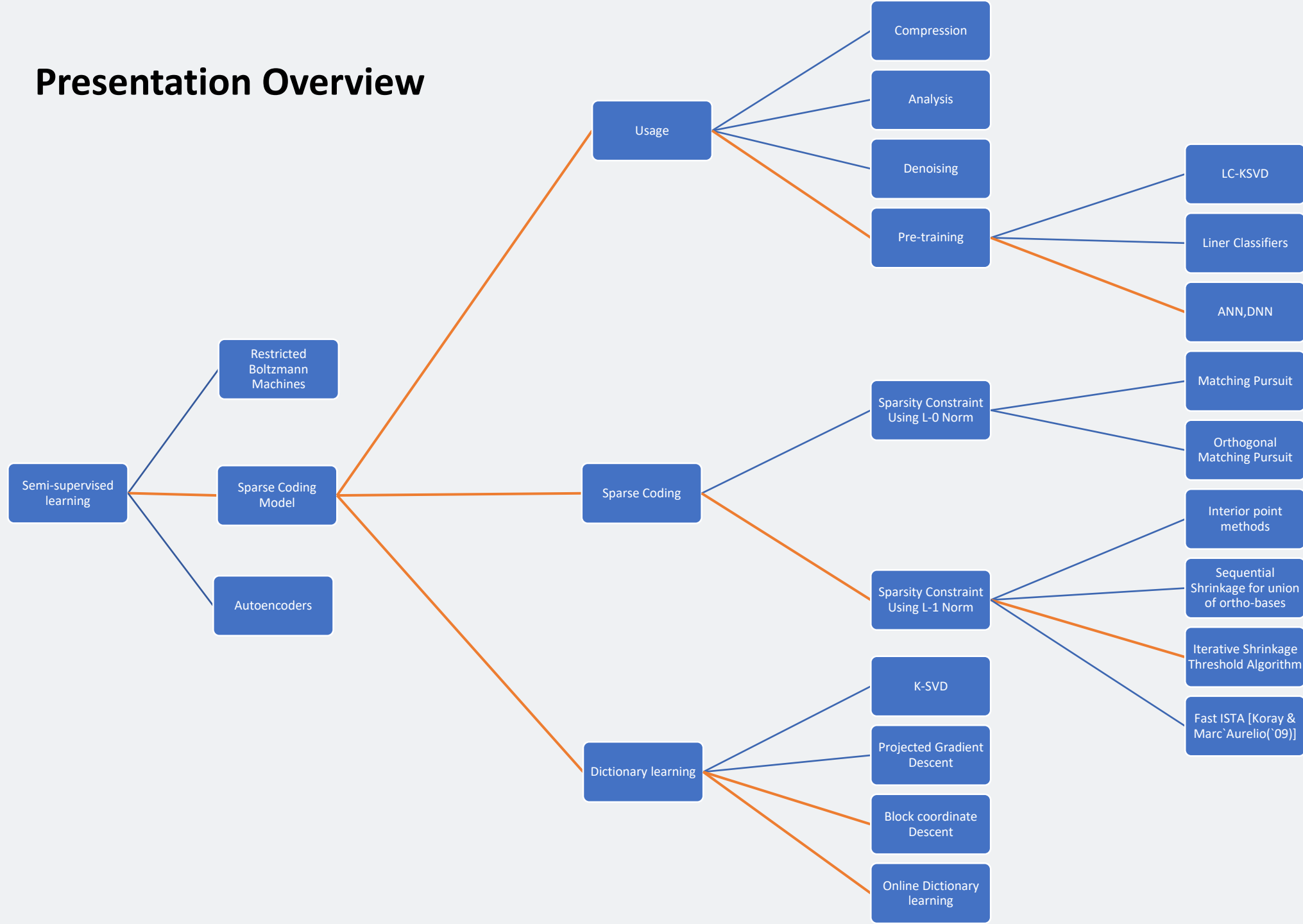
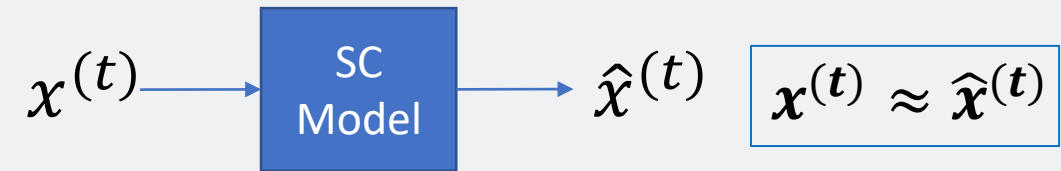


