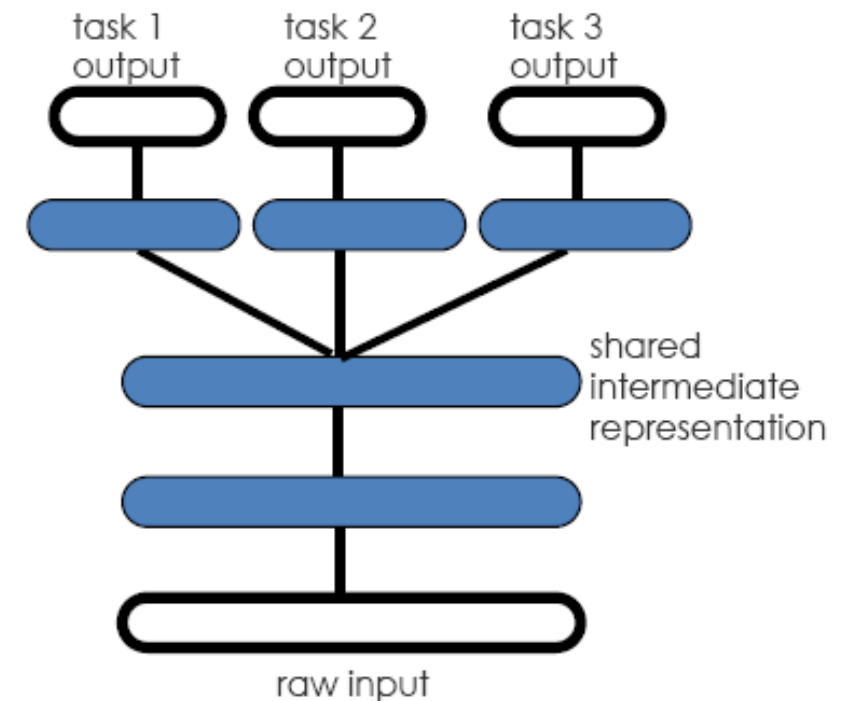


- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Transfer and multi-task learning

## 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
- Transfer learning and multi-task learning is supported by hierarchical and distributed representations

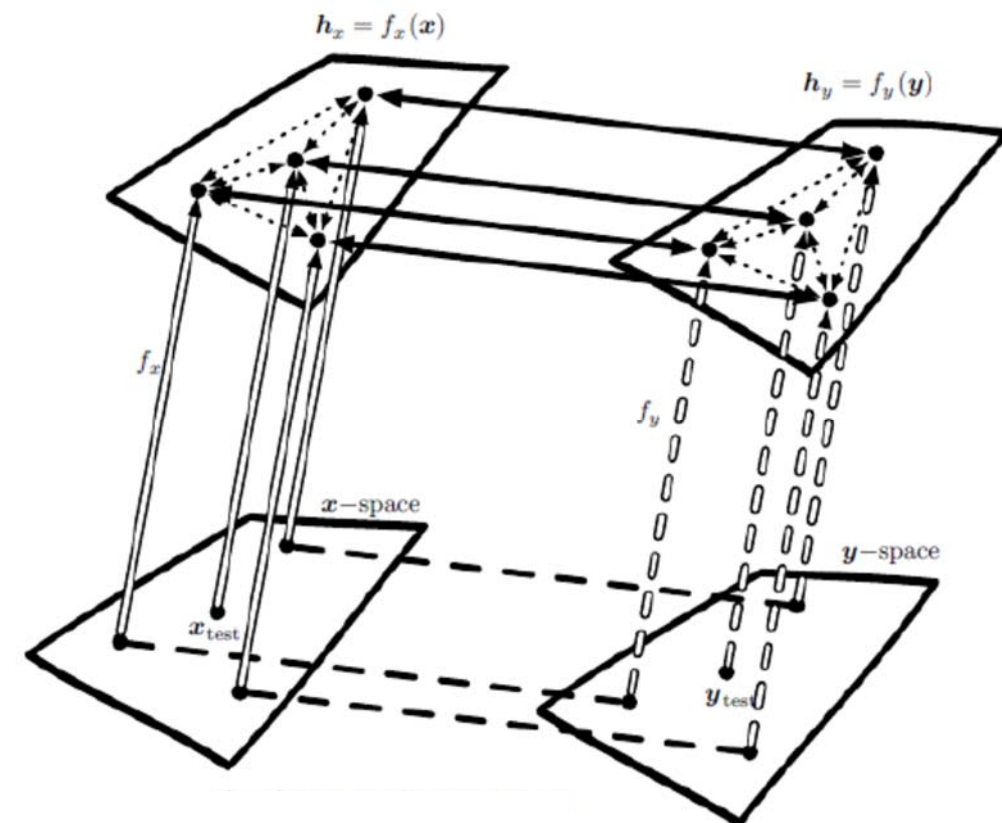


Multi-task learning

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

# Transfer, multi-task/modal, one/zero-shot learning

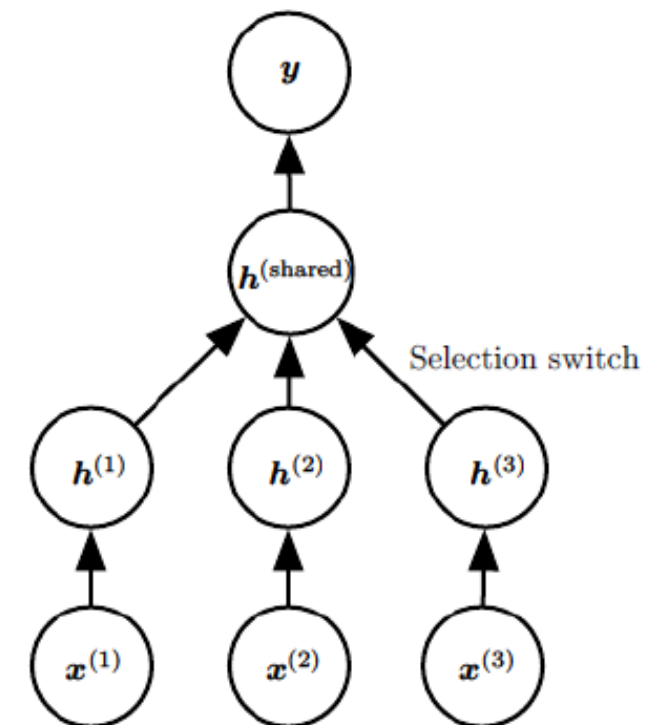
- The concept of *one-shot* learning
- Zero-shot learning as a specific form of *multi-modal learning* (capturing the relationship between representations in different modalities)



