

# Machine Learning for Signal Processing

## Sparse and Overcomplete Representations

Bhiksha Raj  
(slides from Sourish Chaudhuri)

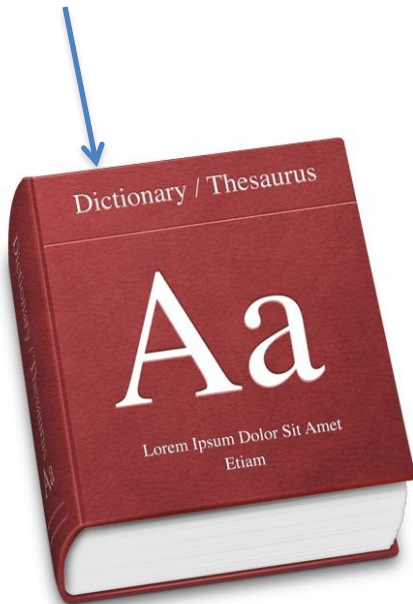
Oct 22, 2015

# Key Topics in this Lecture

- Basics – Component-based representations
  - Overcomplete and Sparse Representations,
  - Dictionaries
- Pursuit Algorithms
- How to learn a dictionary
- Why is an overcomplete representation powerful?

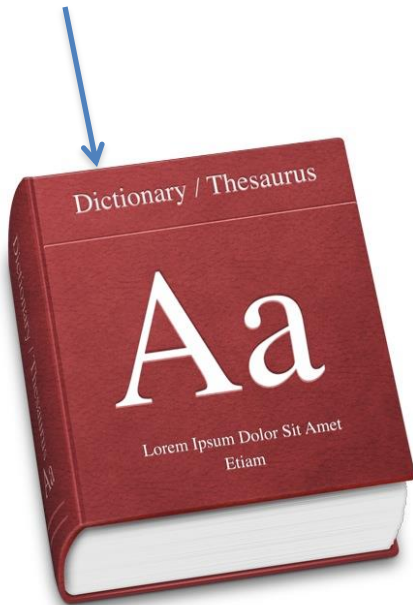
# Representing Data

## Dictionary (codebook)

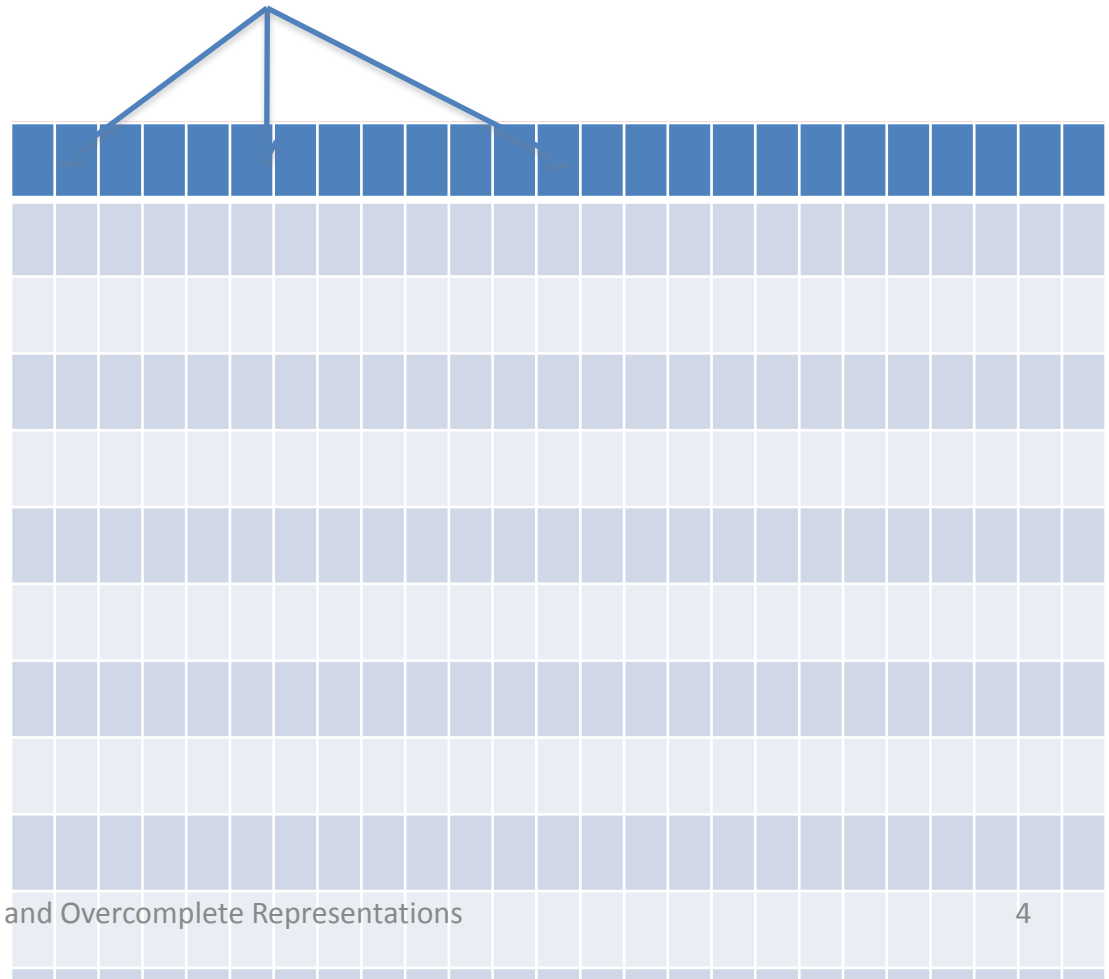


# Representing Data

Dictionary

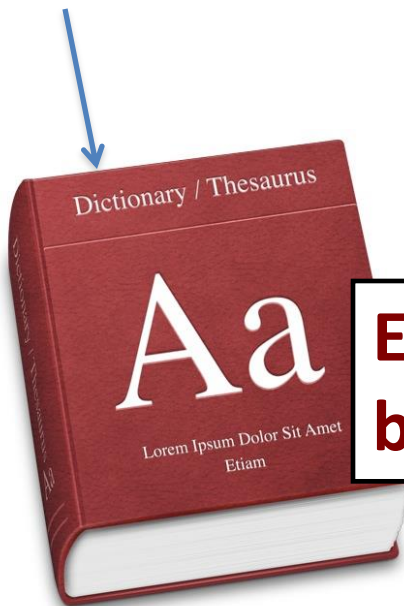


Atoms

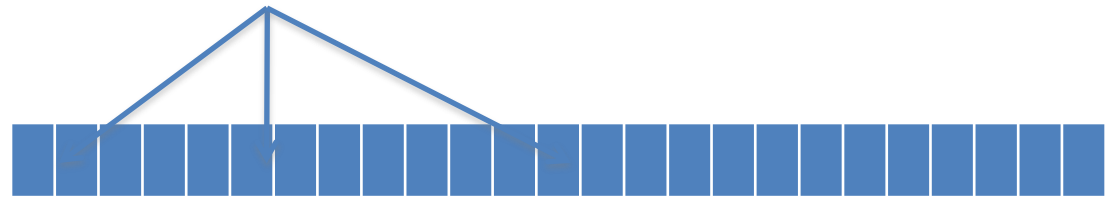


# Representing Data

Dictionary



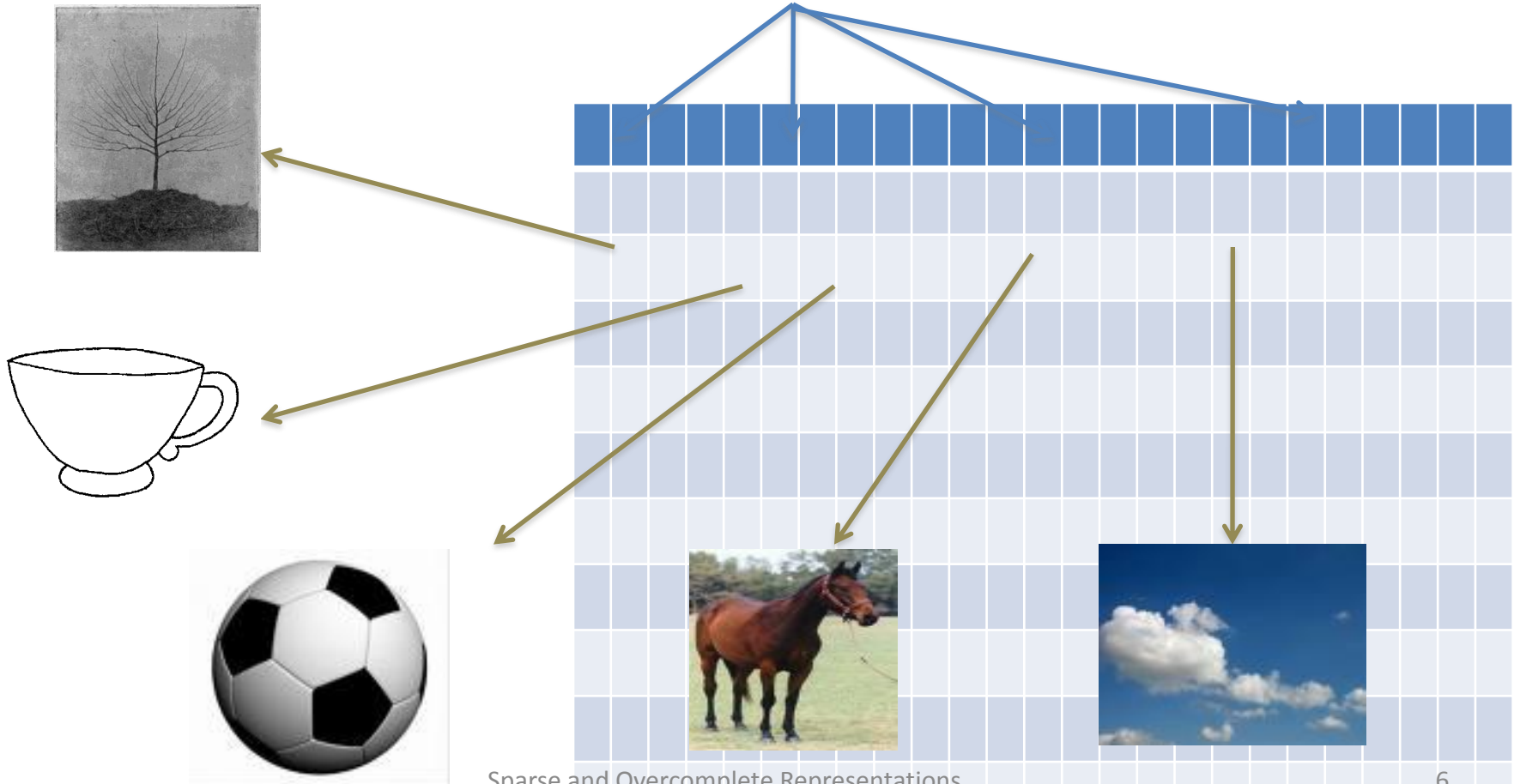
Atoms



Each atom is a basic unit that can be used to “compose” larger units.