Sparse coding and dictionary learning

Presentation Overview



INTRODUCTION

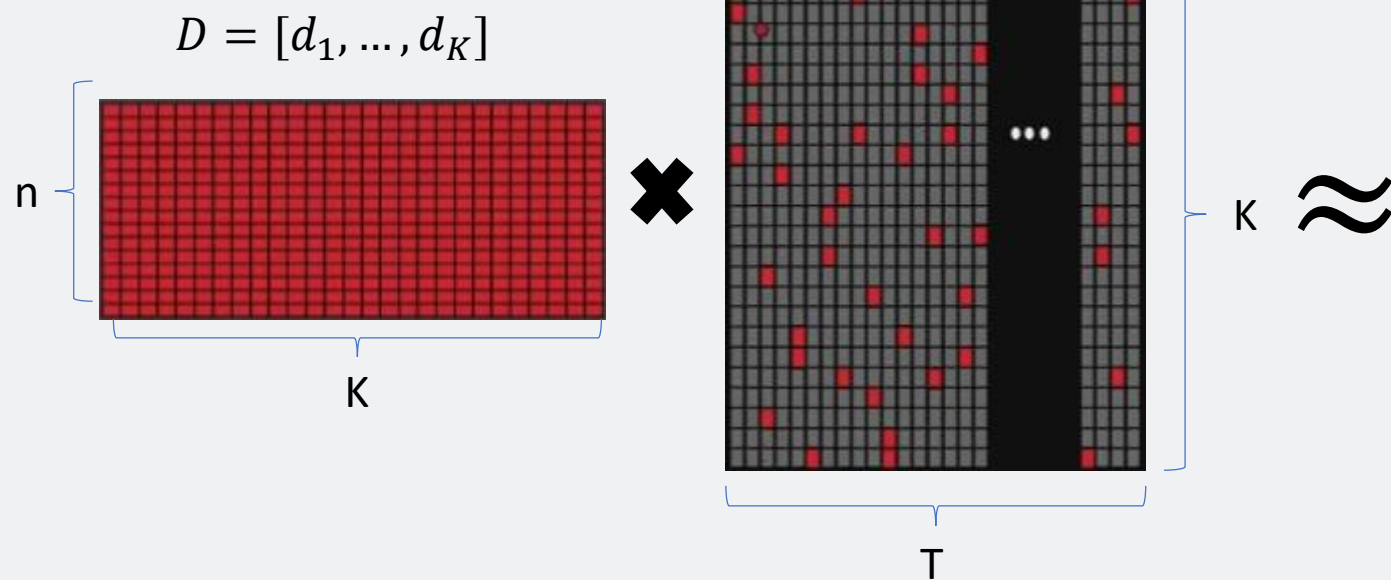


- An Unsupervised training procedure, only need to give the inputs
- Goal is to recreate the given input at the output
- While doing that, model will extract a sparse representation of the input vector
- Sparse representation of an input vector is very unique with respect to a given Dictionary.
- Dictionary consists of most dominant and principle components which will explain entire training set
- Because of that we can use this sparse representation for classification instead of using raw inputs.
- Due to that reason people use Sparse coding with Supervised learning and it's a very robust and efficient method in classification

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \underbrace{\|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2}_{\text{reconstruction error}} + \underbrace{\lambda \|\mathbf{h}^{(t)}\|_1}_{\text{sparsity penalty}}$$

$\underbrace{\mathbf{D} \mathbf{h}^{(t)}}_{\text{reconstruction } \hat{\mathbf{x}}^{(t)}} \quad \underbrace{\lambda \|\mathbf{h}^{(t)}\|_1}_{\text{reconstruction vs. sparsity control}}$

