

Sparse Coding and Dictionary Learning

Vahid Tarokh

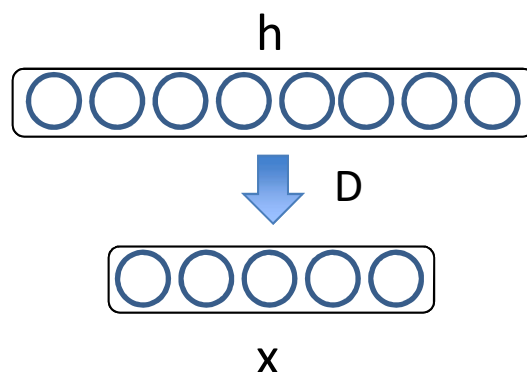
ECE 685D, Fall 2025

Unsupervised Learning

- Unsupervised learning: we only use the inputs $\mathbf{x}^{(t)}$ for learning
 - automatically extract meaningful features for your data
 - leverage the availability of unlabeled data
 - add a data-dependent regularizer to training (in some cases)
- We will consider 3 models for unsupervised learning that will form the basic building blocks for deeper models:
 - Autoencoders
 - Sparse Coding
 - Restricted Boltzmann Machines

Sparse Coding

- Sparse coding was originally developed in AI literature to explain early visual processing in the brain (edge detection) [Olshausen & Field, 1996].
- For each input $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - **The representation is sparse:** the vector $\mathbf{h}^{(t)}$ has many zeros
 - It can be used to well **reconstruct** the original input $\mathbf{x}^{(t)}$



Sparse Coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - **The representation is sparse**: the vector $\mathbf{h}^{(t)}$ has many zeroes
 - It can be used to well **reconstruct** the original input $\mathbf{x}^{(t)}$
- In other words:

Reconstruction: $\hat{\mathbf{x}}^{(t)}$

Sparsity vs. reconstruction control

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

Sparse Coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - The representation is sparse: the vector $\mathbf{h}^{(t)}$ has many zeroes
 - It can be used to well reconstruct the original input $\mathbf{x}^{(t)}$

- In other words:

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

- we also constrain the columns of \mathbf{D} to be of norm 1
- otherwise, \mathbf{D} could grow big while \mathbf{h} becomes small to satisfy the L_1 constraint

Sparse Coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - The representation is sparse: the vector $\mathbf{h}^{(t)}$ has many zeroes
 - It can be used to well reconstruct the original input $\mathbf{x}^{(t)}$
- In other words:

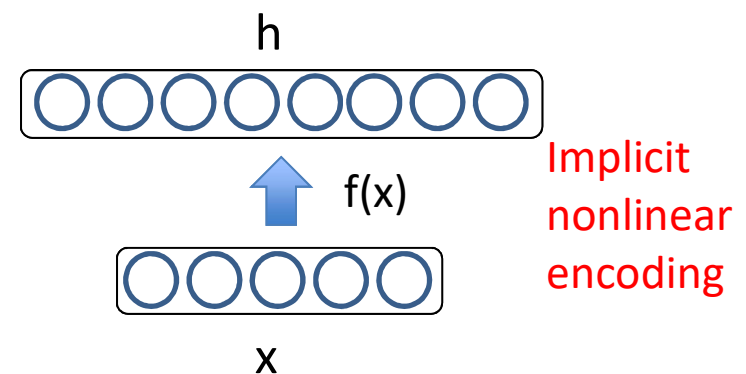
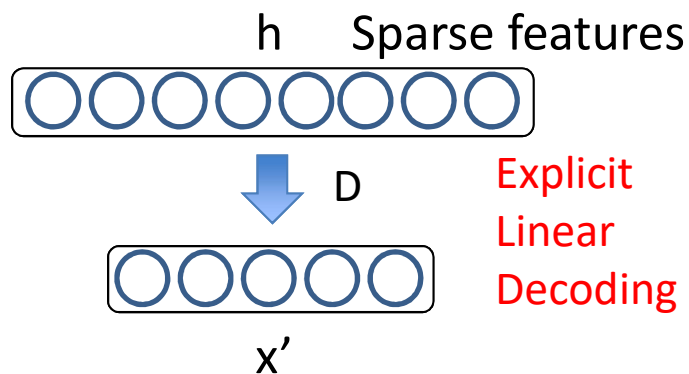
$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

- \mathbf{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 problem:

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

Interpreting Sparse Coding

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



- Sparse, over-complete representation \mathbf{h} .
- Encoding $\mathbf{h} = f(\mathbf{x})$ is implicit and nonlinear function of \mathbf{x} .
- Reconstruction (or decoding) $\mathbf{x}' = \mathbf{D}\mathbf{h}$ is linear and explicit.