# Representing Data

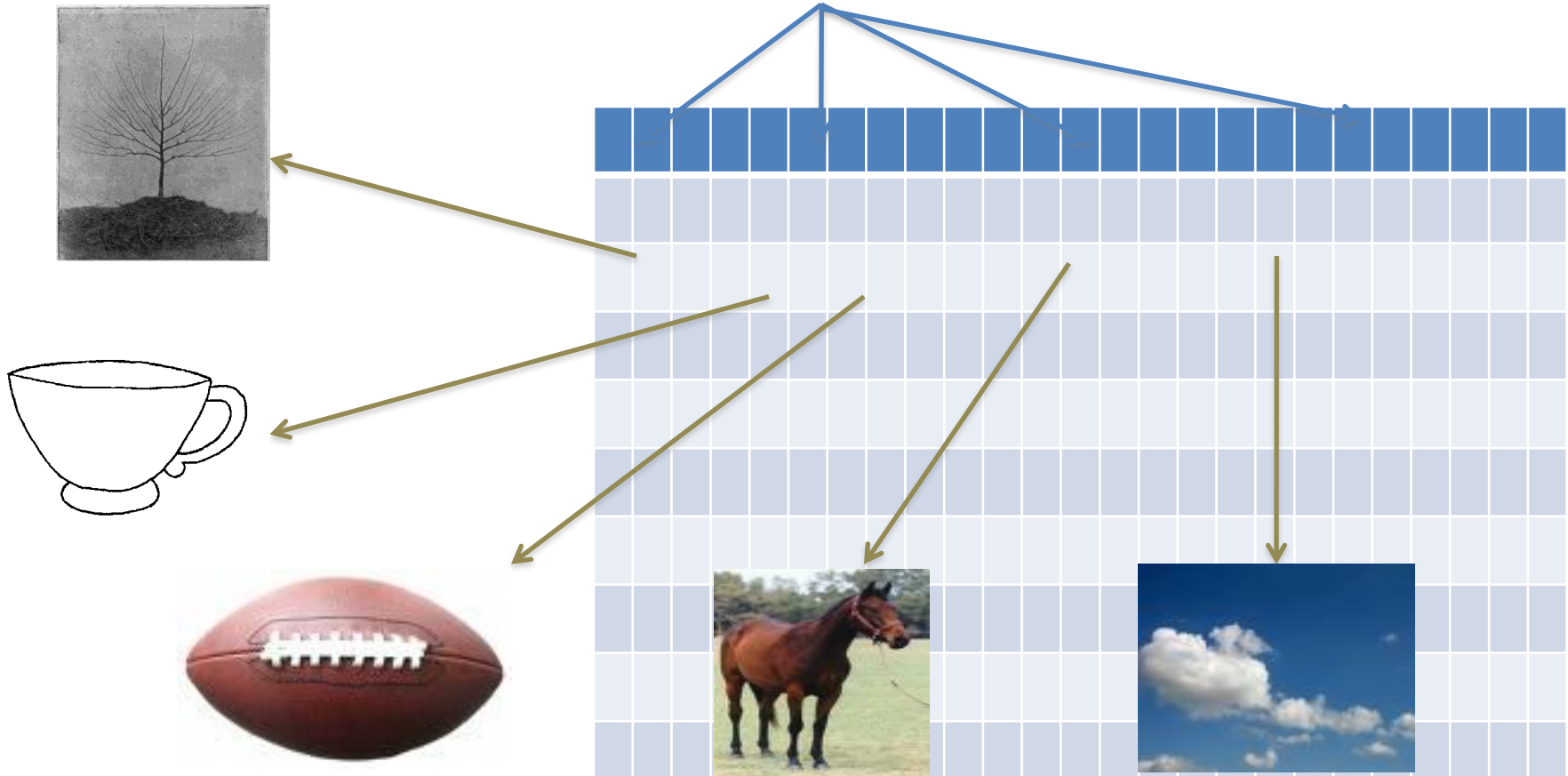
Atoms



Sparse and Overcomplete Representations

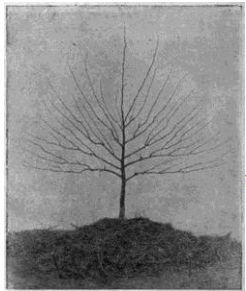
# Representing Data

Atoms



# Representing Data

Atoms



**Many such bases  
(concepts)**





# Representing Data



sparse and Overcomplete Representations

# Representing Data

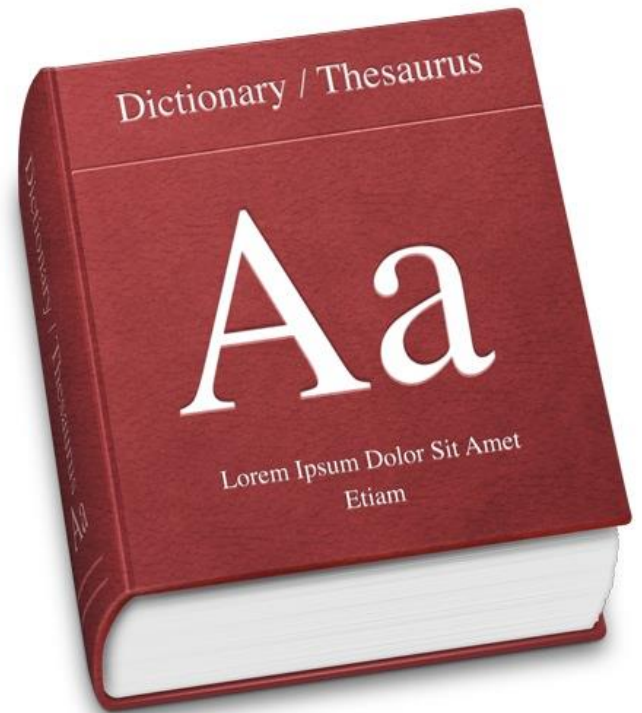


**Using concepts that we know...**

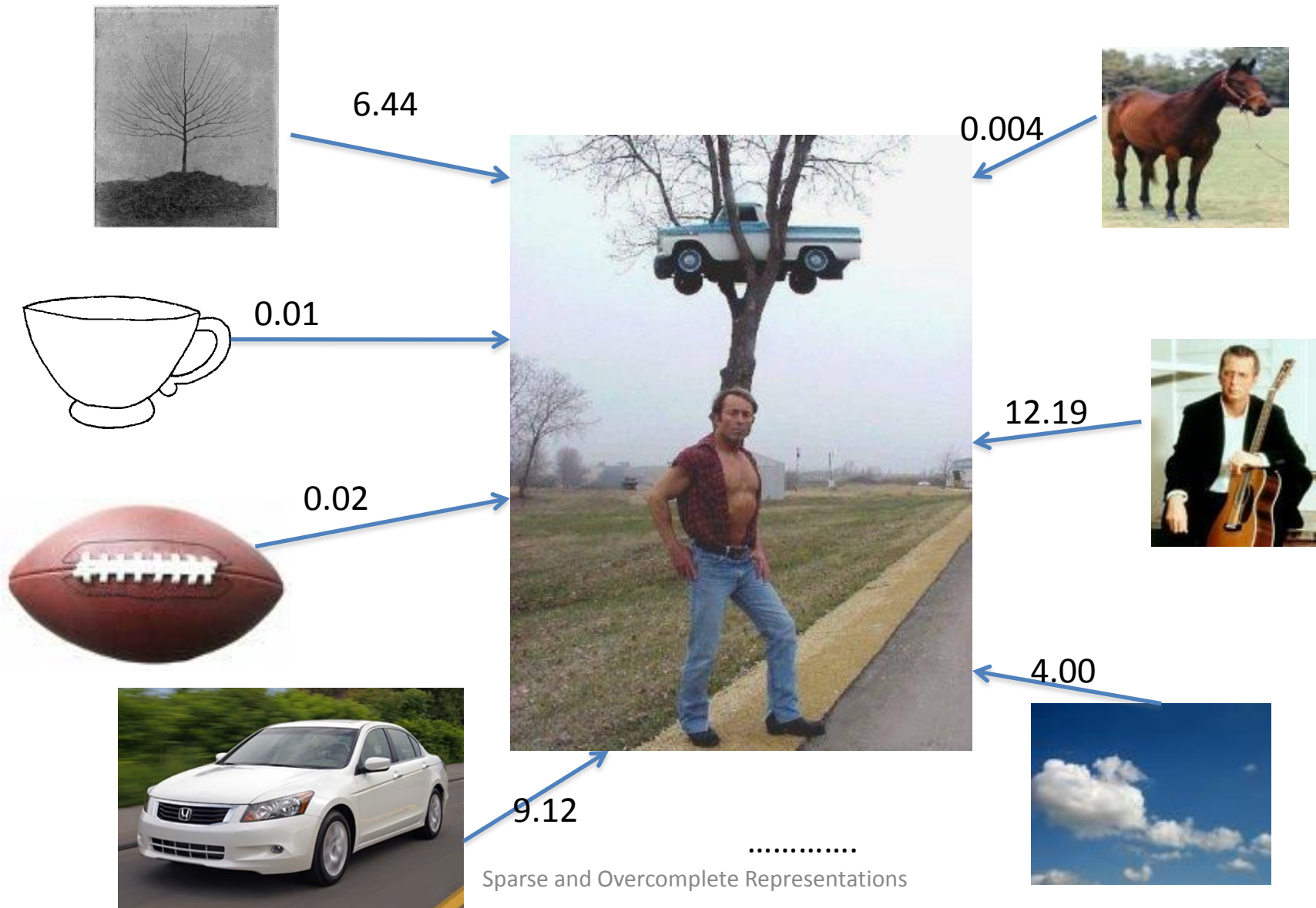
# Representing Data



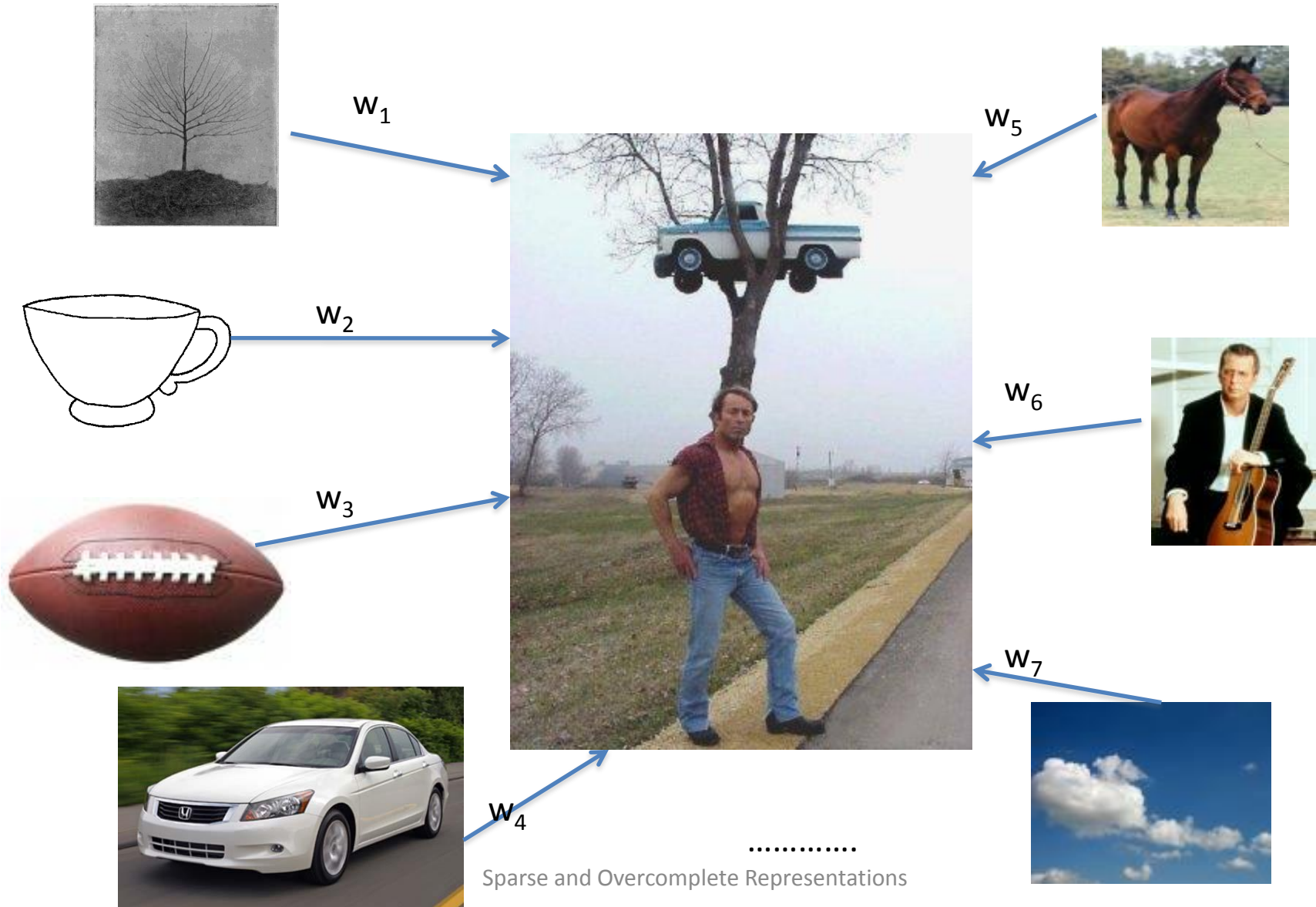
Using concepts that we know...



