

Auto-encoder Based Data Clustering

Chunfeng Song¹, Feng Liu², Yongzhen Huang¹,
Liang Wang¹ and Tieniu Tan¹

1. NLPR, Chinese Academy of Sciences
2. Southeast University



● Data clustering

Previous clustering method

- ✓ K-means
- ✓ Spectral clustering
- ✓ N-cut

→ Linear mapping →

Perform bad for
bad distributed
data.

● Data clustering

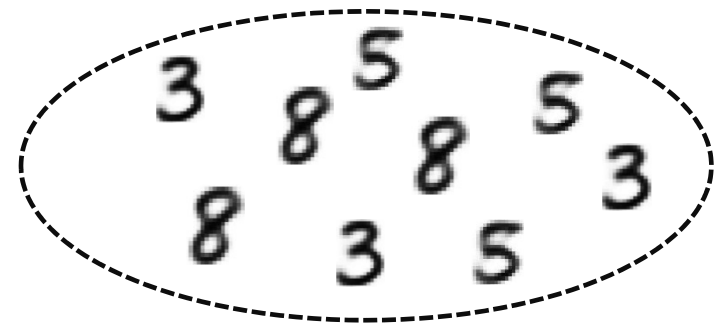
Previous clustering method

- ✓ K-means
- ✓ Spectral clustering
- ✓ N-cut

→ Linear mapping →

Perform bad for
bad distributed
data.

Can **not** deal with this similar images



Original Data Space

● Data clustering

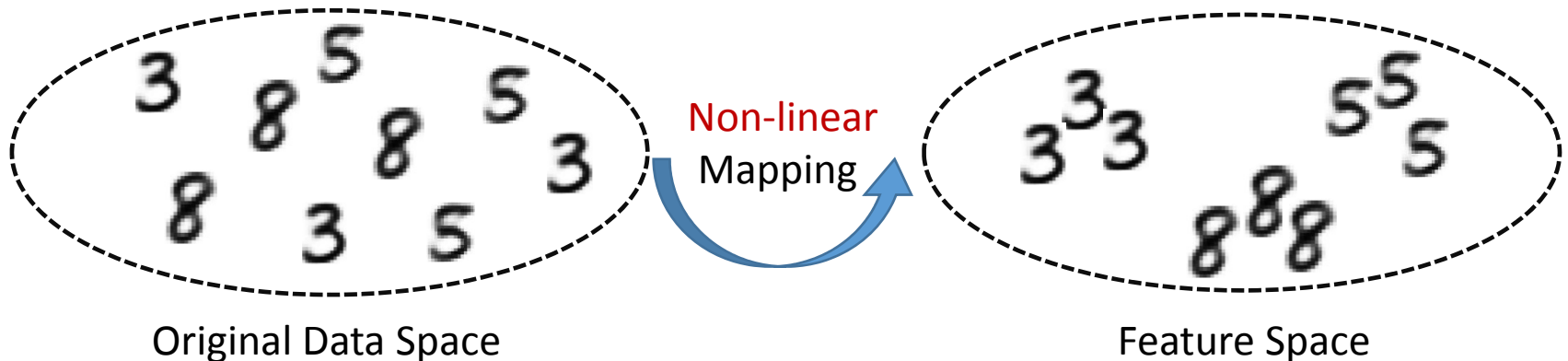
Previous clustering method

- ✓ K-means
- ✓ Spectral clustering
- ✓ N-cut

Linear mapping

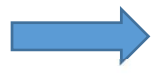
Perform bad for bad distributed data.

✓ Auto-encoder based clustering can provide **non-linear** mapping.