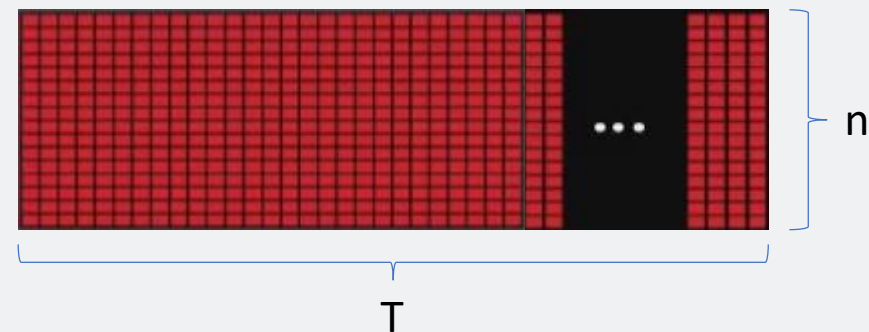
$$H = [h^{(1)}, \dots, h^{(T)}]$$



Reconstruction

\hat{X}

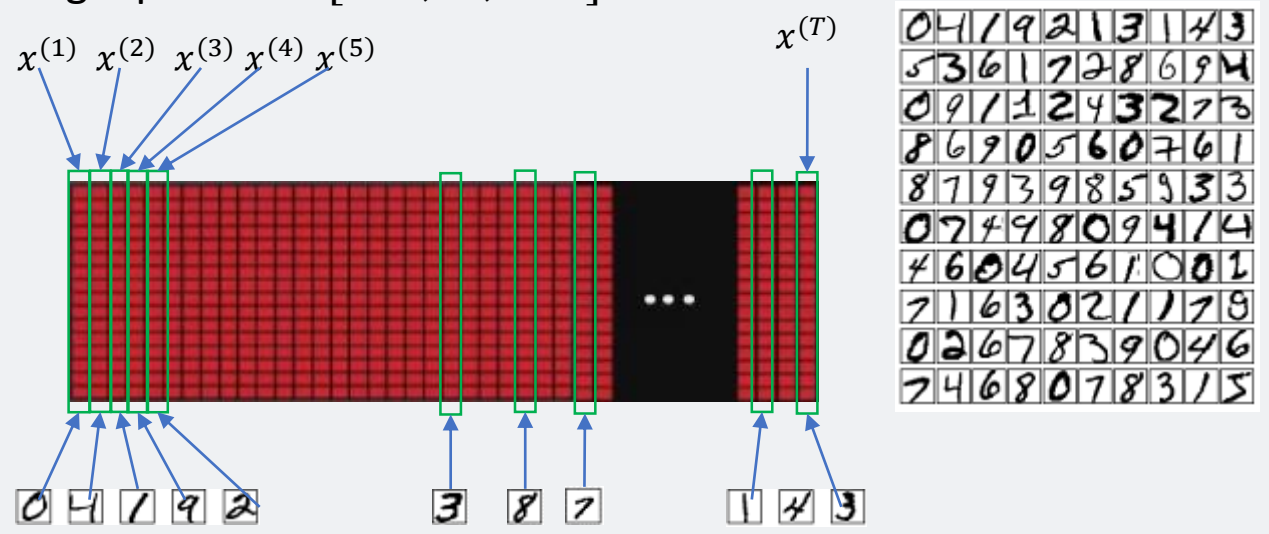
$$X = [x^{(1)}, \dots, x^{(T)}]$$