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

# Transfer, multi-task/modal, one/zero-shot learning

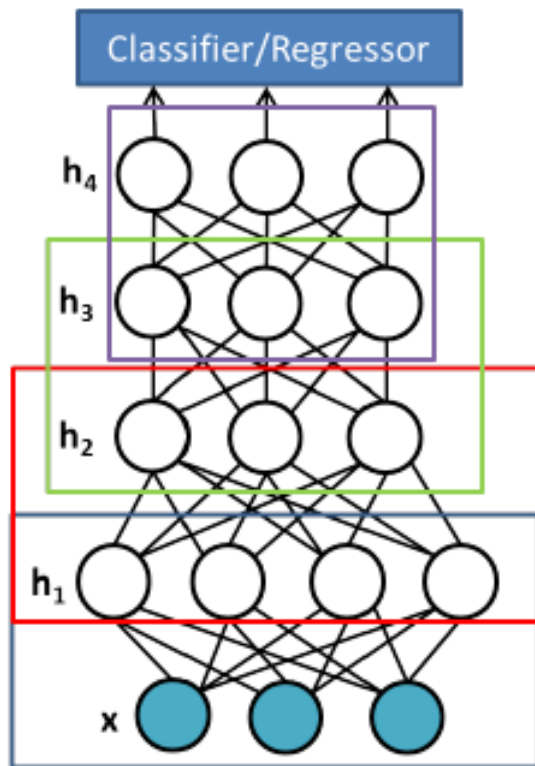
- The concept of *one-shot* learning
- Zero-shot learning as a specific form of *multi-modal learning* (capturing the relationship between representations in different modalities)
- However, sometimes the semantics of the output is shared instead, which requires domain adaptation



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

# Deep belief nets

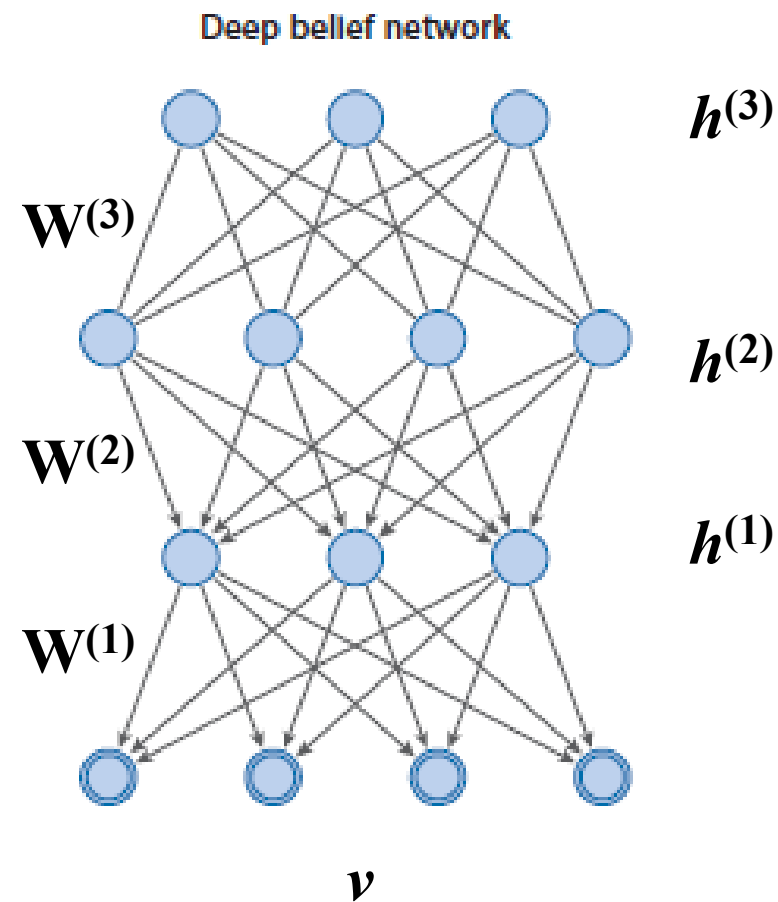
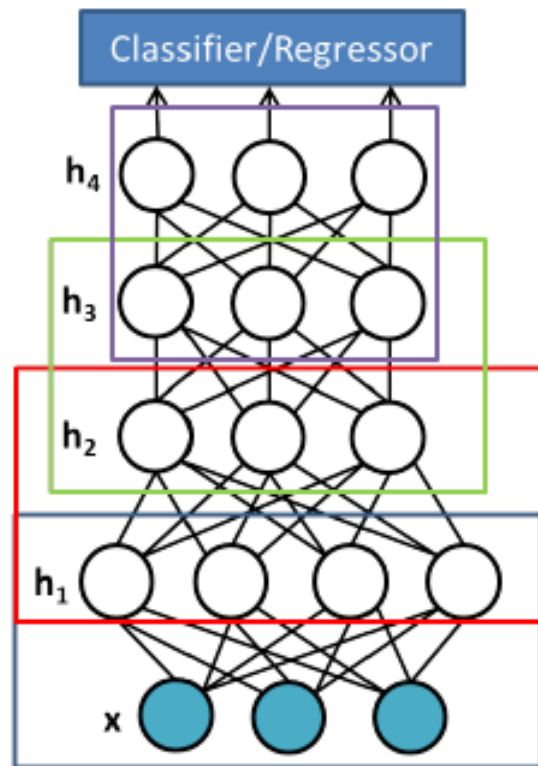
Greedy layer-wise training approach  
with the use of RBMs



