

SHALLOW AUTOENCODER

Ke Chen

Department of Computer Science, The University of Manchester

Ke.Chen@manchester.ac.uk

OVERVIEW

History and main properties

TRADITIONAL AND DENOISING AUTOENCODER

Architecture, Loss function, learning algorithm, illustrative example, manifold perspective

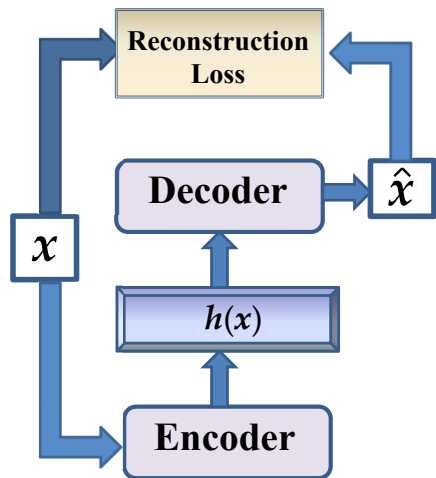
SPARSE AUTOENCODER

Architecture, Loss function, learning algorithm, illustrative example

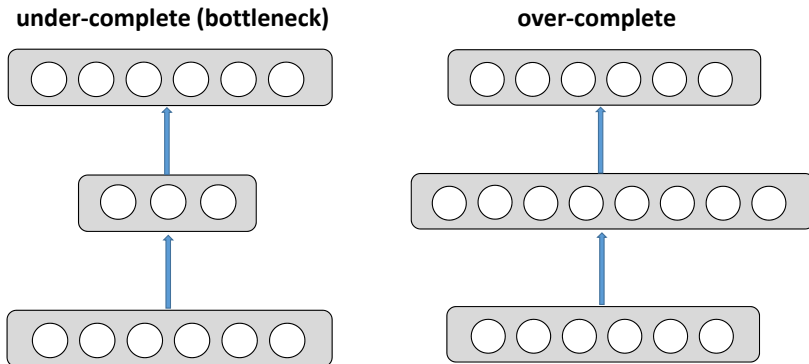
RESTRICTED BOLTZMANN MACHINE

Architecture, energy function, learning algorithm, continuous input, illustrative example

- **Autoencoder (AE)** is a specific type of neural networks, originally named **auto-associator** in 1980s for **dimension reduction** and **feature extraction**.
- In general, **AE** consists of two components, **encoder** and **decoder**, to learn reconstruction of data itself in a **self-supervised learning** manner.
- There are a variety of **AEs**; **discriminative vs. probabilistic**, **static vs. dynamic**, **structured vs. unstructured**, ...
- Nowadays, **AE** is a centre of **representation learning** and closely related to many areas ranging from **manifold learning** to **generative modelling**.



- In general, **AEs** may generate two different types of representations.
 - **Under-complete (bottleneck)**: the dimension of learned representation is **lower** than that of data for dimension reduction
 - **Over-complete**: the dimension of learned representation is higher than that of data to discover and capture intrinsic structures underlying data



Traditional Autoencoder (AE)

- **Architecture:** a MLP of a single **bottleneck** hidden layer ($|\mathbf{h}| < |\mathbf{x}|$) and **tied weight matrix** to generate an **under-complete** representation
- The **hidden layer** is often named **coding** layer.
- **Encoder**

$$\mathbf{h}(\mathbf{x}) = \mathbf{f}(\mathbf{a}_h), \mathbf{a}_h = \mathbf{W}\mathbf{x} + \mathbf{b}_h$$

$$\mathbf{f}(\mathbf{a}_h) = \left\{ f(a_{h,j}) \right\}_{j=1}^{|\mathbf{h}|}$$

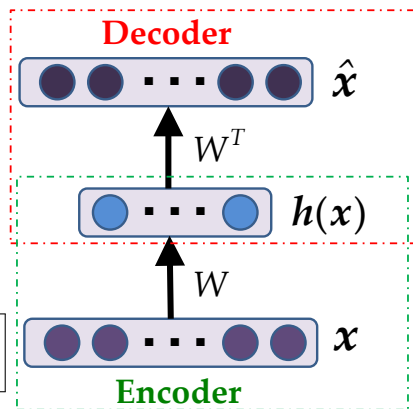
- **Decoder**

$$\hat{\mathbf{x}} = \mathbf{g}(\mathbf{a}_o), \mathbf{a}_o = \mathbf{W}^T \mathbf{h}(\mathbf{x}) + \mathbf{b}_o$$

$$\mathbf{g}(\mathbf{a}_o) = \left\{ g(a_{o,j}) \right\}_{j=1}^{|\mathbf{x}|}$$

$\mathbf{W}_{|\mathbf{h}| \times |\mathbf{x}|}$: (tied) weight matrix

$\mathbf{b}_h, \mathbf{b}_o$: biases for hidden and output layers



Fact: PCA is a special case of the traditional AE when $f(\cdot)$ and $g(\cdot)$ are linear activation functions.

