

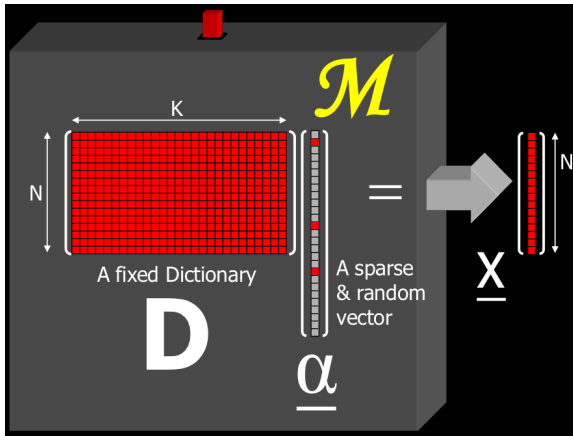
Sparse Coding: An Overview

Brian Booth

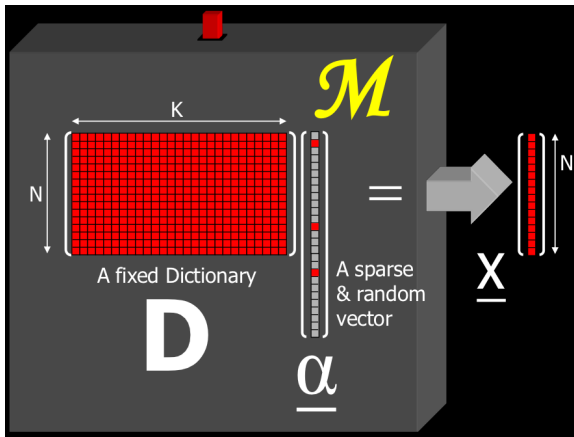
SFU Machine Learning Reading Group

November 12, 2013

The aim of sparse coding

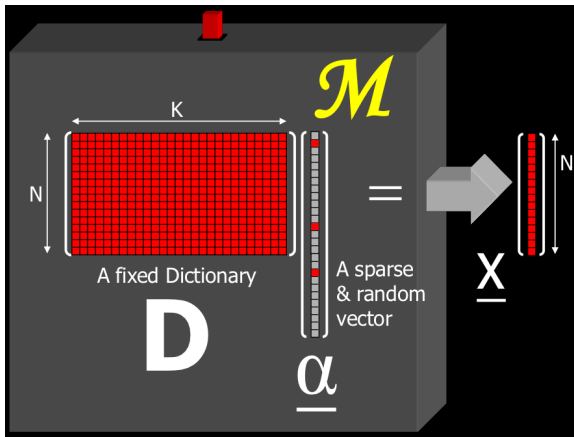


The aim of sparse coding



- Every column of D is a prototype

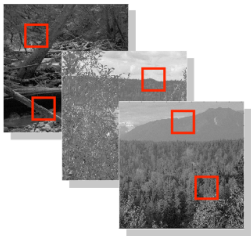
The aim of sparse coding



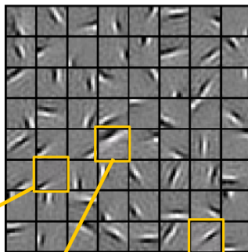
- Every column of D is a prototype
- Similar to, but more general than, PCA

Example: Sparse Coding of Images

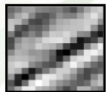
Natural Images



Learned bases ($\phi_1, \phi_2, \dots, \phi_{64}$): "Edges"



Test Example



\approx

$0.8 \times$



ϕ_{36}

$+ 0.3 \times$



ϕ_{42}

$+ 0.5 \times$

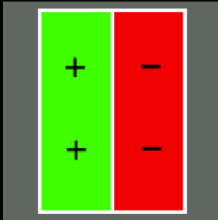


ϕ_{63}

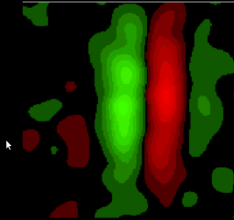
$$[\alpha_1, \dots, \alpha_{64}] = [0, \dots, 0.8, \dots, 0.3, \dots, 0.5, \dots, 0]$$

Sparse Coding in V1

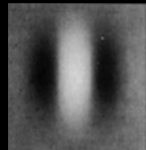
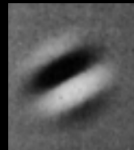
The first stage of visual processing in the brain (V1) does “edge detection.”