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

# Deep belief nets

Greedy layer-wise training approach  
with the use of RBMs

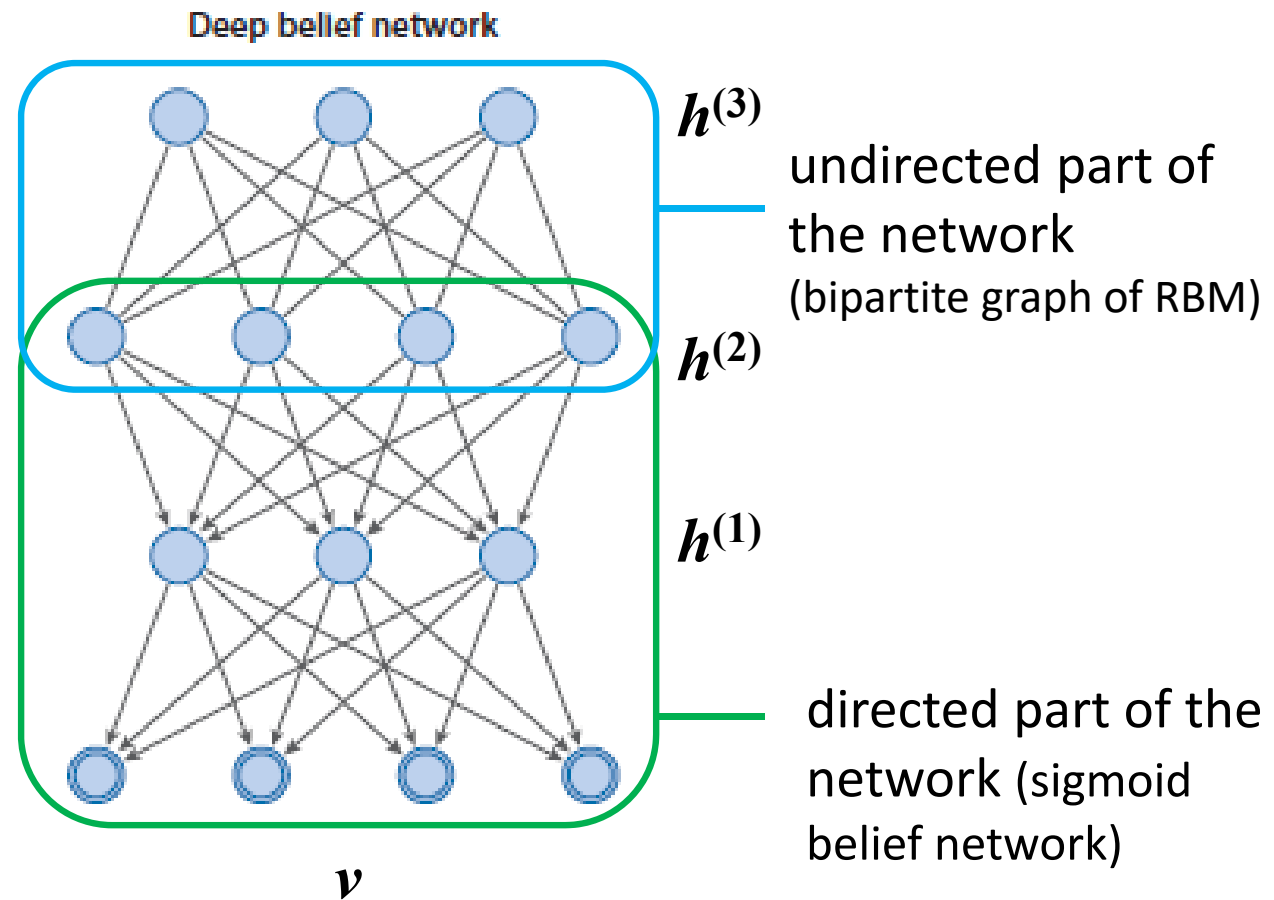
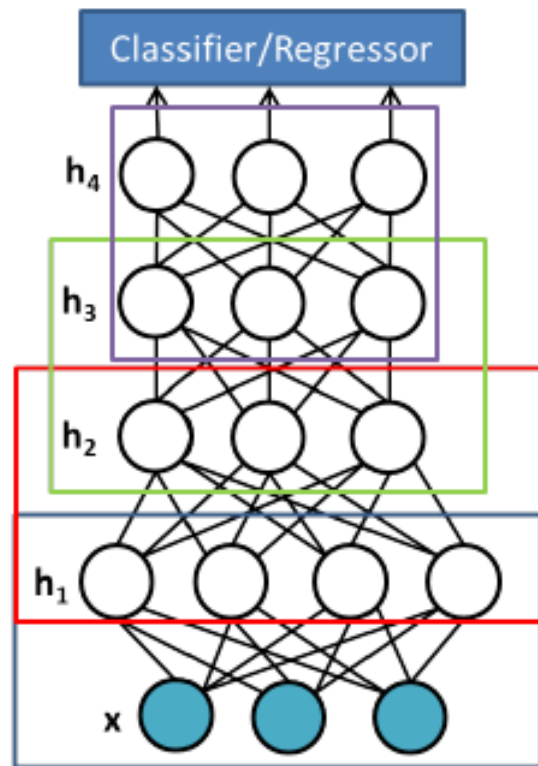


Salakhutdinov, 2015

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

# Deep belief nets

Greedy layer-wise training approach  
with the use of RBMs



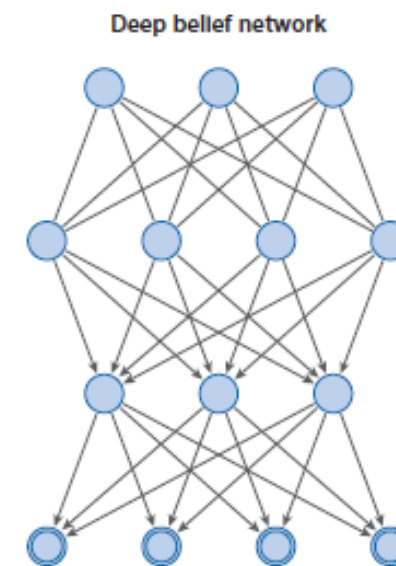
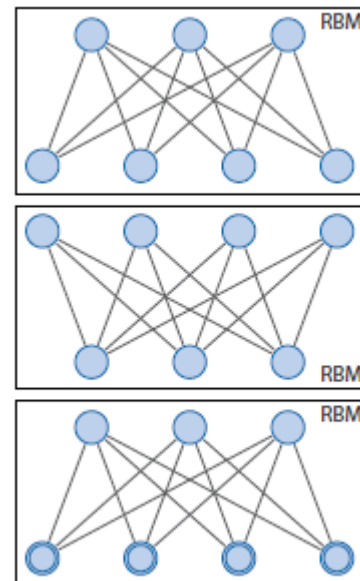
Salakhutdinov, 2015