Denoising Autoencoder (DAE)

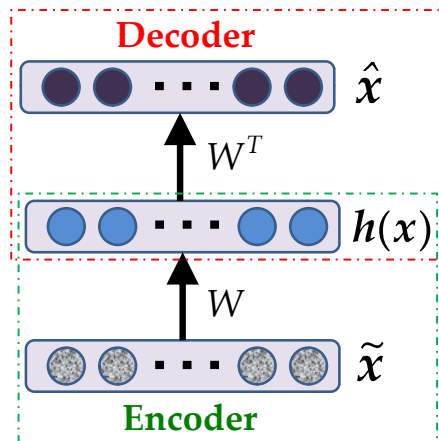
- **Architecture:** a MLP of a hidden layer and **tied weight matrix** to generate either **under-complete** or **over-complete** representation.
- Learn **denoising** by **recovering** a data point, \mathbf{x} , from its **corrupted noisy** version, $\tilde{\mathbf{x}}$
- In **deployment** phase, it can generate a proper representation directly from **test** data.

Encoder

$$h(\tilde{\mathbf{x}}) = f(\mathbf{a}_h), \quad \mathbf{a}_h = W\tilde{\mathbf{x}} + \mathbf{b}_h.$$

Decoder

$$\hat{\mathbf{x}} = g(\mathbf{a}_o), \quad \mathbf{a}_o = W^T h(\mathbf{x}) + \mathbf{b}_o.$$



Fact: In DAE, untied weight matrices are used occasionally in applications.

Noisy Training Data Generation

- Given a training dataset, $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^{|\mathcal{D}|}$, generate noisy data, $\tilde{\mathbf{x}}_i$, to train DAE by corrupting \mathbf{x}_i with
 - Gaussian noise**: for real-valued \mathbf{x}_i , $\tilde{\mathbf{x}}_i = \mathbf{x}_i + \epsilon$, ϵ randomly drawn from $N(\mathbf{0}, \sigma^2 I)$ and the amount of noise controlled by σ
 - Salt-and-peper noise**: for discrete-valued \mathbf{x}_i , generate $\tilde{\mathbf{x}}_i$ by flipping some randomly chosen elements' value of \mathbf{x}_i to either maximum or minimum of the domain range in \mathcal{D}
 - Masking noise**: for discrete-valued \mathbf{x}_i , generate $\tilde{\mathbf{x}}_i$ by setting some randomly chosen elements of \mathbf{x}_i to zero

Loss Function

Given a training dataset, $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{x}_i)\}_{i=1}^{|\mathcal{D}|}$ (AE) or $\mathcal{D} = \{(\tilde{\mathbf{x}}_i, \mathbf{x}_i)\}_{i=1}^{|\mathcal{D}|}$ (DAE), loss functions for AE/DAE are defined based on **real-valued** or **binary-valued** input

- Mean squared error (MSE) loss for real-valued or categorical-valued input (**loss-1**)

$$\mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o; \mathcal{D}) = \frac{1}{2|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2. \quad (1)$$

$\hat{\mathbf{x}}_i$: output of AE/DAE for input \mathbf{x}_i

$g(\cdot)$: linear activation function in output layer (see Slides 5 and 6)

- Cross-entropy loss for binary-valued input (**loss-2**)

$$\mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o; \mathcal{D}) = -\frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} \sum_{j=1}^{|\mathbf{x}_i|} \left(x_{ij} \log \hat{x}_{ij} + (1 - x_{ij}) \log(1 - \hat{x}_{ij}) \right), \quad x_{ij} \in \{0, 1\}. \quad (2)$$

$\hat{\mathbf{x}}_i = (\hat{x}_{i1}, \dots, \hat{x}_{ij}, \dots, \hat{x}_{i|\mathbf{x}_i|})$: output of AE/DAE for input $\mathbf{x}_i = (x_{i1}, \dots, x_{ij}, \dots, x_{i|\mathbf{x}_i|})$

$g(\cdot)$: sigmoid activation function in output layer (see Slides 5 and 6)

Learning Algorithm

Input a training set, $\mathcal{D} = \{(\mathbf{z}_i, \mathbf{x}_i)\}_{i=1}^{|\mathcal{D}|}$ where $\mathbf{z}_i = \mathbf{x}_i$ for AE or $\mathbf{z}_i = \tilde{\mathbf{x}}_i$ for DAE
 Randomly initialize W , \mathbf{b}_h and \mathbf{b}_o and pre-set a learning rate η and batch size $|\mathcal{B}|$

• Forward Computation

For the input \mathbf{z}_i ($i = 1, \dots, |\mathcal{B}|$), output of the hidden layer is

$$\mathbf{h}(\mathbf{z}_i) = \mathbf{f}(\mathbf{a}_h(\mathbf{z}_i)), \quad \mathbf{a}_h(\mathbf{z}_i) = W\mathbf{z}_i + \mathbf{b}_h.$$

And output of the output layer is

$$\hat{\mathbf{x}}_i = \mathbf{g}(\mathbf{a}_o(\mathbf{z}_i)), \quad \mathbf{a}_o(\mathbf{z}_i) = W^T \mathbf{h}(\mathbf{z}_i) + \mathbf{b}_o.$$

where $g(\cdot)$ is linear and sigmoid activation function for real-valued/categorical-valued and binary-valued input, respectively.