Schematic of simple cell

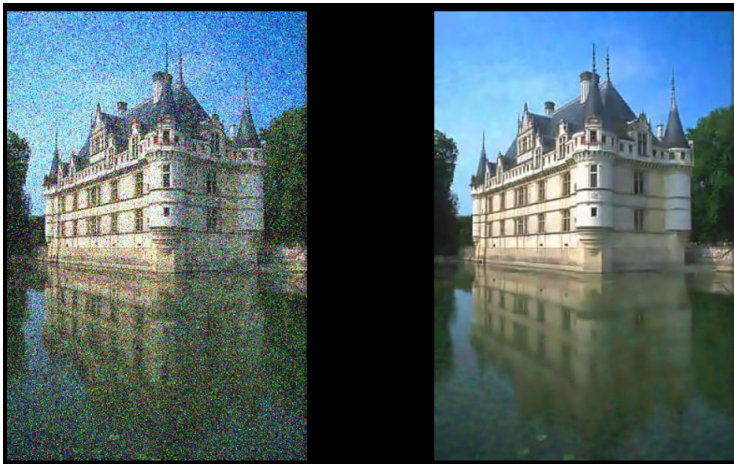


Actual simple cell

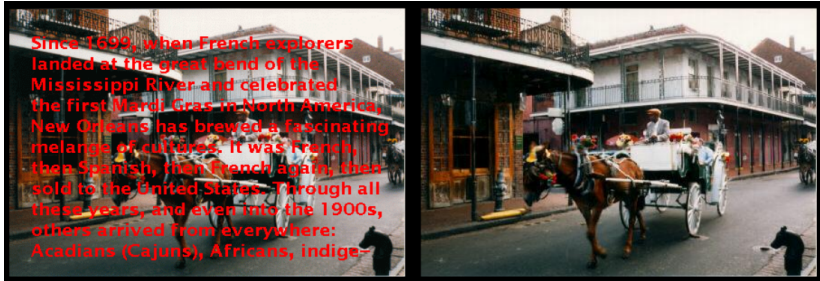


“Gabor functions.”

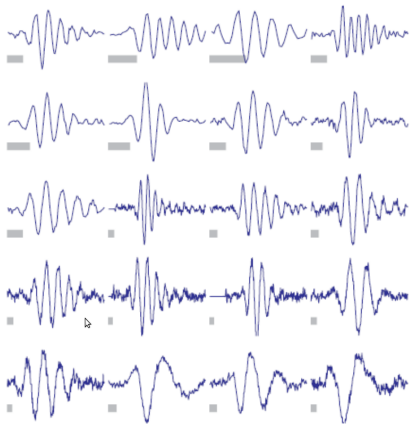
Example: Image Denoising



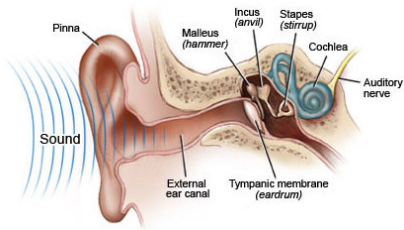
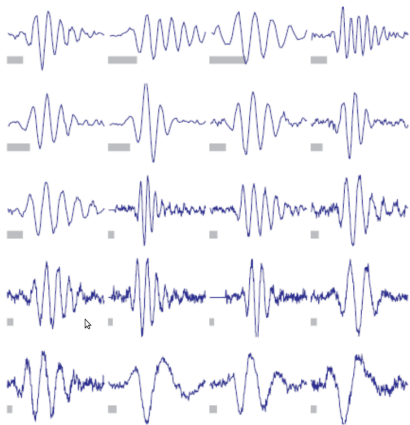
Example: Image Restoration



Sparse Coding and Acoustics



Sparse Coding and Acoustics



- Inner ear (cochlea) also does sparse coding of frequencies

Sparse Coding and Natural Language Processing

a	1	0.1	1.8	0.1
aardvark	0	.2	0	0
aardwolf	0	0	.1	0
able	0	0	0	0.8
account	1	.1	0	0
acid	0	.1	.1	0
across	0	0	0	.1
...	0	0	.2	.2
baby	1	1.2	0	.1
back	0	0	0	.2
...	0	0	.2	0
cradle	1	1	.1	.3
...	0	0	0	0
...	0	.1	.1	0
...	0	0	.1	0
zylophone	1	.1	0	1.2

$\approx \quad 0.7 \times \quad 0.4 \times \quad 0.1 \times$

$\phi_{36} \quad \phi_{42} \quad \phi_{63}$

$$\text{Document} \approx 0.7 \times \text{Topic36} + 0.4 \times \text{Topic42} + 0.1 \times \text{Topic63}$$

Outline

- 1 Introduction: Why Sparse Coding?

Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge

Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 Conclusions

Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 Conclusions

Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 Conclusions

The aim of sparse coding, revisited

We assume our data \mathbf{x} satisfies

$$\mathbf{x} \approx \sum_{i=1}^n \alpha_i \mathbf{d}_i = \alpha \mathbf{D}$$

The aim of sparse coding, revisited

We assume our data \mathbf{x} satisfies

$$\mathbf{x} \approx \sum_{i=1}^n \alpha_i \mathbf{d}_i = \alpha \mathbf{D}$$

Learning:

- Given training data $\mathbf{x}^j, j \in \{1, \dots, m\}$
- Learn dictionary \mathbf{D} and sparse code α

The aim of sparse coding, revisited

We assume our data \mathbf{x} satisfies

$$\mathbf{x} \approx \sum_{i=1}^n \alpha_i \mathbf{d}_i = \alpha \mathbf{D}$$

Learning:

- Given training data $\mathbf{x}^j, j \in \{1, \dots, m\}$
- Learn dictionary \mathbf{D} and sparse code α

Encoding:

- Given test data \mathbf{x} , dictionary \mathbf{D}
- Learn sparse code α

Learning: The Objective Function

Dictionary learning involves optimizing:

$$\arg \min_{\{\mathbf{d}_i\}, \{\alpha^j\}} \sum_{j=1}^m \|\mathbf{x}^j - \sum_{i=1}^n \alpha_i^j \mathbf{d}_i\|^2$$

Learning: The Objective Function

Dictionary learning involves optimizing:

$$\arg \min_{\{\mathbf{d}_i\}, \{\alpha^j\}} \sum_{j=1}^m \|\mathbf{x}^j - \sum_{i=1}^n \alpha_i^j \mathbf{d}_i\|^2 + \beta \sum_{j=1}^m \sum_{i=1}^n |\alpha_i^j|$$

Learning: The Objective Function

Dictionary learning involves optimizing:

$$\arg \min_{\{\mathbf{d}_i\}, \{\alpha^j\}} \sum_{j=1}^m \left\| \mathbf{x}^j - \sum_{i=1}^n \alpha_i^j \mathbf{d}_i \right\|^2 + \beta \sum_{j=1}^m \sum_{i=1}^n |\alpha_i^j|$$

subject to $\|\mathbf{d}_i\|^2 \leq c, \quad \forall i = 1, \dots, n.$

Learning: The Objective Function

Dictionary learning involves optimizing:

$$\arg \min_{\{\mathbf{d}_i\}, \{\alpha_i^j\}} \sum_{j=1}^m \left\| \mathbf{x}^j - \sum_{i=1}^n \alpha_i^j \mathbf{d}_i \right\|^2 + \beta \sum_{j=1}^m \sum_{i=1}^n |\alpha_i^j|$$

subject to $\|\mathbf{d}_i\|^2 \leq c, \quad \forall i = 1, \dots, n.$

In matrix notation:

$$\arg \min_{\mathbf{D}, \mathbf{A}} \|\mathbf{X} - \mathbf{AD}\|_F^2 + \beta \sum_{i,j} |\alpha_{i,j}|$$

subject to $\sum_i \mathbf{D}_{i,j}^2 \leq c, \quad \forall i = 1, \dots, n.$

Learning: The Objective Function

Dictionary learning involves optimizing:

$$\arg \min_{\{\mathbf{d}_i\}, \{\alpha_i^j\}} \sum_{j=1}^m \left\| \mathbf{x}^j - \sum_{i=1}^n \alpha_i^j \mathbf{d}_i \right\|^2 + \beta \sum_{j=1}^m \sum_{i=1}^n |\alpha_i^j|$$

subject to $\|\mathbf{d}_i\|^2 \leq c, \quad \forall i = 1, \dots, n.$

In matrix notation:

$$\arg \min_{\mathbf{D}, \mathbf{A}} \|\mathbf{X} - \mathbf{AD}\|_F^2 + \beta \sum_{i,j} |\alpha_{i,j}|$$

subject to $\sum_i \mathbf{D}_{i,j}^2 \leq c, \quad \forall i = 1, \dots, n.$

Split the optimization over \mathbf{D} and \mathbf{A} in two.