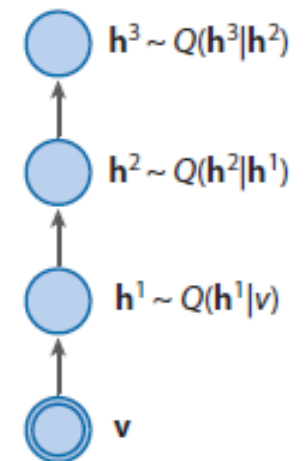
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Deep belief nets

## Approach 1



Bottom-up pass by  
stochastically activating  
higher layers in time

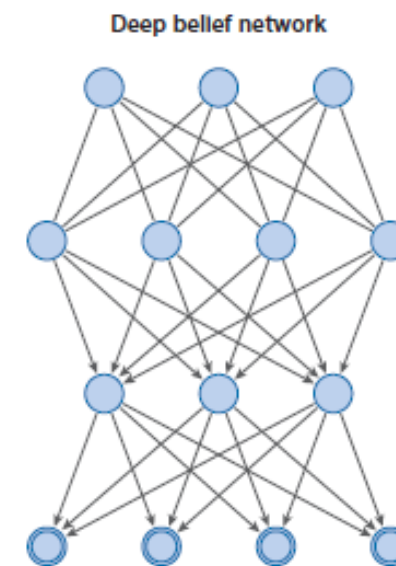
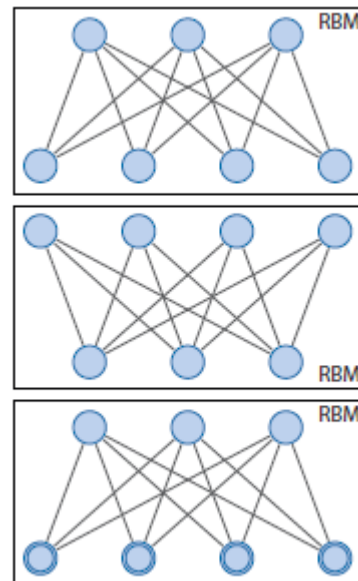


- 1: Fit the parameters  $W^{(1)}$  of the first-layer RBM to data.
- 2: Fix the parameter vector  $W^{(1)}$ , and use samples  $h^{(1)}$  from  $Q(h^{(1)}|v) = P(h^{(1)}|v, W^{(1)})$  as the data for training the next layer of binary features with an RBM.
- 3: Fix the parameters  $W^{(2)}$  that define the second layer of features, and use the samples  $h^{(2)}$  from  $Q(h^{(2)}|h^{(1)}) = P(h^{(2)}|h^{(1)}, W^{(2)})$  as the data for training the third layer of binary features.
- 4: Proceed recursively for the next layers.

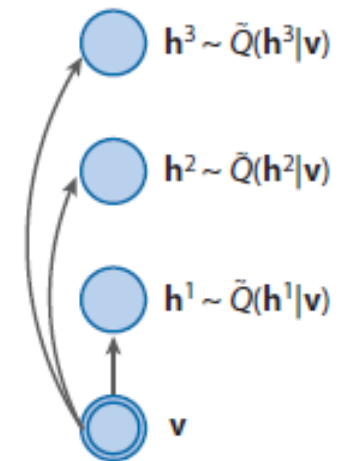
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# Deep belief nets

## Approach 2



Bottom-up pass by stochastically activating higher layers in time



Assumption about fully factorised approximating distribution

$$\tilde{Q}(\mathbf{h}^{(1)}, \dots, \mathbf{h}^{(L)} | \mathbf{v}) = \prod_{l=1}^L \tilde{Q}(\mathbf{h}^{(l)} | \mathbf{v})$$

$$\tilde{Q}(\mathbf{h}^{(1)} | \mathbf{v}) = \prod_j q(h_j^{(1)} | \mathbf{v}), \text{ where:}$$

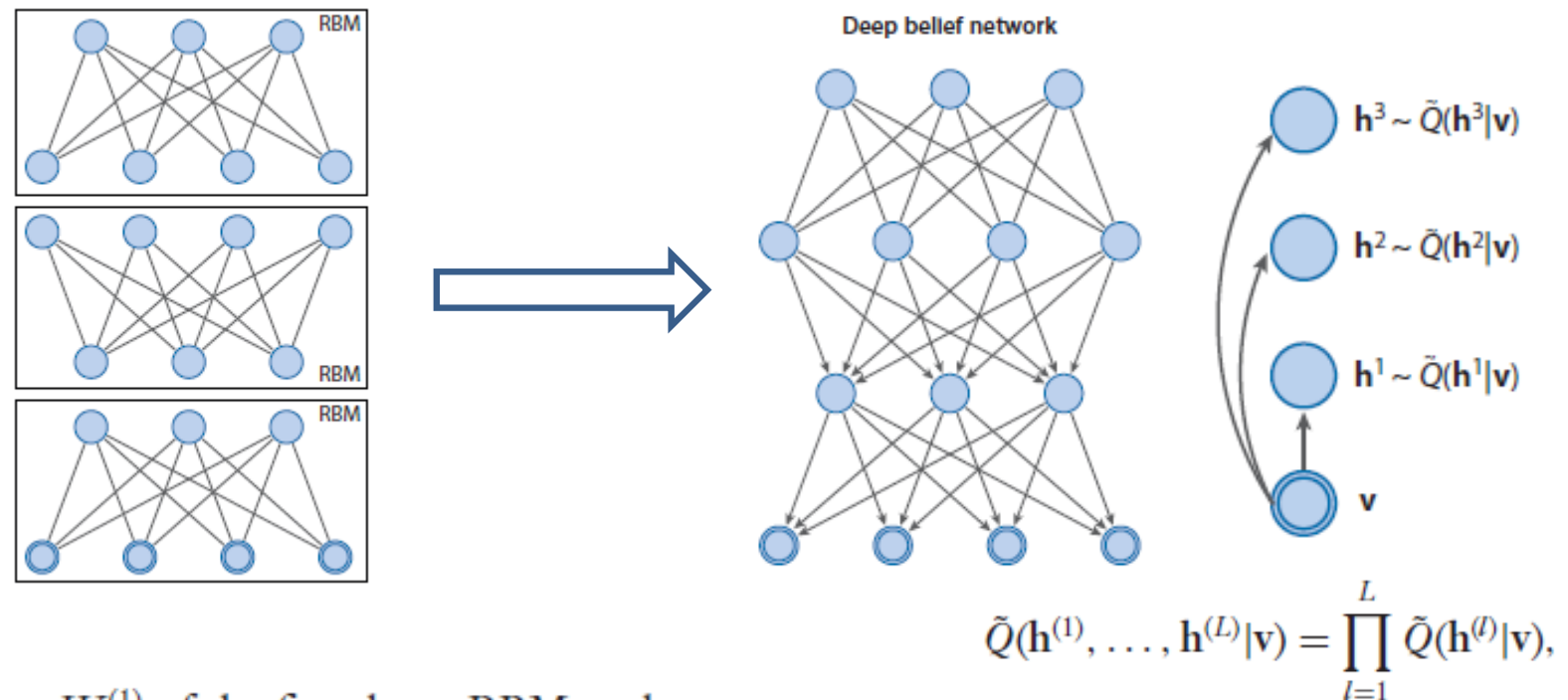
$$q(h_j^{(1)} = 1 | \mathbf{v}) = g \left( \sum_i W_{ij}^{(1)} v_i + a_j^{(1)} \right), \text{ and}$$