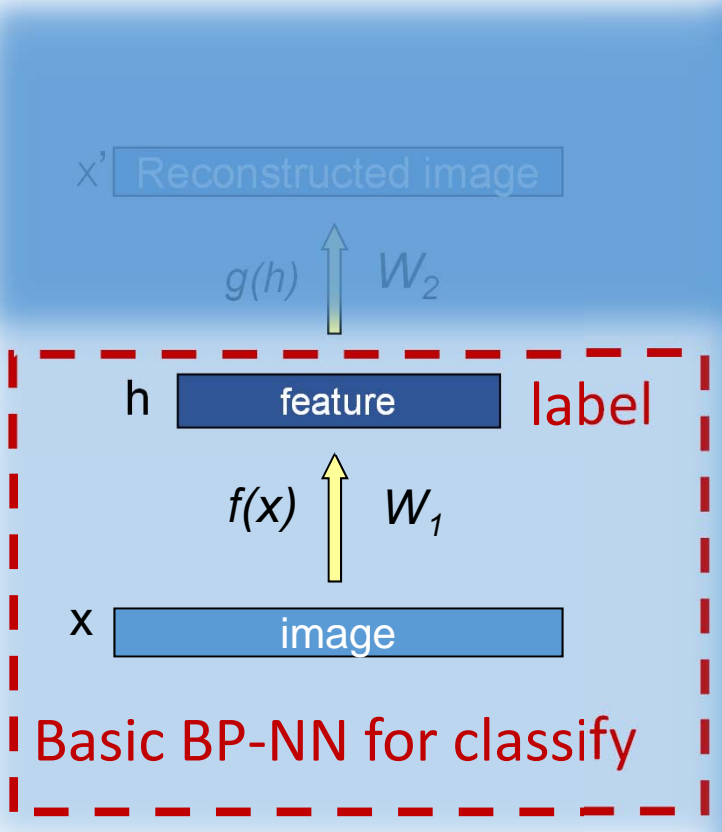


● Auto-encoder

Basic single-layer auto-encoder



Is a kind of **BP-NN**



supervised

Encoder function

Sigmoid-type

$$h_i = f(x_i) = \frac{1}{1 + \exp(-(W_1 x_i + b_1))}$$

Decoder function

$$x'_i = g(h_i) = \frac{1}{1 + \exp(-(W_2 h_i + b_2))}$$

Obj. function

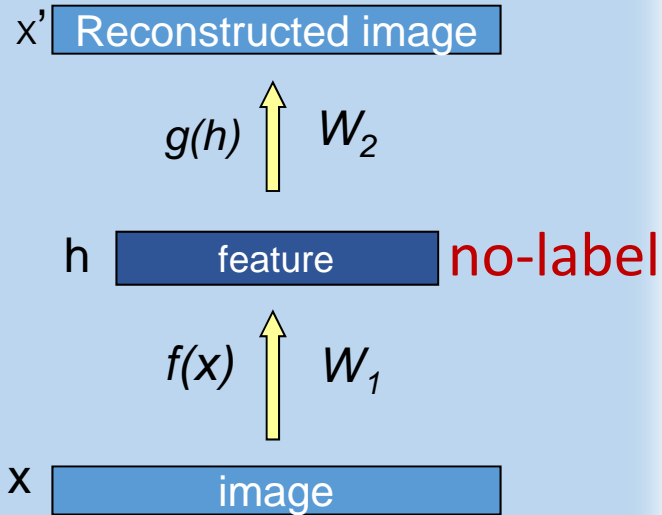
$$\min \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2$$

● Auto-encoder

Basic single-layer auto-encoder



Is a kind of BP-NN



Encoder function

Sigmoid-type

$$h_i = f(x_i) = \frac{1}{1 + \exp(-(W_1 x_i + b_1))}$$

Decoder function

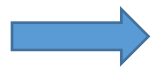
$$x'_i = g(h_i) = \frac{1}{1 + \exp(-(W_2 h_i + b_2))}$$

Obj. function

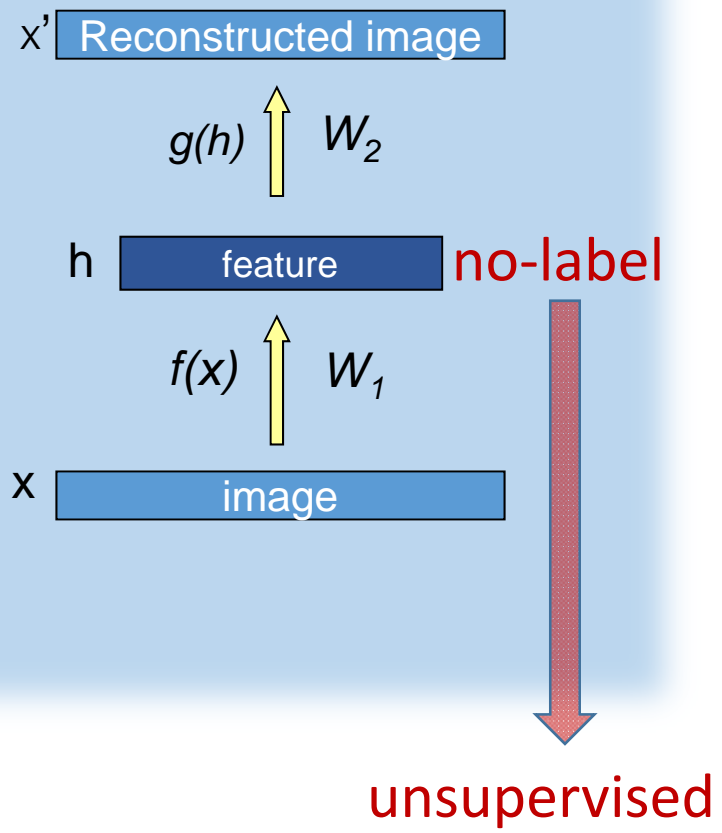
$$\min \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2$$

● Auto-encoder

Basic single-layer auto-encoder



Is a kind of BP-NN



Encoder function

Sigmoid-type

$$h_i = f(x_i) = \frac{1}{1 + \exp(-(W_1 x_i + b_1))}$$

Decoder function

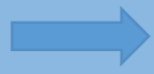
$$x'_i = g(h_i) = \frac{1}{1 + \exp(-(W_2 h_i + b_2))}$$

