

# Probabilistic Models with Latent Variables

Goodfellow 13, 15, 20 Summary

박찬연

# Probabilistic models with latent variables

- Generative model:  $p_{\text{model}}(\boldsymbol{x}) = \mathbb{E}_{\boldsymbol{h}} p_{\text{model}}(\boldsymbol{x} \mid \boldsymbol{h})$  (Ch 20)
- Representation:  $p_{\text{model}}(\boldsymbol{h} \mid \boldsymbol{x})$  (Ch 15)
- A linear factor model is the simplest probabilistic models with latent variables. (Ch 13)

# 13 Linear Factor Model

- Defined by the use of a stochastic, linear decoder function that generates  $\mathbf{x}$  by adding noise to a linear transformation of  $\mathbf{h}$ .
- Data generation process
  - Sample  $\mathbf{h}$  from  $p(\mathbf{h})$ ,

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

- Sample the real-valued observables given the factors.

$$\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \text{noise}$$

# 13.1 Factor Analysis

- The latent variable prior is the unit variance Gaussian.

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$$

- Noise is drawn from a diagonal covariance Gaussian.

$$\boldsymbol{\psi} = \text{diag}(\boldsymbol{\sigma}^2), \text{ with } \boldsymbol{\sigma}^2 = [\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2]^\top$$

- Latent variables capture the dependencies between  $x_i$ .
- $\mathbf{x}$  is a multivariate normal random variable.

$$\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^\top + \boldsymbol{\psi})$$

## 13.1 Probabilistic PCA

- Make the conditional variances  $\sigma_i^2$  equal to each other.
- Conditional distribution:  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^\top + \sigma^2\mathbf{I})$
- Equivalently:  $\mathbf{x} = \mathbf{W}\mathbf{h} + \mathbf{b} + \sigma\mathbf{z}$ ,  $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$
- Becomes PCA as  $\sigma \rightarrow 0$ .

## 13.2 Independent Component Analysis (ICA)

- Require  $p(\mathbf{h})$  to be non-Gaussian.
- From a user-fixed  $p(\mathbf{h})$ , generates  $\mathbf{x} = \mathbf{W}\mathbf{h}$  by nonlinear change of variables to determine  $p(\mathbf{x})$ .
- Example
  - Separating an observed signal into many underlying signals that are scaled and added together to form the observed data.
  - The signals are intended to be fully independent.

## 13.3 Slow Feature Analysis (SFA)

- Idea: important features change very slowly.
- Slowness principle

$$\lambda \sum_t L(f(\mathbf{x}^{(t+1)}), f(\mathbf{x}^{(t)}))$$

- $f$  is the feature extractor.
- $L$  is a loss function measuring the distance, usually the mean squared difference.

## 13.3 Slow Feature Analysis (SFA)

- SFA algorithm

- Define  $f(\mathbf{x}; \boldsymbol{\theta})$  to be linear and solve

$$\min_{\boldsymbol{\theta}} \mathbb{E}_t (f(\mathbf{x}^{(t+1)})_i - f(\mathbf{x}^{(t)})_i)^2 \quad \mathbb{E}_t f(\mathbf{x}^{(t)})_i = 0 \quad \mathbb{E}_t [f(\mathbf{x}^{(t)})_i^2] = 1$$

- May be solved in closed form by a linear algebra package.
- So far the slowness principle has not become the basis of any state of the art applications.
  - Perhaps the slowness prior is too strong.
  - Slowness principle encourages the model to ignore the position of objects that have high velocity.



## 13.4 Sparse Coding

- Assume that the linear factors have Gaussian noise with isotropic precision  $\beta$ .

$$p(\mathbf{x} \mid \mathbf{h}) = \mathcal{N}(\mathbf{x}; \mathbf{W}\mathbf{h} + \mathbf{b}, \frac{1}{\beta}\mathbf{I})$$

- $p(\mathbf{h})$  is chosen to be one with sharp peaks near 0.
- The encoder is an optimization algorithm.

$$\mathbf{h}^* = f(\mathbf{x}) = \arg \max_{\mathbf{h}} p(\mathbf{h} \mid \mathbf{x})$$

$$= \arg \max_{\mathbf{h}} \log p(\mathbf{h} \mid \mathbf{x}) \quad p(h_i) = \text{Laplace}(h_i; 0, \frac{2}{\lambda}) = \frac{\lambda}{4} e^{-\frac{1}{2} \lambda |h_i|}$$

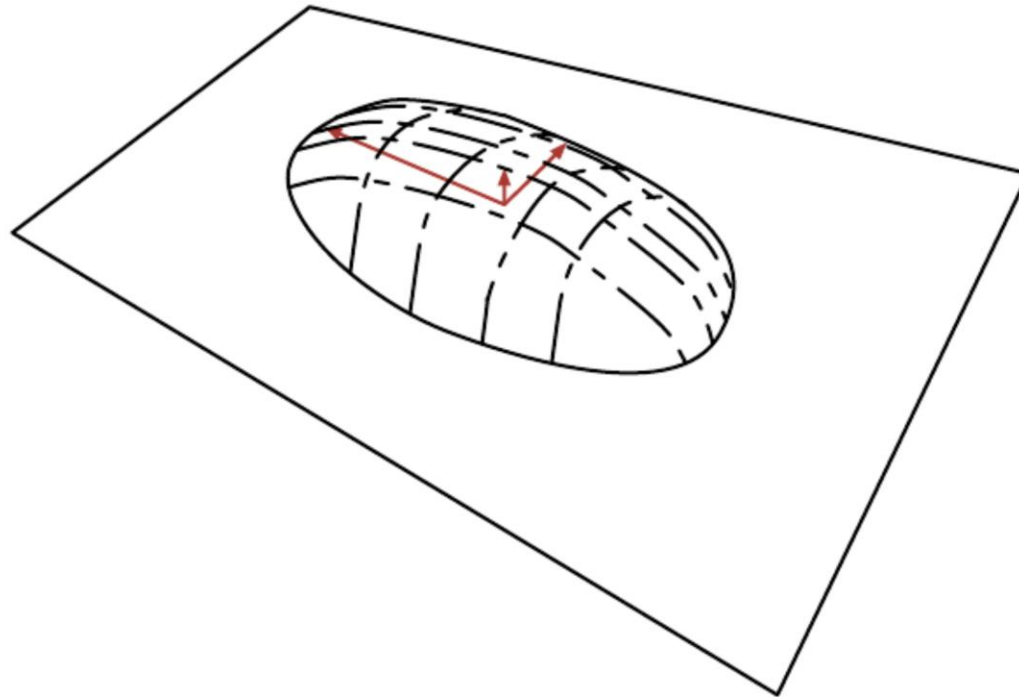
$$= \arg \min_{\mathbf{h}} \lambda \|\mathbf{h}\|_1 + \beta \|\mathbf{x} - \mathbf{W}\mathbf{h}\|_2^2$$

## 13.4 Sparse Coding

- Due to the  $L^1$  norm on  $\mathbf{h}$ , this procedure yield a sparse  $\mathbf{h}^*$ .
- To train the model, alternate between
  - minimization w.r.t  $\mathbf{h}$
  - and minimization w.r.t.  $\mathbf{W}$ ,
  - where we treat  $\beta$  as a hyperparameter.