Learning Algorithm

• Backward Gradient Computation

For $i = 1, 2, \dots, |\mathcal{B}|$, compute gradients of loss function with respect to parameters

- Gradients for the output layer depend on loss functions.

$$\text{loss-1 : } \delta_o(\mathbf{z}_i, \mathbf{x}_i) = \frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{a}_o(\mathbf{z}_i)} = \hat{\mathbf{x}}_i - \mathbf{x}_i$$

$$\text{loss-2 : } \delta_o(\mathbf{z}_i, \mathbf{x}_i) = \frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{a}_o(\mathbf{z}_i)} = \mathbf{x}_i - \hat{\mathbf{x}}_i$$

- Gradients for hidden layer can be computed with gradient of output layer.

$$\frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{h}(\mathbf{z}_i)} = W \delta_o(\mathbf{z}_i), \quad \delta_h(\mathbf{z}_i, \mathbf{x}_i) = \frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{a}_h(\mathbf{z}_i)} = \left(\mathbf{f}'(\mathbf{a}_h(\mathbf{z}_i)) \right) \odot \frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{h}(\mathbf{z}_i)}$$

- Gradients for biases are computed as follows:

$$\frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{b}_o} = \delta_o(\mathbf{z}_i, \mathbf{x}_i), \quad \frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{b}_h} = \delta_h(\mathbf{z}_i, \mathbf{x}_i).$$

Learning Algorithm

- **Parameter Update**

Applying the **gradient descent** method and **tied weight matrix** leads to update rules:

$$W \leftarrow W - \frac{\eta}{2|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left(\underbrace{\delta_h(\mathbf{z}_i, \mathbf{x}_i) \mathbf{z}_i^T}_{\text{encoder}} + \underbrace{\mathbf{h}(\mathbf{z}_i) (\delta_o(\mathbf{z}_i, \mathbf{x}_i))^T}_{\text{decoder}} \right),$$

$$\mathbf{b}_o \leftarrow \mathbf{b}_o - \frac{\eta}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \delta_o(\mathbf{z}_i, \mathbf{x}_i),$$

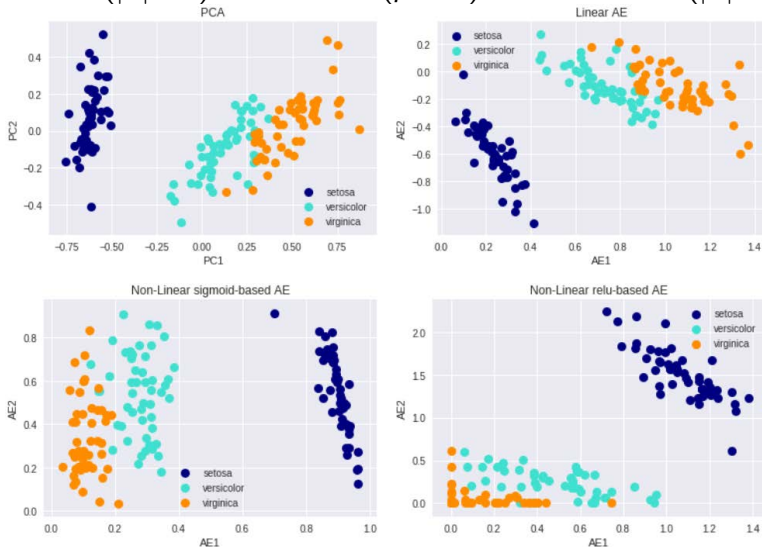
$$\mathbf{b}_h \leftarrow \mathbf{b}_h - \frac{\eta}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \delta_h(\mathbf{z}_i, \mathbf{x}_i).$$

The above three steps repeat for all mini-batches until terminated with **early stopping**.

Fact: While the traditional AE always uses the tied weights, DAE may use untied weights in applications.

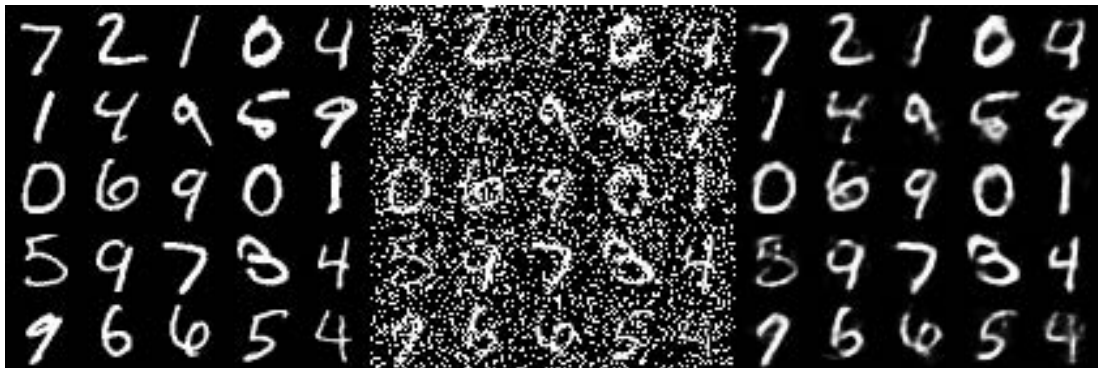
Illustrative Example

- Visualisation: AE ($|h|=2$) versus PCA ($p=2$) on Iris dataset ($|\mathbf{x}|=4$)