Obj. function

$$\min \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2$$

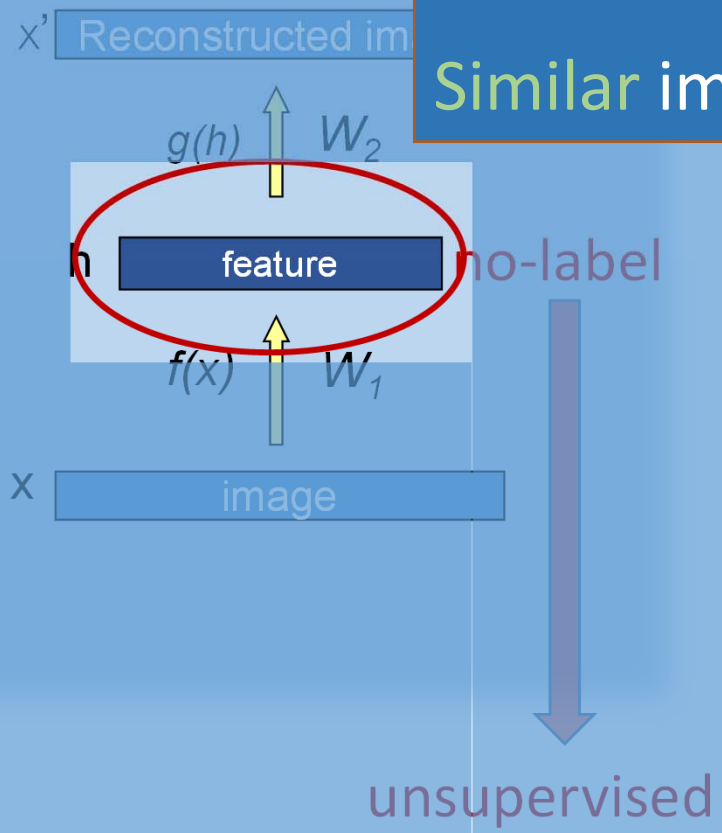
● Auto-encoder

Basic single-layer auto-encoder



Is a kind of BP-NN

Different images may have different feature,
Similar images have similar feature!



Decoder function

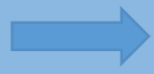
$$x'_i = g(h_i) = \frac{1}{1 + \exp(-(W_2 h_i + b_2))}$$

Obj. function

$$\min \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2$$

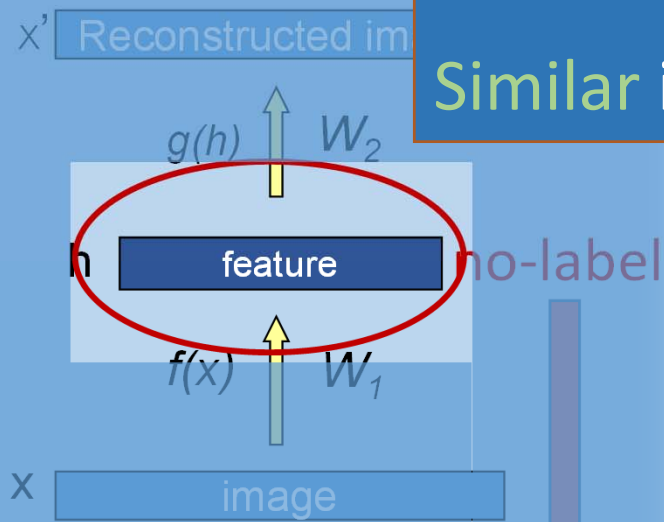
● Auto-encoder

Basic single-layer auto-encoder

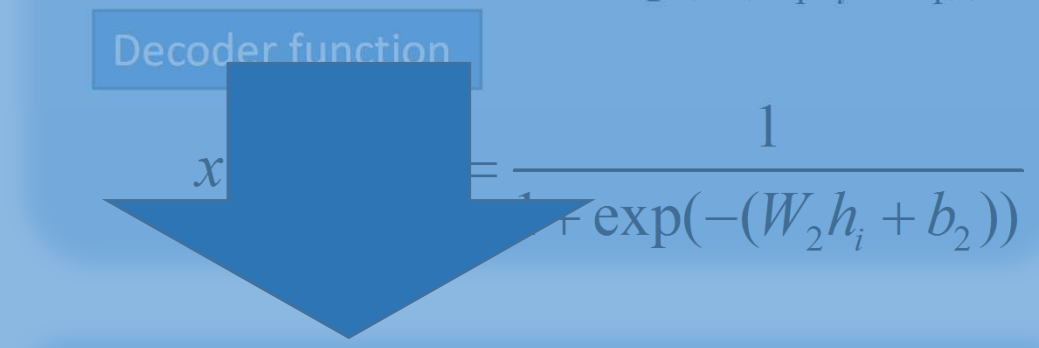


Is a kind of BP-NN

Different images may have different feature,
Similar images have similar feature!