# Representing Data

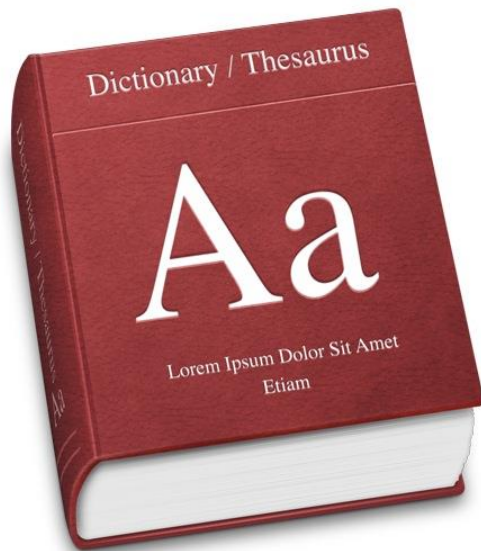


# Representing Data



.....  
Sparse and Overcomplete Representations

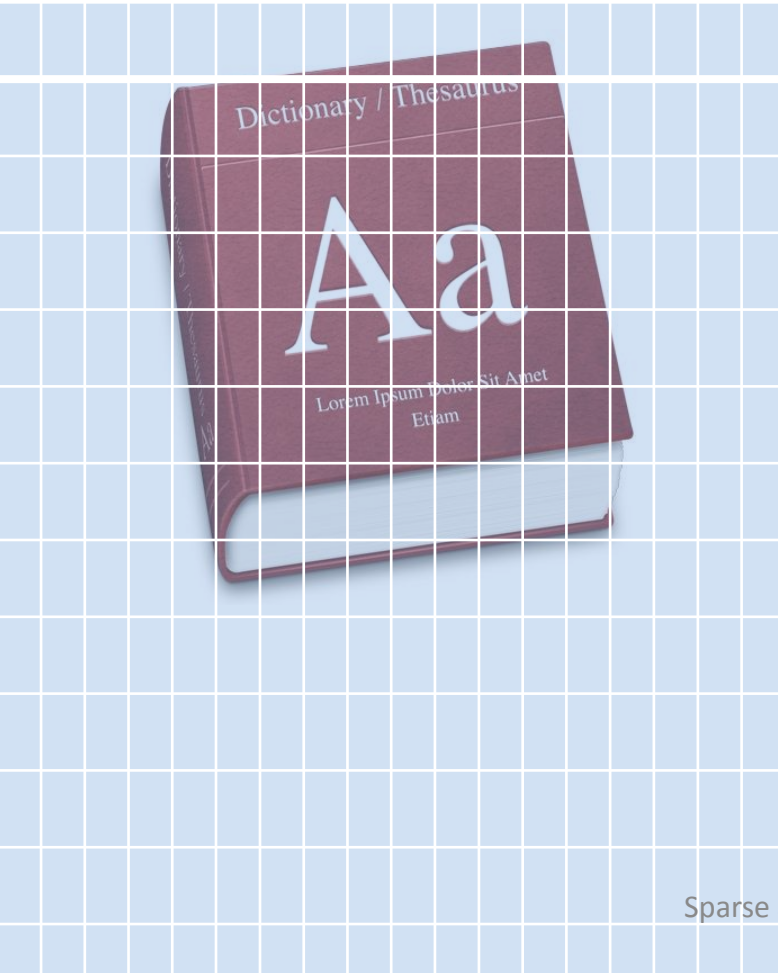
# Representing Data



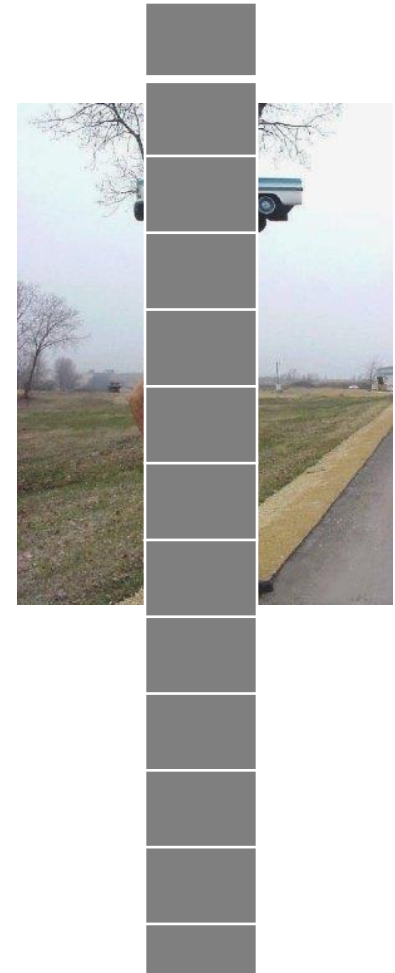
=



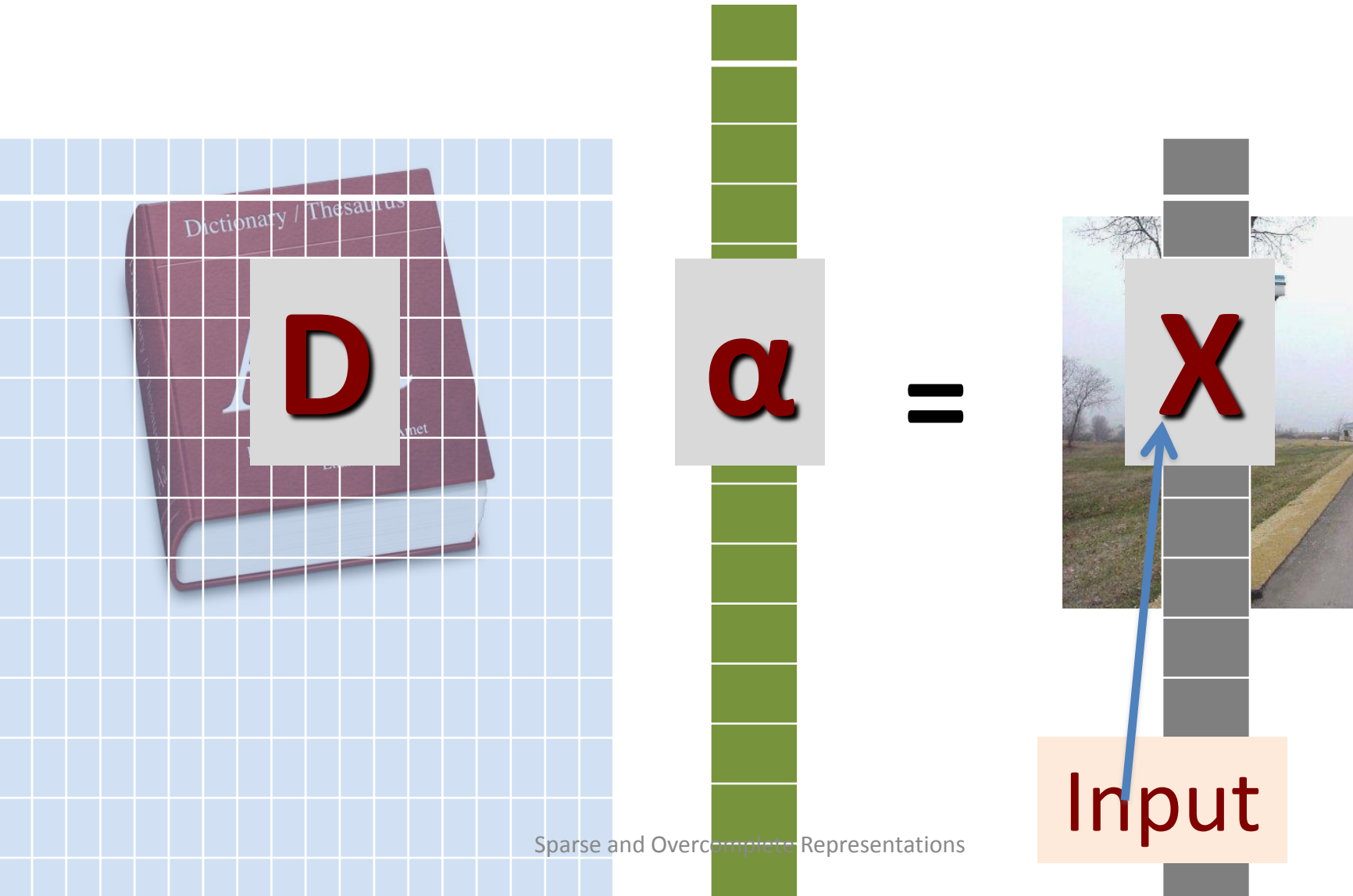
# Representing Data



=

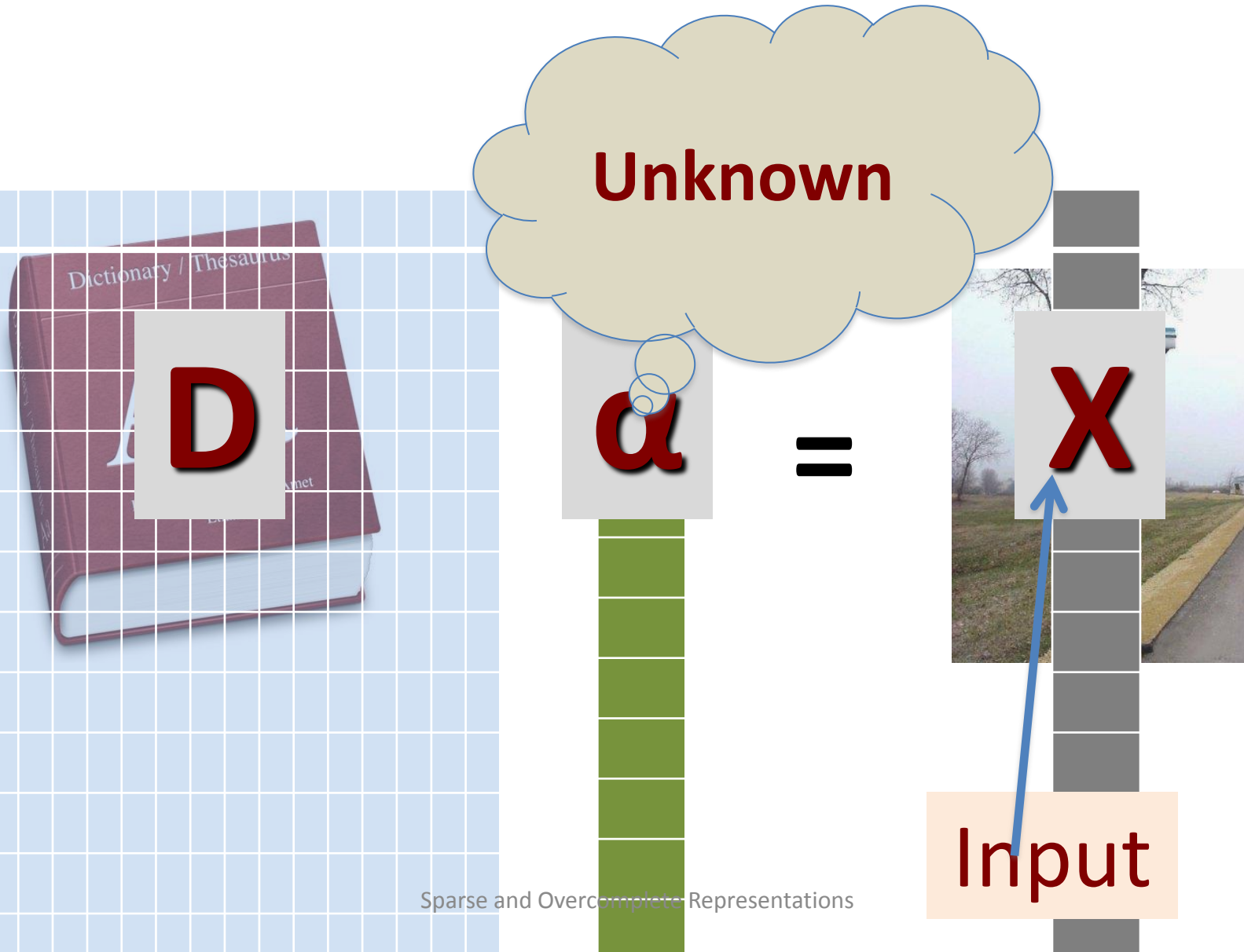


# Representing Data



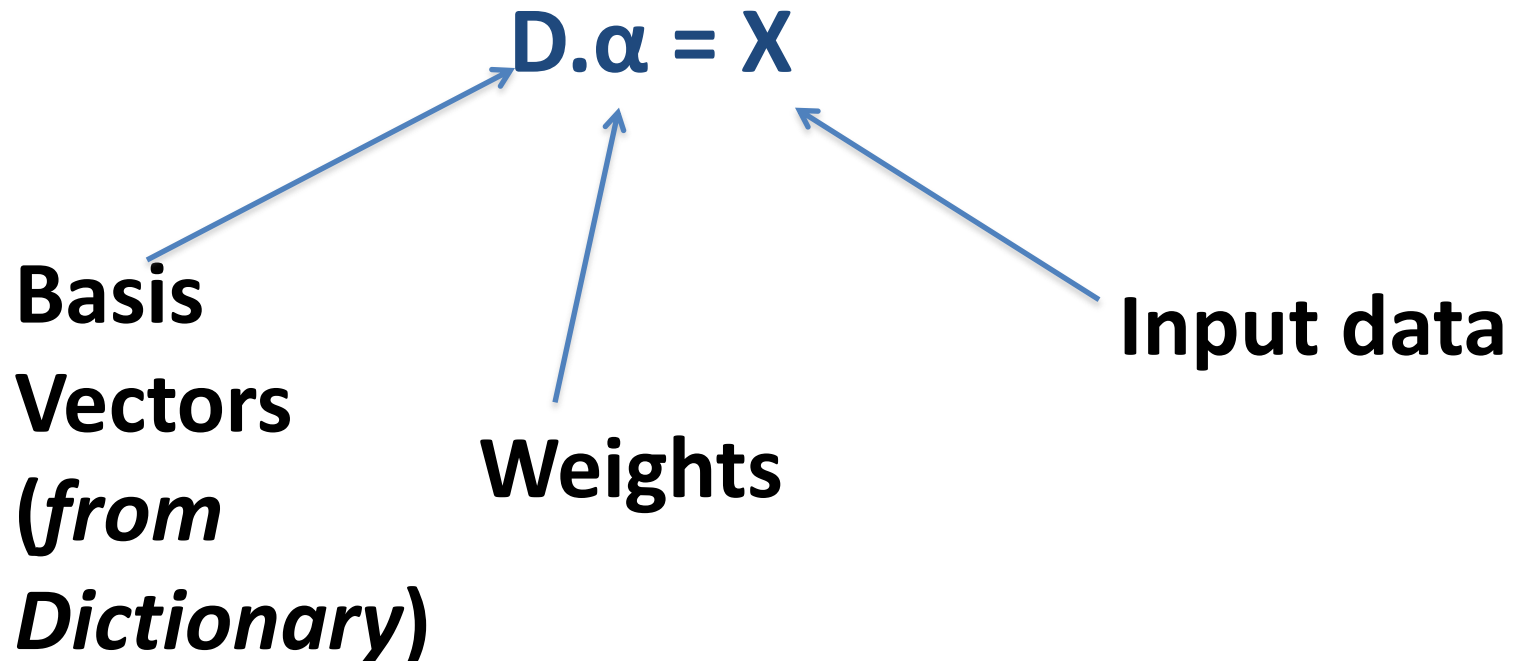


# Representing Data



# Quick Linear Algebra Refresher

- Remember, #(Basis Vectors)= #unknowns



# Overcomplete Representations

- What is the dimensionality of the input image? (say 64x64 image)
  - **4096**
- What is the dimensionality of the dictionary? (each image = 64x64 pixels)
  - **4096 x N**

# Overcomplete Representations

- What is the dimensionality of the input image? (say 64x64 image)

➤ **4096**

- What is the dimensionality of the dictionary? (each image = 64x64 pixels)

➤ **4096 x N**

???