Step 1: Learning the Dictionary

Reduced optimization problem:

$$\begin{aligned} & \arg \min_{\mathbf{D}} \|\mathbf{X} - \mathbf{AD}\|_F^2 \\ & \text{subject to } \sum_i \mathbf{D}_{i,j}^2 \leq c, \quad \forall i = 1, \dots, n. \end{aligned}$$

Step 1: Learning the Dictionary

Reduced optimization problem:

$$\begin{aligned} & \arg \min_{\mathbf{D}} \|\mathbf{X} - \mathbf{AD}\|_F^2 \\ & \text{subject to } \sum_i \mathbf{D}_{i,j}^2 \leq c, \quad \forall i = 1, \dots, n. \end{aligned}$$

Introduce Lagrange multipliers:

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{AD})^T (\mathbf{X} - \mathbf{AD}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j}^2 - c \right)$$

Step 1: Learning the Dictionary

Reduced optimization problem:

$$\begin{aligned} & \arg \min_{\mathbf{D}} \|\mathbf{X} - \mathbf{AD}\|_F^2 \\ & \text{subject to } \sum_i \mathbf{D}_{i,j}^2 \leq c, \quad \forall i = 1, \dots, n. \end{aligned}$$

Introduce Lagrange multipliers:

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{AD})^T (\mathbf{X} - \mathbf{AD}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j} - c \right)$$

where each $\lambda_j \geq 0$ is a dual variable...

Step 1: Moving to the dual

From the Lagrangian

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{AD})^T (\mathbf{X} - \mathbf{AD}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j}^2 - c \right)$$

Step 1: Moving to the dual

From the Lagrangian

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{AD})^T (\mathbf{X} - \mathbf{AD}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j}^2 - c \right)$$

minimize over \mathbf{D} to obtain Lagrange dual

$$\mathbf{D}(\lambda) = \min_{\mathbf{D}} \mathcal{L}(\mathbf{D}, \lambda) =$$

Step 1: Moving to the dual

From the Lagrangian

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{AD})^T (\mathbf{X} - \mathbf{AD}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j}^2 - c \right)$$

minimize over \mathbf{D} to obtain Lagrange dual

$$\mathbf{D}(\lambda) = \min_{\mathbf{D}} \mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left(\mathbf{X}^T \mathbf{X} - \mathbf{XA}^T (\mathbf{AA}^T + \Lambda)^{-1} (\mathbf{XA}^T)^T - c\Lambda \right)$$

Step 1: Moving to the dual

From the Lagrangian

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{AD})^T (\mathbf{X} - \mathbf{AD}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j}^2 - c \right)$$

minimize over \mathbf{D} to obtain Lagrange dual

$$\mathbf{D}(\lambda) = \min_{\mathbf{D}} \mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left(\mathbf{X}^T \mathbf{X} - \mathbf{XA}^T (\mathbf{AA}^T + \Lambda)^{-1} (\mathbf{XA}^T)^T - c\Lambda \right)$$

- The dual can be optimized using conjugate gradient

Step 1: Moving to the dual

From the Lagrangian

$$\mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left((\mathbf{X} - \mathbf{A}\mathbf{D})^T (\mathbf{X} - \mathbf{A}\mathbf{D}) \right) + \sum_{j=1}^n \lambda_j \left(\sum_i \mathbf{D}_{i,j}^2 - c \right)$$

minimize over \mathbf{D} to obtain Lagrange dual

$$\mathbf{D}(\lambda) = \min_{\mathbf{D}} \mathcal{L}(\mathbf{D}, \lambda) = \text{tr} \left(\mathbf{X}^T \mathbf{X} - \mathbf{X}\mathbf{A}^T \left(\mathbf{A}\mathbf{A}^T + \Lambda \right)^{-1} \left(\mathbf{X}\mathbf{A}^T \right)^T - c\Lambda \right)$$