$$q(h_j^{(l)} = 1 | \mathbf{v}) = g \left( \sum_i W_{ij}^{(l)} q(h_i^{(l-1)} = 1 | \mathbf{v}) + a_j^{(l)} \right),$$

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

# Deep belief nets

## Approach 2

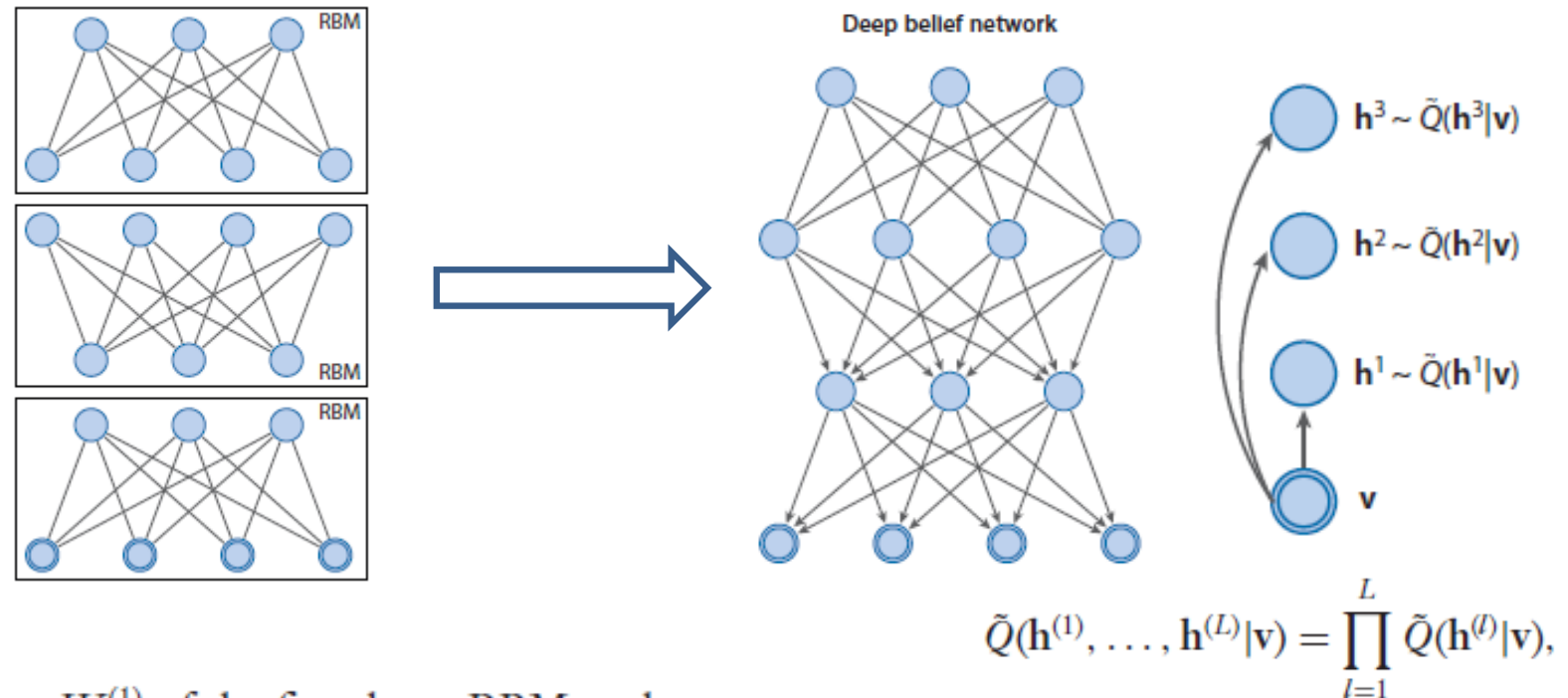


- 1: Fit the parameters  $W^{(1)}$  of the first-layer RBM to data.
- 2: Fix the parameter vector  $W^{(1)}$ , and use samples  $\mathbf{h}^{(1)}$  from  $\tilde{Q}(\mathbf{h}^{(1)} | \mathbf{v}) = P(\mathbf{h}^{(1)} | \mathbf{v}, W^{(1)})$  as the data for training the next layer of binary features with an RBM.
- 3: Fix the parameters  $W^{(2)}$  that define the second layer of features, and use the samples  $\mathbf{h}^{(2)}$  from  $\tilde{Q}(\mathbf{h}^{(2)} | \mathbf{v})$  as the data for training the third layer of binary features.
- 4: Proceed recursively for the next layers.