# Overcomplete Representations

- What is the dimensionality of the input image? (say 64x64 image)

➤ **4096**

- What is the dimensionality of the dictionary? (each image = 64x64 pixels)

➤ **4096 x N**

**VERY LARGE!!!**

# Overcomplete Representations

- What is the dimensionality of the input image? (say  $64 \times 64$  image)

If  $N > 4096$  (as it likely is)

we have an **overcomplete** representation

- What is the dimensionality of the dictionary? (each image =  $64 \times 64$  pixels)

➤  $4096 \times N$

**VERY LARGE!!!**

# Overcomplete Representations

- What is the dimensionality of the input image? (say 64x64 image)

More generally:

If #(basis vectors) > dimensions of input

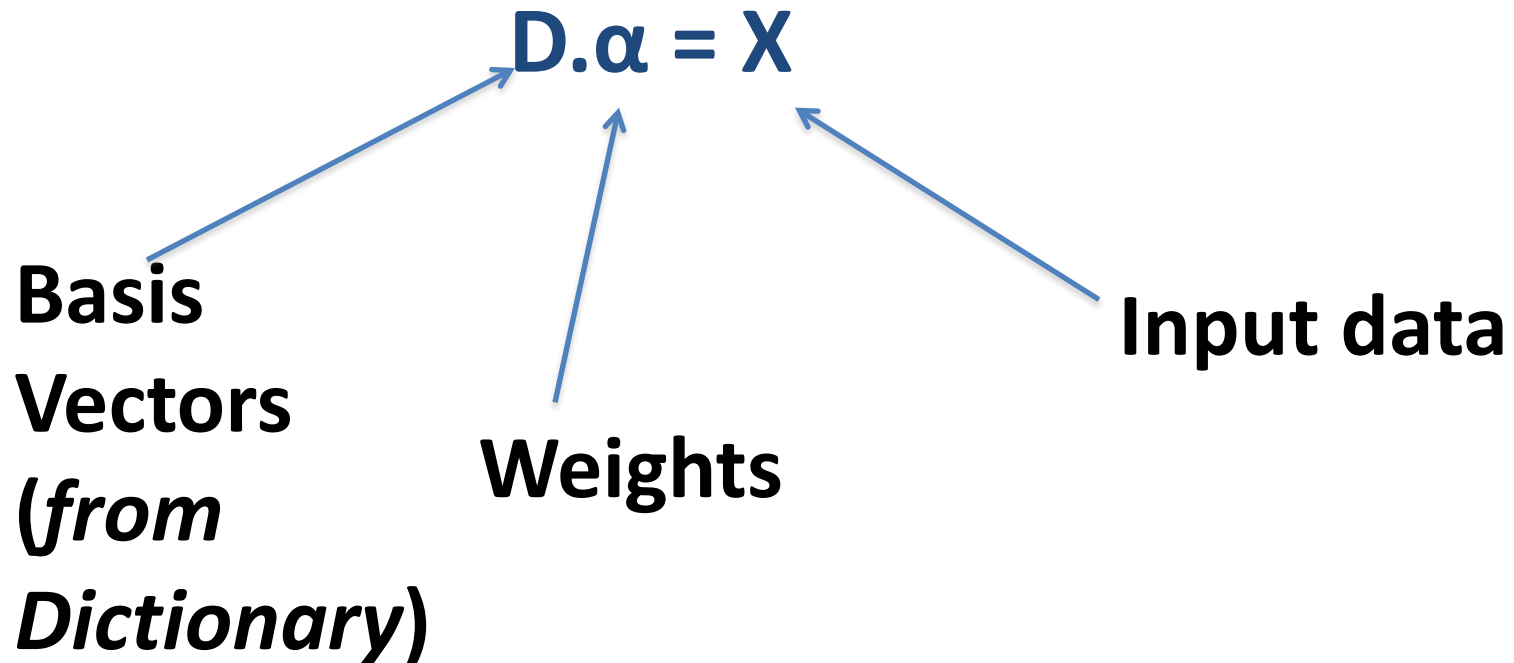
we have an **overcomplete** representation

➤ 4096 x **N**

**VERY LARGE!!!**

# Quick Linear Algebra Refresher

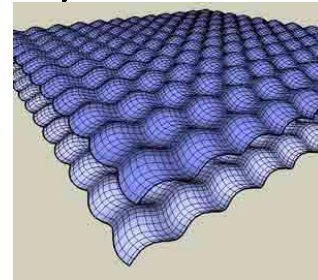
- Remember, #(Basis Vectors)= #unknowns





# Recap – Conventional : Images

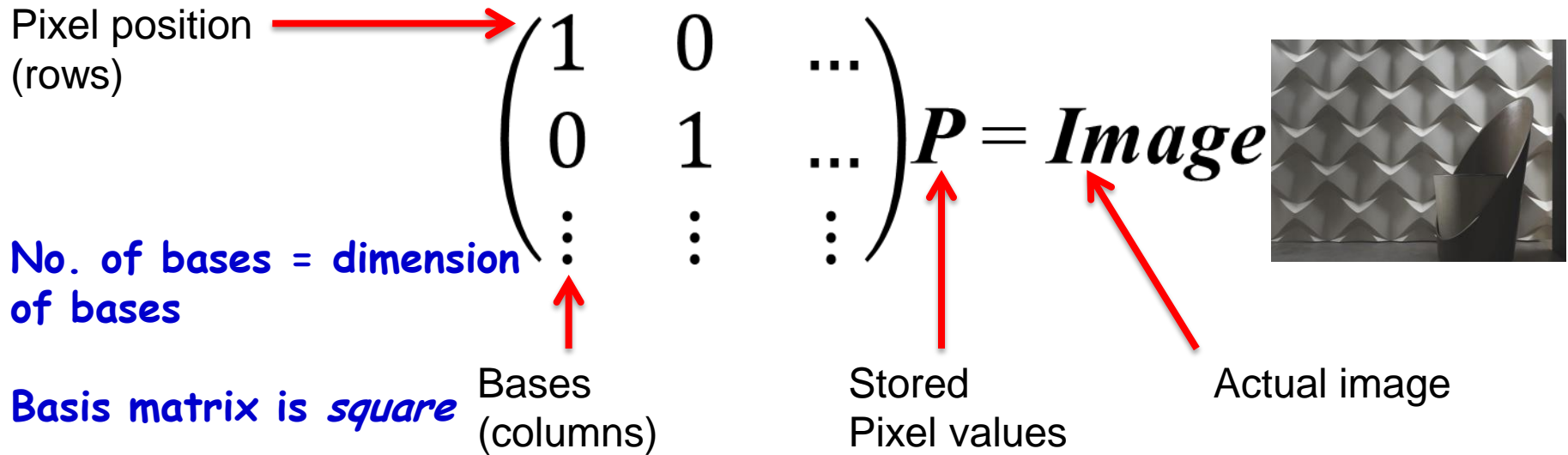
- Conventional characterization: Images
  - Store Pixel values (positions implicit)



- Image:  $128 \times \text{pixel.at.}(1,1) + 134 \times \text{pixel.at.}(1,2) + \dots + 127 \times \text{pixel.at.}(2,1) \dots$
- Store only the numbers (128,134, ..., 127)
- *Bases* are “ $\text{pixel.at.}(1,1)$ ”, “ $\text{pixel.at.}(1,2)$ ” etc..
  - Or rather  $[1 \ 0 \ 0 \ 0 \dots]$ ,  $[0 \ 1 \ 0 \ 0 \dots]$
  - Pixel positions are implicit

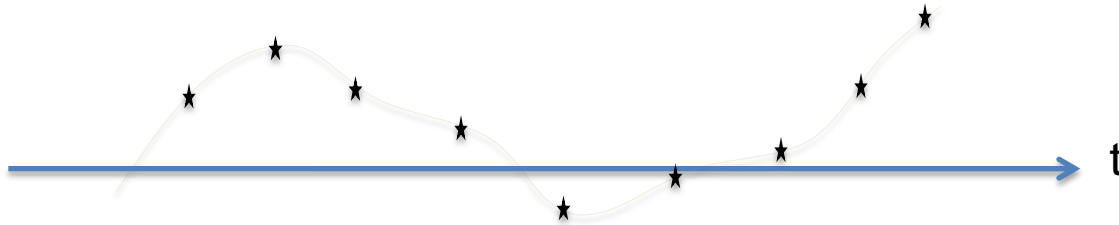
# Recap – Conventional : Images

- Storing an Image  $B \cdot P = \text{Image}$



- “Bases” are unit-valued pixels at specific locations
  - Only weights are stored
  - Basis matrix is *implicit* (everyone knows what it is)

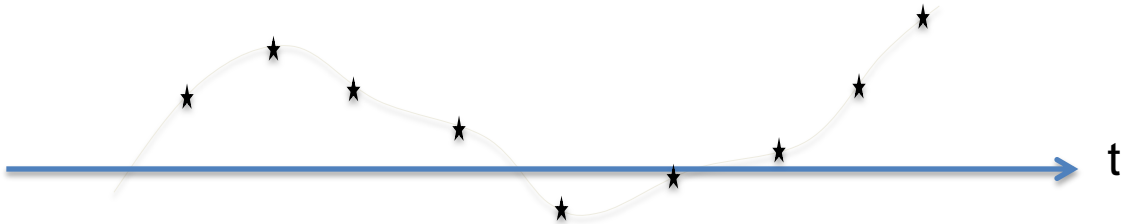
# Recap – Conventional: Sound



- Signal = 3 x sample.at.t=1 + 4 x sample.at.t=2 + 3.5 x sample.at.t=3 ....
- Store only the numbers [3, 4, 3.5...]
- Bases are “sample.at.t=1”, “sample.at.t=2”, ...
  - Or rather [..0 1 0 0 0 ...], [..0 0 1 0 0 0 ...], ....
  - “Time” is implicit

# Recap – Conventional : Sound

- Storing a sound **B.S = Recording**



Sample position  
(rows)

$$\begin{pmatrix} 1 & 0 & \dots \\ 0 & 1 & \dots \\ \vdots & \vdots & \vdots \end{pmatrix} S = \text{Signal}$$

No. of bases = dimension  
of bases

Basis matrix is *square*

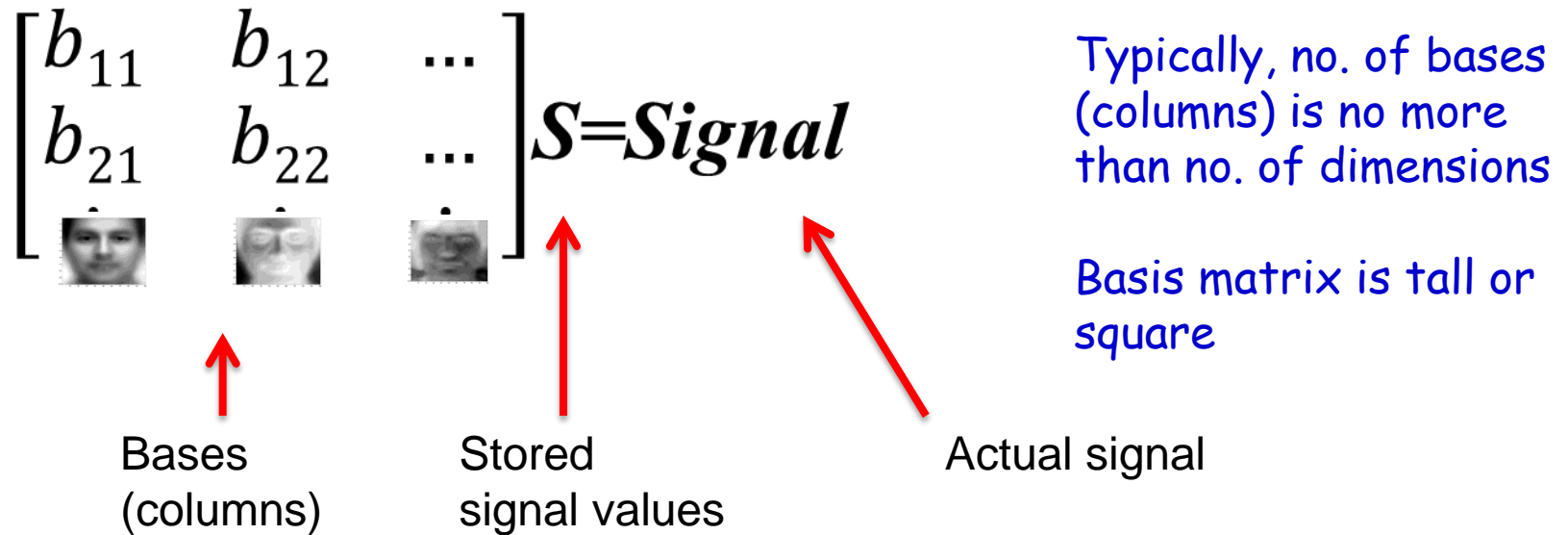
Bases  
(columns)