# 13.5 Manifold Interpretation of PCA

- Linear factor models can be interpreted as learning a manifold.
- Think of PCA as learning a very crude approximation of the data manifold.



# 15 Representation Learning

- Shared representations are useful
  - to handle multiple modalities or domains,
  - or to transfer learned knowledge.
- A feed-forward network trained by supervised learning can be thought of as performing representation learning.
  - Layers except the last one, which is a linear classifier, learns to provide a representation to the classifier.
- We can design an objective function that encourages the elements of the representation vector  $\mathbf{h}$  to be independent.
  - There can be a tradeoff between preserving as much information as possible and attaining nice properties such as independence.
- Representation learning via semi- or unsupervised learning can utilize unlabeled training data and resolve overfitting.

# 15.1 Greedy Layer-Wise Unsupervised Pretraining

- A canonical example of how a representation learned for one task can sometimes be useful for another task.
  - For example from unsupervised learning, trying to capture the shape of the input distribution, to supervised learning with the same input domain.
- Relies on a single-layer representation learning algorithm.
  - Such as RBM, a single-layer autoencoder, a sparse coding model.
- Now we know that greedy layer-wise pretraining is not required to train fully connected architectures.
  - But the unsupervised pretraining approach was the first method to succeed in training deep network architectures.

# 15.1 Greedy Layer-Wise Unsupervised Pretraining

- Greedy because it optimizes each piece independently, one piece at a time, rather than jointly optimizing all pieces.
- Pretraining because It is only a first step before a joint training algorithm is applied to fine-tune all the layers together.
- Can be used as initialization for
  - other unsupervised learning algorithms such as deep autoencoders,
  - and probabilistic models such as deep belief networks and deep Boltzmann machines.

## 15.1.1 Why Does Unsupervised Pretraining Work?

- There are two ideas underlying unsupervised pretraining.
- Regularizing effect on initial parameters
  - It is least well understood.
  - It is possible that pretraining initializes the model in a location that would otherwise be inaccessible.
  - Can expect the unsupervised training to be more effective when the initial representation is poor.
    - Especially useful when processing words, less useful when processing images.
  - Likely to be most useful when the function to be learned is extremely complicated.

## 15.1.1 Why Does Unsupervised Pretraining Work?

- Learning about input distribution can help learning the mapping from inputs to outputs.
  - This is better understood.
  - But not yet at a mathematical, theoretical level.
- Pretraining reduces the variance of the estimation process.
- Less is known about the effect of unsupervised pretraining in conjunction with contemporary approaches.



## 15.1.1 When Does Unsupervised Pretraining Work?

- For Boltzmann machines, greedy layer-wise unsupervised pretraining can yield substantial improvements in test error for classification tasks.
- On many other tasks, however, greedy layer-wise unsupervised pretraining either does not confer a benefit or even causes noticeable harm.
- Unsupervised learning has the disadvantage that it operates with two separate training phases, pretraining and fine-tuning, which makes choosing hyperparameters difficult.

## 15.1.1 When Does Unsupervised Pretraining Work?

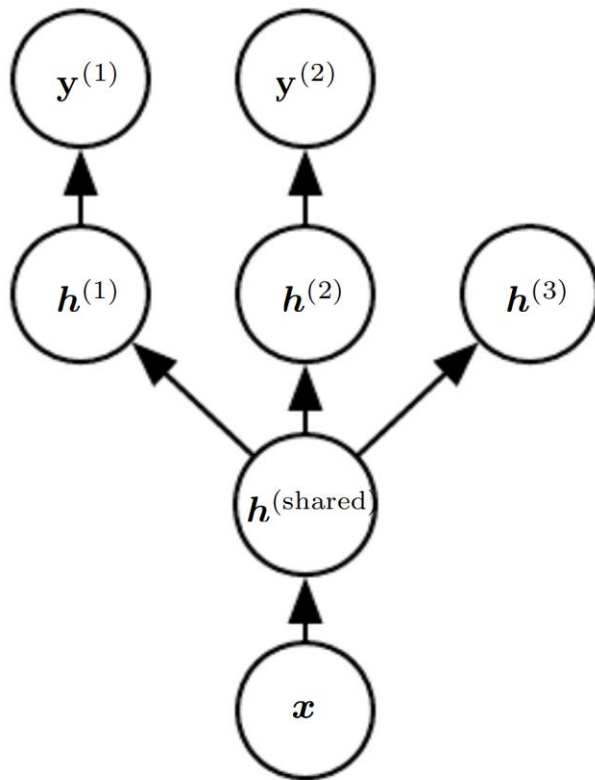
- Today, unsupervised pretraining has been largely abandoned, except in the field of natural language processing.
- Supervised learning, regularized with dropout or batchnorm, outperform unsupervised pretraining on medium-sized datasets.
  - CIFAR-10, MNIST: have ~5k labeled examples per class
- On extremely small datasets Bayesian methods outperform unsupervised pretraining.

## 15.2 Transfer Learning

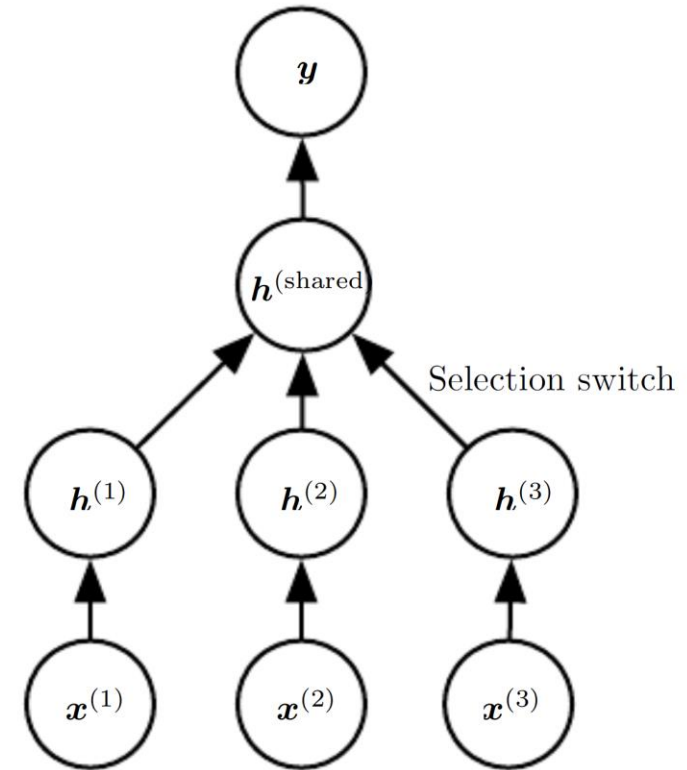
- The learner must perform two or more different tasks  $P_i$ .
- Assume that many of the factors that explain the variations in  $P_1$  are relevant to the variations that need to be captured for learning  $P_2$ .
  - Many visual categories share **low-level** notions such as
    - edges and visual shapes,
    - the effects of geometric changes,
    - changes in lighting, etc.
  - Sometimes what is shared among the different tasks is not the semantics of the input but the semantics of the output, when it makes more sense to share the **upper layers** (near the output) of the neural network, and have a task-specific preprocessing.