↓
unsupervised

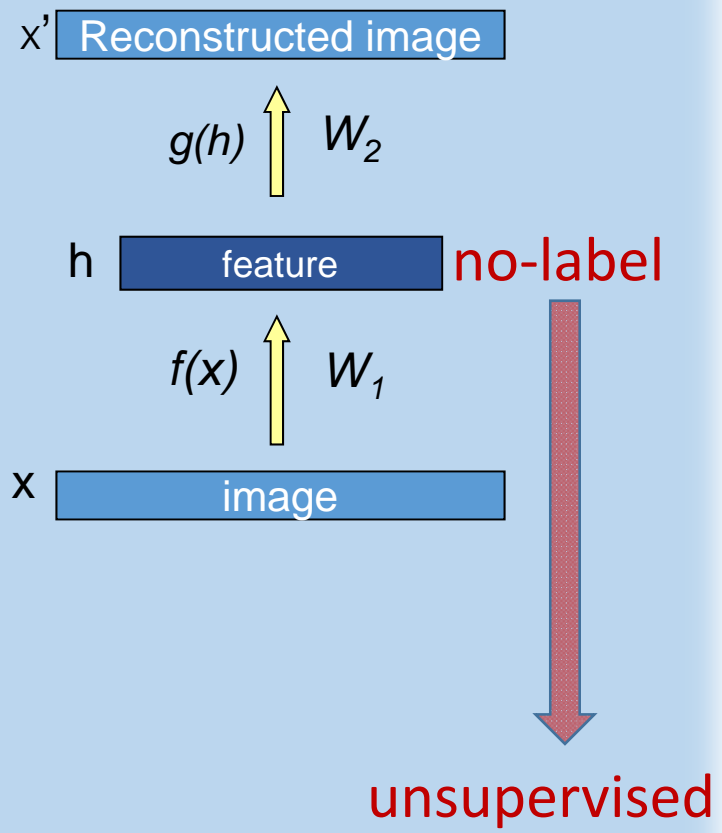


Be fit for clustering

$$\min \frac{1}{N} \sum_{i=1}^N \|x_i - x_i'\|^2$$

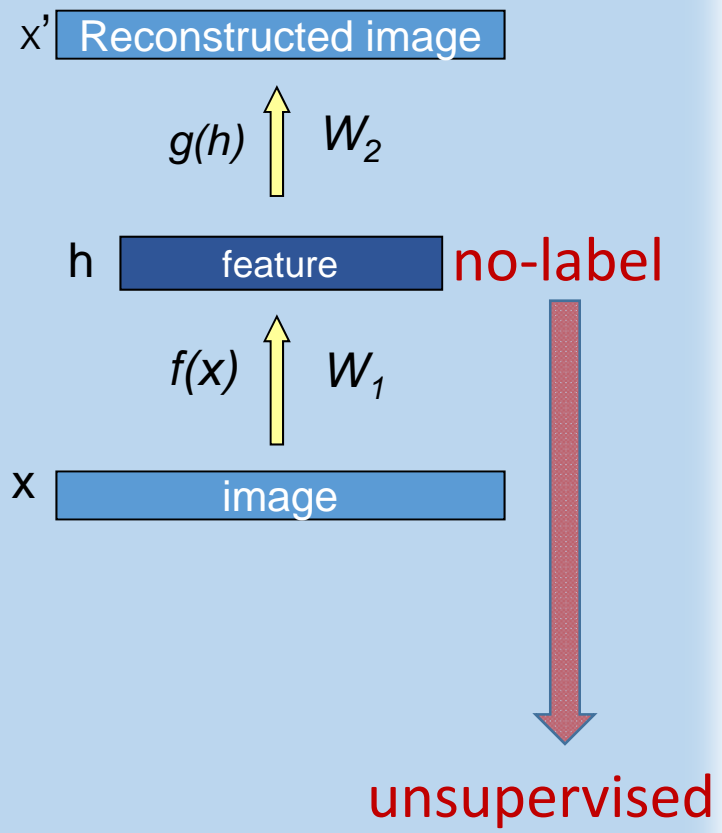
● Auto-encoder

Basic **single**-layer auto-encoder

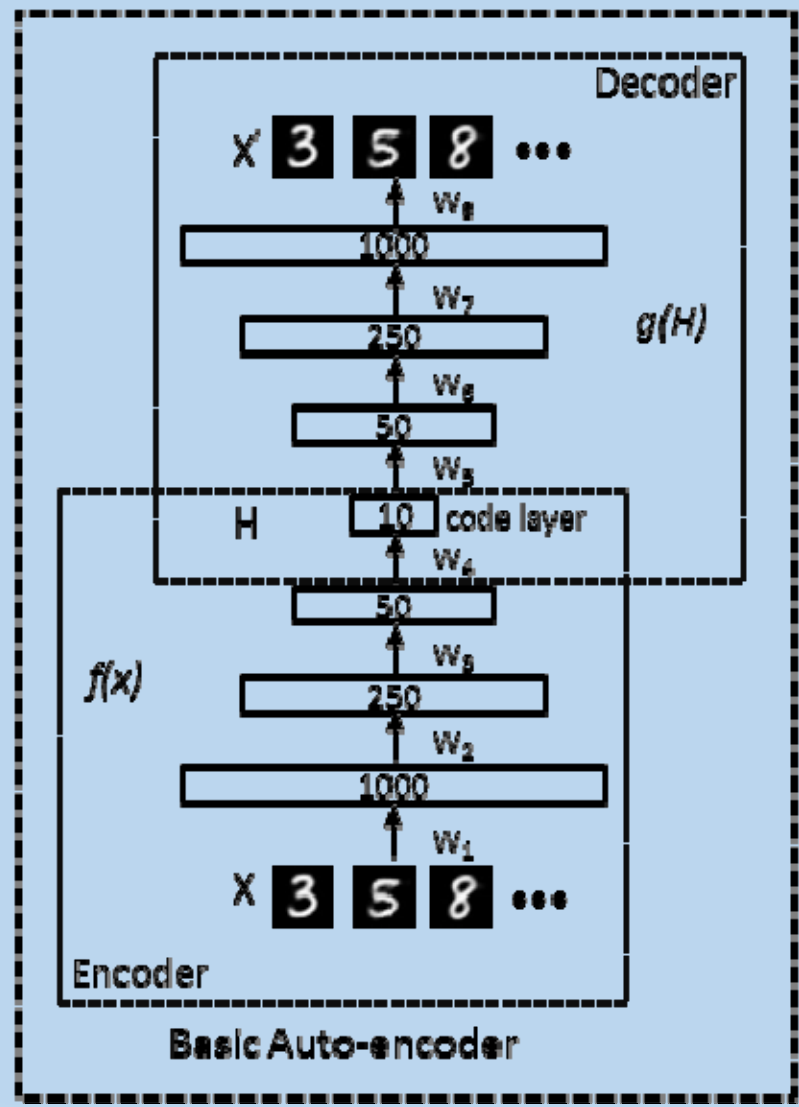


● Auto-encoder