Illustrative Example

- **Denoising:** bottleneck DAE; original (left), noisy (middle), recovered (right)



Illustrative Example

- **Sparse representation:** over-complete DAE (200 hidden neurons, Gaussian noise $\sigma = 0.5$; 12×12 image patches (left), 12×12 weights associated with each of first 100 hidden neurons (right) that learns **Gabor-like local oriented edge detectors**

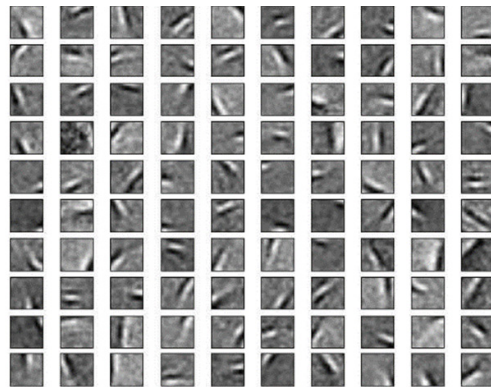
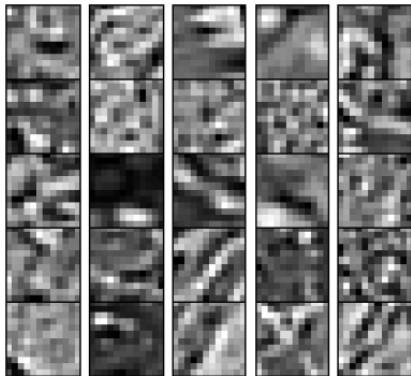
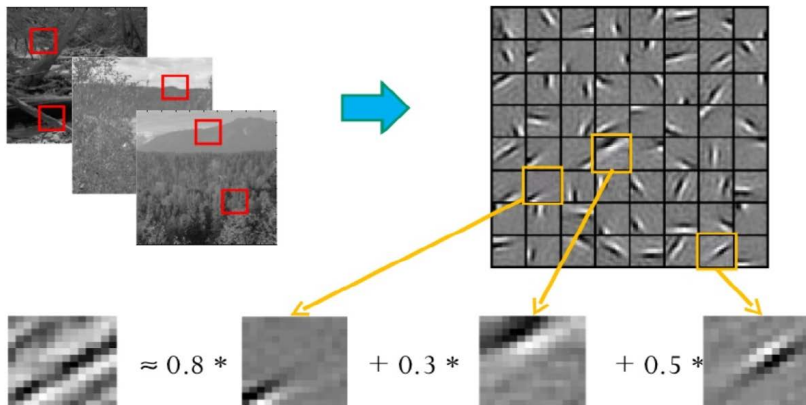


Illustration of Sparse Representation

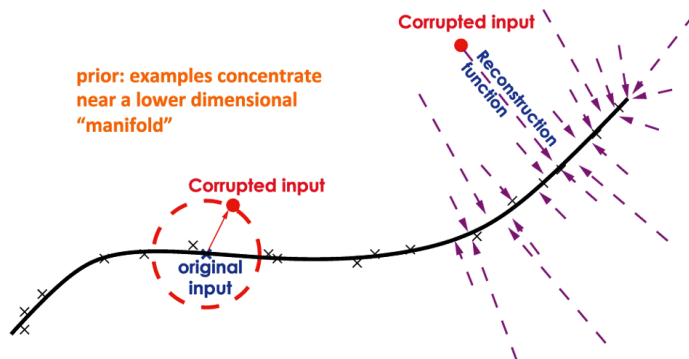
- **Sparse representation**: compact, interpretable, biologically plausible



$[a_1, \dots, a_{64}] = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, 0]$
(feature representation)

Manifold perspective of DAE

- **Manifold assumption**: natural high-dimensional data often concentrated close to a nonlinear low-dimensional **manifold**
- **DAE**: learn **modelling** manifold and **capture main variations** along the manifold
- **Output of encoder**: interpreted as a **coordinate system** on the manifold



Sparse Autoencoder (SAE)

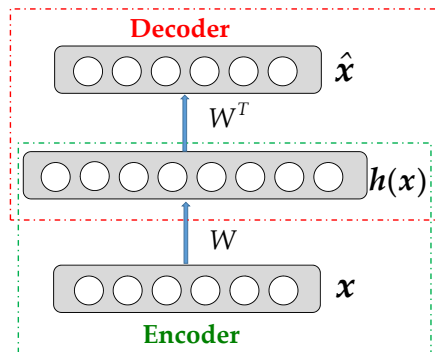
- **Architecture:** NNs of encoder and decoder that generate the **over-complete** representation
- AE extensible to SAE by adding a **regularisation penalty** to avoid learning **identity** function
- Learning with **regularised loss** leads to **sparse representation** reflecting intrinsic structure underlying data.

- **Encoder**

$$h(\mathbf{x}) = f(\mathbf{a}_h), \quad \mathbf{a}_h = W\mathbf{x} + \mathbf{b}_h.$$

- **Decoder**

$$\hat{\mathbf{x}} = g(\mathbf{a}_o), \quad \mathbf{a}_o = W^T h(\mathbf{x}) + \mathbf{b}_o.$$



Fact: In SAE, untied weight matrices may be used in applications.