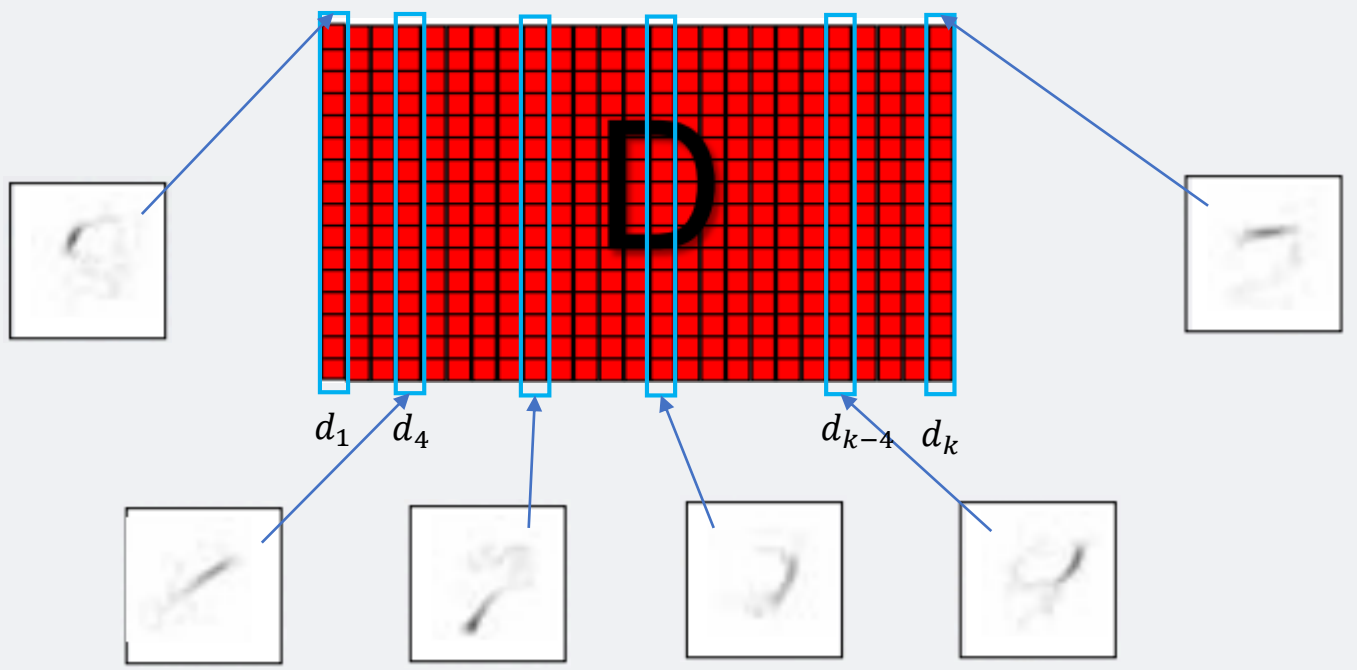
Training inputs

X

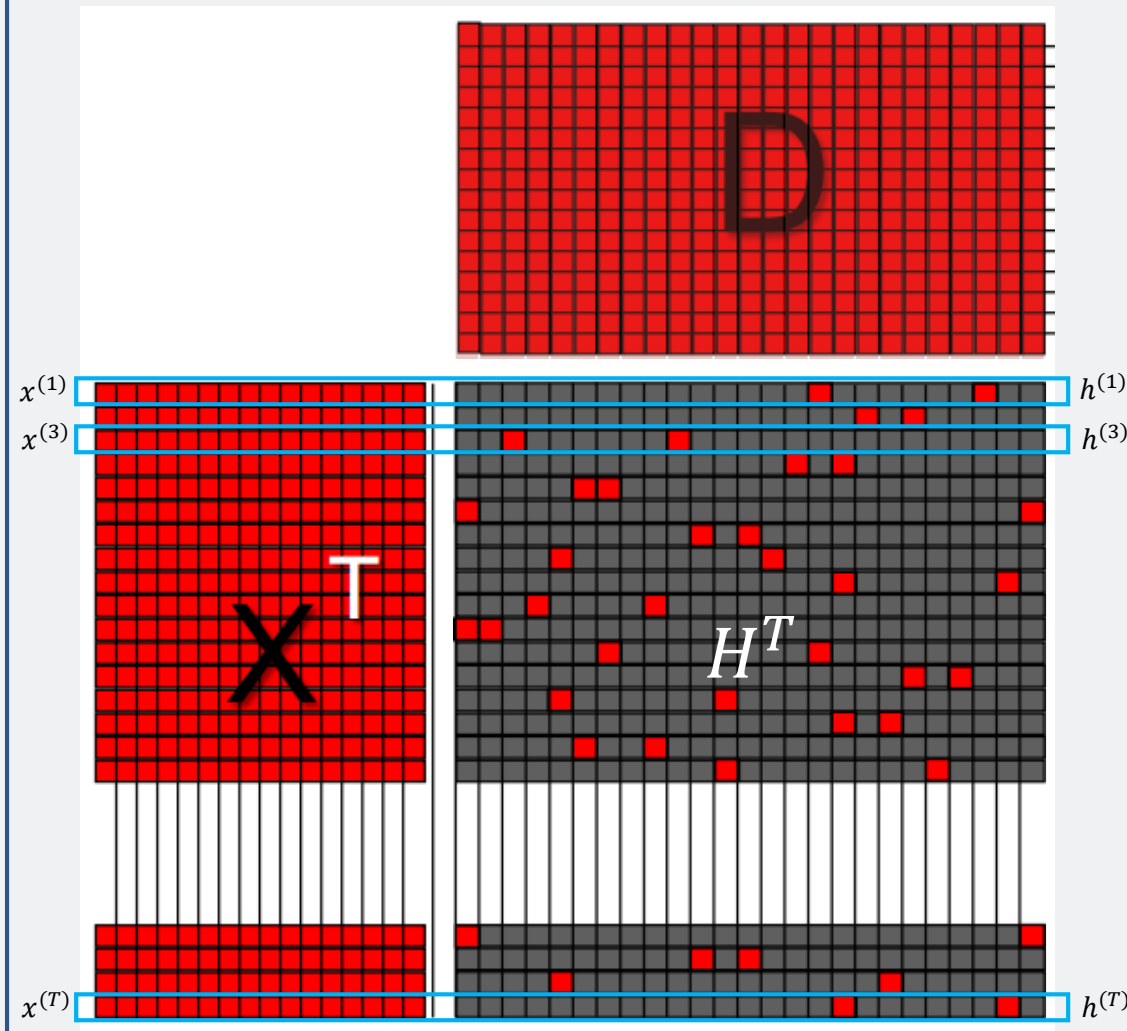
Training inputs $X = [x^{(1)}, \dots, x^{(T)}]$



