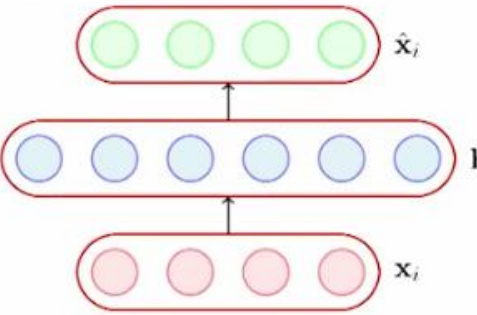


The diagram shows a neural network layer with three layers of nodes. The bottom layer (input) has 4 red nodes labeled \mathbf{x}_i . The middle layer (hidden) has 6 blue nodes labeled \mathbf{h} . The top layer (output) has 4 green nodes labeled $\hat{\mathbf{x}}_i$. Arrows indicate the flow from input to hidden and from hidden to output.

- A hidden neuron with sigmoid activation will have values between 0 and 1
- We say that the neuron is activated when its output is close to 1 and not activated when its output is close to 0.
- A sparse autoencoder tries to ensure the neuron is inactive most of the times.

0:58 / 9:11



The diagram shows a neural network layer with three layers of nodes. The bottom layer (input) has 4 red nodes labeled \mathbf{x}_i . The middle layer (hidden) has 6 blue nodes labeled \mathbf{h} . The top layer (output) has 4 green nodes labeled $\hat{\mathbf{x}}_i$. Arrows indicate the flow from input to hidden and from hidden to output.

- If the neuron l is sparse (i.e. mostly inactive) then $\hat{\rho}_l \rightarrow 0$
- A sparse autoencoder uses a sparsity parameter ρ (typically very close to 0, say, 0.005) and tries to enforce the constraint $\hat{\rho}_l = \rho$
- One way of ensuring this is to add the following term to the objective function

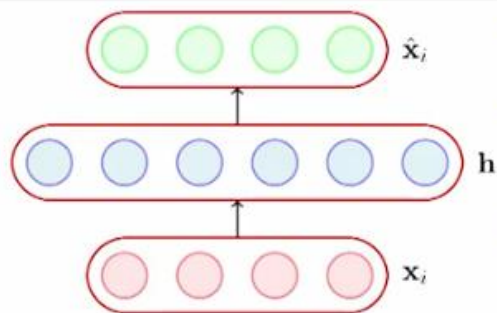
The average value of the activation of a neuron l is given by

$$\hat{\rho}_l = \frac{1}{m} \sum_{i=1}^m h(\mathbf{x}_i)_l$$

$$\Omega(\theta) = \sum_{l=1}^k \rho \log \frac{\rho}{\hat{\rho}_l} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_l}$$

$\mathcal{L}(\theta) = \mathcal{L}'(\theta) + \Omega(\theta)$

4:44 / 9:11



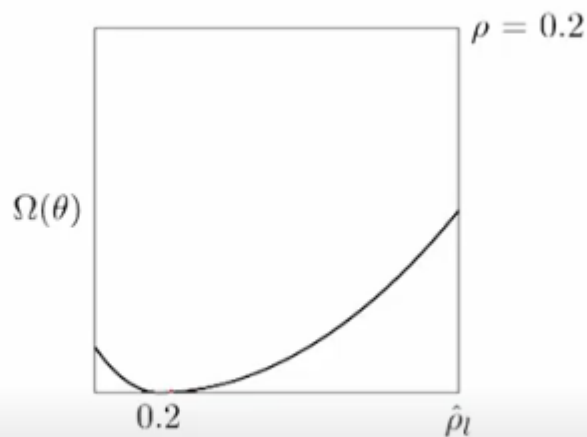
- If the neuron l is sparse (i.e. mostly inactive) then $\hat{\rho}_l \rightarrow 0$
- A sparse autoencoder uses a sparsity parameter ρ (typically very close to 0, say, 0.005) and tries to enforce the constraint $\hat{\rho}_l = \rho$
- One way of ensuring this is to add the following term to the objective function

The average value of the activation of a neuron l is given by

$$\hat{\rho}_l = \frac{1}{m} \sum_{i=1}^m h(\mathbf{x}_i)_l$$

$$\Omega(\theta) = \sum_{l=1}^k \rho \log \frac{\rho}{\hat{\rho}_l} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_l}$$

- When will this term reach its minimum value and what is the minimum value? Let us plot it and check.