Sparse Coding

- We can also write:

$$\hat{\mathbf{x}}^{(t)} = \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)}) = \sum_{\substack{k \text{ s.t.} \\ h(\mathbf{x}^{(t)})_k \neq 0}} \mathbf{D}_{\cdot, k} h(\mathbf{x}^{(t)})_k$$

$$\begin{aligned} \boxed{7} &= 1 \boxed{\text{basis 1}} + 1 \boxed{\text{basis 2}} + 1 \boxed{\text{basis 3}} + 1 \boxed{\text{basis 4}} + 1 \boxed{\text{basis 5}} + 1 \boxed{\text{basis 6}} \\ &\quad + 1 \boxed{\text{basis 7}} + 1 \boxed{\text{basis 8}} + 0.8 \boxed{\text{basis 9}} + 0.8 \boxed{\text{basis 10}} \end{aligned}$$

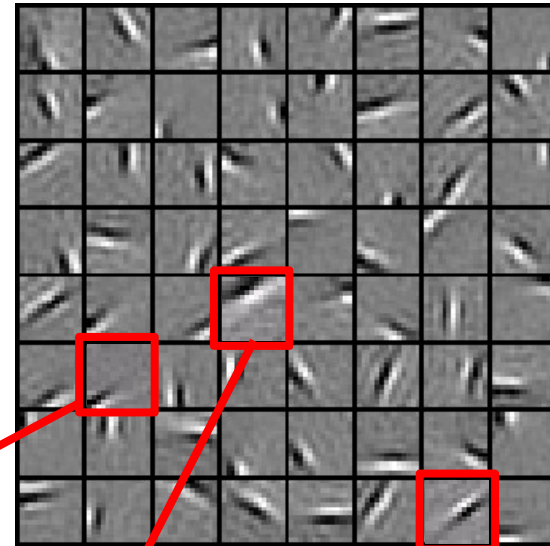
- D is often referred to as **Dictionary**
- In certain applications, we know what dictionary matrix to use
- In many cases, we have to learn it

Sparse Coding

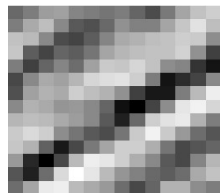
Natural Images



Learned bases: “Edges”



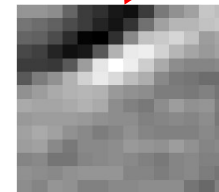
New example



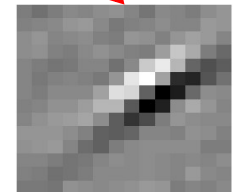
$= 0.8 *$



$+ 0.3 *$



$+ 0.5 *$



x

$= 0.8 *$

ϕ_{36}

$+ 0.3 *$

ϕ_{42}

$+ 0.5 *$

ϕ_{65}

$[0, 0, \dots, \mathbf{0.8}, \dots, \mathbf{0.3}, \dots, \mathbf{0.5}, \dots]$ = coefficients (feature representation)

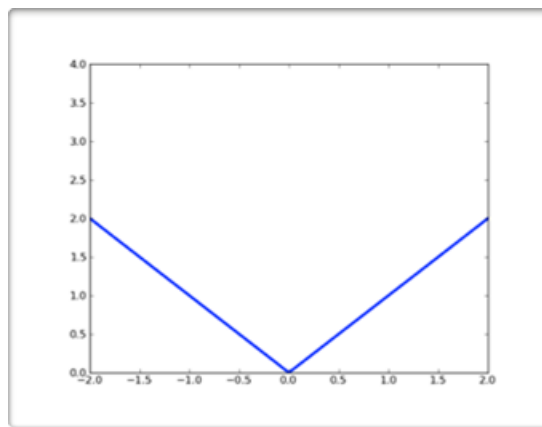
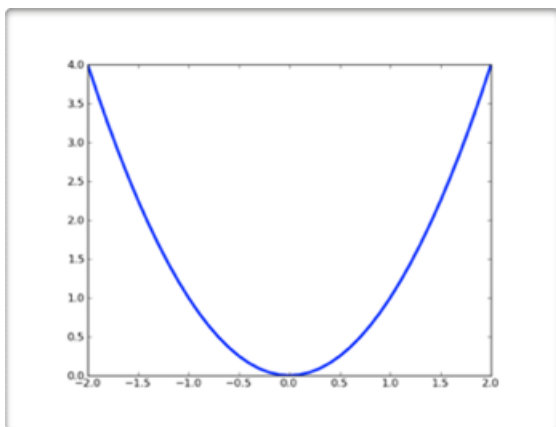
Inference

- Given dictionary \mathbf{D} , how do we compute $\mathbf{h}(\mathbf{x}^{(t)})$?

- We need to optimize:

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

- This is Lasso.



