- A denoising encoder simply corrupts the input data using a probabilistic process ($P(\tilde{x}_{ij}|x_{ij})$) before feeding it to the network
- A simple $P(\tilde{x}_{ij}|x_{ij})$ used in practice is the following
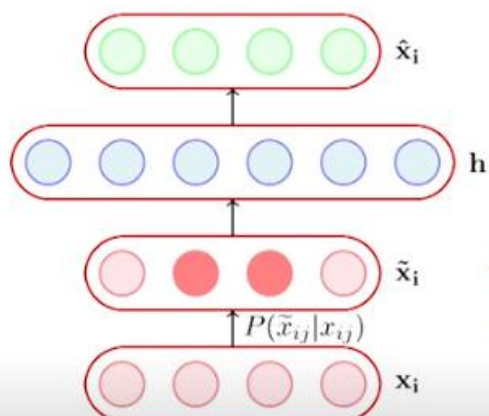
$$P(\tilde{x}_{ij} = 0|x_{ij}) = q$$
$$P(\tilde{x}_{ij} = x_{ij}|x_{ij}) = 1 - q$$

- In other words, with probability $q$ the input is flipped to 0 and with probability $(1-q)$ it is retained as it is

- How does this help ?
- This helps because the objective is still to reconstruct the original (uncorrupted) $x_i$

$$\arg\min_{\theta} \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{n} (\hat{x}_{ij} - x_{ij})^2$$

- It no longer makes sense for the model to copy the corrupted $\tilde{x}_i$ into $h(\tilde{x}_i)$ and then into $\hat{x}_i$ (the objective function will not be minimized by doing so)
- Instead the model will now have to capture the characteristics of the data correctly.

For example, it will have to learn to reconstruct a corrupted $x_{ij}$ correctly by relying on its interactions with other elements of $x_i$

We will now see a practical application in which AEs are used and then compare Denoising Autoencoders with regular autoencoders

## Task: Hand-written digit recognition



Figure: MNIST Data



$$|\mathbf{x}_i| = 784 = 28 \times 28$$

**3**

28*28

Figure: Basic approach(we use raw data as input features)

7:24 / 26:17

---

## Task: Hand-written digit recognition



Figure: MNIST Data



$$\hat{\mathbf{x}}_i \in \mathbb{R}^{784}$$

$$\mathbf{h} \in \mathbb{R}^d$$

$$|\mathbf{x}_i| = 784 = 28 \times 28$$

**3**

Figure: AE approach (first learn important characteristics of data)

7:25 / 26:17

## Task: Hand-written digit recognition



Figure: MNIST Data



$$\mathbf{h} \in \mathbb{R}^d$$
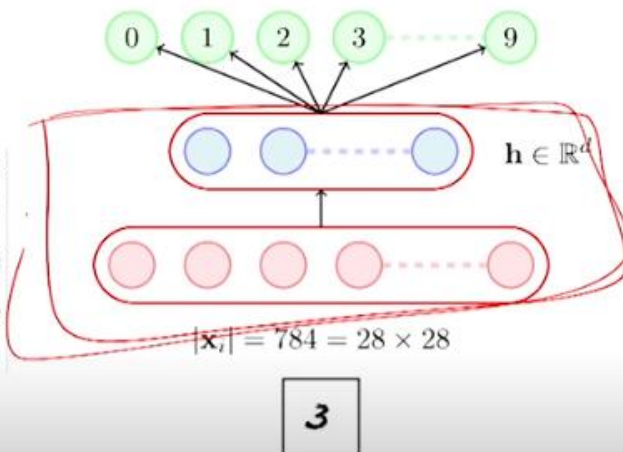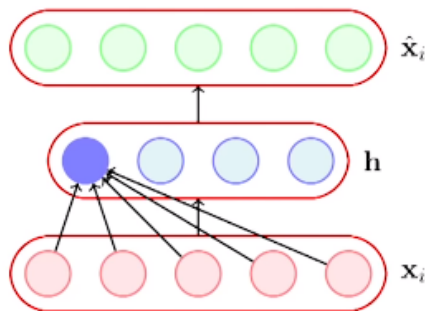
$$|\mathbf{x}_i| = 784 = 28 \times 28$$