KL-sparsity Penalty

- **Motivation**: make **hidden** neurons **inactive in most of time**
- Given a training dataset, $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^{|\mathcal{D}|}$, the **averaged** activation of hidden neuron j :

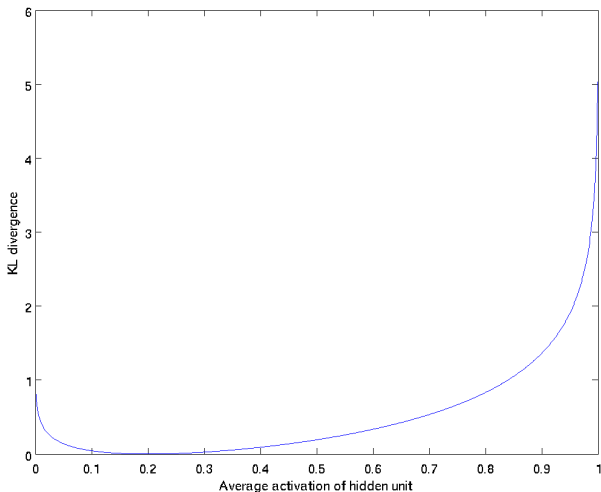
$$\hat{\rho}_j = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} h_j(\mathbf{x}_i) = \frac{1}{|\mathcal{D}|} \sum_{i=1}^{|\mathcal{D}|} f(a_{h,j}(\mathbf{x}_i)), \quad j = 1, 2, \dots, |\mathbf{h}|.$$

- To ensure the sparsity, set up a **constraint**: $\hat{\rho}_j = \rho$, where ρ (**sparsity degree**) is a constant of small value close to zero, e.g., 0.05.
- **KL-sparsity penalty**

$$\mathcal{R}_{KL}(W, \mathbf{b}_h) = \sum_{j=1}^{|\mathbf{h}|} KL(\rho || \hat{\rho}_j) = \sum_{j=1}^{|\mathbf{h}|} \left(\rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j} \right).$$

KL-sparsity Penalty

- **Property:** $KL(\rho || \hat{\rho}_j) = 0$ if $\hat{\rho}_j = \rho$ or increases monotonically as $\hat{\rho}_j$ diverges from ρ .
For instance, set $\rho = 0.2$ for one **sigmoid** hidden neuron



SAE Learning

- Based on **loss-1** or **loss-2**, the **regularised loss for SAE** is

$$\mathcal{L}_R(W, \mathbf{b}_h, \mathbf{b}_o; \mathcal{D}, \rho) = \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o; \mathcal{D}) + \lambda \mathcal{R}_{KL}(W, \mathbf{b}_h),$$

where λ is a trade-off coefficient (hyper-parameter to be tuned during learning).

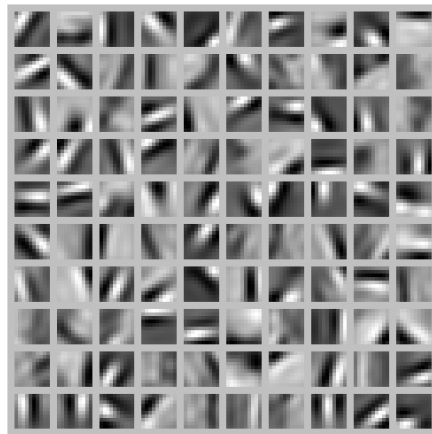
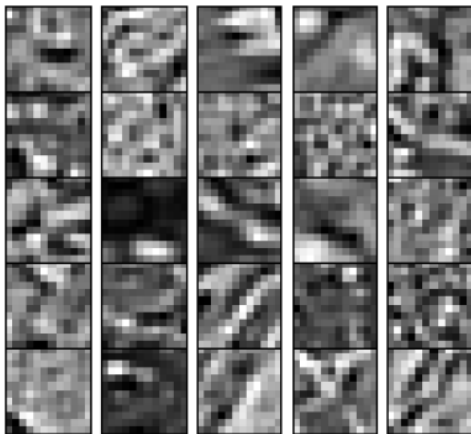
- Adapt the **learning algorithm** for AE to SAE with the following modification:
 - Forward computation:** further compute the averaged activations of all $|\mathbf{h}|$ hidden neurons, $\hat{\rho}_1, \hat{\rho}_2, \dots, \hat{\rho}_{|\mathbf{h}|}$, on training dataset \mathcal{D}
 - Backward gradient computation:** add the **gradient of $\mathcal{R}_{KL}(W, \mathbf{b}_h)$** to that associated with the **hidden layer**

$$\delta_h(\mathbf{x}_i, \mathbf{x}_i) = \frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{a}_h(\mathbf{x}_i)} = \left(\mathbf{f}'(\mathbf{a}_h(\mathbf{x}_i)) \right) \odot \left(\frac{\partial \mathcal{L}(W, \mathbf{b}_h, \mathbf{b}_o)}{\partial \mathbf{h}(\mathbf{x}_i)} + \lambda \boldsymbol{\delta}_{KL} \right),$$

where $\boldsymbol{\delta}_{KL} = \{\delta_{KL}^{(j)}\}_{j=1}^{|\mathbf{h}|}$, $\delta_{KL}^{(j)} = \frac{\partial \mathcal{R}(W, \mathbf{b}_h)}{\partial a_{h,j}(\mathbf{x}_i)} = -\frac{\rho}{\hat{\rho}_j} + \frac{1-\rho}{1-\hat{\rho}_j}$ and $\mathbf{a}_h(\mathbf{x}_i) = \{a_{h,j}\}_{j=1}^{|\mathbf{h}|}$.