- We could use a **gradient descent** method:

$$\nabla_{\mathbf{h}^{(t)}} l(\mathbf{x}^{(t)}) = \mathbf{D}^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)}) + \lambda \text{sign}(\mathbf{h}^{(t)})$$

Inference

- For a single hidden unit:

$$\frac{\partial}{\partial h_k^{(t)}} l(\mathbf{x}^{(t)}) = (\mathbf{D}_{\cdot, k})^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)}) + \lambda \operatorname{sign}(h_k^{(t)})$$

- **issue:** L_1 norm **not differentiable** at 0
- very unlikely for gradient descent to “land” on $h_k^{(t)} = 0$ (even if it’s the solution)
- **Solution:** if $h_k^{(t)}$ changes sign because of L_1 norm gradient, clamp to 0.

Inference

- For a single hidden unit:

$$\frac{\partial}{\partial h_k^{(t)}} l(\mathbf{x}^{(t)}) = (\mathbf{D}_{\cdot, k})^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)}) + \lambda \text{sign}(h_k^{(t)})$$

- **Solution:** if $h_k^{(t)}$ changes sign because of L1 norm gradient, clamp to 0.

- Each hidden unit update would be performed as follows:

➤ $h_k^{(t)} \Leftarrow h_k^{(t)} - \alpha (\mathbf{D}_{\cdot, k})^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)})$

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 $h_k^{(t)} \Leftarrow h_k^{(t)} - \alpha \lambda \text{sign}(h_k^{(t)})$

Update sparsity
term

ISTA Algorithm

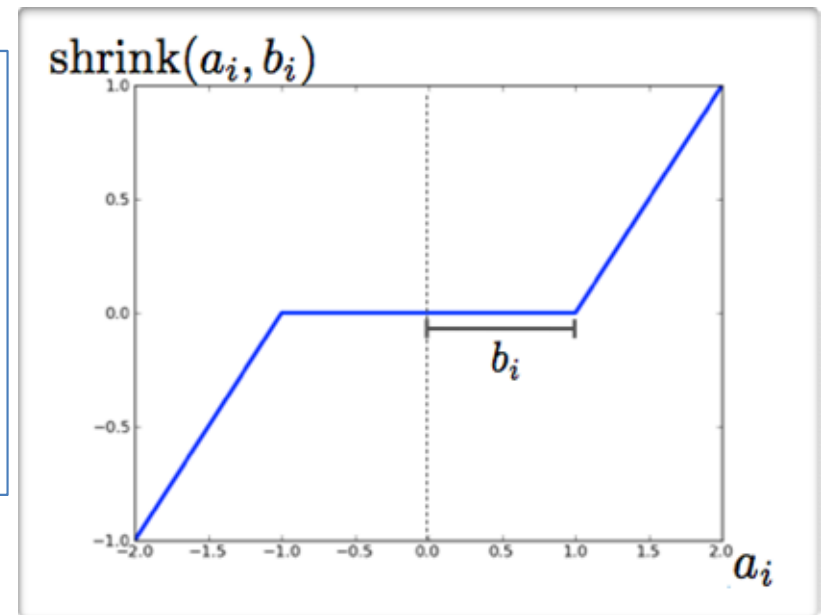
- This process corresponds to the **ISTA** (Iterative Shrinkage and Thresholding) Algorithm:

- Initialize $\mathbf{h}^{(t)}$ (for example to 0)
- While $\mathbf{h}^{(t)}$ has not converged
$$\mathbf{h}^{(t)} \Leftarrow \mathbf{h}^{(t)} - \alpha \mathbf{D}^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)})$$
$$\mathbf{h}^{(t)} \Leftarrow \text{shrink}(\mathbf{h}^{(t)}, \alpha \lambda)$$

where

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

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



Dictionary Learning I

- Remember our **optimization problem**:

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

- Let us first assume that $\mathbf{h}(\mathbf{x}^{(t)})$ does not depend on \mathbf{D}

- We then minimize:

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

- we must also constrain the columns of \mathbf{D} to be of unit norm

column norm constraints: $\|\mathbf{D}_{:,j}\|_2 = 1, \forall j$

Dictionary Learning I

- We can use **projected gradient descent** algorithm.

- While D has not converged:

- Perform gradient update of D

$$\mathbf{D} \Leftarrow \mathbf{D} + \alpha \frac{1}{T} \sum_{t=1}^T (\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)})) \mathbf{h}(\mathbf{x}^{(t)})^\top$$

- Renormalize the columns of D

- For each column of D:

$$\mathbf{D}_{\cdot,j} \Leftarrow \frac{\mathbf{D}_{\cdot,j}}{\|\mathbf{D}_{\cdot,j}\|_2}$$