Basic **single**-layer auto-encoder

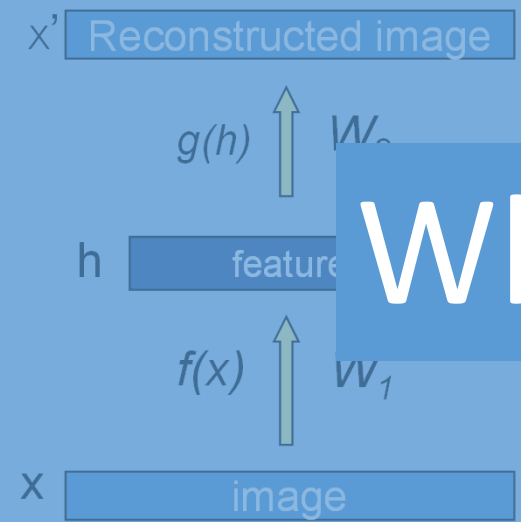


multi-layers auto-encoders



● Auto-encoder

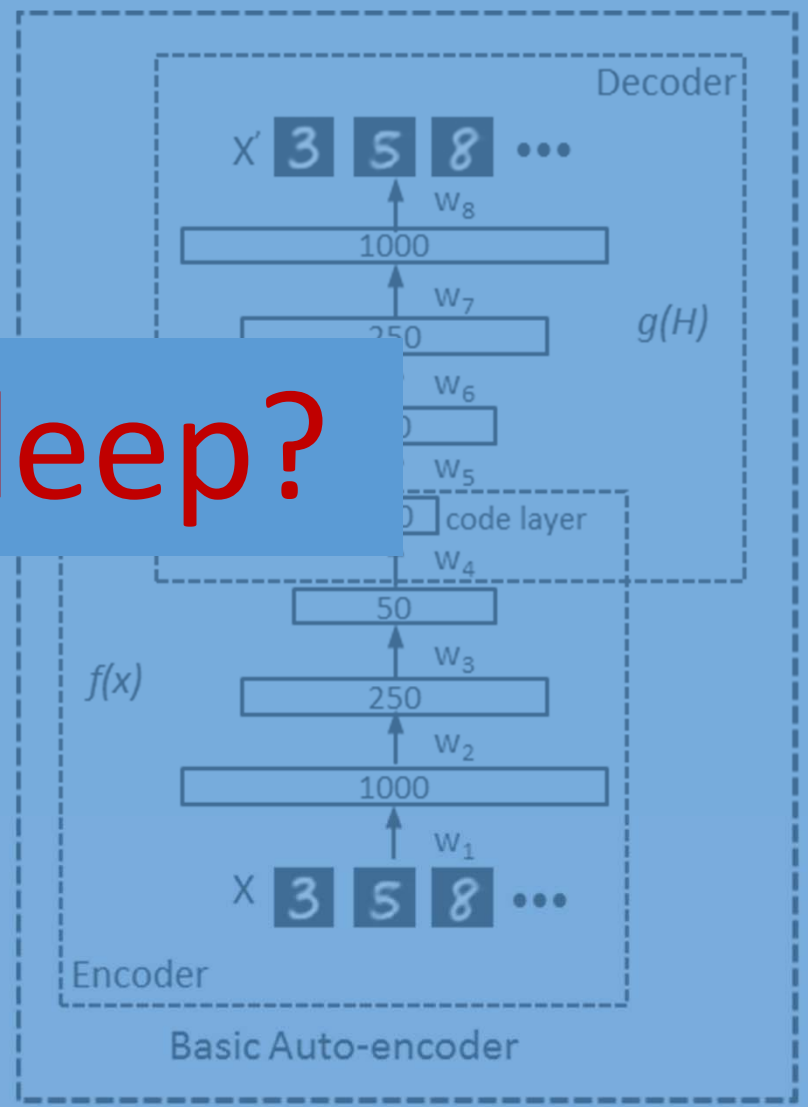
Basic single-layer auto-encoder



Why so deep?

unsupervised

multi-layers auto-encoders



Why so deep?