# 15.2 Transfer Learning

- Many visual categories share low-level notions such as edges and visual shapes, the effects of geometric changes, changes in lighting, etc.



- Sometimes what is shared among the different tasks is not the semantics of the input but the semantics of the output, when it makes more sense to share the upper layers (near the output) of the neural network, and have a task-specific preprocessing.



## 15.2 Domain Adaptation

- The task remains the same, but the input distribution is slightly different.
- For example, a sentiment predictor trained on **customer reviews of media content** such as books, videos, and music is later used to analyze **comments about consumer electronics** such as televisions or smartphones.

# 15.2 Other Forms of Transfer Learning

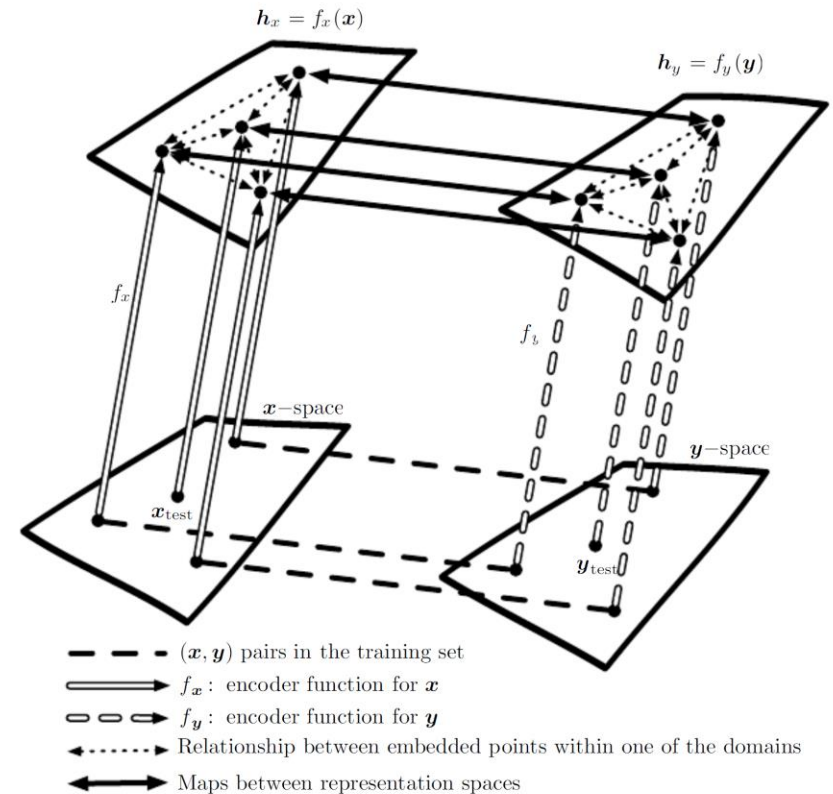
- Concept Drift
  - A form of transfer learning due to gradual changes in the data distribution over time.
- One-shot learning
  - Only one labeled example of the transfer task is given.
  - Possible because the representation is learned to cleanly **separate** the underlying classes **during the first stage**.
  - Then during **the transfer learning stage** only one labeled example is needed to **infer the label** of many possible test examples that all cluster around the same point in representation space.

## 15.2 Other Forms of Transfer Learning

- Zero-shot learning, zero-data learning
  - No labeled examples are given at all.
  - Only possible because additional information has been exploited during training.
  - Includes three random variables
    - The traditional inputs,  $\mathbf{x}$ ,
    - The traditional outputs or targets,  $\mathbf{y}$ ,
    - An additional random variable describing the **task**,  $T$ .
      - For example, a representation of a question such as “Is there a cat in this image?”.
  - The model is trained to estimate  $p(\mathbf{y}|\mathbf{x}, T)$ .

# 15.2 Zero-shot learning & Zero-data Learning

- Zero-shot learning between similar input domains.
  - A word  $\mathbf{x}_{\text{test}}$  and another word  $\mathbf{y}_{\text{test}}$  from two different languages.
- Zero-data learning between different domains.
  - An image  $\mathbf{x}_{\text{test}}$  and a word  $\mathbf{y}_{\text{test}}$ .





## 15.3 Semi-Supervised Disentangling of Causal Factors

- Ideal representation is the one whose features or **directions** in the feature space correspond to different **causes**, thereby disentangling the causes.
- Two main strategies for dealing with a large number of causes
  - Do supervised learning and unsupervised learning at the same time.
  - Use much larger representations if using purely unsupervised learning.
- An emerging strategy is designing a model that **learns** how to determine what is salient, such as generative adversarial networks (GANs).
- The benefit of learning causal factors is that the causal mechanisms remain **invariant** when considering changes in distribution of the cause due to different domains, temporal non-stationarity, or changes in the nature of the task.

# 15.4 Distributed Representation

- Representations composed of many elements that can be set separately from each other.
  - $n$  features with  $k$  values describe  $k^n$  different concepts,
  - whereas symbolic representations such as one-hot representation can only divide  $n$  different regions in input space.
- Examples of learning algorithms based on non-distributed representations
  - Clustering methods including  $k$ -means algorithm
  - $k$ -nearest neighbors algorithms
  - Decision trees
  - Gaussian mixtures
  - Kernel machines
  - Language models based on  $n$ -grams

# 15.4 Distributed Representation

- Examples of distributed representation
  - Hidden units in a deep convolutional network trained on the ImageNet data.
  - Latent space of a deep convolutional generative adversarial network (DCGAN) trained on faces.
  - Word2Vec

# 15.5 Exponential Gains from Depth

- Deep distributed representations
  - The higher level features (functions of the input) or factors (generative causes) are obtained through the composition of many **nonlinearities**.
  - This can give an exponential boost to statistical efficiency **on top of** the exponential boost given by using a distributed representation.