- The dual can be optimized using conjugate gradient
- Only n, λ values compared to \mathbf{D} being $n \times k$

Step 1: Dual to the Dictionary

With the optimal Λ , our dictionary is

$$\mathbf{D}^T = (\mathbf{A}\mathbf{A}^T + \Lambda)^{-1} (\mathbf{X}\mathbf{A}^T)^T$$

Step 1: Dual to the Dictionary

With the optimal Λ , our dictionary is

$$\mathbf{D}^T = (\mathbf{A}\mathbf{A}^T + \Lambda)^{-1} (\mathbf{X}\mathbf{A}^T)^T$$

Key point: Moving to the dual reduces the number of optimization variables, speeding up the optimization.

Step 2: Learning the Sparse Code

With \mathbf{D} now fixed, optimize for \mathbf{A}

$$\arg \min_{\mathbf{A}} \|\mathbf{X} - \mathbf{AD}\|_F^2 + \beta \sum_{i,j} |\alpha_{i,j}|$$

Step 2: Learning the Sparse Code

With \mathbf{D} now fixed, optimize for \mathbf{A}

$$\arg \min_{\mathbf{A}} \|\mathbf{X} - \mathbf{AD}\|_F^2 + \beta \sum_{i,j} |\alpha_{i,j}|$$

- Unconstrained, convex quadratic optimization

Step 2: Learning the Sparse Code

With **D** now fixed, optimize for **A**

$$\arg \min_{\mathbf{A}} \|\mathbf{X} - \mathbf{AD}\|_F^2 + \beta \sum_{i,j} |\alpha_{i,j}|$$

- Unconstrained, convex quadratic optimization
 - Many solvers for this (e.g. interior point methods, in-crowd algorithm, fixed-point continuation)
-

Step 2: Learning the Sparse Code

With \mathbf{D} now fixed, optimize for \mathbf{A}

$$\arg \min_{\mathbf{A}} \|\mathbf{X} - \mathbf{AD}\|_F^2 + \beta \sum_{i,j} |\alpha_{i,j}|$$

- Unconstrained, convex quadratic optimization
- Many solvers for this (e.g. interior point methods, in-crowd algorithm, fixed-point continuation)

Note:

- Same problem as the encoding problem.
- Runtime of optimization in the encoding stage?

Speeding up the testing phase

Fair amount of work on speeding up the encoding stage:

- H. Lee et al., *Efficient sparse coding algorithms*
<http://ai.stanford.edu/~hllee/nips06-sparsecoding.pdf>
- K. Gregor and Y. LeCun, *Learning Fast Approximations of Sparse Coding*
<http://yann.lecun.com/exdb/publis/pdf/gregor-icml-10.pdf>
- S. Hawe et al., *Separable Dictionary Learning*
<http://arxiv.org/pdf/1303.5244v1.pdf>

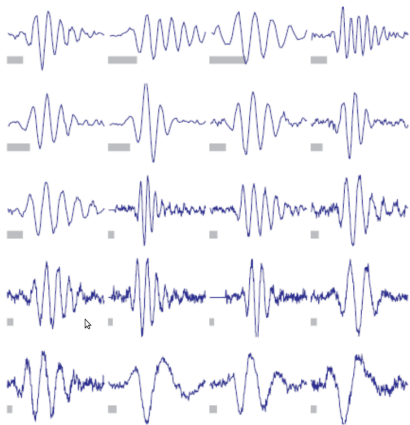
Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 Conclusions

Outline

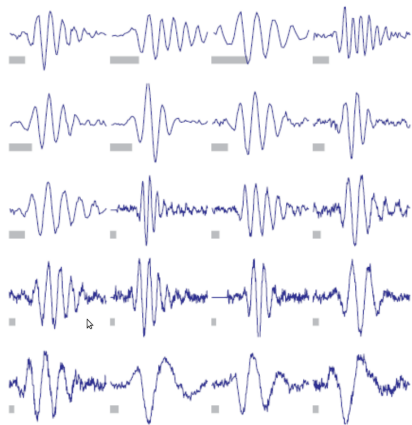
- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 Conclusions