Stored  
signal values

Actual signal

- “Bases” are unit-valued samples at specific time instants
  - Only weights are stored
  - Basis matrix is *implicit* (everyone knows what it is)

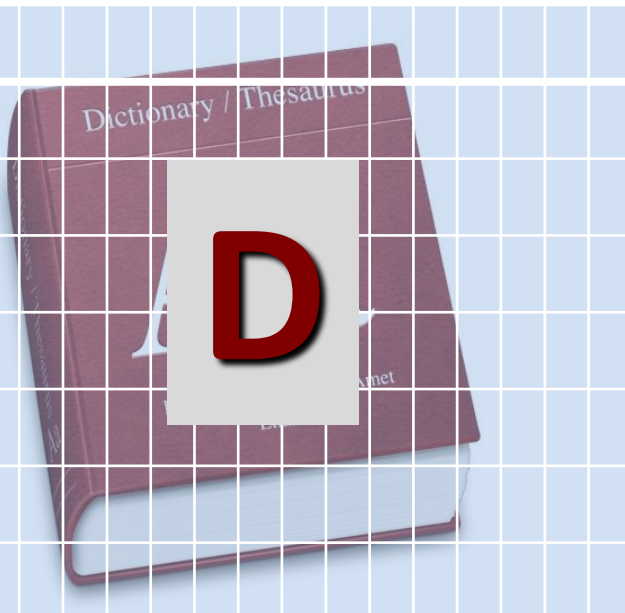
# Recap: *Component-based* representations



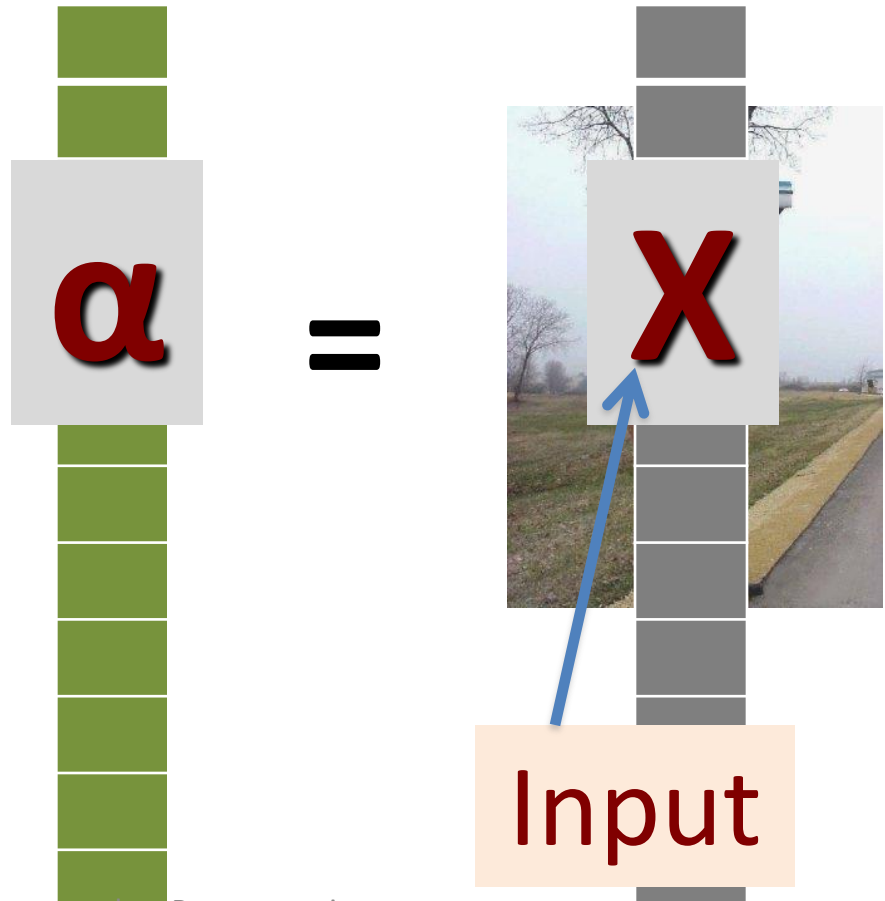
- Bases may be deterministic, e.g. sinusoids/wavelets or derived, e.g. PCA / ICA / NMF bases
- Only store  $w$  to represent individual signals. Bases matrix  $B$  stored separately as a one-time deal

# Dictionary based Representations

- Overcomplete “dictionary”-based representations are composition-based representations with more bases than the dimensionality of the data



Bases matrix is *wide*  
(more bases than dimensions)

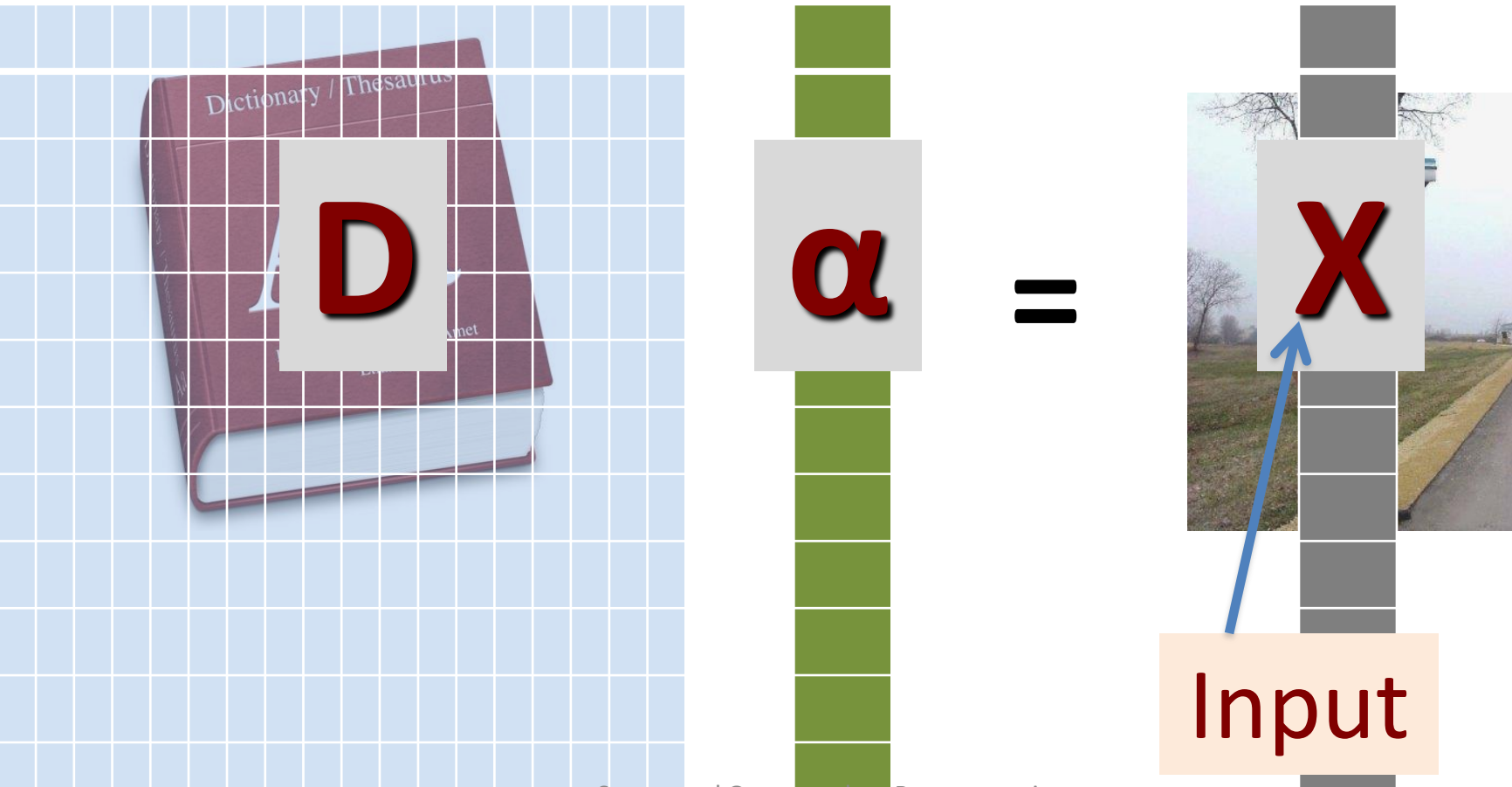


# Why Dictionary-based Representations?

- Dictionary based representations are semantically more meaningful
- Enable content-based description
  - Bases can capture entire structures in data
  - E.g. notes in music
  - E.g. image structures (such as faces) in images
- Enable content-based processing
  - Reconstructing, separating, denoising, manipulating speech/music signals
  - Coding, compression, etc.
- Statistical reasons: We will get to that in an hour..

# Problems

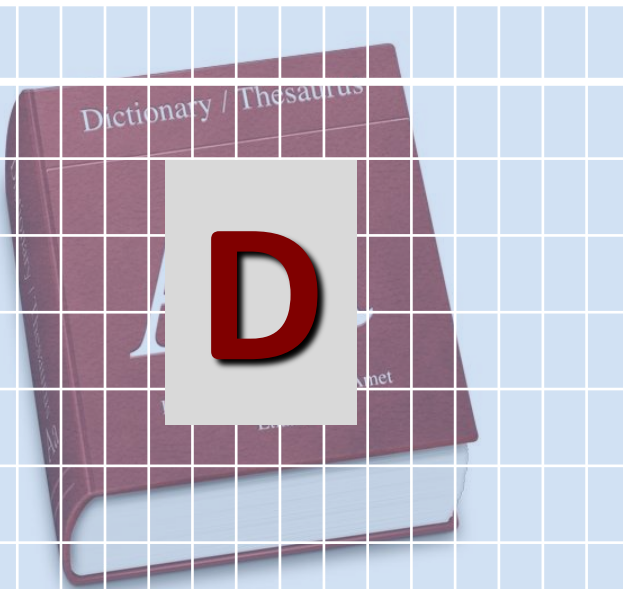
- How to obtain the dictionary
  - Which will give us meaningful representations
- How to compute the weights?



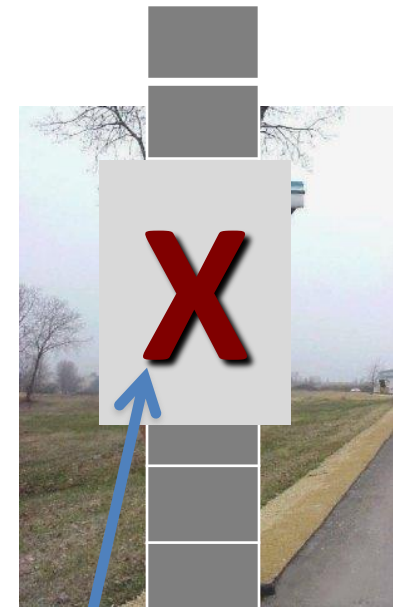


# Problems

- How to obtain the dictionary
  - Which will give us meaningful representations
- How to compute the weights?



