Neural networks Sparse coding - definition



UNSUPERVISED LEARNING

Topics: unsupervised learning

- Unsupervised learning: only use the inputs $\mathbf{x}^{(t)}$ for learning
 - automatically extract meaningful features for your data
 - Ieverage the availability of unlabeled data
 - add a data-dependent regularizer to trainings

- We will see 3 neural networks for unsupervised learning
 - restricted Boltzmann machines
 - autoencoders
 - sparse coding model



- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - lacksimit is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - lacksimwe can reconstruct the original input $\mathbf{x}^{(t)}$ as well as possible
- More formally:

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^{T} \min_{\mathbf{h}^{(t)}} \frac{1}{2} ||\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}||_{2}^{2} + \lambda ||\mathbf{h}^{(t)}||_{1}$$

- ${\scriptstyle \blacktriangleright}$ we also constrain the columns of ${\bf D}\,$ to be of norm [
 - otherwise, ${f D}$ could grow big while ${f h}^{(t)}$ becomes small to satisfy the prior
- sometimes the columns are constrained to be no greater than I

3

1

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - lacksimit is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - lacksimwe can reconstruct the original input $\mathbf{x}^{(t)}$ as well as possible
- More formally: reconstruction error

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^{T} \min_{\mathbf{h}^{(t)}} \frac{1}{2} ||\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}||_{2}^{2} + \lambda ||\mathbf{h}^{(t)}|$$

- ${\scriptstyle \blacktriangleright}$ we also constrain the columns of ${\bf D}\,$ to be of norm [
 - otherwise, ${f D}$ could grow big while ${f h}^{(t)}$ becomes small to satisfy the prior
- sometimes the columns are constrained to be no greater than I

3

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - lacksimit is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - lacksimwe can reconstruct the original input $\mathbf{x}^{(t)}$ as well as possible
- More formally: reconstruction error

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^{T} \min_{\mathbf{h}^{(t)}} \frac{1}{2} ||\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}||_{2}^{2} + \lambda ||\mathbf{h}^{(t)}||_{1}$$
reconstruction $\widehat{\mathbf{x}}^{(t)}$

- ${\scriptstyle \blacktriangleright}$ we also constrain the columns of ${\bf D}\,$ to be of norm [
 - otherwise, ${f D}$ could grow big while ${f h}^{(t)}$ becomes small to satisfy the prior
- sometimes the columns are constrained to be no greater than I

3

Topics: sparse coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - it is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - lacksimwe can reconstruct the original input $\mathbf{x}^{(t)}$ as well as possible
- More formally: reconstruction error sparsity penalty $\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^{T} \min_{\mathbf{h}^{(t)}} \frac{1}{2} ||\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}||_{2}^{2} + \lambda ||\mathbf{h}^{(t)}||_{1}$

• we also constrain the columns of \mathbf{D} to be of norm |

- otherwise, ${f D}$ could grow big while ${f h}^{(t)}$ becomes small to satisfy the prior

reconstruction $\widehat{\mathbf{x}}^{(t)}$

sometimes the columns are constrained to be no greater than I

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - it is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - we can reconstruct the original input $\mathbf{x}^{(t)}$ as well as possible



- we also constrain the columns of \mathbf{D} to be of norm |
 - otherwise, ${f D}$ could grow big while ${f h}^{(t)}$ becomes small to satisfy the prior
- sometimes the columns are constrained to be no greater than I

Topics: sparse coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - it is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - we can reconstruct the original input $\mathbf{x}^{(t)}$ as well as possible



- D is equivalent to the autoencoder output weight matrix
- however, $\mathbf{h}(\mathbf{x}^{(t)})$ is now a complicated function of $\mathbf{x}^{(t)}$

- encoder is the minimization $\mathbf{h}(\mathbf{x}^{(t)}) = \operatorname*{arg\,min}_{\mathbf{h}^{(t)}} \frac{1}{2} ||\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}||_2^2 + \lambda ||\mathbf{h}^{(t)}||_1$

Topics: dictionary • Can also write $\widehat{\mathbf{x}}^{(t)} = \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)}) = \sum \mathbf{D}_{\cdot,k} h(\mathbf{x}^{(t)})_k$ $k \text{ s.t.} \\ h(\mathbf{x}^{(t)})_k \neq 0$

- we also refer to **D** as the dictionary
 - in certain applications, we know what dictionary matrix to use
 - often however, we have to learn it



Topics: dictionary • Can also write $\widehat{\mathbf{x}}^{(t)} = \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)}) = \sum \mathbf{D}_{\cdot,k} h(\mathbf{x}^{(t)})_k$ k s.t. $h(\mathbf{x}^{(t)})_k \neq 0$



- we also refer to **D** as the dictionary
 - in certain applications, we know what dictionary matrix to use
 - often however, we have to learn it



k s.t. $h(\mathbf{x}^{(t)})_k \neq 0$

Topics: dictionary

• Can also write $\widehat{\mathbf{x}}^{(t)} = \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)}) = \sum \mathbf{D}_{\cdot,k} h(\mathbf{x}^{(t)})_k$



- we also refer to **D** as the dictionary
 - in certain applications, we know what dictionary matrix to use
 - often however, we have to learn it