Relationships between Dictionary atoms



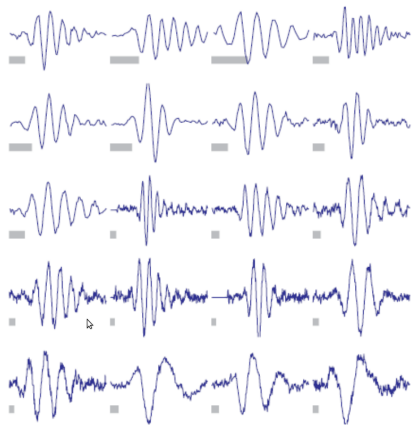
- Dictionaries are over-complete bases

Relationships between Dictionary atoms



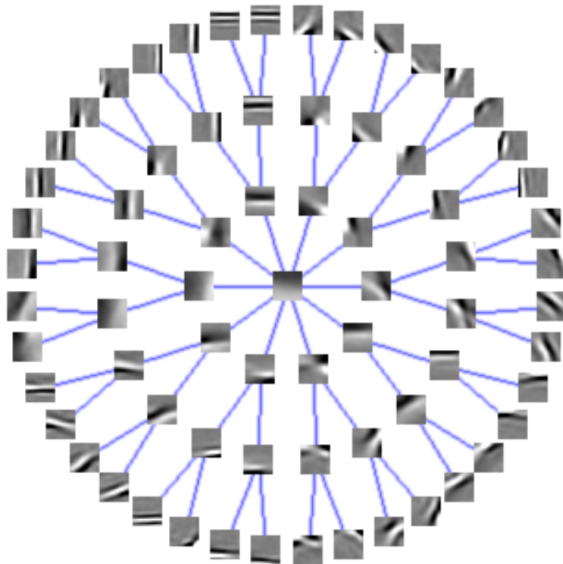
- Dictionaries are over-complete bases
- Dictate relationships between atoms

Relationships between Dictionary atoms

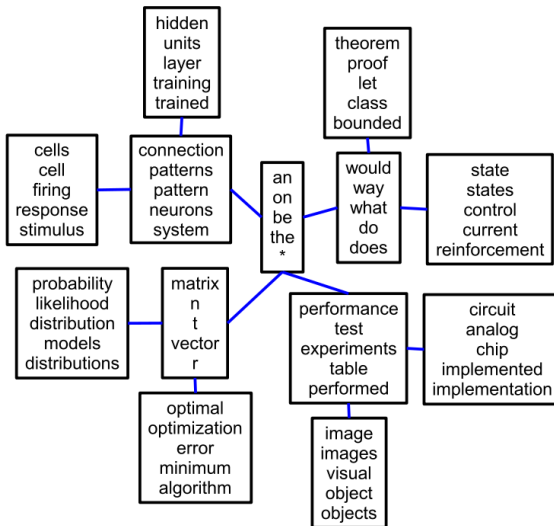


- Dictionaries are over-complete bases
- Dictate relationships between atoms
- Example: Hierarchical dictionaries

Example: Image Patches



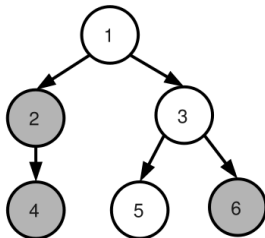
Example: Document Topics



Problem Statement

Goal:

- Have sub-groups of sparse code α all be non-zero (or zero).



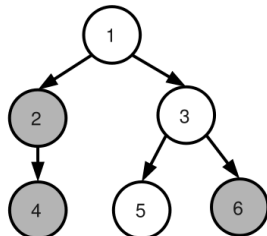
Problem Statement

Goal:

- Have sub-groups of sparse code α all be non-zero (or zero).

Hierarchical:

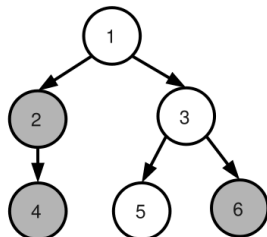
- If a node is non-zero, its parent must be non-zero
- If a node's parent is zero, the node must be zero



Problem Statement

Goal:

- Have sub-groups of sparse code α all be non-zero (or zero).