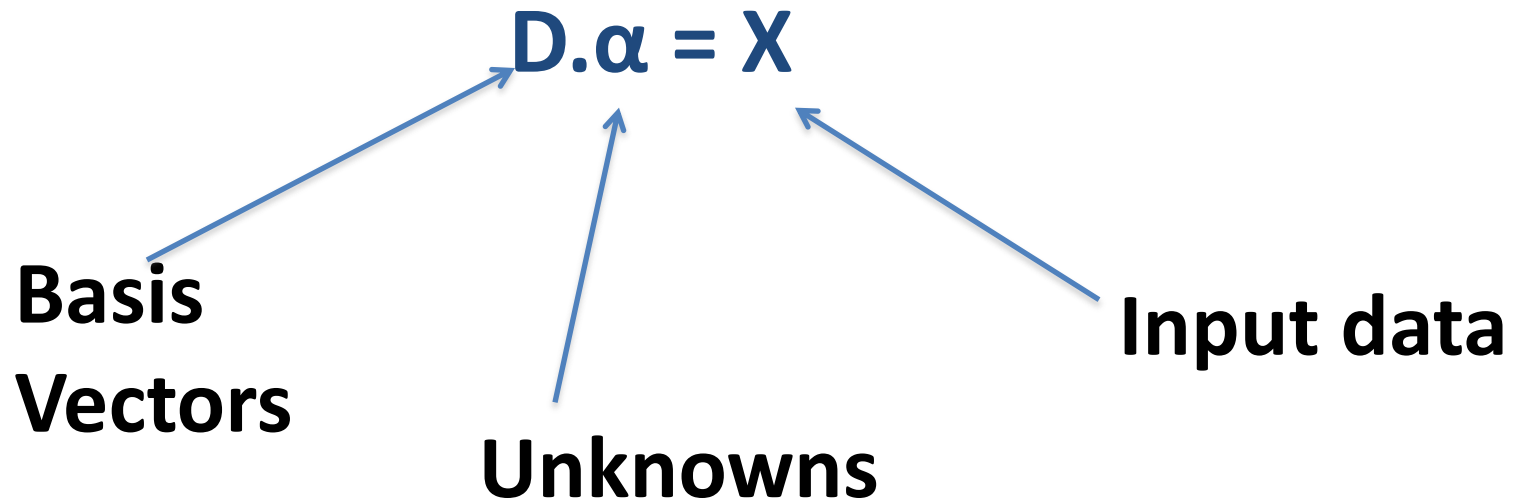
=



Input

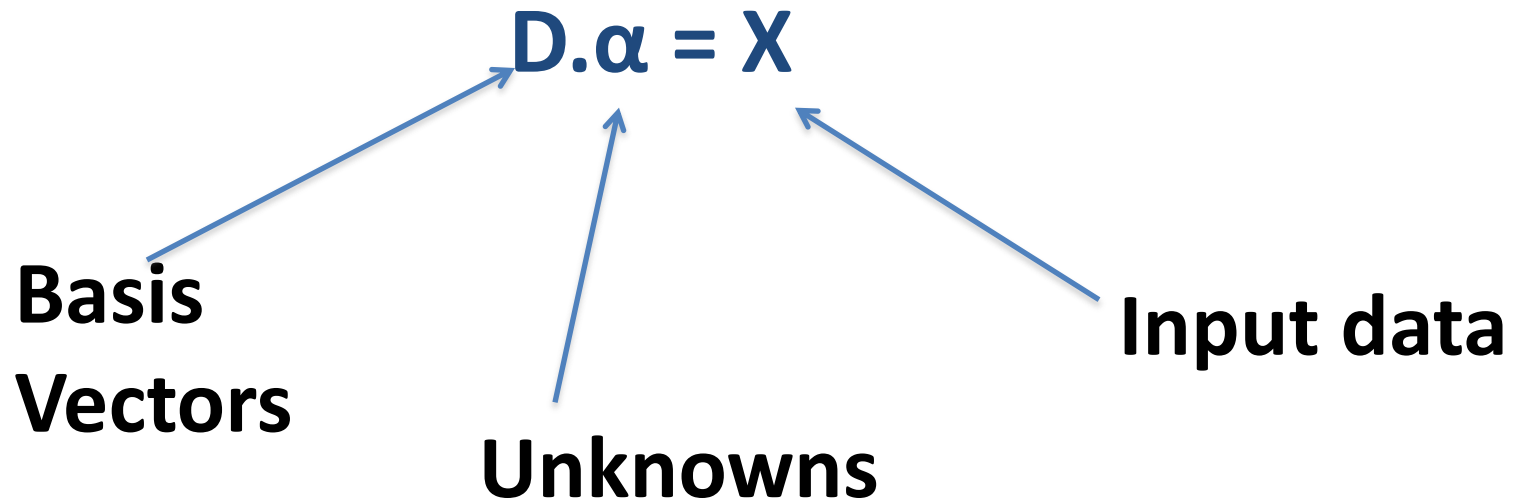
# Quick Linear Algebra Refresher

- Remember, #(Basis Vectors) = #unknowns



# Quick Linear Algebra Refresher

- Remember, #(Basis Vectors) = #unknowns



When can we solve for  $\alpha$ ?

# Quick Linear Algebra Refresher

$$D \cdot \alpha = X$$

- When  $\#(\text{Basis Vectors}) = \dim(\text{Input Data})$ , we have a unique solution
- When  $\#(\text{Basis Vectors}) < \dim(\text{Input Data})$ , we may have no exact solution
- When  $\#(\text{Basis Vectors}) > \dim(\text{Input Data})$ , we have infinitely many solutions

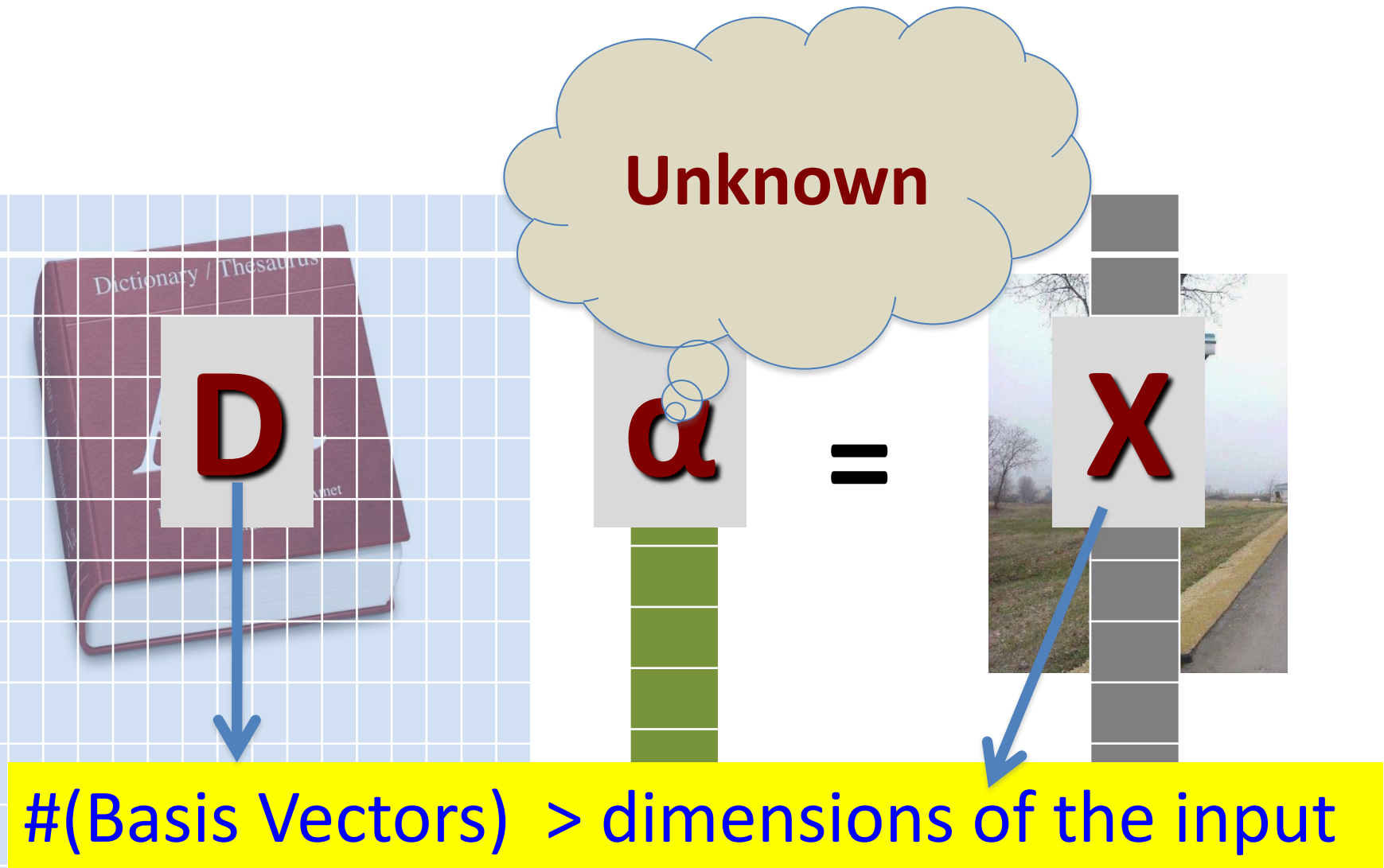
# Quick Linear Algebra Refresher

$$D \cdot \alpha = X$$

- When  $\#(\text{Basis Vectors}) = \dim(\text{Input Data})$ , we have a unique solution
- When  $\#(\text{Basis Vectors}) < \dim(\text{Input Data})$ , we may have no solution
- When  $\#(\text{Basis Vectors}) > \dim(\text{Input Data})$ , we have infinitely many solutions

Our Case

# Overcomplete Representation



# Overcompleteness and Sparsity

- To solve an overcomplete system of the type:

$$\mathbf{D} \cdot \boldsymbol{\alpha} = \mathbf{X}$$

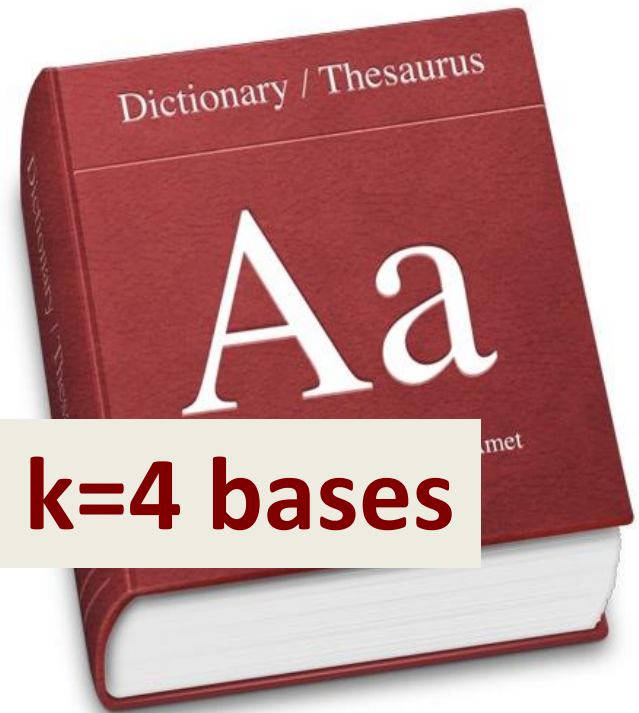
- Make assumptions about the data.
- Suppose, we say that  $\mathbf{X}$  is composed of no more than a fixed number ( $\mathbf{k}$ ) of “bases” from  $\mathbf{D}$  ( $\mathbf{k} \leq \dim(\mathbf{X})$ )
  - The term “bases” is an abuse of terminology..
- Now, we can find the set of  $\mathbf{k}$  bases that best fit the data point,  $\mathbf{X}$ .

# Representing Data



Using bases that we know...