- ✓ Deep makes **more accurate** results.
- ✓ Deep networks can **learn better**.
- ✓ Deep networks can provide **better non-linear** mapping.

Why so deep?

- ✓ Deep makes more accurate results
- ✓ Deep networks can learn better
- ✓ Deep networks can provide better non-linear mapping

Is auto-encoder perfect for clustering?

Why so deep?

- ✓ Deep makes more accurate results
- ✓ Deep networks can learn better
- ✓ Deep networks can provide better non-linear mapping

Is auto-encoder perfect for clustering?

Not enough.

● Clustering Based on Auto-encoder

Basic auto-encoder

Obj.
fun

$$\min_{W,b} \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2$$

● Clustering Based on Auto-encoder

Obj. fun

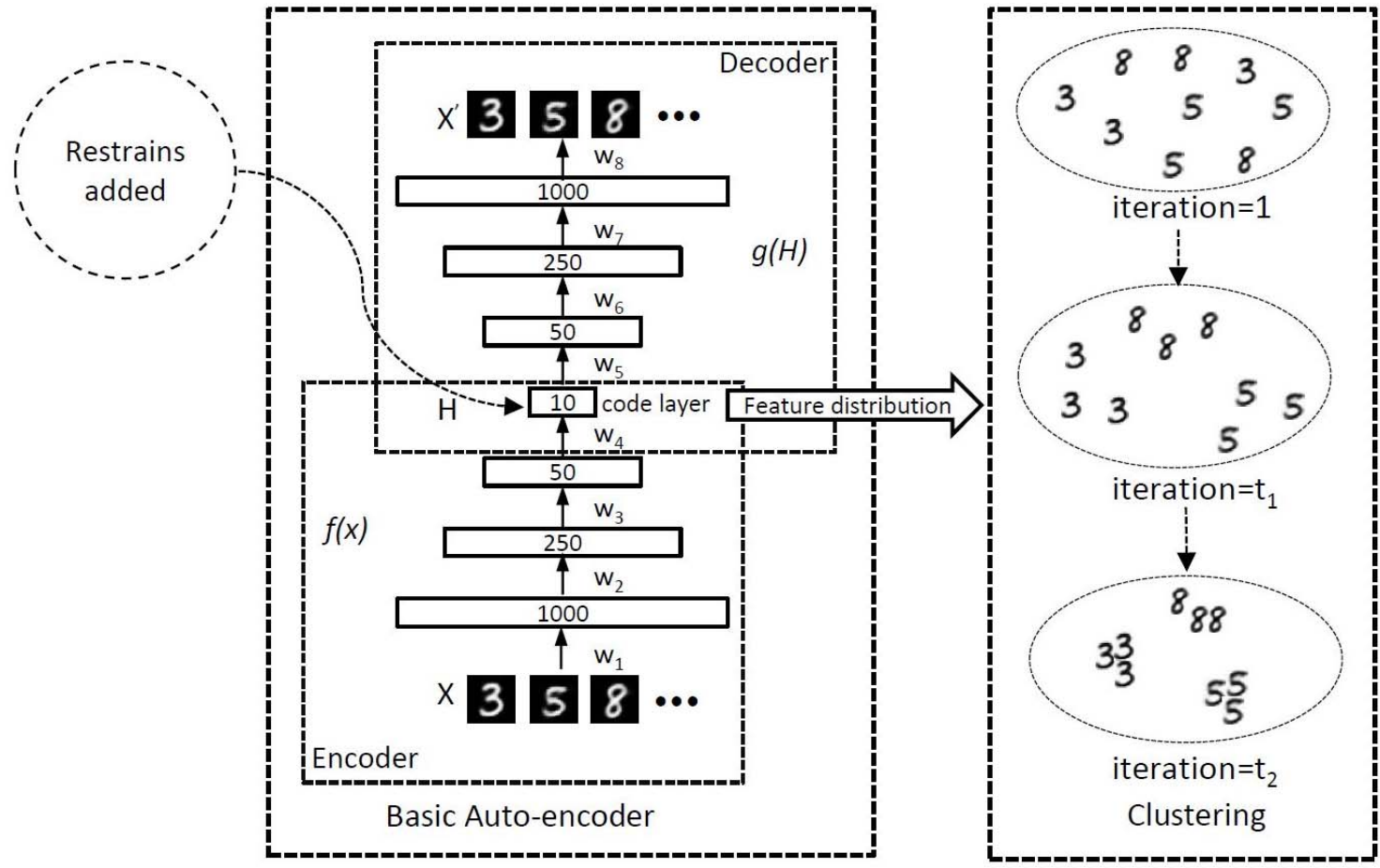
Basic auto-encoder

$$\min_{W,b} \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2$$

$$\lambda \cdot \sum_{i=1}^N \|f^t(x_i) - c_i^*\|^2$$

proposed

✓ **Restrains** added to achieve compact distribution in **feature** layer



● Algorithm

$$\min_{W,b} \frac{1}{N} \sum_{i=1}^N \|x_i - x'_i\|^2 - \lambda \cdot \sum_{i=1}^N \|f^t(x_i) - c_i^*\|^2 \quad (4)$$