## 15.6 Providing Clues to Discover Underlying Causes

- To use unlabeled data, representation learning makes use of implicit prior beliefs about the underlying factors that act as generic regularization strategies.
  - Smoothness and linearity
  - Multiple explanatory factors, causal factors, and hierarchical organization of explanatory factors
  - Shared factors across tasks
  - Manifold of probability mass and natural clustering
  - Temporal and spatial coherence from slowness principle
  - Sparsity in the sense that most features should not be relevant to describing most inputs
  - Simplicity of factor dependencies

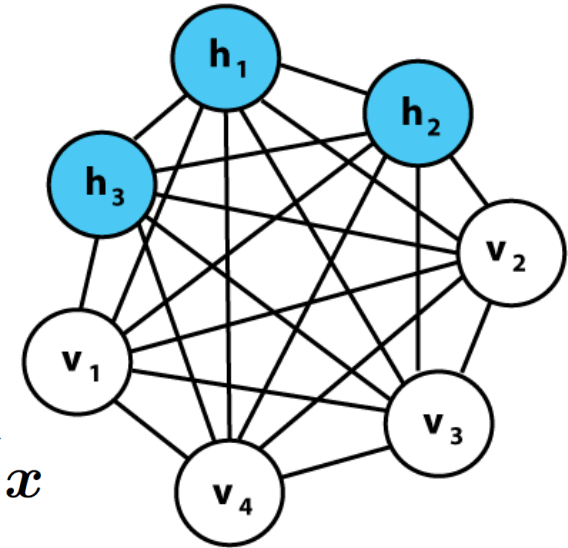
## 20. Deep Generative Models

- All of the models represent probability distributions over multiple variables in some way.
- Some allow the probability distribution function to be evaluated explicitly.
  - Boltzmann machines
- Others do not allow the evaluation of the probability distribution function, but support drawing samples from the distribution.
  - GAN
- Some of these models are structured probabilistic models (graphical models) described in terms of graphs and factors.

# 20.1 Boltzmann Machines

- An energy-based model

$$P(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z} \quad E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{U} \mathbf{x} - \mathbf{b}^\top \mathbf{x}$$



- When all variables are observed, the probability of one unit being on is given by a linear model.
- Becomes more powerful when having non-observed variables, or latent variables.
  - In this case the Boltzmann machine becomes a universal approximator over discrete variables.

## 20.1 Boltzmann machine

- Learning algorithms for Boltzmann machines are usually based on maximum likelihood.
- All Boltzmann machines have an intractable partition function.
  - So the maximum likelihood gradient must be approximated using MCMC algorithms for maximizing the log-likelihood with an intractable partition function using gradient ascent, for example contrastive divergence (CD, or CD- $k$  to indicate CD with  $k$  Gibbs steps)



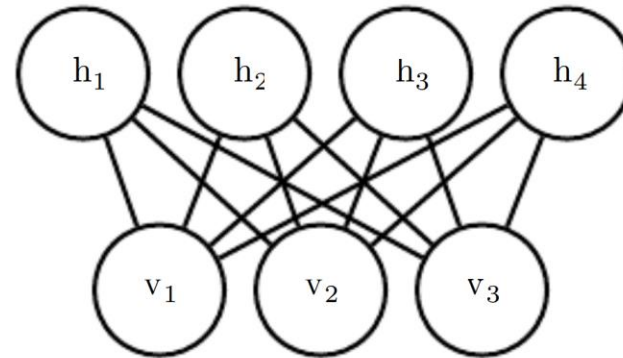
## 20.2 Restricted Boltzmann Machine

- Has a bipartite graph model.

$$P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}$$

$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp\{-E(\mathbf{v}, \mathbf{h})\}$$



- Naïve method of computing  $Z$  is intractable, implying  $P(\mathbf{v})$  is also intractable to evaluate.

## 20.2.1 Conditional Distributions

- The bipartite graph structure of the RBM has the very special property that its conditional distributions  $P(\mathbf{h} \mid \mathbf{v})$  and  $P(\mathbf{v} \mid \mathbf{h})$  are factorial and relatively simple to compute and to sample from.

$$\begin{aligned} P(\mathbf{h} \mid \mathbf{v}) &= \frac{P(\mathbf{h}, \mathbf{v})}{P(\mathbf{v})} \\ &= \frac{1}{P(\mathbf{v})} \frac{1}{Z} \exp \left\{ \mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \\ &= \frac{1}{Z'} \exp \left\{ \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \\ &= \frac{1}{Z'} \exp \left\{ \sum_{j=1}^{n_h} c_j h_j + \sum_{j=1}^{n_h} \mathbf{v}^\top \mathbf{W}_{:,j} h_j \right\} \\ &= \frac{1}{Z'} \prod_{j=1}^{n_h} \exp \left\{ c_j h_j + \mathbf{v}^\top \mathbf{W}_{:,j} h_j \right\} \end{aligned}$$

$$\begin{aligned} P(h_j = 1 \mid \mathbf{v}) &= \frac{\tilde{P}(h_j = 1 \mid \mathbf{v})}{\tilde{P}(h_j = 0 \mid \mathbf{v}) + \tilde{P}(h_j = 1 \mid \mathbf{v})} \\ &= \frac{\exp \{c_j + \mathbf{v}^\top \mathbf{W}_{:,j}\}}{\exp \{0\} + \exp \{c_j + \mathbf{v}^\top \mathbf{W}_{:,j}\}} \\ &= \sigma \left( c_j + \mathbf{v}^\top \mathbf{W}_{:,j} \right). \end{aligned}$$

$$P(\mathbf{h} \mid \mathbf{v}) = \prod_{j=1}^{n_h} \sigma \left( (2\mathbf{h} - 1) \odot (\mathbf{c} + \mathbf{W}^\top \mathbf{v}) \right)_j$$

$$P(\mathbf{v} \mid \mathbf{h}) = \prod_{i=1}^{n_v} \sigma \left( (2\mathbf{v} - 1) \odot (\mathbf{b} + \mathbf{W} \mathbf{h}) \right)_i$$

## 20.2.2 Training RBM

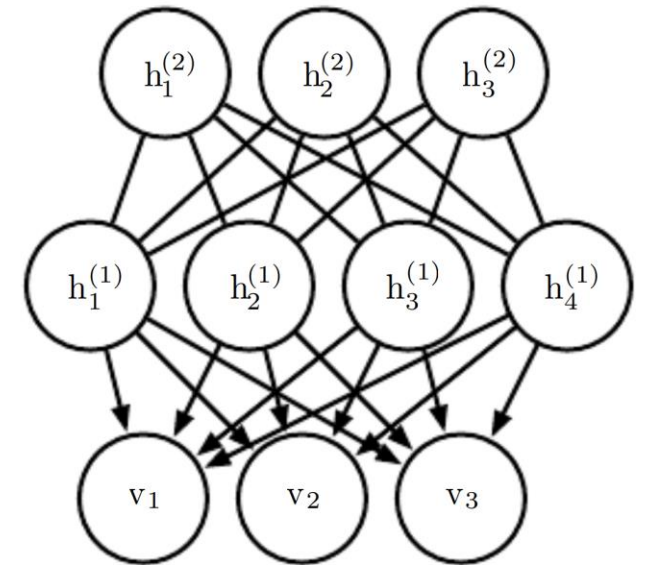
- RBM admits efficient evaluation and differentiation of  $\tilde{P}(\mathbf{v})$  and efficient MCMC sampling in the form of block Gibbs sampling.
- Can be trained with any technique for training models that have intractable partition functions.
  - This includes contrastive divergence (CD), stochastic maximum likelihood (SML, or persistent contrastive divergence (PCD)), ratio matching and so on.
- Compared to other undirected models used in deep learning, the RBM is relatively straightforward to train because we can compute  $P(\mathbf{h} \mid \mathbf{v})$  exactly in closed form.