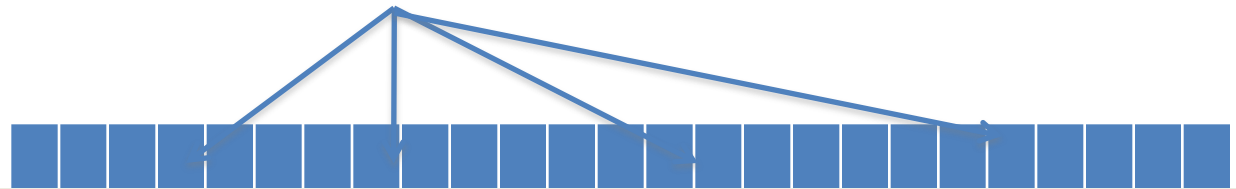
But no more than  $k=4$  bases





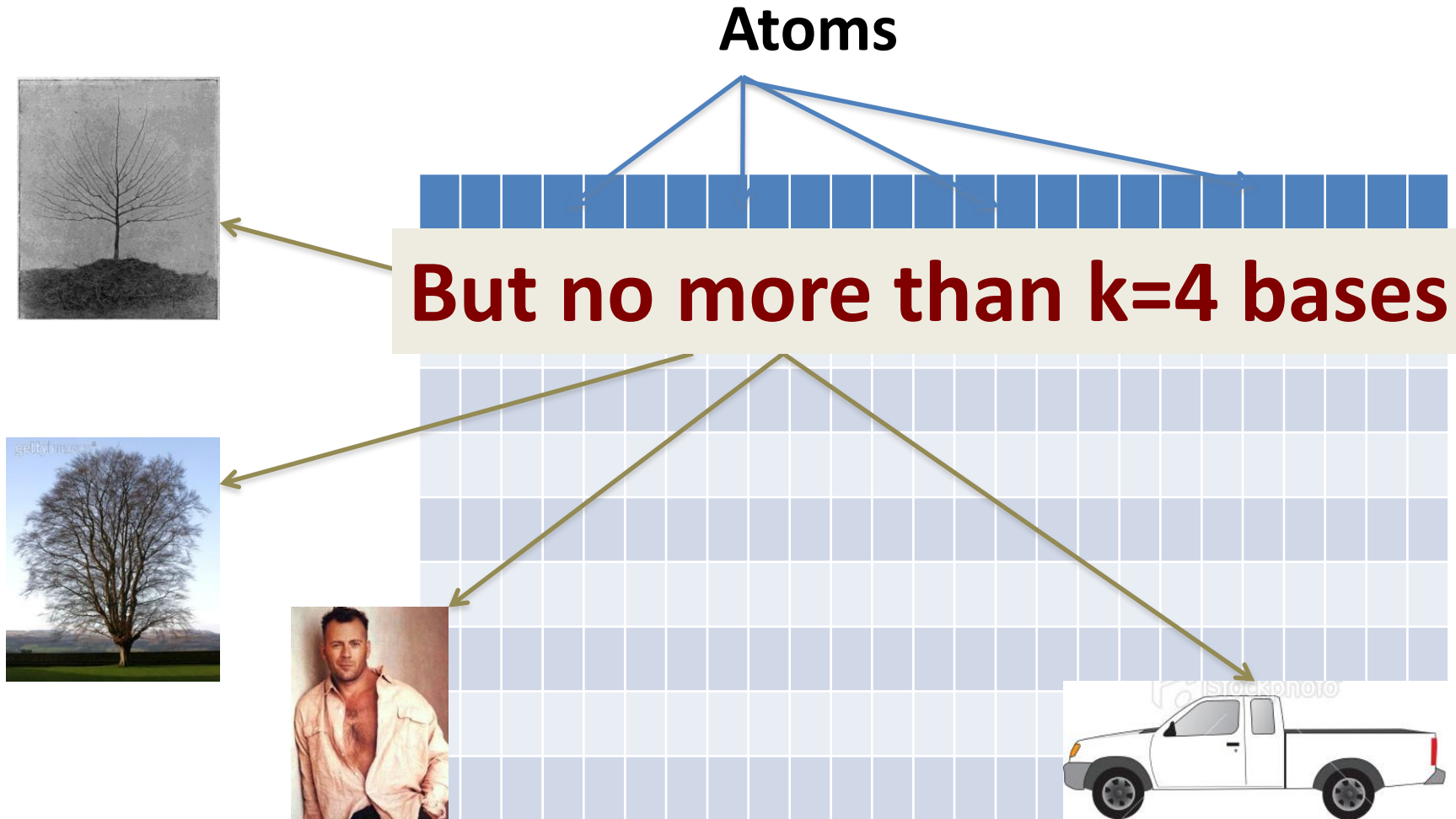
# Overcompleteness and Sparsity

Atoms

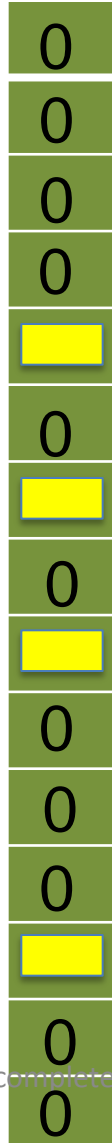
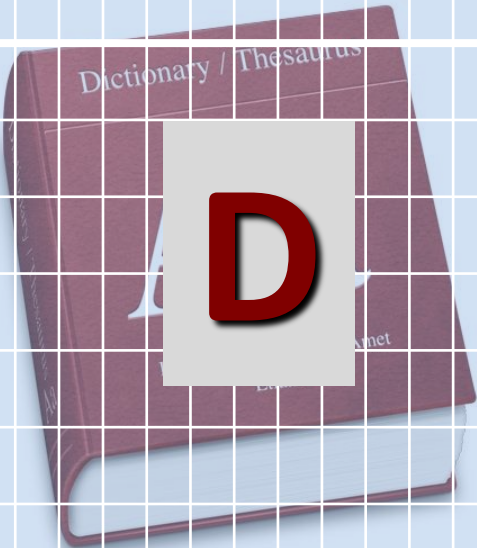
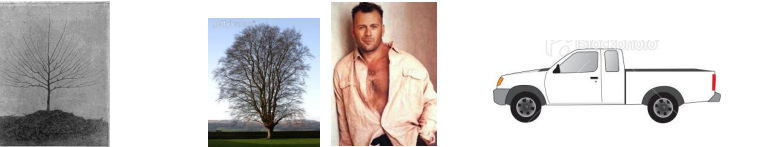


**But no more than  $k=4$  bases are “active”**

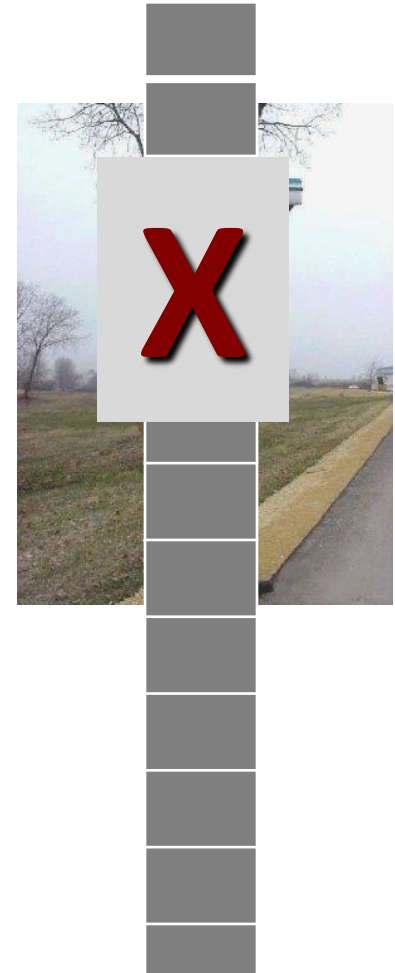
# Overcompleteness and Sparsity



# No more than 4 bases



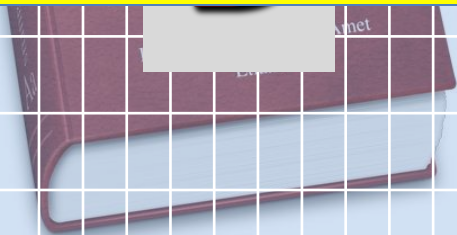
=



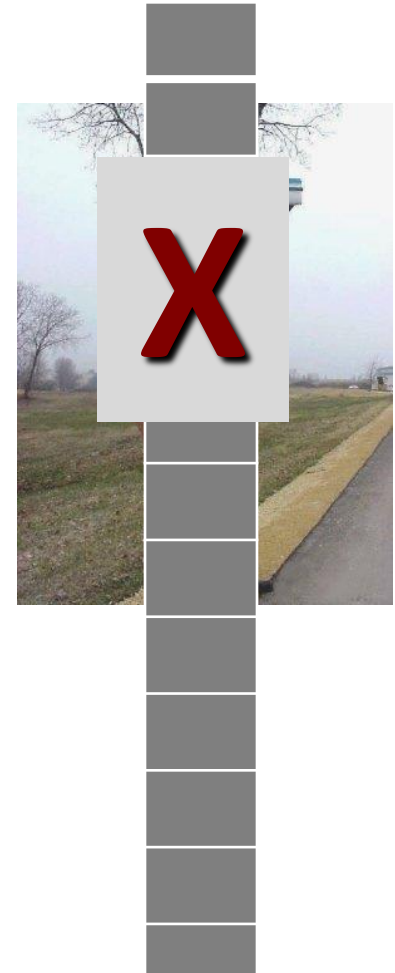
# No more than 4 bases



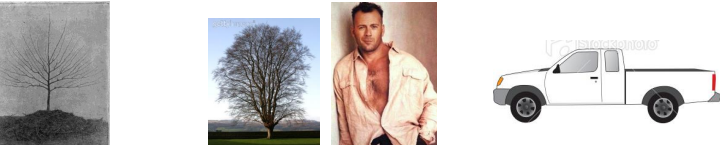
ONLY THE  $\alpha$  COMPONENTS  
CORRESPONDING  
TO THE 4 KEY  
DICTIONARY ENTRIES  
ARE NON-ZERO



=

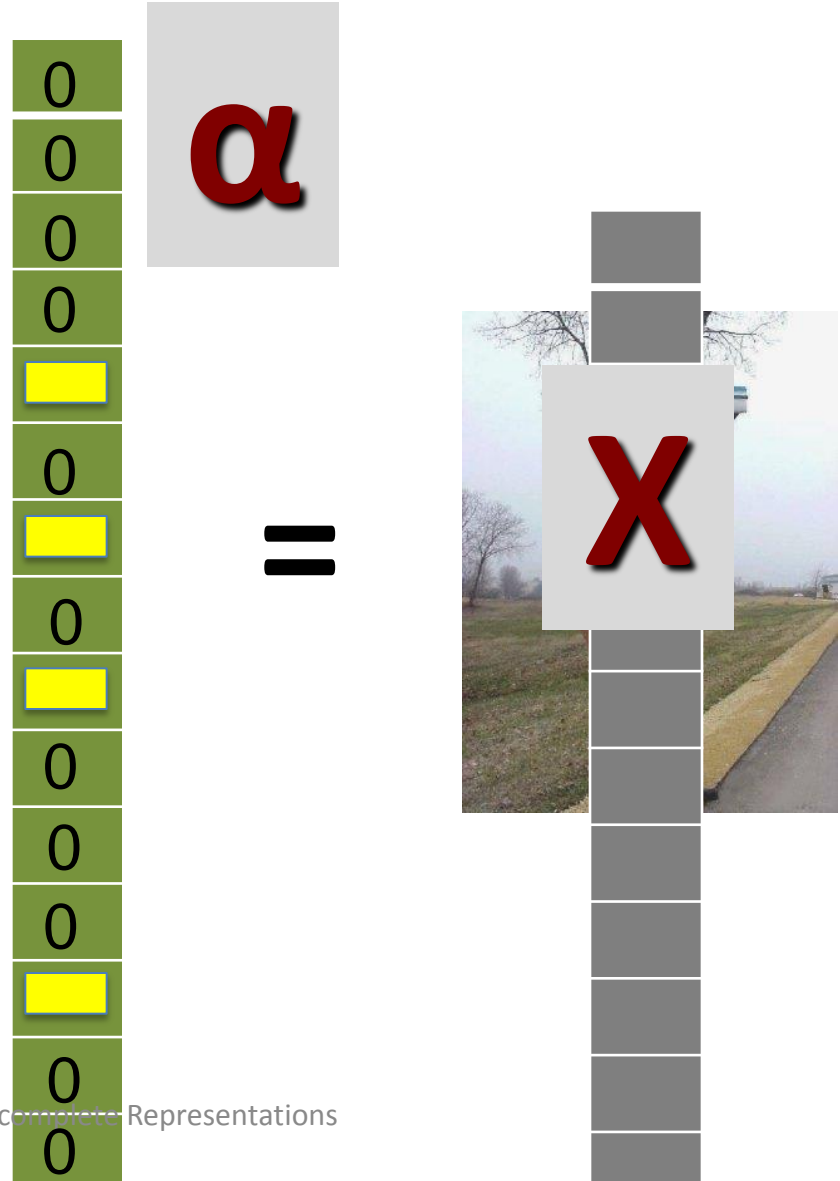


# No more than 4 bases



ONLY THE  $\alpha$  COMPONENTS CORRESPONDING TO THE 4 KEY DICTIONARY ENTRIES ARE NON-ZERO

MOST OF  $\alpha$  IS ZERO!!  
 $\alpha$  IS SPARSE



# Sparsity- Definition

- *Sparse representations* are representations that account for most or all information of a signal with a linear combination of a small number of atoms.


(from: [www.see.ed.ac.uk/~tblumens/Sparse/Sparse.html](http://www.see.ed.ac.uk/~tblumens/Sparse/Sparse.html))

# The Sparsity Problem

- We don't really know  $\mathbf{k}$
- You are given a signal  $\mathbf{X}$
- Assuming  $\mathbf{X}$  was generated using the dictionary, can we find  $\alpha$  that generated it?

# The Sparsity Problem

- We want to use as few basis vectors as possible to do this.

$$\begin{array}{l} \text{Min}_{\underline{\alpha}} \quad \|\underline{\alpha}\|_0 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$




# The Sparsity Problem

- We want to use as few basis vectors as possible to do this.

$$\begin{array}{l} \text{Min}_{\underline{\alpha}} \|\underline{\alpha}\|_0 \\ \text{s.t. } \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

Counts the number of non-zero elements in  $\alpha$

# The Sparsity Problem

- We want to use **as few basis vectors** as possible to do this
  - Ockham's razor: Choose the simplest explanation invoking the fewest variables

$$\begin{array}{l} \text{Min}_{\underline{\alpha}} \quad \|\underline{\alpha}\|_0 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

# The Sparsity Problem

- We want to use as few basis vectors as possible to do this.

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_0 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

**How can we solve the above?**

# Obtaining Sparse Solutions

- We will look at 2 algorithms:
  - Matching Pursuit (MP)
  - Basis Pursuit (BP)

# Matching Pursuit (MP)


- Greedy algorithm
- Finds an atom in the dictionary that best matches the input signal
- Remove the weighted value of this atom from the signal
- Again, find an atom in the dictionary that best matches the remaining signal.
- Continue till a defined stop condition is satisfied.

# Matching Pursuit

- Find the dictionary atom that best matches the given signal.



Weight =  $w_1$



# Matching Pursuit

- Remove weighted image to obtain updated signal



Find best match for  
this signal from the  
dictionary

# Matching Pursuit

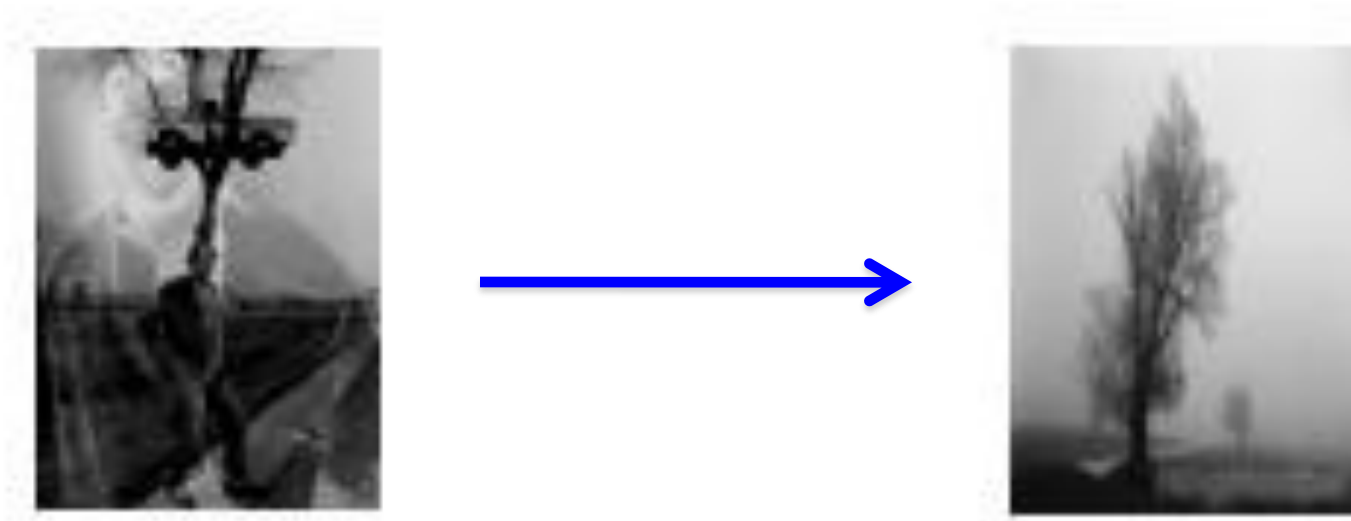
- Find best match for updated signal





# Matching Pursuit

- Find best match for updated signal



Iterate till you reach a stopping condition,  
 **$\text{norm}(\text{ResidualInputSignal}) < \text{threshold}$**

# Matching Pursuit

## Algorithm Matching Pursuit

Input: Signal:  $f(t)$ .