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

# Deep belief nets

## Approach 2



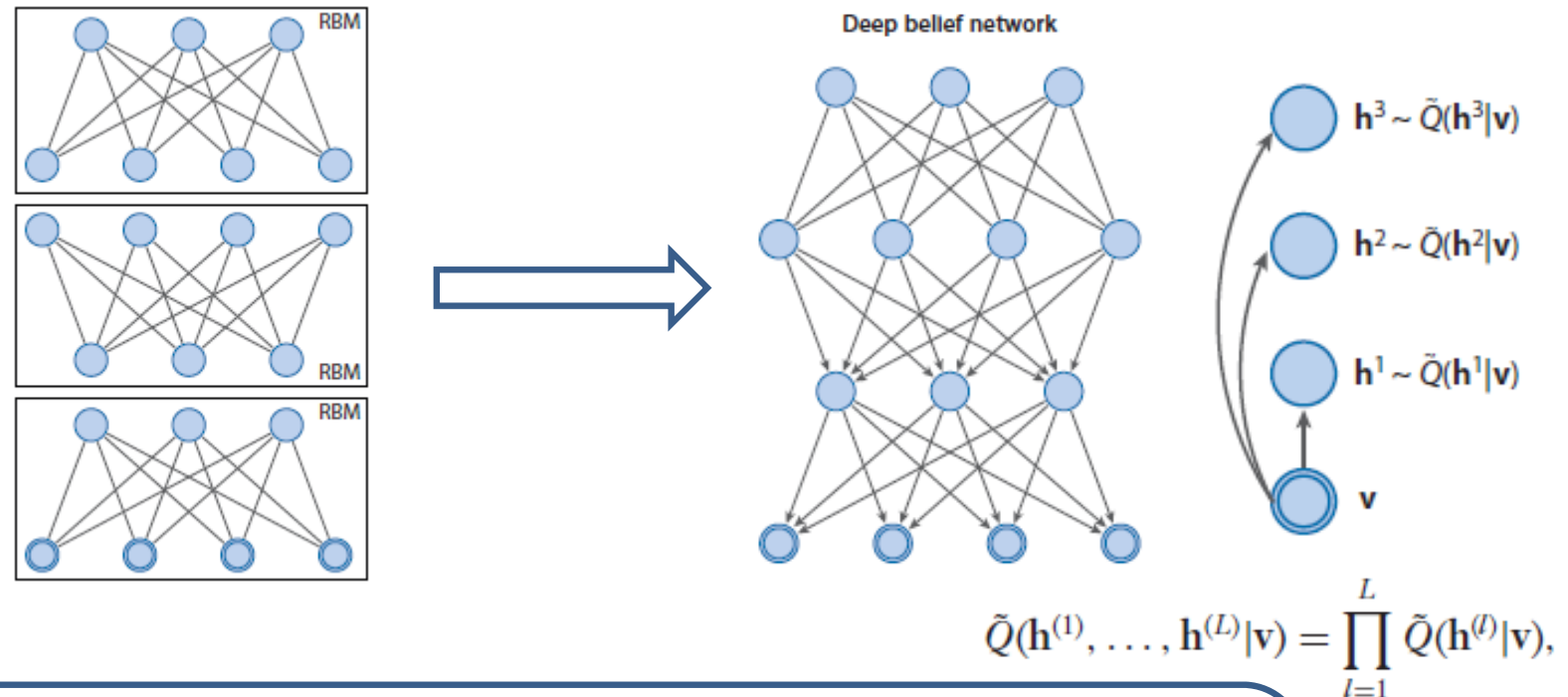
- 1: Fit the parameters  $W^{(1)}$  of the first-layer RBM to data.
- 2: Fix the parameter vector  $W^{(1)}$ , and use samples  $\mathbf{h}^{(1)}$  from  $\tilde{Q}(\mathbf{h}^{(1)}|\mathbf{v}) = P(\mathbf{h}^{(1)}|\mathbf{v}, W^{(1)})$  as the data for training the next layer of binary features with an RBM.

The difference lies in how we obtain sampling distribution to generate input for greedy training another layer.

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

# Deep belief nets

## Approach 2



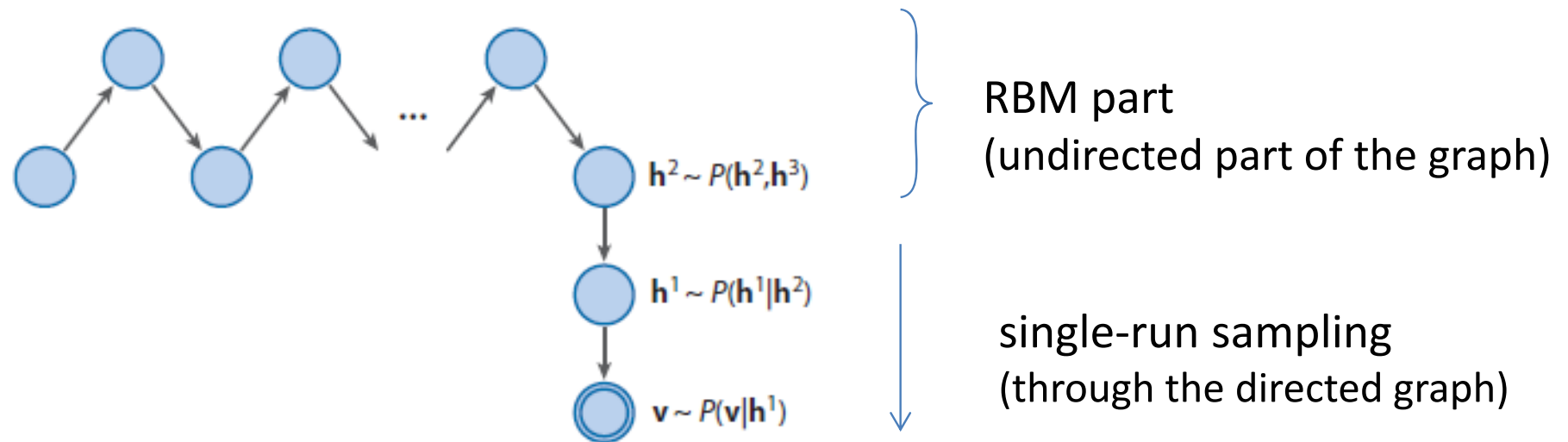
For a fully factorised model,  $\tilde{Q}(\mathbf{h}^{(l)} | \mathbf{v})$ , needed for sampling data at the next level, a single deterministic bottom-up pass can be executed on real-valued probabilities.

$$q(b_j^{(l)} = 1 | \mathbf{v}) = g \left( \sum_i W_{ij}^{(l)} q(b_i^{(l-1)} = 1 | \mathbf{v}) + a_j^{(l)} \right),$$