## 20.3 Deep Belief Networks

- The introduction of deep belief networks (DBNs) in 2006 began the current deep learning renaissance.
- Prior to the introduction of DBNs, deep models were considered too difficult to optimize, kernel machines with convex objective functions dominated the research landscape.
- DBNs demonstrated that deep architectures can be successful, by outperforming kernelized SVMs on the MNIST dataset.
- Today, DBNs are rarely used, even compared to other unsupervised or generative learning algorithms.
- But they are still deservedly recognized for their important role in deep learning history.

## 20.3 DBNs

- Generative models with several layers of latent variables.
- The latent variables are typically binary, while the visible units may be binary or real.
- There are no intra-layer connections.
- The connections between the top two layers are undirected.
- The connections between all other layers are directed.



## 20.3 DBNs

- The probability distribution represented by the DBN

$$P(\mathbf{h}^{(l)}, \mathbf{h}^{(l-1)}) \propto \exp \left( \mathbf{b}^{(l)\top} \mathbf{h}^{(l)} + \mathbf{b}^{(l-1)\top} \mathbf{h}^{(l-1)} + \mathbf{h}^{(l-1)\top} \mathbf{W}^{(l)} \mathbf{h}^{(l)} \right),$$

$$P(h_i^{(k)} = 1 \mid \mathbf{h}^{(k+1)}) = \sigma \left( b_i^{(k)} + \mathbf{W}_{:,i}^{(k+1)\top} \mathbf{h}^{(k+1)} \right) \forall i, \forall k \in 1, \dots, l-2,$$

$$P(v_i = 1 \mid \mathbf{h}^{(1)}) = \sigma \left( b_i^{(0)} + \mathbf{W}_{:,i}^{(1)\top} \mathbf{h}^{(1)} \right) \forall i.$$

- In the case of real-valued visible units,

$$\mathbf{v} \sim \mathcal{N} \left( \mathbf{v}; \mathbf{b}^{(0)} + \mathbf{W}^{(1)\top} \mathbf{h}^{(1)}, \beta^{-1} \right)$$

## 20.3 DBNs

- To generate a sample,
  - first run several steps of Gibbs sampling on the top two hidden layers,
  - then use a single pass of ancestral sampling through the rest of the model to draw a sample from the visible units.

## 20.3 DBNs

- Have many of the problems associated with both directed and undirected models.
- Inference is intractable
  - due to the explaining away effect within each directed layer,
  - and due to the interaction between the two hidden layers with undirected connections.
- Evaluating or maximizing the standard ELBO on the log-likelihood is also intractable
  - because the ELBO takes the expectation of cliques of size equal to the network width.
- Evaluating or maximizing the log-likelihood requires
  - not just confronting the problem of intractable inference to marginalize out the latent variables,
  - but also the problem of an intractable partition function within the undirected model of the top two layers.



## 20.3 DBNs

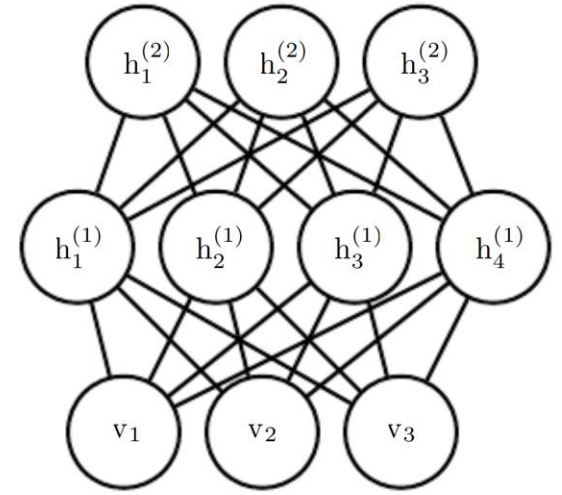
- To train DBNs,
  - begin by training an RBM to maximize  $\mathbb{E}_{\mathbf{v} \sim p_{\text{data}}} \log p(\mathbf{v})$  using CD or SML,
  - then the RBM defines the first layer of the DBN.
  - Next a second RBM is trained to approximately maximize
$$\mathbb{E}_{\mathbf{v} \sim p_{\text{data}}} \mathbb{E}_{\mathbf{h}^{(1)} \sim p^{(1)}(\mathbf{h}^{(1)}|\mathbf{v})} \log p^{(2)}(\mathbf{h}^{(1)})$$
  - And repeat adding as many as layers to the DBN as desired.
  - This can be justified as increasing a variational lower bound on the log-likelihood of the data under the DBN.
- In most applications, no effort is made to jointly train the DBN after the greedy layer-wise procedure is complete.

## 20.3 DBNs

- The trained DBN may be used directly as a generative model.
- But most of the interest in DBNs arose from their ability to improve classification models by taking the weights from the DBN and using them to define an MLP.
- The term “deep belief network” may cause some confusion because the term “belief network” is sometimes used to refer to purely directed models, while DBNs contain an undirected layer.

## 20.4 Deep Boltzmann Machines

- Entirely undirected model unlike DBN.
- Has several layers of latent variables unlike RBM.
- Within each layer, like RBM,
  - each variable is mutually independent,
  - and is conditioned on the variables in the neighboring layers.
- Typically contain only binary units like RBMs and DBNs, but it is straightforward to include real-valued visible units.



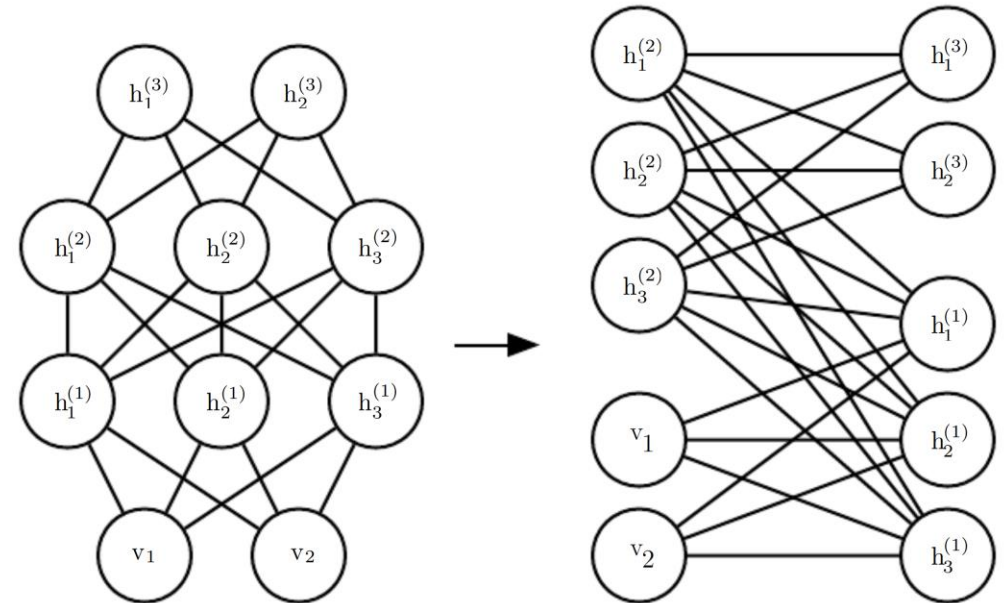
## 20.4 DBMs

- An energy-based model

$$P\left(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}\right) = \frac{1}{Z(\boldsymbol{\theta})} \exp\left(-E\left(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \boldsymbol{\theta}\right)\right)$$