Illustrative Example

- **Sparse representation:** SAE (200 hidden neurons, $\rho = 0.05$); 12×12 image patches (left), 12×12 weights associated with each of first 100 hidden neurons (right) that learns **Gabor-like local oriented edge detectors**



Restricted Boltzmann Machine (RBM)

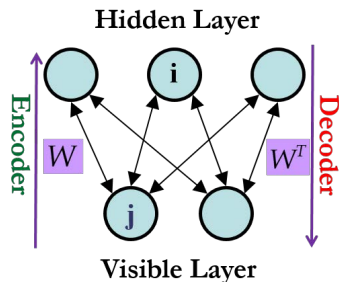
- **Architecture:** 2-layer probabilistic NN of encoder and decoder with **bi-directional** weight matrix: $W_{ij} = W_{ji}$
- **Connections:** only between visible and hidden layer
- **Probabilistic neuron:** output probability for **hidden state** and “**reconstruction**” for **binary-valued hidden** and **visible** units, $h_i \in \{1, 0\}$, $v_j \in \{1, 0\}$, $P(h_i = 1|\mathbf{v})$, $P(v_j = 1|\mathbf{h})$

- **Encoder**

$$P(h_i|\mathbf{v}) = \phi\left(\sum_{j=1}^{|\mathbf{v}|} W_{ij} v_j + b_{h,i}\right), \quad P(\mathbf{h}|\mathbf{v}) = \prod_{i=1}^{|\mathbf{h}|} P(h_i|\mathbf{v}).$$

- **Decoder**

$$P(v_j|\mathbf{h}) = \phi\left(\sum_{i=1}^{|\mathbf{h}|} W_{ji} h_i + b_{v,j}\right), \quad P(\mathbf{h}|\mathbf{v}) = \prod_{j=1}^{|\mathbf{v}|} P(v_j|\mathbf{h}).$$



Sigmoid activation function

$$\phi(x) = \frac{1}{1 + e^{-x}}$$

Fact: RBM inspired by Physics was invented by P. Smolensky in 1986 and rose to prominence after G. Hinton in 2002.

RESTRICTED BOLTZMMAN MACHINE

- **Energy function**: associate a scalar energy to each configuration of states in system

$$E(\mathbf{v}, \mathbf{h}; \Theta) = -\mathbf{v}^T W \mathbf{h} - \mathbf{b}_v^T \mathbf{v} - \mathbf{b}_h^T \mathbf{h} = -\sum_{i=1}^{|\mathbf{h}|} \sum_{j=1}^{|\mathbf{v}|} W_{ij} h_i v_j - \sum_{j=1}^{|\mathbf{v}|} v_j b_{v,j} - \sum_{i=1}^{|\mathbf{h}|} h_i b_{h,i},$$

where $\Theta = \{W, \mathbf{b}_h, \mathbf{b}_v\}$.

- **Boltzmann distribution**: joint probability of all random variables with energy function

$$P(\mathbf{v}, \mathbf{h}) = \frac{e^{-E(\mathbf{v}, \mathbf{h}; \Theta)}}{Z}, \quad \text{Partition: } Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h}; \Theta)}$$

- **Loss function**: learn “reconstructing” \mathbf{x} from a training dataset $\mathcal{D} = \{\mathbf{x}_t\}_{t=1}^{|\mathcal{D}|}$

$$\mathcal{L}(\Theta; \mathcal{D}) = -\log P(\mathbf{v}) = -\log \sum_{\mathbf{h}} P(\mathbf{v}, \mathbf{h}) = -(\log \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h}; \Theta)} - \log Z).$$

The optimisation of the loss function is **intractable** due to the partition Z . Different **sampling-based approximation algorithms** were proposed for RBM learning.

Gibbs sampling: one step ($k = 1$)

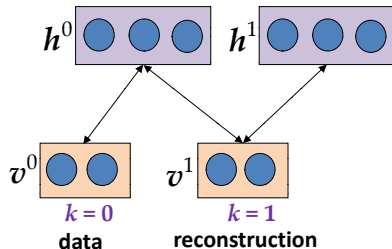
- set $\mathbf{v}^0 = \mathbf{x}$, $\mathbf{x} \in \mathcal{D}$.
- estimate $\{P(h_i^0 | \mathbf{v}^0)\}_{i=1}^{|\mathbf{h}|}$ with **encoder** and then form **realisation** of \mathbf{h}^0 by **sampling** with these probabilities

$$P(h_i^0 = 1 | \mathbf{v}^0) = \phi\left(\sum_{j=1}^{|\mathbf{v}|} W_{ij} v_j^0 + b_{h,i}\right)$$

- use **decoder** to estimate $\{P(v_j^1 | \mathbf{h}^0)\}_{j=1}^{|\mathbf{v}|}$ and then generate “**reconstruction**”, \mathbf{v}^1 , via **sampling**

$$P(v_j^1 = 1 | \mathbf{h}^0) = \phi\left(\sum_{i=1}^{|\mathbf{h}|} W_{ji} h_i^0 + b_{v,j}\right)$$

- with the “**reconstruction**”, \mathbf{v}^1 , use **encoder** to estimate probabilities $\{P(h_i^1 | \mathbf{v}^1)\}_{i=1}^{|\mathbf{h}|}$, ...