- The function will reach its minimum value(s) when $\hat{\rho}_l = \rho$.

$$\Omega(\theta) = \sum_{i=1}^k \rho \log \frac{\rho}{\hat{\rho}_i} + (1-\rho) \log \frac{1-\rho}{1-\hat{\rho}_i}$$

Can be re-written as

$$-\Omega(\theta) = \sum_{i=1}^k \rho \log \rho - \rho \log \hat{\rho}_i + (1-\rho) \log(1-\rho) - (1-\rho) \log(1-\hat{\rho}_i)$$

By Chain rule

$$\frac{\partial \Omega(\theta)}{\partial W} = \frac{\partial \Omega(\theta)}{\partial \hat{\rho}} * \frac{\partial \hat{\rho}}{\partial W}$$

$$\frac{\partial \Omega(\theta)}{\partial \hat{\rho}} = -\frac{\rho}{\hat{\rho}} + \frac{(1-\rho)}{1-\hat{\rho}}$$

For each neuron $l \in 1 \dots k$ in hidden layer, we have

$$\frac{\partial \hat{\rho}_l}{\partial W} = \mathbf{x}_i (g'(W^T \mathbf{x}_i + \mathbf{b}))^T$$

Finally,

$$\frac{\partial \hat{\mathcal{L}}(\theta)}{\partial W} = \frac{\partial \mathcal{L}(\theta)}{\partial W} + \frac{\partial \Omega(\theta)}{\partial W}$$

and we know how to calculate both terms on R.H.S)

• Now,

$$\hat{\mathcal{L}}(\theta) = \mathcal{L}(\theta) + \Omega(\theta)$$

• $\mathcal{L}(\theta)$ is the squared error loss or cross entropy loss and $\Omega(\theta)$ is the sparsity constraint.

• We already know how to calculate $\frac{\partial \mathcal{L}(\theta)}{\partial W}$

• Let us see how to calculate $\frac{\partial \Omega(\theta)}{\partial W}$.

Derivation

$$\frac{\partial \hat{\rho}}{\partial W} = \begin{bmatrix} \frac{\partial \hat{\rho}_1}{\partial W} & \frac{\partial \hat{\rho}_2}{\partial W} & \dots & \frac{\partial \hat{\rho}_k}{\partial W} \end{bmatrix}$$

For each element in the above equation we can calculate $\frac{\partial \hat{\rho}_l}{\partial W}$ (which is the partial derivative of a scalar w.r.t. a matrix = matrix). For a single element of a matrix W_{jl} :-

$$\begin{aligned} \frac{\partial \hat{\rho}_l}{\partial W_{jl}} &= \frac{\partial \left[\frac{1}{m} \sum_{i=1}^m g(W_{:,l}^T \mathbf{x}_i + b_l) \right]}{\partial W_{jl}} \\ &= \frac{1}{m} \sum_{i=1}^m \frac{\partial [g(W_{:,l}^T \mathbf{x}_i + b_l)]}{\partial W_{jl}} \\ &= \frac{1}{m} \sum_{i=1}^m g'(W_{:,l}^T \mathbf{x}_i + b_l) x_{ij} \end{aligned}$$

So in matrix notation we can write it as :

$$\frac{\partial \hat{\rho}_l}{\partial W} = \mathbf{x}_i (g'(W^T \mathbf{x}_i + \mathbf{b}))^T$$