Dictionary Learning II

- An **alternative method** is to solve for each column $\mathbf{D}_{:,j}$ in cycle.
 - setting the gradient for $\mathbf{D}_{:,j}$ to zero, we have

$$\begin{aligned} 0 &= \frac{1}{T} \sum_{t=1}^T (\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)})) h(\mathbf{x}^{(t)})_j \\ 0 &= \sum_{t=1}^T \left(\mathbf{x}^{(t)} - \left(\sum_{i \neq j} \mathbf{D}_{:,i} h(\mathbf{x}^{(t)})_i \right) - \mathbf{D}_{:,j} h(\mathbf{x}^{(t)})_j \right) h(\mathbf{x}^{(t)})_j \\ \sum_{t=1}^T \mathbf{D}_{:,j} h(\mathbf{x}^{(t)})_j^2 &= \sum_{t=1}^T \left(\mathbf{x}^{(t)} - \left(\sum_{i \neq j} \mathbf{D}_{:,i} h(\mathbf{x}^{(t)})_i \right) \right) h(\mathbf{x}^{(t)})_j \\ \mathbf{D}_{:,j} &= \frac{1}{\sum_{t=1}^T h(\mathbf{x}^{(t)})_j^2} \sum_{t=1}^T \left(\mathbf{x}^{(t)} - \left(\sum_{i \neq j} \mathbf{D}_{:,i} h(\mathbf{x}^{(t)})_i \right) \right) h(\mathbf{x}^{(t)})_j \end{aligned}$$

- Note that we don't need to specify a learning rate to update \mathbf{D} .

Dictionary Learning II

- An **alternative method** is to solve for each column $\mathbf{D}_{:,j}$ in cycle.

➤ We can rewrite

$$\begin{aligned}\mathbf{D}_{:,j} &= \frac{1}{\sum_{t=1}^T h(\mathbf{x}^{(t)})_j^2} \sum_{t=1}^T \left(\mathbf{x}^{(t)} - \left(\sum_{i \neq j} \mathbf{D}_{:,i} h(\mathbf{x}^{(t)})_i \right) \right) h(\mathbf{x}^{(t)})_j \\ &= \frac{1}{\sum_{t=1}^T h(\mathbf{x}^{(t)})_j^2} \left(\left(\sum_{t=1}^T \mathbf{x}^{(t)} h(\mathbf{x}^{(t)})_j \right) - \sum_{i \neq j} \mathbf{D}_{:,i} \left(\sum_{t=1}^T h(\mathbf{x}^{(t)})_i h(\mathbf{x}^{(t)})_j \right) \right) \\ &= \frac{1}{A_{j,j}} (\mathbf{B}_{:,j} - \mathbf{D} \mathbf{A}_{:,j} + \mathbf{D}_{:,j} A_{j,j})\end{aligned}$$

➤ this way, we only need to store:

$$\mathbf{A} \Leftarrow \sum_{t=1}^T \mathbf{h}(\mathbf{x}^{(t)}) \mathbf{h}(\mathbf{x}^{(t)})^\top$$

$$\mathbf{B} \Leftarrow \sum_{t=1}^T \mathbf{x}^{(t)} \mathbf{h}(\mathbf{x}^{(t)})^\top$$

Dictionary Learning II

- This leads to the following algorithm

- While \mathbf{D} has not converged:

- for each column $\mathbf{D}_{:,j}$ perform updates

$$\mathbf{D}_{:,j} \leftarrow \frac{1}{A_{j,j}} (\mathbf{B}_{:,j} - \mathbf{D} \mathbf{A}_{:,j} + \mathbf{D}_{:,j} A_{j,j})$$

$$\mathbf{D}_{:,j} \leftarrow \frac{\mathbf{D}_{:,j}}{\|\mathbf{D}_{:,j}\|_2}$$

- This is referred to as a **block-coordinate descent algorithm**
 - a different block of variables are updated at each step
 - the “blocks” are the columns $\mathbf{D}_{:,j}$

Learning Sparse Coding Model

- Putting it all together, we have the following algorithm, where learning alternates between **inference and dictionary learning**.

- While D has not converged:

- find the sparse codes $\mathbf{h}(\mathbf{x}^{(t)})$ for all $\mathbf{x}^{(t)}$ in the training set with ISTA
- Update the dictionary:

$$\mathbf{A} \leftarrow \sum_{t=1}^T \mathbf{x}^{(t)} \mathbf{h}(\mathbf{x}^{(t)})^\top$$

$$\mathbf{B} \leftarrow \sum_{t=1}^T \mathbf{h}(\mathbf{x}^{(t)}) \mathbf{h}(\mathbf{x}^{(t)})^\top$$

$$\mathbf{D}_{\cdot,j} \leftarrow \frac{1}{A_{j,j}} (\mathbf{B}_{\cdot,j} - \mathbf{D} \mathbf{A}_{\cdot,j} + \mathbf{D}_{\cdot,j} A_{j,j})$$

$$\mathbf{D}_{\cdot,j} \leftarrow \frac{\mathbf{D}_{\cdot,j}}{\|\mathbf{D}_{\cdot,j}\|_2}$$

Instead of computing A and B from scratch over all data, we maintain running averages that update with each new sample.