Realisation via sampling

$$h_i = 1/0, \text{ if } P(h_i^0 | \mathbf{v}^0) \geq u \sim U[0, 1]$$

$$P(h_i^1 = 1 | \mathbf{v}^1) = \phi\left(\sum_{j=1}^{|\mathbf{v}|} W_{ij} v_j^1 + b_{h,i}\right)$$

Contrastive Divergence (CD-1) Learning Algorithm

Input: a training dataset, $\mathcal{D} = \{\mathbf{x}_t\}_{t=1}^{|\mathcal{D}|}$, randomly initialise W , \mathbf{b}_h and \mathbf{b}_v , and pre-set a learning rate, η , and a mini-batch size, $|\mathcal{B}|$.

For $t = 1, 2, \dots, |\mathcal{B}|$, do **steps 1 & 2** as follows:

① Data Phase

- present an instance to the visible layer, i.e., $\mathbf{v}_t^0 = \mathbf{x}_t$.
- estimate probabilities with encoder: $P(\mathbf{h}_t^0 | \mathbf{v}_t^0) = \left(P(h_{ti}^0 = 1 | \mathbf{v}_t^0) \right)_{i=1}^{|\mathbf{h}|}$

② Model Phase

- Form a realisation of \mathbf{h}_t^0 by sampling with probabilities $P(\mathbf{h}_t^0 | \mathbf{v}_t^0)$.
- With the realisation of \mathbf{h}_t^0 , apply the decoder to estimate probabilities: $P(\mathbf{v}_t^1 | \mathbf{h}_t^0) = \left(P(v_{tj}^0 = 1 | \mathbf{h}_t^0) \right)_{j=1}^{|\mathbf{v}|}$, and then produce a “reconstruction” \mathbf{v}_t^1 via sampling with probabilities $P(\mathbf{v}_t^1 | \mathbf{h}_t^0)$.
- With the encoder and the “reconstruction”, \mathbf{v}_t^1 , estimate probabilities: $P(\mathbf{h}_t^1 | \mathbf{v}_t^1) = \left(P(h_{ti}^1 = 1 | \mathbf{v}_t^1) \right)_{i=1}^{|\mathbf{h}|}$.

Contrastive Divergence (CD-1) Learning Algorithm (cont.)

• Parameter Update

Based on Gibbs sampling results in the **data** and the **model** phases, parameters are updated as follows:

$$W \leftarrow W + \frac{\eta}{|\mathcal{B}|} \sum_{t=1}^{|\mathcal{B}|} \left(\underbrace{P(\mathbf{h}_t^0 | \mathbf{v}_t^0)}_{\text{data}} (\mathbf{v}_t^0)^T - \underbrace{P(\mathbf{h}_t^1 | \mathbf{v}_t^1)}_{\text{model}} (\mathbf{v}_t^1)^T \right),$$

$$\mathbf{b}_h \leftarrow \mathbf{b}_h + \frac{\eta}{|\mathcal{B}|} \sum_{t=1}^{|\mathcal{B}|} \left(P(\mathbf{h}_t^0 | \mathbf{v}_t^0) - P(\mathbf{h}_t^1 | \mathbf{v}_t^1) \right),$$

$$\mathbf{b}_v \leftarrow \mathbf{b}_v + \frac{\eta}{|\mathcal{B}|} \sum_{t=1}^{|\mathcal{B}|} (\mathbf{v}_t^0 - \mathbf{v}_t^1).$$

The CD-1 learning algorithm runs iteratively (for several epochs) until it converges.

Gaussian–Bernoulli RBM for Continuous Input

- Energy function (\mathbf{v} : real-valued input; \mathbf{h} : binary-valued hidden state)

$$E(\mathbf{v}, \mathbf{h}; \Theta) = \sum_{j=1}^{|\mathbf{v}|} \frac{(v_j - b_{v,j})^2}{2\sigma_j^2} - \sum_{i=1}^{|\mathbf{h}|} \sum_{j=1}^{|\mathbf{v}|} W_{ij} h_i \frac{v_j}{\sigma_j} - \sum_{i=1}^{|\mathbf{h}|} h_i b_{h,i},$$

where $\Theta = \{W, \mathbf{b}_h, \mathbf{b}_v\}$ and σ_j is standard deviation in Gaussian for visible neuron j .