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

# Approximate sampling from DBN

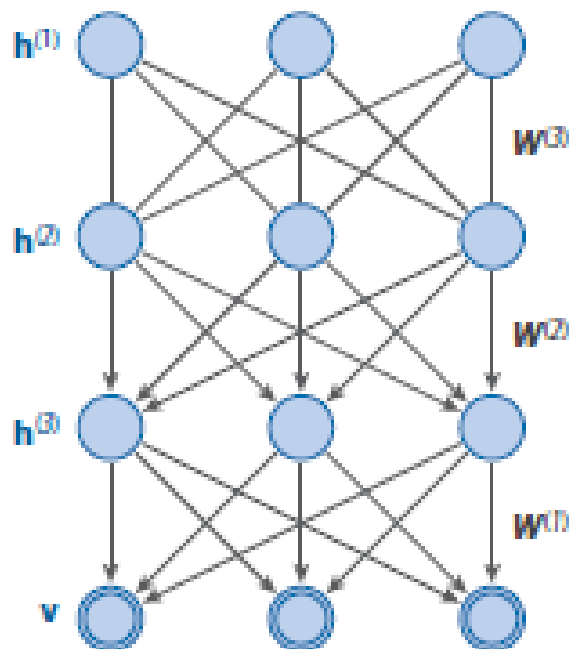
Gibbs sampling chain in the RBM part



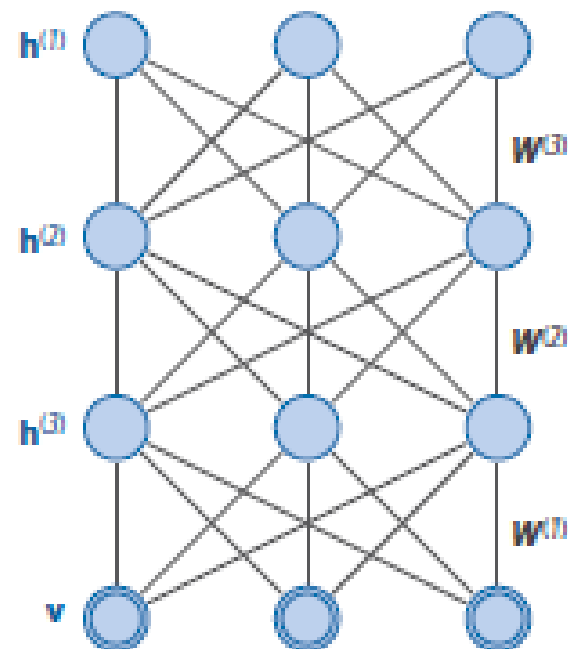
- Recap
- Data representations
- Learning data representations in deep networks
- Deep generative models

# DBN vs DBM

**Deep belief network**



**Deep Boltzmann machine**



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

# DBN vs DBM

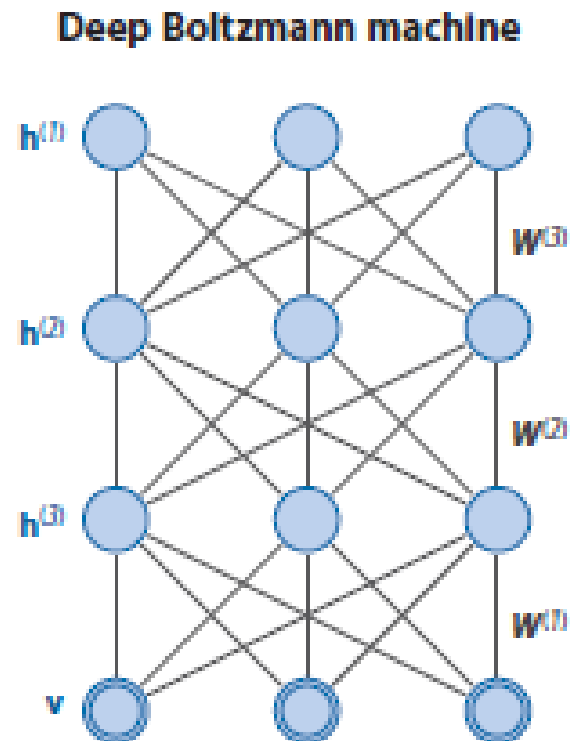
Greedy layer-wise pre-training approach\*

$$p(b_j^{(1)} = 1 | \mathbf{v}, \mathbf{h}^{(2)}) = g \left( \sum_i W_{ij}^{(1)} v_i + \sum_m W_{jm}^{(2)} b_m^{(2)} \right),$$

$$p(b_m^{(2)} = 1 | \mathbf{h}^{(1)}, \mathbf{h}^{(3)}) = g \left( \sum_j W_{jm}^{(2)} b_j^{(1)} + \sum_l W_{ml}^{(3)} b_l^{(3)} \right),$$

$$p(b_l^{(3)} = 1 | \mathbf{h}^{(2)}) = g \left( \sum_m W_{ml}^{(3)} b_m^{(2)} \right),$$

$$p(v_i = 1 | \mathbf{h}^{(1)}) = g \left( \sum_j W_{ij}^{(1)} b_j^{(1)} \right).$$