$$E\left(\mathbf{v}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \boldsymbol{\theta}\right) = -\mathbf{v}^\top \mathbf{W}^{(1)} \mathbf{h}^{(1)} - \mathbf{h}^{(1)\top} \mathbf{W}^{(2)} \mathbf{h}^{(2)} - \mathbf{h}^{(2)\top} \mathbf{W}^{(3)} \mathbf{h}^{(3)}$$

- Can be organized into a bipartite graph.
  - Conditional independence between variables of odd layers and variables of even layers.



## 20.4 DBMs

- The bipartite structure means that we can apply the same equations used for an RBM to determine the conditional distributions in a DBM.
- The units within a layer are conditionally independent from each other given the values of the neighboring layers.

- For example, when a DBM has two hidden layers,

$$P(v_i = 1 \mid \mathbf{h}^{(1)}) = \sigma \left( \mathbf{W}_{i,:}^{(1)} \mathbf{h}^{(1)} \right)$$

$$P(h_i^{(1)} = 1 \mid \mathbf{v}, \mathbf{h}^{(2)}) = \sigma \left( \mathbf{v}^\top \mathbf{W}_{:,i}^{(1)} + \mathbf{W}_{i,:}^{(2)} \mathbf{h}^{(2)} \right)$$

$$P(h_k^{(2)} = 1 \mid \mathbf{h}^{(1)}) = \sigma \left( \mathbf{h}^{(1)\top} \mathbf{W}_{:,k}^{(2)} \right)$$

- The bipartite structure makes Gibbs sampling in a DBM efficient.
  - Gibbs sampling can be divided into two blocks of updates.

## 20.4.1 Interesting Properties

- The posterior distribution  $P(\mathbf{h}|\mathbf{v})$  is simpler compared to DBNs.
  - The simplicity allows richer approximation.
  - DBNs have interactions between hidden units within the same layer because of the explaining-away effect, but a heuristic choice of  $Q(\mathbf{h})$  to obtain a variational lower bound of the posterior ignores those interactions, for example.
  - In DBMs, all of the hidden units within a layer are conditionally independent given the other layers. This lack of intra-layer interaction makes it possible to use fixed point equations to actually optimize the variational lower bound and find the true optimal mean field expectations.
- Sampling from DBMs is relatively difficult.
  - DBNs only need to use MCMC sampling in the top pair of layers.
  - For DBMs it is necessary to use MCMC across all layers.

## 20.4.2 DBM Mean Field Inference

- The conditional distribution over one DBM layer given the neighboring layers is factorial.
- The DBM posterior distribution over their hidden units is easy to approximate with a variational approximation, specifically a mean field approximation.
  - The mean field approximation is a simple form of variational inference, where we restrict the approximating distribution to fully factorial distributions.
  - In the case of the mean field approximation, the approximating family is the set of distributions where the hidden units are conditionally independent.

## 20.4.2 DBM Mean Field Inference

- The mean field approach for the example with two hidden layers
- Let  $Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$  be the approximation of  $P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$ .

- The mean field assumption

$$Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) = \prod_j Q(h_j^{(1)} | \mathbf{v}) \prod_k Q(h_k^{(2)} | \mathbf{v})$$

- Minimize

$$\text{KL}(Q \| P) = \sum_{\mathbf{h}} Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) \log \left( \frac{Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})}{P(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})} \right)$$

- Parametrize Q as a product of Bernoulli distributions.

$$Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v}) = \prod_j (\hat{h}_j^{(1)})^{h_j^{(1)}} (1 - \hat{h}_j^{(1)})^{(1-h_j^{(1)})} \times \prod_k (\hat{h}_k^{(2)})^{h_k^{(2)}} (1 - \hat{h}_k^{(2)})^{(1-h_k^{(2)})}$$



## 20.4.2 DBM Mean Field Inference

- Choosing a member of this  $Q(\mathbf{h}^{(1)}, \mathbf{h}^{(2)} | \mathbf{v})$  family that best fits  $P$ .

$$\tilde{q}(h_i | \mathbf{v}) = \exp \left( \mathbb{E}_{\mathbf{h}_{-i} \sim q(\mathbf{h}_{-i} | \mathbf{v})} \log \tilde{p}(\mathbf{v}, \mathbf{h}) \right)$$

- Applying the above, we obtain the update rule.

$$\hat{h}_j^{(1)} = \sigma \left( \sum_i v_i W_{i,j}^{(1)} + \sum_{k'} W_{j,k'}^{(2)} \hat{h}_{k'}^{(2)} \right), \quad \forall j$$

$$\hat{h}_k^{(2)} = \sigma \left( \sum_{j'} W_{j',k}^{(2)} \hat{h}_{j'}^{(1)} \right), \quad \forall k.$$

- At a fixed point, we have a local maximum of the variational lower bound  $\mathcal{L}(Q)$ .

## 20.4.2 DBM Mean Field Inference

- These fixed point update equations define an iterative algorithm where we alternate updates of  $\hat{h}_j^{(1)}$  and updates of  $\hat{h}_k^{(2)}$ .
- On small problems such as MNIST,  $\sim 10$  iterations can be sufficient.

## 20.4.3 DBM Parameter Learning

- Learning in the DBM must confront both
  - the challenge of an intractable posterior distribution,
  - and the challenge of an intractable partition function.
- Variational inference allows the construction of  $Q(\mathbf{h}|\mathbf{v})$  that approximate the intractable posterior distribution  $P(\mathbf{h}|\mathbf{v})$ .
- Learning then proceeds by maximizing  $\mathcal{L}(\mathbf{v}, Q, \boldsymbol{\theta})$ , the variational lower bound on the intractable log-likelihood,  $\log P(\mathbf{v}; \boldsymbol{\theta})$ , w.r.t  $\boldsymbol{\theta}$ .