Dictionary $D = [d_1, \dots, d_K]$

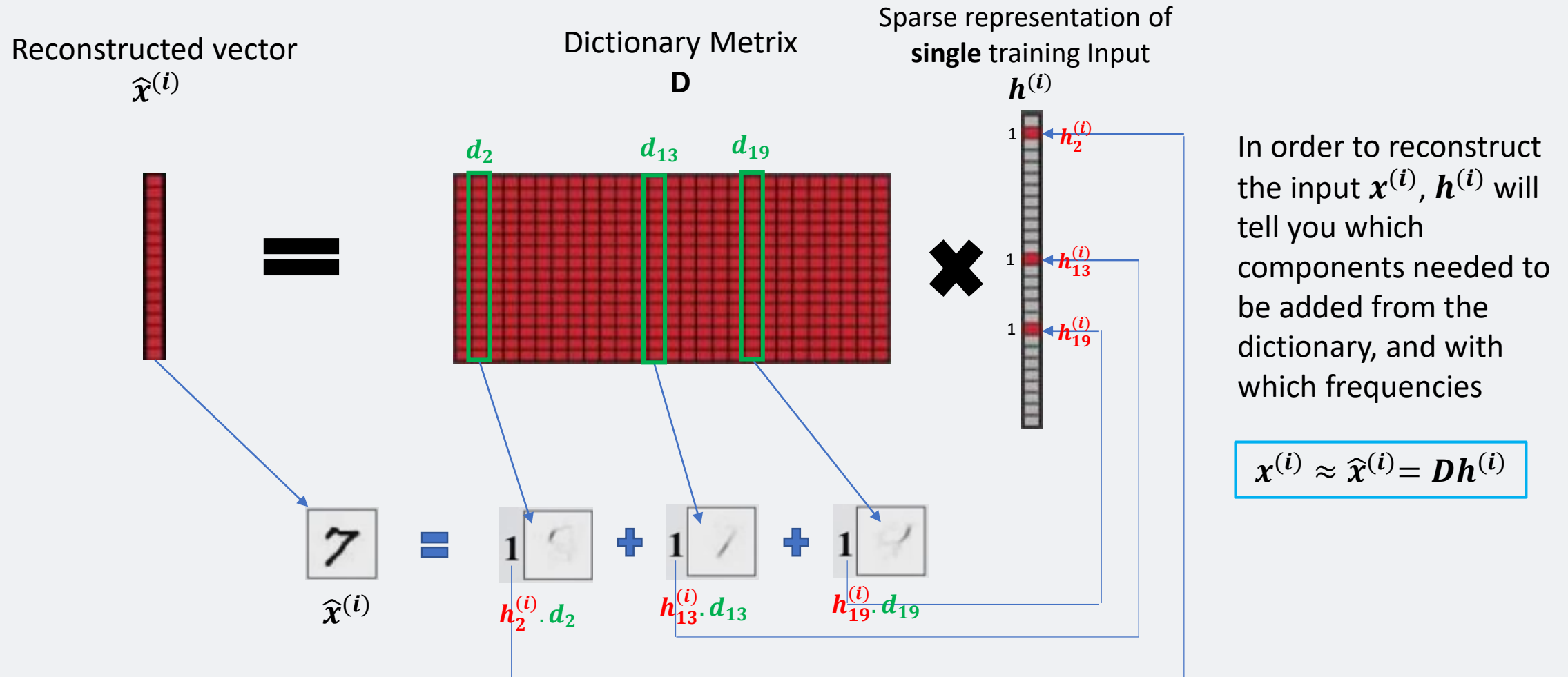


Relationship

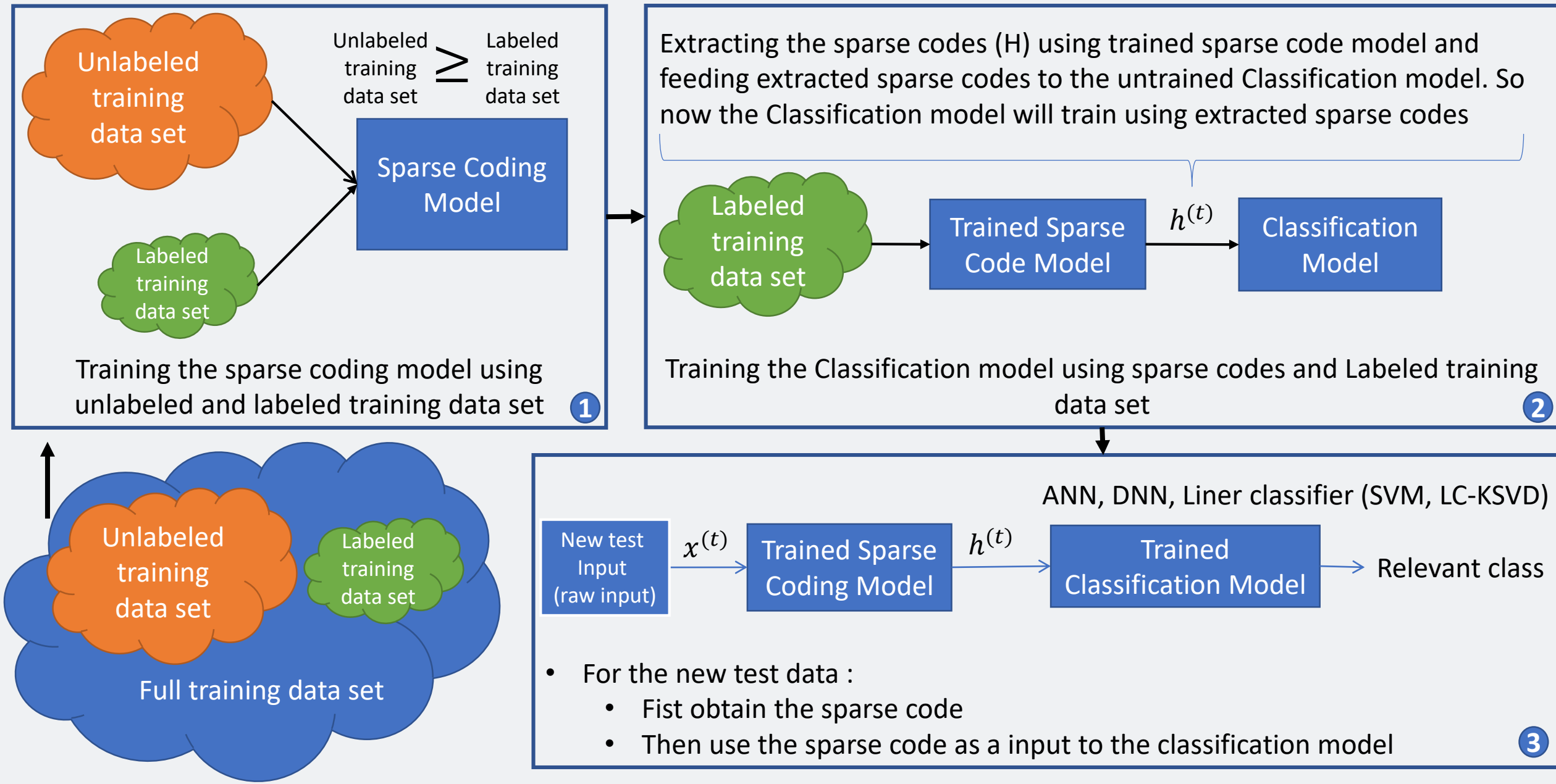


Ex : $h^{(1)}$ is the sparse representation of $x^{(1)}$

Reconstruction Example of one single input $x^{(i)} = \boxed{7}$



For classification – semi supervised learning



Training Process

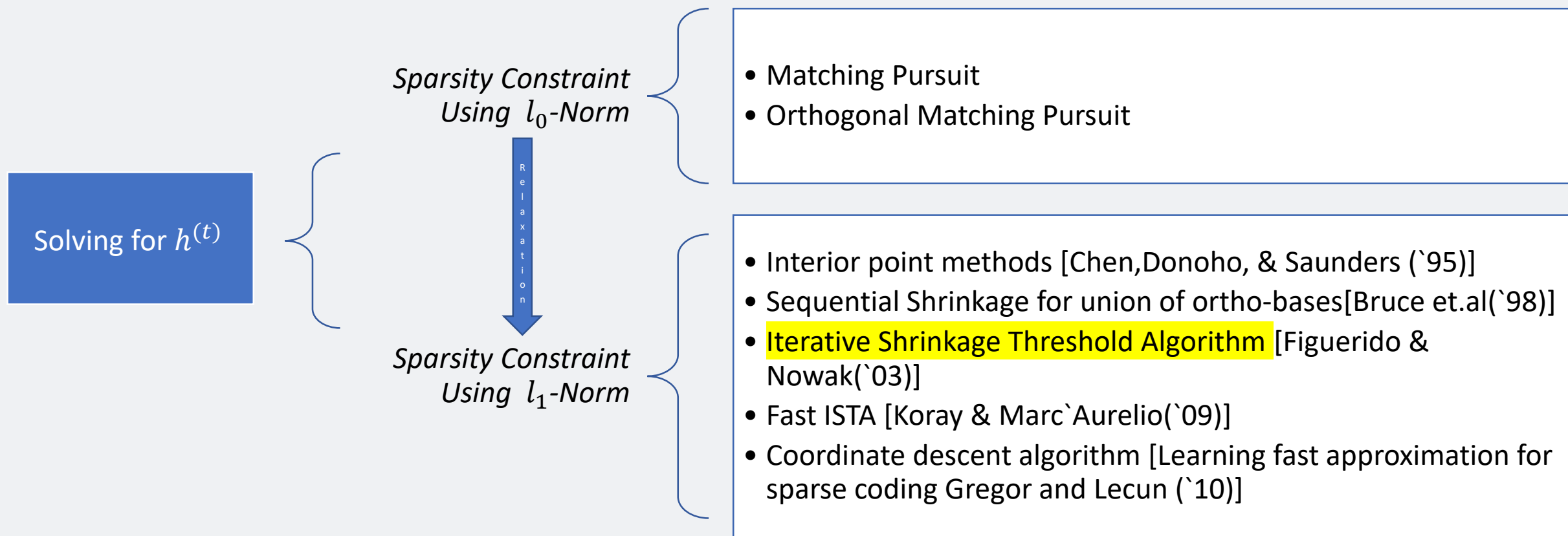
The diagram shows the equation for the training process with several annotations:

- reconstruction error**: Points to the term $\frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2$.
- sparsity penalty**: Points to the term $\lambda \|\mathbf{h}^{(t)}\|_1$.
- reconstruction $\hat{\mathbf{x}}^{(t)}$** : Points to the term $\mathbf{D} \mathbf{h}^{(t)}$.
- reconstruction vs. sparsity control**: Points to the entire equation.

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

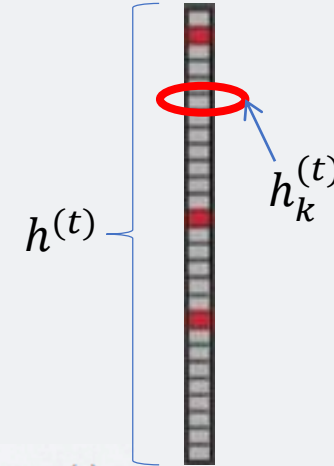
- Learning will be alternating between computing $h^{(t)}$ (sparse codes for each training input) and updating the Dictionary
- This process will go back and forth until it satisfies the convergence requirement
- Constrain the columns of \mathbf{D} to be of unit norm $\|d_i\|_2 = 1$