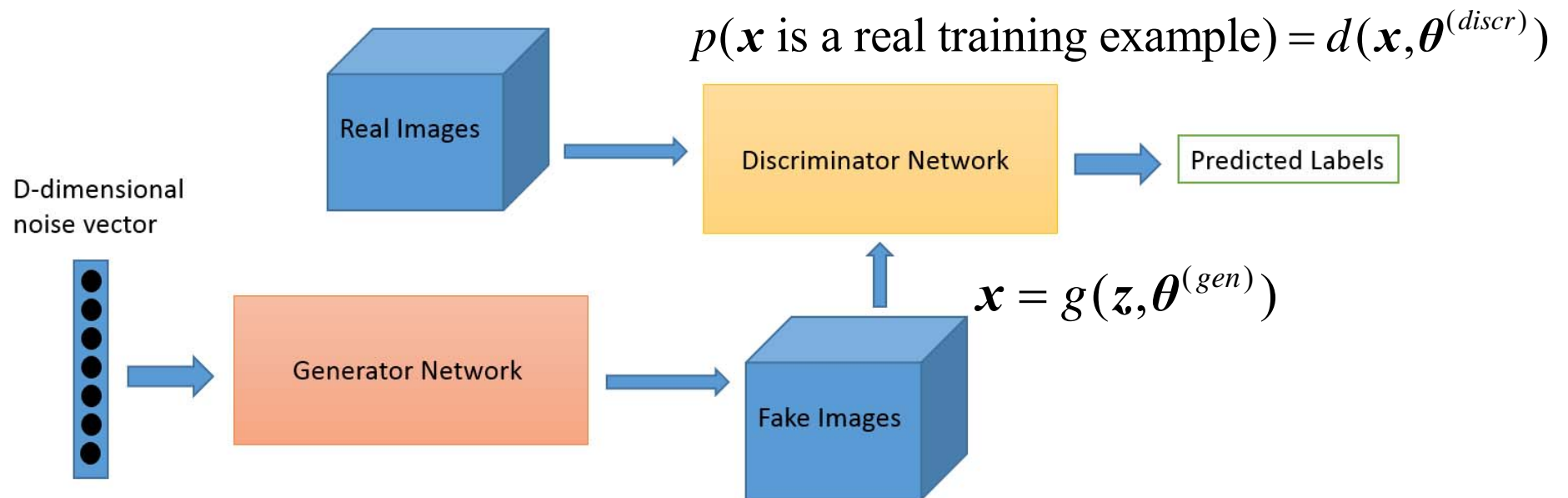


\*Salakhutdinov and Hinton, 2009



- 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:

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

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

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

# Summary

1. Central role of data representations (sparse distributed code in the cortex).
2. Deep learning is about learning features that constitute representations (string link to hierarchical processing in the brain).
3. RBMs and autoencoders are key computational blocks for learning representations and building deep generative models – let them learn/extract features.
4. The earlier popularity of greedy layer-wise pretraining, today we rather rely on dropout and ReLU units (rather than sigmoidal).
5. Generative power of deep models.

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

## Recommended reading

- Goodfellow, I., Bengio, Y., & Courville, A. *Deep learning*, chapters 6, 14, 15, 20 .
- Salakhutdinov, R. (2015) Learning deep generative models. *Annual Reviews of Statistics and Its Application*, 2, p.361–385.
- Bengio, Y., Courville, A., & Vincent, P. (2013). Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35 (8), p.1798-1828.
- Hinton, G. (2010). A Practical Guide to Training Restricted Boltzmann Machines. Technical report UTML TR 2010–003.
- Bengio, Y. (2009) Learning deep architectures for AI. *Foundations and trends® in Machine Learning*, 2.1. p.1-127.
- Salakhutdinov ,R., & Hinton, G. (2009). Deep Boltzmann machines. *Proc. 12th Int. Conf. Artif. Intell. Stat.*, Clearwater Beach, FL, pp. 448–55. Brookline, MA: Microtome.



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

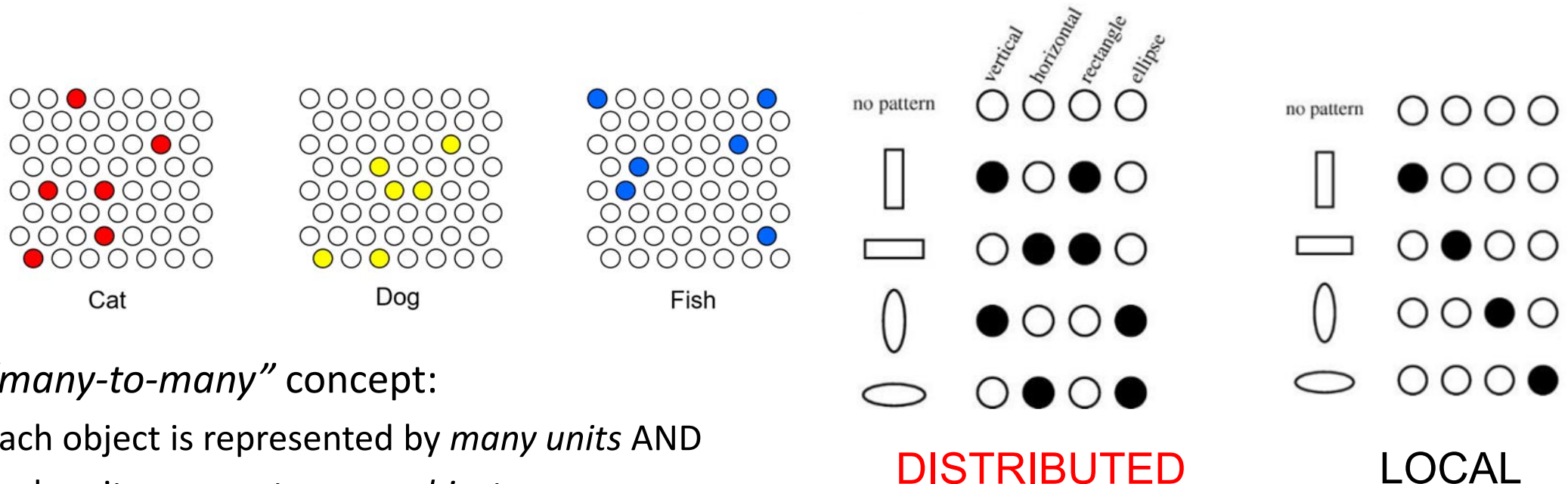
# Recap on deep networks & representation learning

- Learning data representations (what is a good data representation)
  - uncovering *causal (latent)* factors
  - *end-to-end learning* vs hand-crafting features
  - hierarchy of distributed features
  - *distributed* representations as opposed to local or dense codes

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

# Recap on deep networks & representation learning

- Learning data representations (what is a good data representation)
  - uncovering *causal (latent)* factors
  - *end-to-end learning* vs hand-crafting features
  - hierarchy of distributed features
  - *distributed* representations as opposed to local or dense codes



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

# Recap on deep networks & representation learning

- Learning data representations (what is a good data representation)
    - uncovering *causal (latent)* factors
    - *end-to-end learning* vs hand-crafting features
    - hierarchy of distributed features
    - *distributed* representations as opposed to local or dense codes
      - *multi-modal* learning
      - *multi-task* and *transfer* (sharing representations) learning
      - enhanced generalisation power (shared attributes and semantic proximity in distributed representations)
      - expressiveness, fault tolerance, content (feature) addressability
- 