Figure: AE approach (and then train a classifier on top of this hidden representation)

---

We will now see a way of visualizing AEs and use this visualization to compare different AEs

- We can think of each neuron as a filter which will fire (or get maximally) activated for a certain input configuration $\mathbf{x}_i$
- For example,

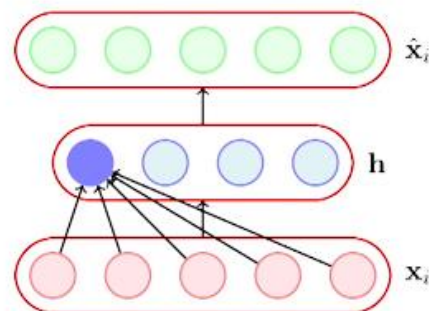$$\mathbf{h}_1 = \sigma(W_1^T \mathbf{x}_i) \ [\text{ignoring bias } b]$$

Where $W_1$ is the trained vector of weights connecting the input to the first hidden neuron

- What values of $\mathbf{x}_i$ will cause $\mathbf{h}_1$ to be maximum (or maximally activated)
- Suppose we assume that our inputs are normalized so that $\|\mathbf{x}_i\| = 1$

$$\max_{\mathbf{x}_i} \ \{w_1^T \mathbf{x}_i\}$$
$$s.t. \ \|\mathbf{x}_i\|^2 = \mathbf{x}_i^T \mathbf{x}_i = 1$$

---



- Thus the inputs

$$\dot{\mathbf{x}}_i = \frac{W_1}{\sqrt{W_1^T W_1}}, \frac{W_2}{\sqrt{W_2^T W_2}}, \dots \frac{W_n}{\sqrt{W_n^T W_n}}$$

will respectively cause hidden neurons 1 to $n$ to maximally fire

- Let us plot these images ($\mathbf{x}_i$'s) which maximally activate the first $k$ neurons of the hidden representations learned by a vanilla autoencoder and different denoising autoencoders
- These $\mathbf{x}_i$'s are computed by the above formula using the weights $(W_1, W_2 \dots W_k)$ learned by the respective autoencoders

$$\max_{\mathbf{x}_i} \ \{w_1^T \mathbf{x}_i\}$$
$$s.t. \ \|\mathbf{x}_i\|^2 = \mathbf{x}_i^T \mathbf{x}_i = 1$$
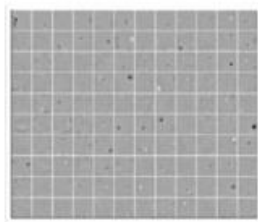$$\text{Solution:} \ \ \mathbf{x}_i = \frac{w_1}{\sqrt{w_1^T w_1}}$$

Figure: Vanilla AE (No noise)
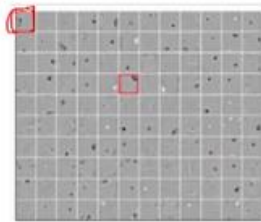
Figure: 25% Denoising AE (q=0.25)

Figure: 50% Denoising AE (q=0.5)

- The vanilla AE does not learn many meaningful patterns
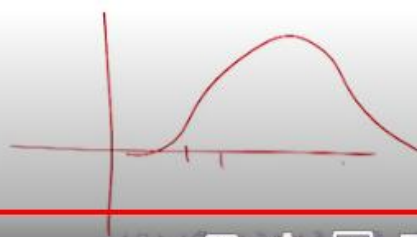
- We saw one form of $P(\widetilde{x}_{ij}|x_{ij})$ which flips a fraction $q$ of the inputs to zero
- Another way of corrupting the inputs is to add a Gaussian noise to the input

$$\widetilde{x}_{ij} = x_{ij} + \mathcal{N}(0,1)$$

- We will now use such a denoising AE on a different dataset and see their performance