- If not \mathbf{D} could grow bigger while $h^{(t)}$ becomes small to satisfy the sparsity constraints
- First, consider optimizing \mathbf{H} and after that optimize \mathbf{D}

Optimizing $h^{(t)}$

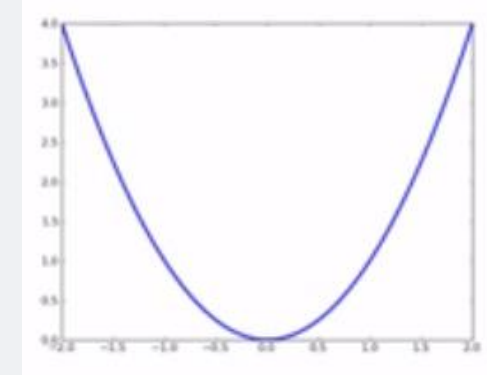


Iterative Shrinkage Threshold Algorithm (ISTA)

$$\mathbf{h}(\mathbf{x}^{(t)}) = \arg \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$



$$\frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2$$



Now the optimization will be $l(\mathbf{x}^{(t)}) = \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$ w.r.t. $\mathbf{h}^{(t)}$

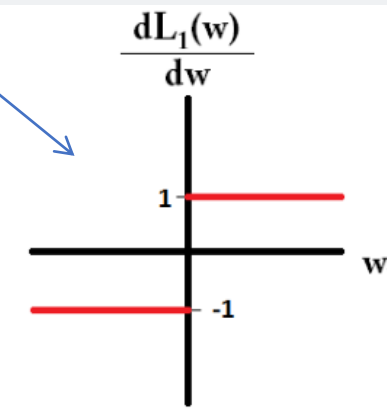
Gradient Descend method can be apply $\nabla_{\mathbf{h}^{(t)}} l(\mathbf{x}^{(t)}) = \mathbf{D}^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)}) + \lambda \text{sign}(\mathbf{h}^{(t)})$

For a single element in $\mathbf{h}^{(t)}$ $\frac{\partial}{\partial h_k^{(t)}} l(\mathbf{x}^{(t)}) = (\mathbf{D}_{:,k})^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)}) + \lambda \text{sign}(h_k^{(t)})$

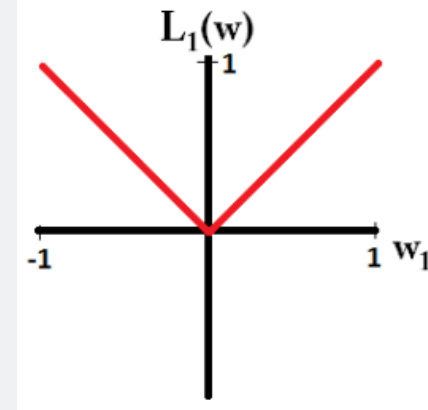
Algorithm try to minimize $h_k^{(t)}$ to zero in order to get a sparse representation