## 20.4.3 DBM Parameter Learning

- For a DBM with two hidden layers,

$$\mathcal{L}(Q, \theta) = \sum_i \sum_{j'} v_i W_{i,j'}^{(1)} \hat{h}_{j'}^{(1)} + \sum_{j'} \sum_{k'} \hat{h}_{j'}^{(1)} W_{j',k'}^{(2)} \hat{h}_{k'}^{(2)} - \log Z(\theta) + \mathcal{H}(Q)$$

- This contains the log partition function  $\log Z(\theta)$ .
- Evaluating the probability of a DBM requires approximate methods such as annealed importance sampling (AIS).
- Training the model requires approximations to the gradient of the log partition function, typically stochastic maximum likelihood/persistent contrastive divergence (SML/PCD).
  - Use variational stochastic maximum likelihood as applied to the DBM, Algorithm 20.1.

## 20.4.3 DBM Parameter Learning

- Methods not applicable to DBM
  - Techniques such as pseudolikelihood require the ability to evaluate the unnormalized probabilities, rather than merely obtain a variational lower bound on them.
  - Contrastive divergence is slow for deep Boltzmann machines because they do not allow efficient sampling of the hidden units given the visible units — instead, contrastive divergence would require burning in a Markov chain every time a new negative phase sample is needed.

## 20.4.4 Layer-Wise Pretraining

- Unfortunately, training a DBM using SML/PCD from a random initialization usually results in failure.
  - In some cases, the model fails to learn to represent the distribution adequately.
  - In other cases, the DBM may represent the distribution well, but with no higher likelihood than could be obtained with just an RBM.
- The original and most popular method for overcoming the joint training problem of DBMs is greedy layer-wise pretraining.
  - In this method, each layer of the DBM is trained in isolation as an RBM.
- The DBM may then be trained with PCD.
  - Typically PCD training will make only a small change in the model's parameters
  - and its performance as measured by the log-likelihood it assigns to the data, or its ability to classify inputs.

## 20.4.5 Jointly Training DBMs

- Classic DBMs require greedy unsupervised pretraining and a separate MLP-based classifier on top of the hidden features they extract.
- This has some undesirable properties.
  - It is hard to track performance during training because we cannot evaluate properties of the full DBM while training the first RBM.
  - Implementations of DBMs need to have many different components
    - for CD training of individual RBMs,
    - PCD training of the full DBM,
    - and training based on back-propagation through the MLP.
  - Thus, it is hard to tell how well our hyperparameters are working until quite late in the training process.

## 20.4.5 Jointly Training DBMs

- Two main ways to resolve the joint training problem
  - Centered DBM
    - Unable to compete with appropriately regularized MLPs as a classifier.
  - Multi-prediction DBM
    - Does not lead to good likelihood or samples.
    - Does lead to superior classification and ability to reason well about missing inputs.



## 20.5 Boltzmann Machines for Real-Valued Data

- Gaussian-Bernoulli RBMs
  - RBM with binary hidden units and real-valued visible units
- Undirected Models of Conditional Covariance
  - Gaussian RBM inductive bias is not well suited to some types of data.
    - Information in natural images is embedded in the covariance between pixels rather than in the raw pixel values, the Gaussian RBM only models the conditional mean.
  - These models include the mean and covariance RBM (mcRBM), the mean-product of t-distribution (mPoT) model and the spike and slab RBM (ssRBM).

## 20.5 Boltzmann Machines for Real-Valued Data

- mcRBM

- Uses its hidden units to independently encode the conditional mean and covariance of all observed units.

$$E_{\text{mc}}(\mathbf{x}, \mathbf{h}^{(m)}, \mathbf{h}^{(c)}) = E_{\text{m}}(\mathbf{x}, \mathbf{h}^{(m)}) + E_{\text{c}}(\mathbf{x}, \mathbf{h}^{(c)})$$

$$E_{\text{m}}(\mathbf{x}, \mathbf{h}^{(m)}) = \frac{1}{2} \mathbf{x}^\top \mathbf{x} - \sum_j \mathbf{x}^\top \mathbf{W}_{:,j} h_j^{(m)} - \sum_j b_j^{(m)} h_j^{(m)} \quad E_{\text{c}}(\mathbf{x}, \mathbf{h}^{(c)}) = \frac{1}{2} \sum_j h_j^{(c)} \left( \mathbf{x}^\top \mathbf{r}^{(j)} \right)^2 - \sum_j b_j^{(c)} h_j^{(c)}$$

$$p_{\text{mc}}(\mathbf{x} \mid \mathbf{h}^{(m)}, h^{(c)}) = \mathcal{N} \left( \mathbf{x}; \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{mc}} \left( \sum_j \mathbf{W}_{:,j} h_j^{(m)} \right), \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{mc}} \right) \quad \mathbf{C}_{\mathbf{x}|\mathbf{h}}^{\text{mc}} = \left( \sum_j h_j^{(c)} \mathbf{r}^{(j)} \mathbf{r}^{(j)\top} + \mathbf{I} \right)^{-1}$$

- Sampling from  $p_{\text{mc}}(\mathbf{x} \mid \mathbf{h}^{(m)}, \mathbf{h}^{(c)})$  requires computing  $(\mathbf{C}^{\text{mc}})^{-1}$  at every iteration of learning, which can be impractical.
- Avoid it by sampling directly from the marginal  $p_{\text{mc}}(\mathbf{x})$  using Hamiltonian (hybrid) Monte Carlo.

## 20.6 Convolutional Boltzmann Machines

- Unclear how to generalize a pooling operation of a CNN to the setting of energy-based models.
- CNN work well with inputs of many different sizes; for Boltzmann machines, it is difficult to change the input size.
- Pixels at the boundary of the image pose some difficulty.
  - Both doing zero-padding and not doing zero-padding have issues.

## 20.7 Boltzmann Machines for Structured or Sequential Outputs

- A natural way to represent the relationships between the entries in a structured output  $\mathbf{y}$  is to use a probability distribution  $p(\mathbf{y}|\mathbf{x})$ .
- Boltzmann machines extended to model conditional distributions, can supply this probabilistic model.
- For sequence modeling, conditional Boltzmann machines can represent factors of the form  $p(\mathbf{x}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t-1)})$ .
  - Conditional RBM modeling  $p(\mathbf{x}^{(t)} | \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(t-m)})$  for small  $m$ .
  - RNN-RBM sequence model that consists of an RNN that emits the RBM parameters for each time step.

## 20.8 Other Boltzmann Machines

- Discriminative RBM to maximize  $p(y|\boldsymbol{v})$ 
  - Not as powerful as MLP.
- Higher-order Boltzmann machines in energy functions

## 20.9 Backprop through Random Operations

- Reparametrization trick
  - $y \sim \mathcal{N}(\mu, \sigma^2) \rightarrow y = \mu + \sigma \cdot z, z \sim \mathcal{N}(0, 1)$
- More generally,  $\mathbf{y} \sim p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) \rightarrow \mathbf{y} = f(\mathbf{z}; \mathbf{x}, \boldsymbol{\theta})$ 
  - $f$  is continuous and differentiable everywhere.
  - $\boldsymbol{\theta}$  must not be a function of  $\mathbf{z}$  and vice versa.
  - $\mathbf{z}$  from unit uniform or unit Gaussian.