- Probabilistic neurons

$$P_B(h_i = 1 | \mathbf{v}) = \phi \left(\sum_{j=1}^{|\mathbf{v}|} W_{ij} \frac{v_j}{\sigma_j} + b_{h,i} \right)$$

where $\phi(\cdot)$ is the sigmoid activation function.

$$P_G(v_j = x | \mathbf{h}) = \frac{1}{\sqrt{2\pi}\sigma_j} \exp \left(-\frac{\left(x - b_{v,j} - \sigma_j \sum_{i=1}^{|\mathbf{h}|} W_{ji} h_i \right)^2}{2\sigma_j^2} \right)$$

Gaussian–Bernoulli RBM for Continuous Input

- Contrastive Divergence (CD-1) Learning Algorithm

For $t = 1, 2, \dots, |\mathcal{B}|$, $i = 1, 2, \dots, |\mathbf{h}|$ and $j = 1, 2, \dots, |\mathbf{v}|$ do steps 1 & 2

① **Data Phase**: estimate probabilities with $P_B(h_{ti}^0 = 1 | \mathbf{v}_t^0)$ where $\mathbf{v}_t^0 = \mathbf{x}_t$

② **Model Phase**

- Form a realisation of h_{ti}^0 by sampling with $P_B(h_{ti}^0 = 1 | \mathbf{v}_t^0)$
- With h_{ti}^0 , estimate the “reconstruction”: $v_{tj}^1 = \sigma_j \sum_{i=1}^{|\mathbf{h}|} W_{ji} h_{ti}^0 + b_{v,j}$
- With v_{tj}^1 , estimate $P_B(h_{ti}^1 = 1 | \mathbf{v}_t^1)$.

Parameter Update: $i = 1, 2, \dots, |\mathbf{h}|$ and $j = 1, 2, \dots, |\mathbf{v}|$ do

$$W_{ij} \leftarrow W_{ij} + \frac{\eta}{|\mathcal{B}|} \sum_{t=1}^{|\mathcal{B}|} \left(P_B(h_{ti}^0 = 1 | \mathbf{v}_t^0) \frac{v_{tj}^0}{\sigma_j} - P_B(h_{ti}^1 = 1 | \mathbf{v}_t^1) \frac{v_{tj}^1}{\sigma_j} \right)$$

$$b_{h,i} \leftarrow b_{h,i} + \frac{\eta}{|\mathcal{B}|} \sum_{t=1}^{|\mathcal{B}|} \left(P_B(h_{ti}^0 = 1 | \mathbf{v}_t^0) - P_B(h_{ti}^1 = 1 | \mathbf{v}_t^1) \right)$$

$$b_{v,j} \leftarrow b_{v,j} + \frac{\eta}{|\mathcal{B}|} \sum_{t=1}^{|\mathcal{B}|} \left(v_{tj}^0 - v_{tj}^1 \right)$$

Tip: Each input feature is often standardised to $N(0, 1)$ to avoid setting σ_j .

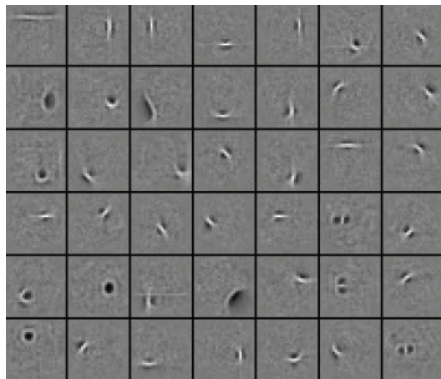
Illustrative Example

- **Handwritten characters:** binary-valued RBM of 200 hidden neurons
 28×28 image size (left), 28×28 weights associated with each of 42 randomly chosen hidden neurons (right) that learns **receptive fields**

Training samples



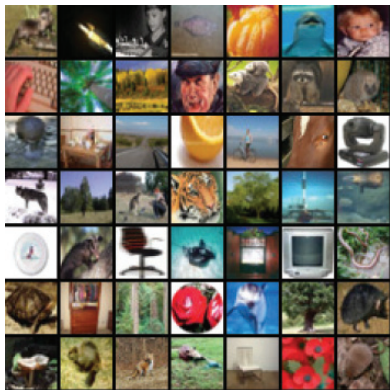
Learned receptive fields



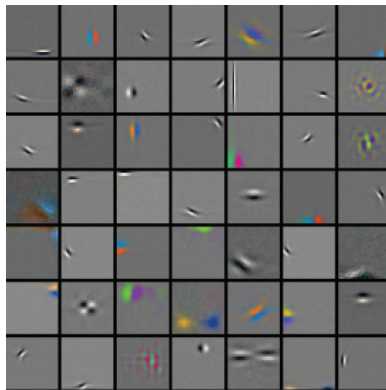
Illustrative Example

- **CIFAR-100 image dataset:** real-valued RBM of 300 hidden neurons
 32×32 image size (left), 32×32 weights associated with each of 49 randomly chosen hidden neurons (right) that learns **receptive fields**

Training samples



Learned receptive fields



If you want to deepen your understanding and learn something beyond this lecture, you can self-study the optional references below.

[Goodfellow et al., 2016] Goodfellow I., Bengio Y., and Courville A. (2016): *Deep Learning*, MIT Press. (Chapter 14, Sections 18.1-18.2 & 20.2)

[Hinton, 2010] Hinton G. (2010): A practical guide to training restricted Boltzmann machines. *Technical Report: UTML TR 2010(003)*, Department of Computer Science, University of Toronto. Online: <https://www.cs.toronto.edu/~hinton/absps/guideTR.pdf>

[Chen, 2015] Chen K. (2015): Deep and modular neural networks. In *Springer Handbook of Computational Intelligence*, Chapter 28, pp. 473-492. (Sections 28.1-28.2)