But L_1 norm not differentiable at 0 so its very unlikely for gradient descent to land on zero even if it is the solution

Solution is if $h_k^{(t)}$ changes the sign because of L_1 norm gradient, then it should be clamp to zero



Assume $w_1 = h_k^{(t)}$



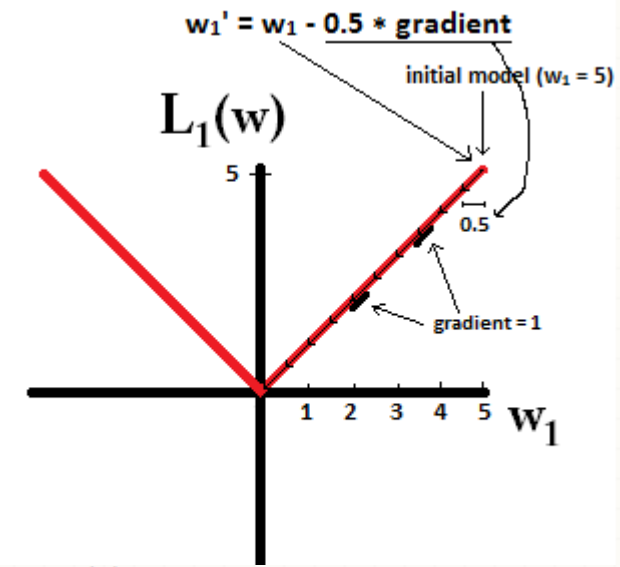
$\|\mathbf{h}^{(t)}\|_1$

▪ Each element $h_k^{(t)}$ update :

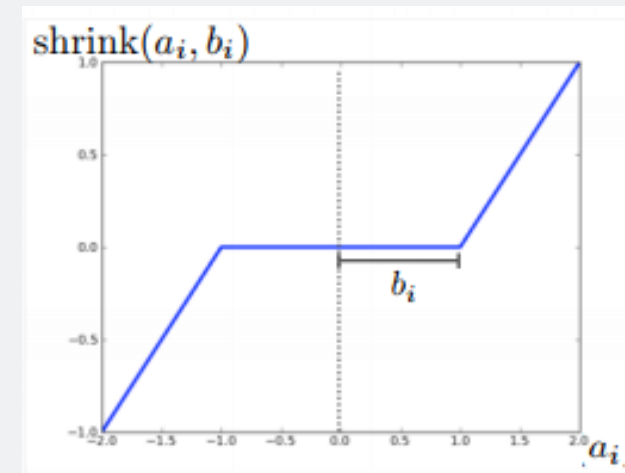
- $h_k^{(t)} \leftarrow h_k^{(t)} - \alpha (D_{:,k})^T (Dh^{(t)} - x^{(t)})$ } Update from Reconstruction
- if $\text{sign}(h_k^{(t)}) \neq \text{sign}(h_k^{(t)} - \alpha \lambda \text{sign}(h_k^{(t)}))$
 - then: $h_k^{(t)} \leftarrow 0$
- else } Update from sparsity
 - $h_k^{(t)} \leftarrow h_k^{(t)} - \alpha \lambda \text{sign}(h_k^{(t)})$

▪ Algorithm

- initialize $h^{(t)}$ (for instance to 0)
- while $h^{(t)}$ has not converged
 - $h^{(t)} \leftarrow h^{(t)} - \alpha D^T (Dh^{(t)} - x^{(t)})$
 - $h^{(t)} \leftarrow \text{shrink}(h^{(t)}, \alpha \lambda)$
- return $h^{(t)}$



$$w_1 := w_1 - \eta \cdot \frac{dL_1(w)}{dw} = w_1 - \frac{1}{2} \cdot 1, \text{ until reaching a model with } w_1 = 0$$



$$\text{shrink}(a, b) = [\text{sign}(a_i) \max(|a_i| - b_i, 0)]$$

Will converge if $\frac{1}{\alpha}$ is bigger than the largest eigenvalue of $D^T D$

Dictionary learning

- K-SVD
 - The K-SVD: An algorithm for designing of overcomplete dictionaries for sparse representations [M. Aharon, M. Elad, and A. M. Bruckstein (`06)]
- Projected Gradient Descent
- Block coordinate Descent
- Online Dictionary learning
 - Online Dictionary Learning for Sparse Coding [Mairal, Bach, Ponce and Sapiro (`09)]

Block Coordinate Descent

- Now the optimization will be

- $$\min_D \frac{1}{T} \sum_{t=1}^T \min_{h^{(t)}} l(x^{(t)}) = \min_D \frac{1}{T} \sum_{t=1}^T \frac{1}{2} \|x^{(t)} - Dh^{(t)}\|_2^2 + \lambda \|h^{(t)}\|_1$$

- For a given $h^{(t)}$

- $$\min_D \frac{1}{T} \sum_{t=1}^T \frac{1}{2} \|x^{(t)} - Dh^{(t)}\|_2^2$$

- we must also constrain the columns of D to be of unit norm ($\|D_{:,j}\|_2 = 1$)
- This algorithm doesn't need a learning rate
- Idea of this algorithms is to minimize one column at a time