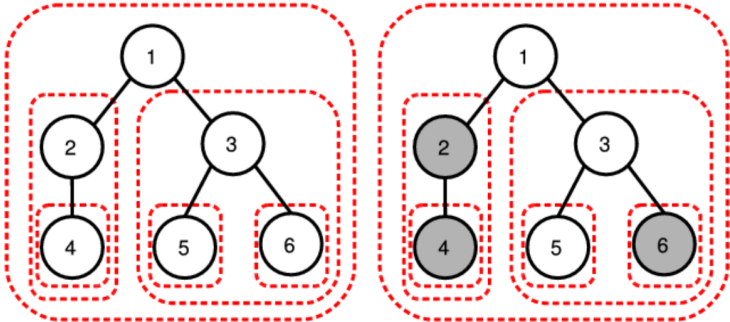
Hierarchical:

- If a node is non-zero, its parent must be non-zero
- If a node's parent is zero, the node must be zero

Implementation:

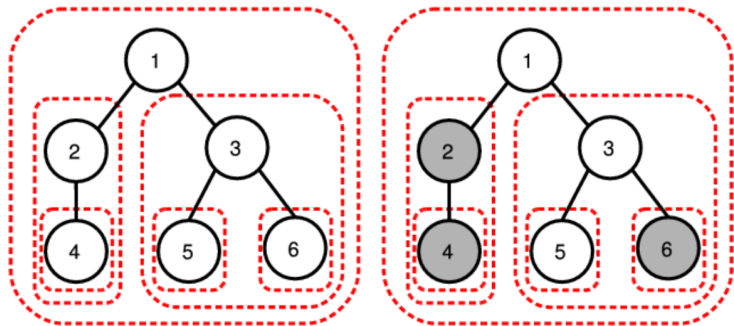
- Change the regularization
- Enforce sparsity differently...

Grouping Code Entries



- Level k included in $k + 1$ groups

Grouping Code Entries



- Level k included in $k + 1$ groups
- Add $|\alpha_j|$ to objective function once for each group

Group Regularization

Updated objective function:

$$\arg \min_{\mathbf{D}, \{\alpha^j\}} \sum_{j=1}^m \left[\|\mathbf{x}^j - \mathbf{D}\alpha^j\|^2 \right]$$

Group Regularization

Updated objective function:

$$\arg \min_{\mathbf{D}, \{\alpha^j\}} \sum_{j=1}^m \left[\|\mathbf{x}^j - \mathbf{D}\alpha^j\|^2 + \beta \Omega(\alpha^j) \right]$$

Group Regularization

Updated objective function:

$$\arg \min_{\mathbf{D}, \{\alpha^j\}} \sum_{j=1}^m \left[\|\mathbf{x}^j - \mathbf{D}\alpha^j\|^2 + \beta \Omega(\alpha^j) \right]$$

where

$$\Omega(\alpha) = \sum_{g \in \mathcal{P}} w_g \|\alpha_{|g}\|$$

Group Regularization

Updated objective function:

$$\arg \min_{\mathbf{D}, \{\alpha^j\}} \sum_{j=1}^m \left[\|\mathbf{x}^j - \mathbf{D}\alpha^j\|^2 + \beta \Omega(\alpha^j) \right]$$

where

$$\Omega(\alpha) = \sum_{g \in \mathcal{P}} w_g \|\alpha_{|g}\|$$

- $\alpha_{|g}$ are the code values for group g .

Group Regularization

Updated objective function:

$$\arg \min_{\mathbf{D}, \{\alpha^j\}} \sum_{j=1}^m \left[\|\mathbf{x}^j - \mathbf{D}\alpha^j\|^2 + \beta \Omega(\alpha^j) \right]$$

where

$$\Omega(\alpha) = \sum_{g \in \mathcal{P}} w_g \|\alpha_{|g}\|$$

- $\alpha_{|g}$ are the code values for group g .
- w_g weights the enforcement of the hierarchy

Group Regularization

Updated objective function:

$$\arg \min_{\mathbf{D}, \{\alpha^j\}} \sum_{j=1}^m \left[\|\mathbf{x}^j - \mathbf{D}\alpha^j\|^2 + \beta \Omega(\alpha^j) \right]$$

where

$$\Omega(\alpha) = \sum_{g \in \mathcal{P}} w_g \|\alpha_{|g}\|$$

- $\alpha_{|g}$ are the code values for group g .
- w_g weights the enforcement of the hierarchy
- Solve using proximal methods.