Online Learning

- This algorithm is “**batch**” (i.e. not online)
 - single update of the dictionary per pass on the training set
 - for large datasets, we’d like to update D after visiting each $\mathbf{x}^{(t)}$
- **Solution**: for each input $\mathbf{x}^{(t)}$
 - perform inference of $\mathbf{h}(\mathbf{x}^{(t)})$ for the current input
 - update running averages of the quantities required to update D:

$$\mathbf{B} \leftarrow \beta \mathbf{B} + (1 - \beta) \mathbf{x}^{(t)} \mathbf{h}(\mathbf{x}^{(t)})^\top$$

$$\mathbf{A} \leftarrow \beta \mathbf{A} + (1 - \beta) \mathbf{h}(\mathbf{x}^{(t)}) \mathbf{h}(\mathbf{x}^{(t)})^\top$$

Online Learning

- While D has not converged:

- For each $\mathbf{x}^{(t)}$

- Infer code $\mathbf{h}(\mathbf{x}^{(t)})$

- Update the dictionary:

$$\mathbf{A} \Leftarrow \beta \mathbf{A} + (1 - \beta) \mathbf{h}(\mathbf{x}^{(T+1)}) \mathbf{h}(\mathbf{x}^{(T+1)})^\top$$

$$\mathbf{B} \Leftarrow \beta \mathbf{B} + (1 - \beta) \mathbf{x}^{(t)} \mathbf{h}(\mathbf{x}^{(t)})^\top$$

- while D hasn't converged

- for each column of D perform gradient update

$$\mathbf{D}_{:,j} \Leftarrow \frac{1}{A_{j,j}} (\mathbf{B}_{:,j} - \mathbf{D} \mathbf{A}_{:,j} + \mathbf{D}_{:,j} A_{j,j})$$

$$\mathbf{D}_{:,j} \Leftarrow \frac{\mathbf{D}_{:,j}}{\|\mathbf{D}_{:,j}\|_2}$$

Online Dictionary Learning for Sparse
Coding. [1] Mairal, Bach, Ponce and Sapiro, 2009.

Whitening

- Before running a sparse coding algorithm, it is beneficial to remove “obvious” structure from the data

- normalize such that mean is 0 and covariance is the identity (whitening)
- this will remove 1st and 2nd order statistical structure

- preprocessing

- let the empirical mean be $\hat{\mu}$ and the empirical covariance matrix be $\hat{\Sigma} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$ (in its eigenvalue/eigenvector representation)
- ZCA transforms each input as follows:

$$\mathbf{x} \longleftarrow \mathbf{U} \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^\top (\mathbf{x} - \hat{\mu})$$

Whitening

- After this transformation
 - the empirical mean is 0

$$\begin{aligned} & \frac{1}{T} \sum_t \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top (\mathbf{x}^{(t)} - \hat{\boldsymbol{\mu}}) \\ &= \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \left(\left(\frac{1}{T} \sum_t \mathbf{x}^{(t)} \right) - \hat{\boldsymbol{\mu}} \right) \\ &= \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top (\hat{\boldsymbol{\mu}} - \hat{\boldsymbol{\mu}}) \\ &= 0 \end{aligned}$$

Whitening

- After this transformation
 - the empirical covariance matrix is the identity

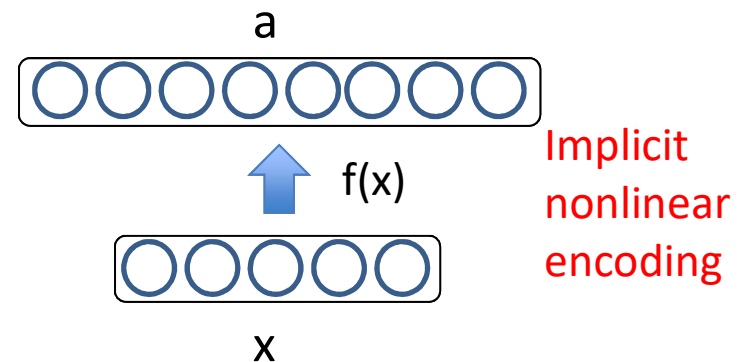
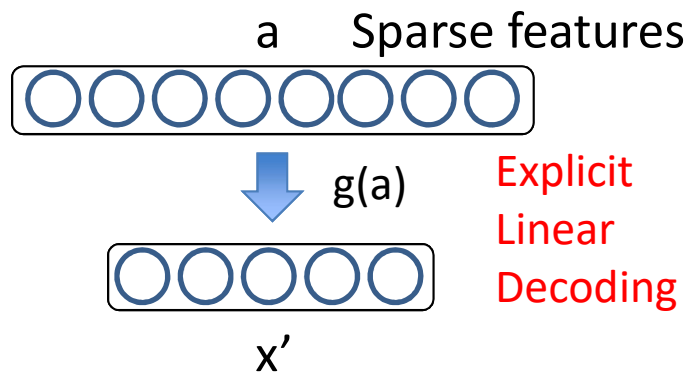
$$\begin{aligned} & \frac{1}{T-1} \sum_t \left(\mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top (\mathbf{x}^{(t)} - \hat{\boldsymbol{\mu}}) \right) \left(\mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top (\mathbf{x}^{(t)} - \hat{\boldsymbol{\mu}}) \right)^\top \\ &= \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \left(\frac{1}{T-1} \sum_t (\mathbf{x}^{(t)} - \hat{\boldsymbol{\mu}})(\mathbf{x}^{(t)} - \hat{\boldsymbol{\mu}})^\top \right) \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \\ &= \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \hat{\boldsymbol{\Sigma}} \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \\ &= \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \mathbf{U} \Lambda \mathbf{U}^\top \mathbf{U} \Lambda^{-\frac{1}{2}} \mathbf{U}^\top \\ &= \mathbf{I} \end{aligned}$$