$$c_i^* = \arg \min_{c_j^{t-1}} \|f^t(x_i) - c_j^{t-1}\|^2, \quad (5)$$

$$c_j^t = \frac{\sum_{x_i \in C_j^{t-1}} f^t(x_i)}{|C_j^{t-1}|}, \quad (6)$$

Algorithm 1 Auto-encoder based data clustering algorithm

- 1: **Input:** Dataset X , the number of clusters K , hyper-parameter λ , the maximum number of iterations T .
 - 2: **Initialize** sample assignment C^0 randomly.
 - 3: **Set** t to 1.
 - 4: **repeat**
 - 5: Update the mapping network by minimizing Eqn. (4) with stochastic gradient descent for one epoch.
 - 6: Update cluster center c^t via Eqn. (6).
 - 7: Partition X into K clusters and update the sample assignment C^t via Eqn. (5).
 - 8: $t = t + 1$.
 - 9: **until** $t > T$
 - 10: **Output:** Final sample assignment C .
-

● Iteration

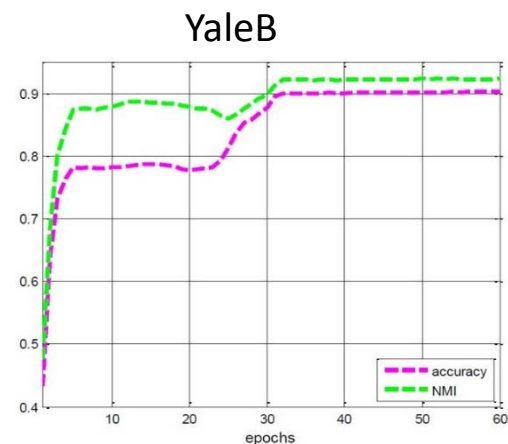
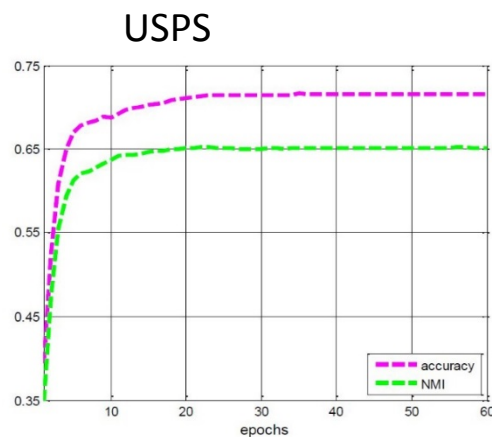
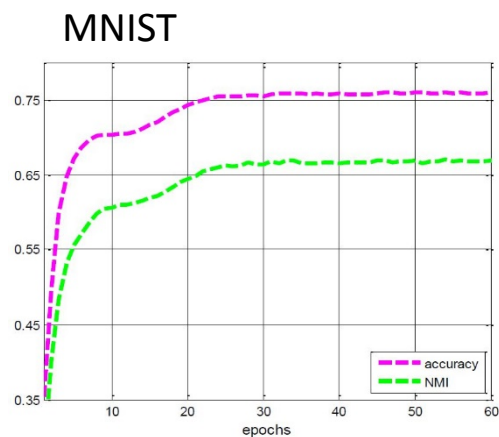
ACC: the cluster accuracy. **Distance**: the sum of distances between 10 clusters in **feature layer**.

| epoch | ACC | Distance | Visualization of 10 cluster centers(the reconstruction of feature) | | | | | | | | | |
|-------|------|----------|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.30 | 0.003 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0.46 | 0.296 | 9 | 9 | 0 | 6 | 3 | 7 | 9 | 0 | 1 | 0 |
| 3 | 0.53 | 0.432 | 9 | 9 | 0 | 6 | 3 | 7 | 9 | 0 | 1 | 0 |
| 4 | 0.56 | 0.493 | 9 | 9 | 0 | 6 | 3 | 7 | 9 | 0 | 1 | 0 |
| 5 | 0.59 | 0.515 | 9 | 9 | 0 | 6 | 3 | 7 | 2 | 0 | 1 | 8 |
| 6 | 0.61 | 0.526 | 9 | 9 | 5 | 6 | 3 | 7 | 2 | 0 | 1 | 8 |
| 7 | 0.63 | 0.534 | 9 | 9 | 5 | 6 | 3 | 7 | 2 | 0 | 1 | 8 |
| 8 | 0.65 | 0.537 | 9 | 9 | 5 | 6 | 3 | 7 | 2 | 0 | 1 | 8 |
| 9 | 0.67 | 0.538 | 9 | 9 | 5 | 6 | 3 | 7 | 2 | 0 | 1 | 8 |
| 10 | 0.68 | 0.539 | 9 | 9 | 5 | 6 | 3 | 7 | 2 | 0 | 1 | 8 |