## 20.9.1 Backprop through Discrete Stochastic Operations

- $\mathbf{y} = f(\mathbf{z}; \mathbf{x}, \boldsymbol{\theta})$ 
  - When  $\mathbf{y}$  is discrete,  $f$  must be a step function.
  - Then the derivative of any cost function  $J(\mathbf{y})$  do not give any information.
- REINFORCE algorithm
  - The expected cost  $\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} J(f(\mathbf{z}; \mathbf{x}, \boldsymbol{\theta}))$  is often a smooth function.
  - It can be estimated without bias using a Monte Carlo average.

$$\begin{aligned}\mathbb{E}_{\mathbf{z}}[J(\mathbf{y})] &= \sum_{\mathbf{y}} J(\mathbf{y}) p(\mathbf{y}) & \frac{\partial \mathbb{E}[J(\mathbf{y})]}{\partial \boldsymbol{\omega}} &= \sum_{\mathbf{y}} J(\mathbf{y}) \frac{\partial p(\mathbf{y})}{\partial \boldsymbol{\omega}} \\ & & &= \sum_{\mathbf{y}} J(\mathbf{y}) p(\mathbf{y}) \frac{\partial \log p(\mathbf{y})}{\partial \boldsymbol{\omega}} \\ & & &\approx \frac{1}{m} \sum_{\mathbf{y}^{(i)} \sim p(\mathbf{y}), i=1}^m J(\mathbf{y}^{(i)}) \frac{\partial \log p(\mathbf{y}^{(i)})}{\partial \boldsymbol{\omega}}\end{aligned}$$

## 20.9.1 Backprop through Discrete Stochastic Operations

- One issue with the simple REINFORCE estimator is that it has a high variance, therefore requires many samples.
- Use variance reduction by using a baseline  $b(\boldsymbol{\omega})$  to offset  $J(\mathbf{y})$ .
- Optimal baseline  $b^*(\boldsymbol{\omega})$  and the changed gradient estimator are

$$b^*(\boldsymbol{\omega})_i = \frac{E_{p(\mathbf{y})} \left[ J(\mathbf{y}) \frac{\partial \log p(\mathbf{y})^2}{\partial \omega_i} \right]}{E_{p(\mathbf{y})} \left[ \frac{\partial \log p(\mathbf{y})^2}{\partial \omega_i} \right]} \quad (J(\mathbf{y}) - b(\boldsymbol{\omega})_i) \frac{\partial \log p(\mathbf{y})}{\partial \omega_i}$$

- $b(\boldsymbol{\omega})$  is obtained by training a neural network to estimate the numerator and the denominator of  $b^*(\boldsymbol{\omega})$  for each element of  $\boldsymbol{\omega}$ .



## 20.10 Directed Generative Nets

- DBNs are a partially directed model.
- Sparse coding models can be thought of as shallow directed generative models.
- Sigmoid belief networks
  - Very similar to DBNs, but the beginning layer units are independent, rather than being a RBM.
- Differentiable generator networks
  - Transforms samples of latent variables  $\mathbf{z}$  to samples  $\mathbf{x}$  or to distributions over samples using a differentiable function  $g(\mathbf{z}; \boldsymbol{\theta})$ , which is a neural network.
  - Includes VAEs (generator network + inference network) and GANs (generator network + discriminator network)

## 20.10.3 Variational Autoencoders

- Sampling
  - Draws a sample  $\mathbf{z}$  from  $p_{\text{model}}(\mathbf{z})$ .
  - Run  $\mathbf{z}$  through a differentiable generator network  $g(\mathbf{z})$ .
  - Sample  $\mathbf{x}$  from  $p_{\text{model}}(\mathbf{x}; g(\mathbf{z})) = p_{\text{model}}(\mathbf{x}|\mathbf{z})$ .
- Training
  - Approximate inference network (i.e. encoder)  $q(\mathbf{z}|\mathbf{x})$  is used to obtain  $\mathbf{z}$ .
  - $p_{\text{model}}(\mathbf{x}|\mathbf{z})$  works as a decoder network

## 20.10.3 Variational Autoencoders

- VAE are trained by maximizing the variational lower bound (ELBO).

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} | \mathbf{x})) \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) || p_{\text{model}}(\mathbf{z}))\end{aligned}$$

- The first equality
  - The first term is the joint log-likelihood of the visible and hidden variables under the approximate posterior over the latent variables.
    - Just like with EM, except that we use an approximate rather than the exact posterior.
  - The second term is the entropy of the approximate posterior.
    - This term encourages  $q(\mathbf{z}|\mathbf{x})$  to place high probability mass on many  $\mathbf{z}$  that could have generated  $\mathbf{x}$ , rather than collapsing to a single point estimate of the most likely value.
- The second equality
  - The first term is the reconstruction log-likelihood found in other autoencoders.
  - The second term tries to make the approximate posterior distribution  $q(\mathbf{z}|\mathbf{x})$  and the model prior  $p_{\text{model}}(\mathbf{z})$  approach each other.

## 20.10.3 Variational Autoencoders

- Train a parametric encoder (inference network)  $f(\mathbf{x}; \theta)$  that produces the parameters of  $q(\mathbf{z}|\mathbf{x})$ .
- Backprop through samples of  $\mathbf{z}$  drawn from  $q(\mathbf{z}|\mathbf{x}) = q(\mathbf{z}; f(\mathbf{x}; \theta))$
- Maximize  $\mathcal{L}$  w.r.t  $\theta$  and the parameters of  $p_{\text{model}}(\mathbf{x}|\mathbf{z})$  by approximating the expectations in  $\mathcal{L}$  by sampling.
- Main drawback is that samples are blurry.
  - Intrinsic effect of maximum likelihood that minimizes  $D_{\text{KL}}(p_{\text{data}} \parallel p_{\text{model}})$ .
  - Practice of using a Gaussian distribution for  $p_{\text{model}}(\mathbf{x}; g(\mathbf{z}))$ .

# 20.14 Evaluating Generative Models

- Evaluation metrics for generative models are often hard research problems.
- Often evaluate generative models by visually inspecting the samples.
  - A generative model trained on data may ignore some modes (for example dropping cats and only generating dogs),
  - and a human observer would not easily be able to inspect or remember enough images to detect the missing variation.
- Inception score
  - Shane Barratt, Rishi Sharma, *A Note on the Inception Score*, <https://arxiv.org/abs/1801.01973>

$$\text{IS}(G) = \exp \left( \mathbb{E}_{\mathbf{x} \sim p_g} D_{KL}(p(y|\mathbf{x}) \parallel p(y)) \right)$$

- IS aimed to codify two desirable qualities of a generative model into a metric
  - The images generated should contain clear objects (i.e. the images are sharp rather than blurry), or  $p(y|x)$  should be low entropy. In other words, the Inception Network should be highly confident there is a single object in the image.
  - The generative algorithm should output a high diversity of images from all the different classes in ImageNet, or  $p(y)$  should be high entropy.