Feature Learning

- A sparse coding model can be used to extract features
 - given a **labeled** training set $\{(\mathbf{x}^{(t)}, y^{(t)})\}$
 - train sparse coding dictionary only on training inputs $\{\mathbf{x}^{(t)}\}$
 - this yields a **dictionary** from which to infer sparse codes
 - train your favorite classifier on transformed training set $\{(\mathbf{h}(\mathbf{x}^{(t)}), y^{(t)})\}$
- When classifying test input \mathbf{x}
 - infer its sparse representation: $\mathbf{h}(\mathbf{x})$
 - feed it to the classifier

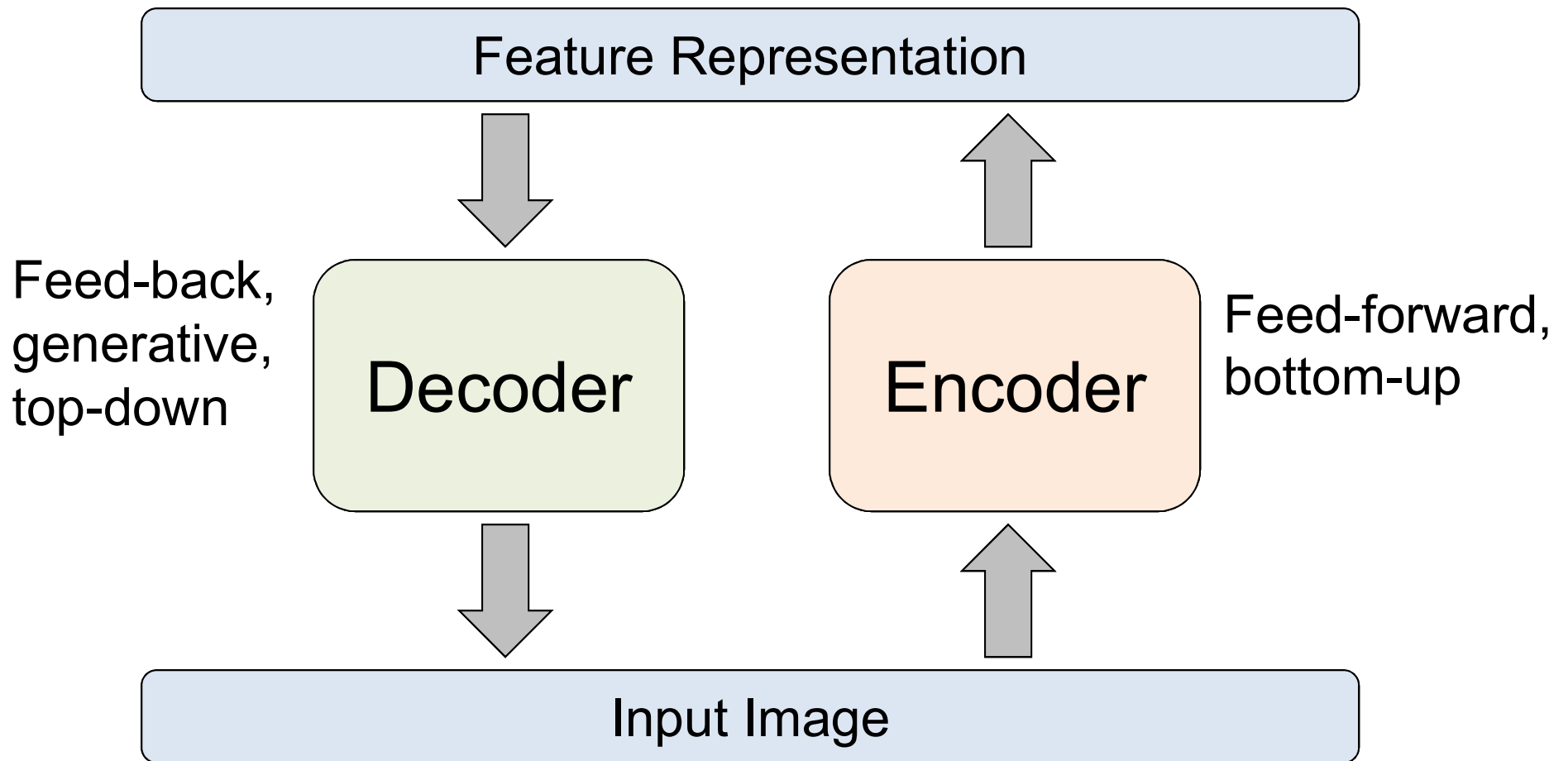
Interpreting Sparse Coding

$$\min_{\mathbf{a}, \phi} \sum_{n=1}^N \left\| \mathbf{x}_n - \sum_{k=1}^K a_{nk} \phi_k \right\|_2^2 + \lambda \sum_{n=1}^N \sum_{k=1}^K |a_{nk}|$$



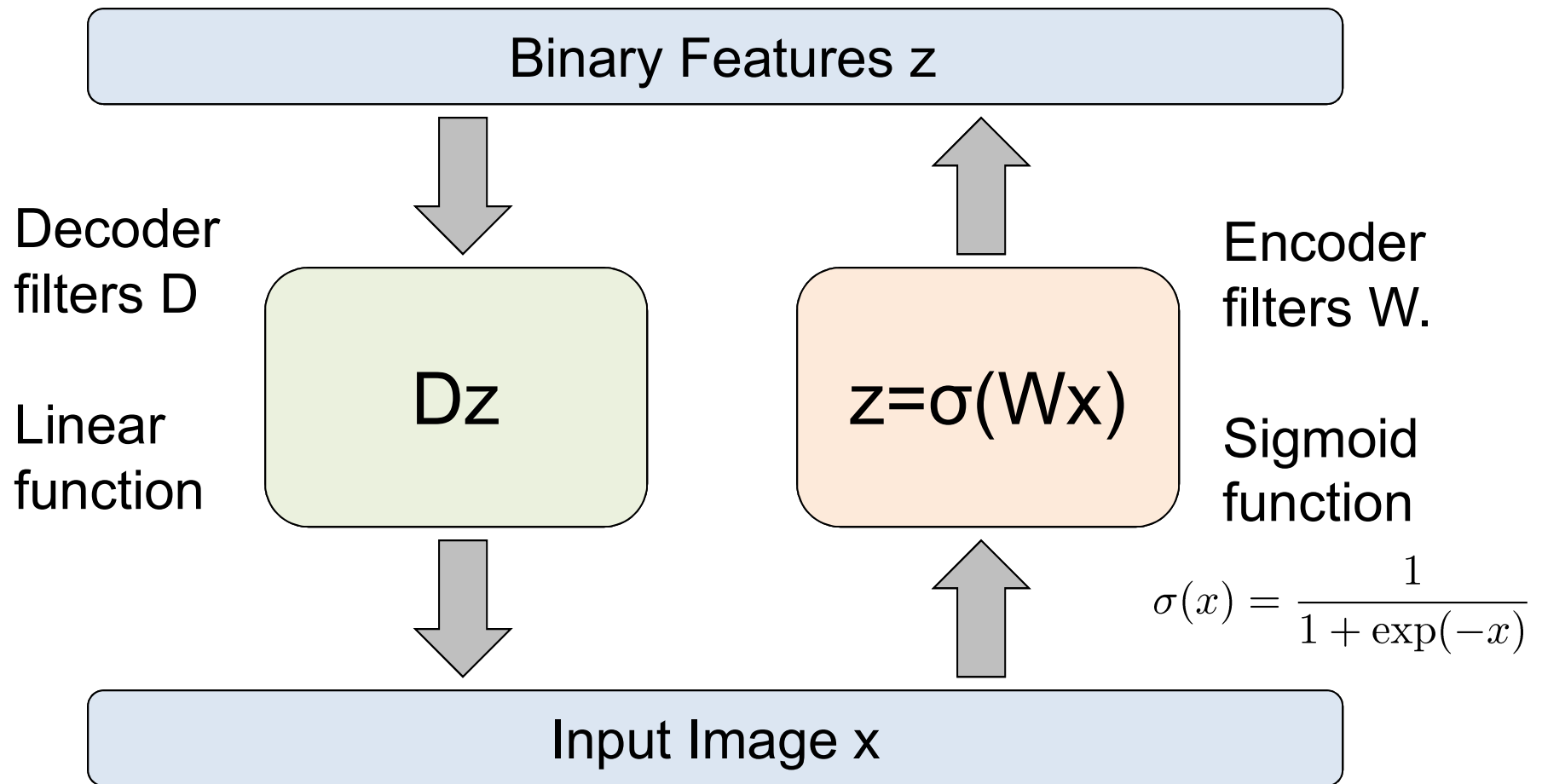
- Sparse, over-complete representation \mathbf{a} .
- Encoding $\mathbf{a} = f(\mathbf{x})$ is implicit and nonlinear function of \mathbf{x} .
- Reconstruction (or decoding) $\mathbf{x}' = g(\mathbf{a})$ is linear and explicit.