Test on MNIST datasets (including 60000 images with 28*28 resolution)

● Experiments

✓ Influence of the iteration number on three databases

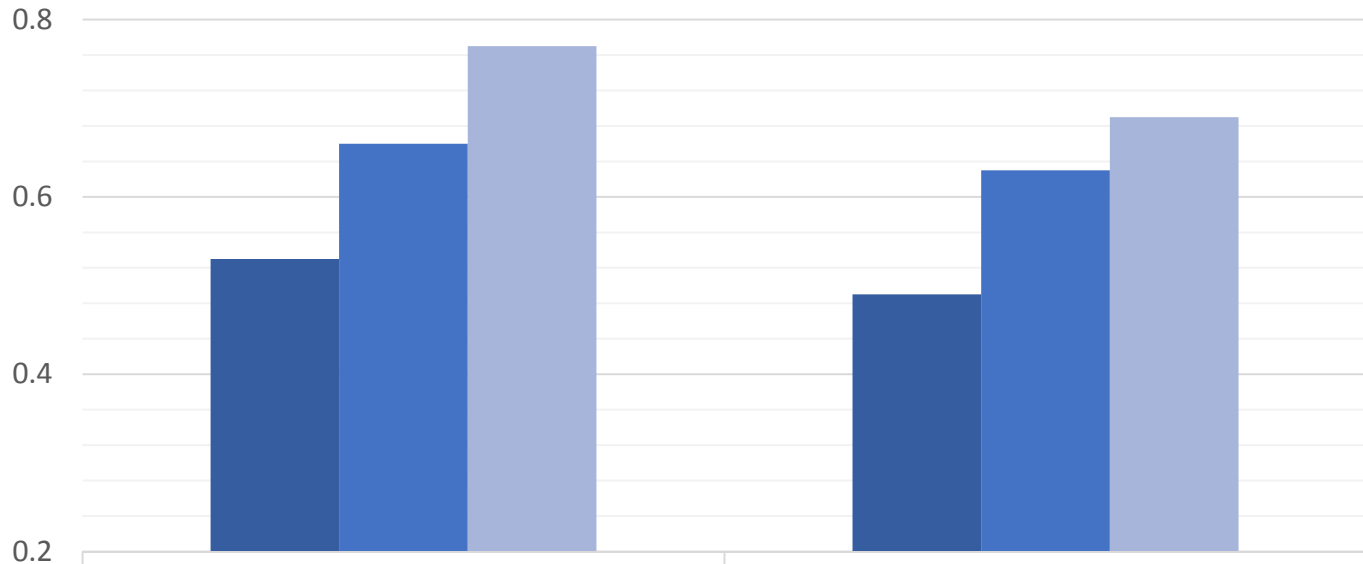


✓ Performance comparison of clustering algorithms on three databases

| Datasets | MNIST | | USPS | | YaleB | |
|-----------|--------------|--------------|--------------|--------------|--------------|--------------|
| Criterion | NMI | ACC | NMI | ACC | NMI | ACC |
| K-means | 0.494 | 0.535 | 0.615 | 0.674 | 0.866 | 0.793 |
| Spectral | 0.482 | 0.556 | 0.662 | 0.693 | 0.881 | 0.851 |
| N-cut | 0.507 | 0.543 | 0.657 | 0.696 | 0.883 | 0.821 |
| Proposed | 0.669 | 0.760 | 0.651 | 0.715 | 0.923 | 0.902 |

● Experiments

✓ Performance comparison in three different spaces with **k-means**



| | NMI | ACC |
|--------------|------|------|
| Original | 0.53 | 0.49 |
| Auto-encoder | 0.66 | 0.63 |
| Proposed | 0.77 | 0.69 |

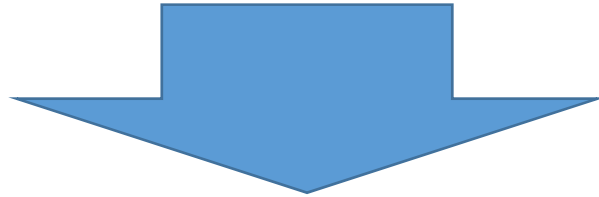
*Original means the **images(pixel) space**.
Auto-encoder means the **feature space** trained by auto-encoder nets.
Proposed means the **feature space** trained by **restrains added** auto-encoder nets.

● Conclusions

- ✓ Auto-encoder can provide good **non-linear** mapping.
- ✓ Auto-encoder nets can provide **data-stable** network.
- ✓ Restrains added can ensure **compact**.

● Conclusions

- ✓ Auto-encoder can provide good **non-linear** mapping.
- ✓ Auto-encoder nets can provide **data-stable** network.
- ✓ Restrains added can ensure **compact**.



Is **good** for **clustering**.

Thank you!

Any questions?