Set the gradient of optimization function to zero with respect to $D_{.,j}$

$$0 = \nabla_{D_{.,j}} \left(\frac{1}{T} \sum_{t=1}^T \frac{1}{2} \|x^{(t)} - Dh^{(t)}\|_2^2 \right)$$

$$0 = \frac{1}{T} \sum_{t=1}^T (x^{(t)} - Dh^{(t)}) h_j^{(t)}$$

$$0 = \sum_{t=1}^T \left(x^{(t)} - \left(\sum_{i \neq j} D_{.,i} h_i^{(t)} \right) - D_{.,j} h_j^{(t)} \right) h_j^{(t)}$$

$$\sum_{t=1}^T D_{.,j} h_j^{(t)^2} = \sum_{t=1}^T \left(x^{(t)} - \left(\sum_{i \neq j} D_{.,i} h_i^{(t)} \right) \right) h_j^{(t)}$$

$$D_{.,j} = \frac{1}{\sum_1^T h_j^{(t)^2}} \sum_{t=1}^T \left(x^{(t)} - \left(\sum_{i \neq j} D_{.,i} h_i^{(t)} \right) \right) h_j^{(t)}$$

If $A = \sum_{t=1}^T h^{(t)} h^{(t)T}$ and $B = \sum_{t=1}^T x^{(t)} h^{(t)T}$

$$D_{.,j} = \frac{1}{\underbrace{\sum_1^T h_j^{(t)^2}}_{A_{j,j}}} \left(\underbrace{\left(\sum_{t=1}^T x^{(t)} h_j^{(t)} \right)}_{B_{.,j}} - \sum_{i \neq j} D_{.,i} \underbrace{\left(\sum_{t=1}^T h_i^{(t)} h_j^{(t)} \right)}_{A_{i,j}} \right)$$

$$D_{.,j} = \frac{1}{A_{j,j}} (B_{.,j} - \sum_{i \neq j} D_{.,i} A_{i,j})$$

Find the best value for j^{th} column with respect to all the other parameters (including fixing of other columns of **D** matrix)

Then iterate cycles over each columns
 $1^{st} \rightarrow 2^{nd} \rightarrow 3^{rd} \rightarrow 4^{th} \dots \rightarrow 1^{st} \dots$
 until we get an stable values for **D**

Pseudocode for block coordinate descent algorithm

while \mathbf{D} hasn't converged

for each column $D_{:,j}$ perform updates

$$D_{:,j} = \frac{1}{A_{j,j}} (B_{:,j} - \sum_{i \neq j} D_{:,i} A_{i,j}) \leftarrow \text{coordinate descent update}$$

$$D_{:,j} = \frac{D_{:,j}}{\|D_{:,j}\|_2} \leftarrow \text{projection ; normalizing to satisfy the column constrain of D (unit norm)}$$

Implementation note :-

Since \mathbf{A} and \mathbf{B} not depend of \mathbf{D} , we can calculate \mathbf{A} and \mathbf{B} before doing any optimization to \mathbf{D} and keep it in memory. It will be efficient and save lot of resources

Complete Pseudo-code For Sparse Coding And Dictionary Learning

- *Initialize D and normalize the columns*
- *while D has not converged :*
 - *find the sparse codes $h^{(t)} \forall x^{(t)}$ in training set using **ISTA** for a given D*
 - *update the dictionary*
 - *compute A and B for calculated $h^{(t)}$*
 - *run **block coordinate descend** algorithm*

- As we can see no learning rate will be required by Block coordinate descent
- Learning will be alternating between computing $h^{(t)}$ (sparse codes for each training input) and updating the Dictionary.
- This process will go back and forth until it satisfies the convergence requirement
- This method is related to batch learning category
- It means for a single update, it will pass through the whole given training set
- This method could be inefficient for a big amount of training data
- The solution is an online learning algorithm and this means algorithm update while data is coming
- This can apply context where constantly obtaining new data in online fashion (like downloading images, solar irradiance)

Online Dictionary Learning Algorithm

- This method required an update of **running averages** for **A** and **B**
- $B_{new} \Leftarrow \beta B_{pre} + (1 - \beta)x^{(t)}h^{(t)T}$
- $A_{new} \Leftarrow \beta A_{pre} + (1 - \beta)h^{(t)}h^{(t)T}$
- Now the previously mentioned A and B has changed with a running average
- β is a hyper parameter and if β close to 1 then model will give more priority to previous **B** and give less priority to incoming new data

Initializing D

- Not to 0
 - Because columns of D vectors should have unit norm
 - It is impossible to reconstruct and do the learning
- Use random matrix and then normalize each of the column
- Normalize first few training vectors and use those as columns in matrix D

Online Dictionary Learning Algorithm

- *Initialize \mathbf{D}*
- *While \mathbf{D} hasn't converged*
 - *for each $x^{(t)}$*
 - *infer sparse codes $h^{(t)}$ using **ISTA*** } Run ISTA to extract the sparse code for new data
 - *update Dictionary*
 - $B_{new} \Leftarrow \beta B_{prev} + (1 - \beta)x^{(t)}h^{(t)T}$
 - $A_{new} \Leftarrow \beta A_{prev} + (1 - \beta)h^{(t)}h^{(t)T}$} Calculate the running average
 - *while \mathbf{D} hasn't converged*
 - *for each $D_{:,j}$ column perform Gradient update*
 - $D_{:,j} = \frac{1}{A_{j,j}}(B_{:,j} - \sum_{i \neq j} D_{:,i} A_{i,j})$
 - $D_{:,j} = \frac{D_{:,j}}{\|D_{:,j}\|_2}$} Run Block Coordinate Descent

Thank You