Autoencoder

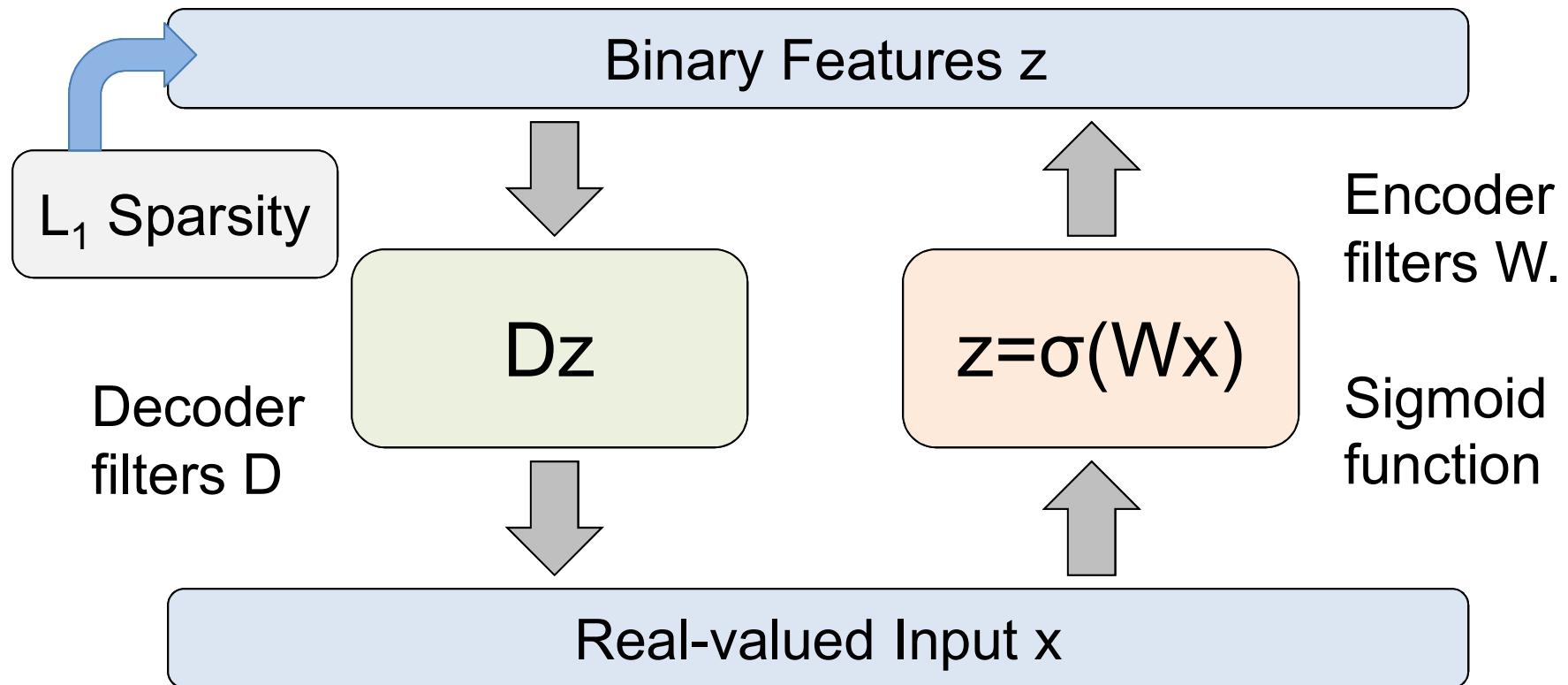


- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity.

Autoencoder



Predictive Sparse Decomposition



At training time

$$\min_{D, W, \mathbf{z}} \underbrace{\|D\mathbf{z} - \mathbf{x}\|_2^2 + \lambda \|\mathbf{z}\|_1}_{\text{Decoder}} + \underbrace{\|\sigma(W\mathbf{x}) - \mathbf{z}\|_2^2}_{\text{Encoder}}$$

Stacked Sparse Coding?