Other Examples

Other examples of structured sparsity:

- M. Stojnic et al., *On the Reconstruction of Block-Sparse Signals With an Optimal Number of Measurements*, <http://dx.doi.org/10.1109/TSP.2009.2020754>
- J. Mairal et al., *Convex and Network Flow Optimization for Structured Sparsity*, <http://jmlr.org/papers/volume12/mairal11a/mairal11a.pdf>

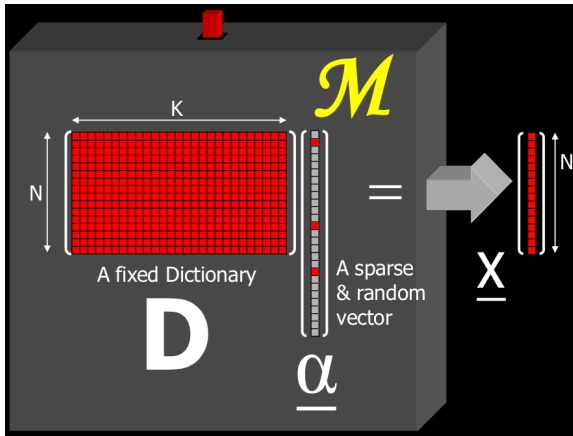
Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 Conclusions

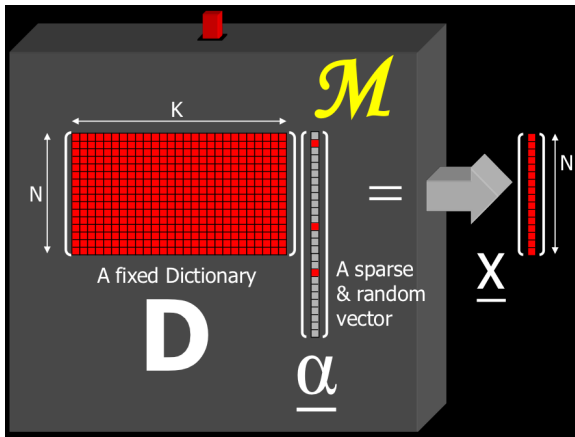
Outline

- 1 Introduction: Why Sparse Coding?
- 2 Sparse Coding: The Basics
- 3 Adding Prior Knowledge
- 4 **Conclusions**

Summary

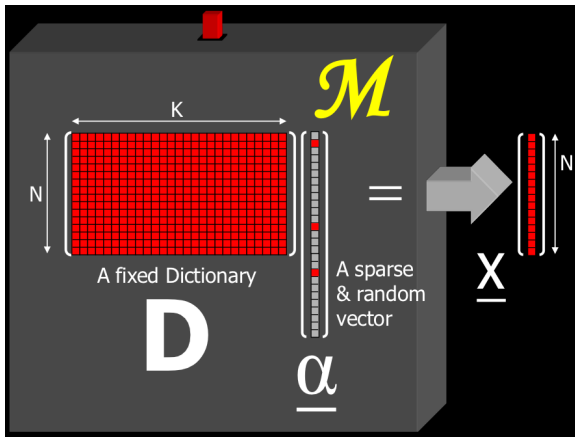


Summary



Two interesting directions:

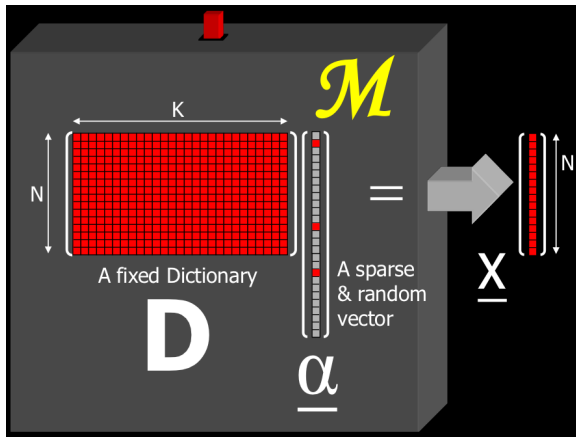
Summary



Two interesting directions:

- Increasing speed of the testing phase

Summary



Two interesting directions:

- Increasing speed of the testing phase
- Optimizing dictionary structure