Output: List of coefficients:  $(a_n, g_{\gamma_n})$ .

Initialization:

$$Rf_1 \leftarrow f(t);$$

Repeat

find  $g_{\gamma_n} \in D$  with maximum inner product  $\langle Rf_n, g_{\gamma_n} \rangle$ ;

$$a_n \leftarrow \langle Rf_n, g_{\gamma_n} \rangle;$$

$$Rf_{n+1} \leftarrow Rf_n - a_n g_{\gamma_n};$$

$$n \leftarrow n + 1;$$

Until stop condition (for example:  $\|Rf_n\| < \text{threshold}$ )

From [http://en.wikipedia.org/wiki/Matching\\_pursuit](http://en.wikipedia.org/wiki/Matching_pursuit)

# Matching Pursuit

- Problems ???

# Matching Pursuit

- Main Problem
  - Computational complexity
  - The entire dictionary has to be searched at every iteration

# Comparing MP and BP

Matching Pursuit	Basis Pursuit
Hard thresholding  (remember the equations)	
Greedy optimization at each step	
Weights obtained using greedy rules	

# Basis Pursuit (BP)

- Remember,

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_0 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

# Basis Pursuit

- Remember,

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_0 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

In the general case, this is intractable

# Basis Pursuit

- Remember,

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_0 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

In the general case, this is intractable

Requires combinatorial optimization



# Basis Pursuit

- Replace the intractable expression by an expression that is solvable

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_1 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

# Basis Pursuit

- Replace the intractable expression by an expression that is solvable

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_1 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

This will provide identical solutions when  $\mathbf{D}$  obeys the *Restricted Isometry Property*.

# Basis Pursuit

- Replace the intractable expression by an expression that is solvable

$$\begin{array}{l} \underset{\underline{\alpha}}{\text{Min}} \quad \|\underline{\alpha}\|_1 \\ \text{s.t.} \quad \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

Objective

Constraint

# Basis Pursuit

- We can formulate the optimization term as:

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

Constraint

Objective

# Basis Pursuit


- We can formulate the optimization term as:

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

$\lambda$  is a penalty term on the non-zero elements and promotes sparsity

# Basis Pursuit

Equivalent to *LASSO*; for more details, see [this paper by Tibshirani](#)


$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

$\lambda$  is a penalty term on the non-zero elements and promotes sparsity

# Basis Pursuit

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

$$\frac{\partial \|\underline{\alpha}\|_1}{\partial \alpha_j} = \begin{cases} +1 & \text{at } \alpha_j > 0 \\ [-1, 1] & \text{at } \alpha_j = 0 \\ -1 & \text{at } \alpha_j < 0 \end{cases}$$

- $\|\alpha\|_1$  is not differentiable at  $\alpha_j = 0$
- Gradient of  $\|\alpha\|_1$  for gradient descent update
- At optimum, following conditions hold

$$\nabla_j \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \text{sign}(\alpha_j) = 0, \quad \text{if } |\alpha_j| > 0$$

$$\nabla_j \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 \leq \lambda, \quad \text{if } \alpha_j = 0$$

# Basis Pursuit

- There are efficient ways to solve the LASSO formulation. [Link to [Matlab code](#)]
- Simplest solution: Coordinate descent algorithms
  - On webpage..

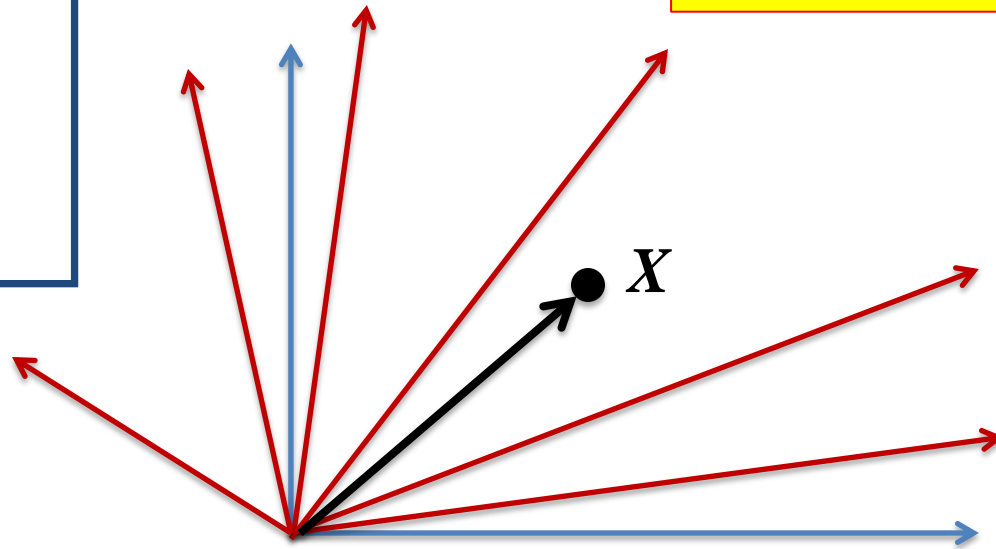


# $L_1$ VS $L_0$

$$\text{Min}_{\underline{\alpha}} \|\underline{\alpha}\|_0$$

$$\text{s.t. } \underline{X} = \mathbf{D}\underline{\alpha}$$

Overcomplete set  
of 6 bases



- **$L_0$  minimization**

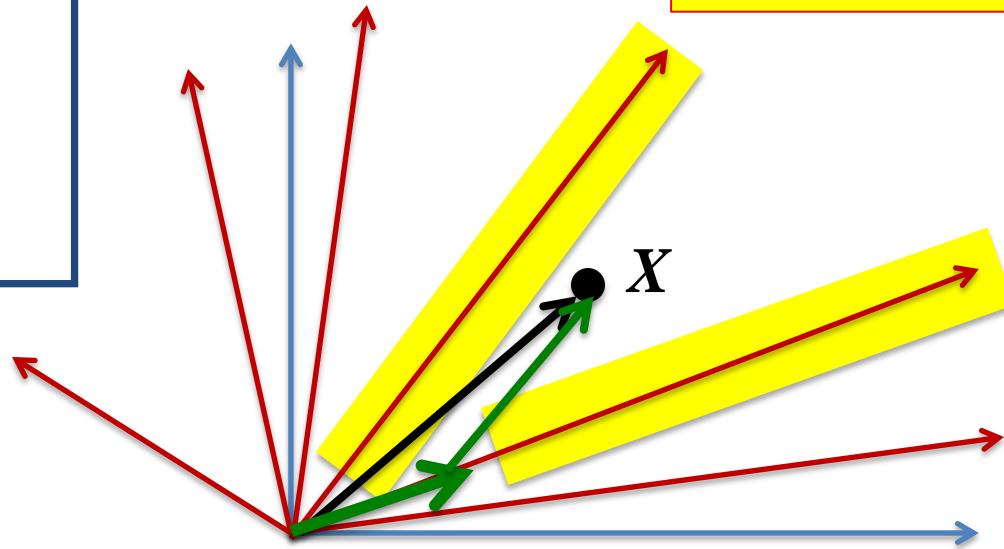
- Two-sparse solution

- ANY pair of bases can explain  $X$  with 0 error

# $L_1$ VS $L_0$

$$\begin{array}{l} \text{Min}_{\underline{\alpha}} \|\underline{\alpha}\|_1 \\ \text{s.t. } \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

Overcomplete set  
of 6 bases



- **$L_1$  minimization**

- Two-sparse solution
- All else being equal, the two closest bases are chosen

# Comparing MP and BP

Matching Pursuit	Basis Pursuit
Hard thresholding	Soft thresholding
(remember the equations)	
Greedy optimization at each step	Global optimization
Weights obtained using greedy rules	Can force N-sparsity with appropriately chosen weights

# General Formalisms

- $L_0$  minimization
- $L_0$  constrained optimization

$$\underset{\underline{\alpha}}{\text{Min}} \|\underline{\alpha}\|_0$$

$$\text{s.t. } \underline{X} = \mathbf{D}\underline{\alpha}$$

$$\underset{\underline{\alpha}}{\text{Min}} \|\underline{X} - \mathbf{D}\underline{\alpha}\|_2^2$$

$$\text{s.t. } \|\underline{\alpha}\|_0 < C$$

- $L_1$  minimization
- $L_1$  constrained optimization

$$\underset{\underline{\alpha}}{\text{Min}} \|\underline{\alpha}\|_1$$

$$\text{s.t. } \underline{X} = \mathbf{D}\underline{\alpha}$$

$$\underset{\underline{\alpha}}{\text{Min}} \|\underline{X} - \mathbf{D}\underline{\alpha}\|_2^2$$

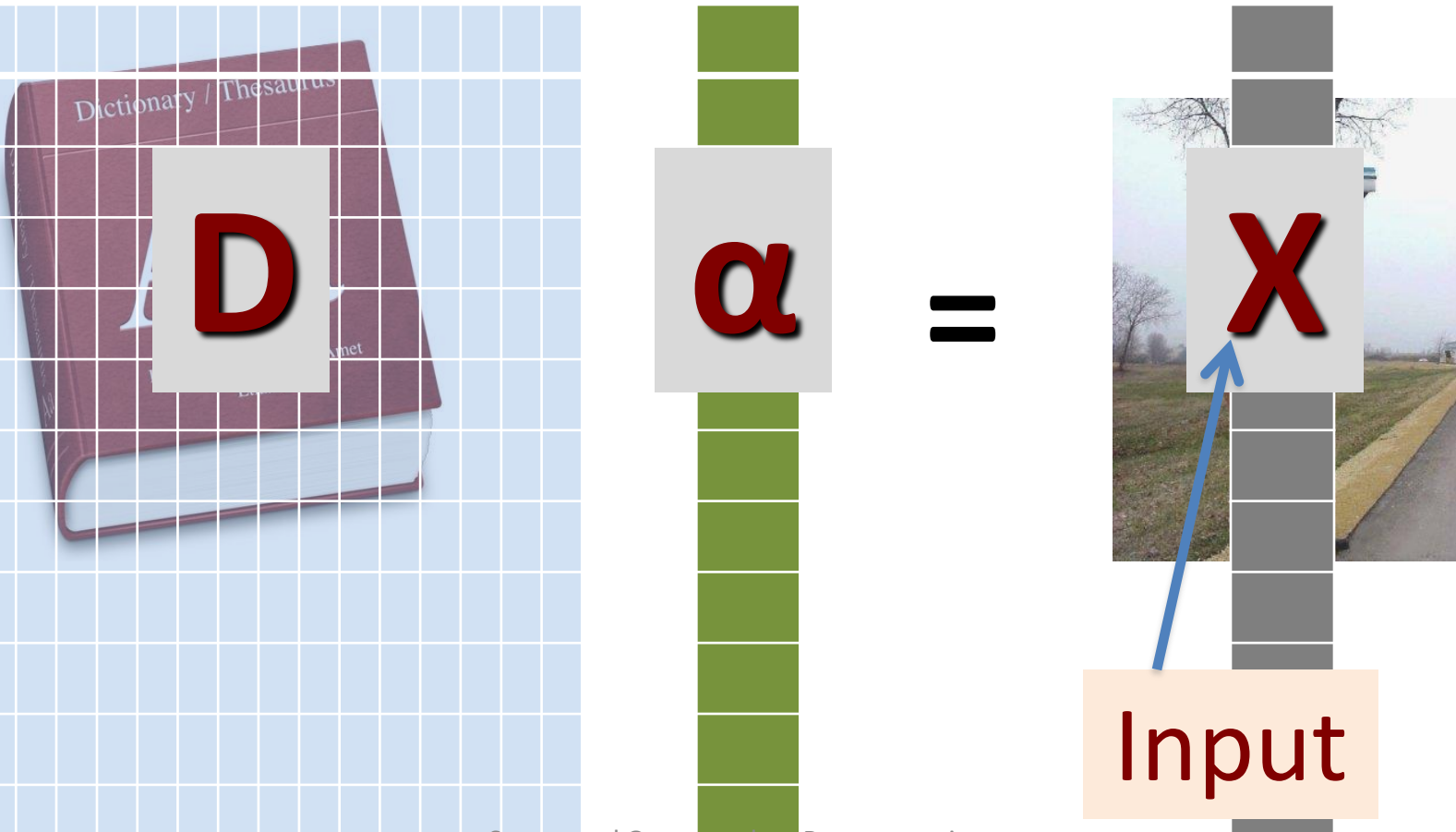
$$\text{s.t. } \|\underline{\alpha}\|_1 < C$$

# Many Other Methods..

- Iterative Hard Thresholding (IHT)
- CoSAMP
- OMP
- ...

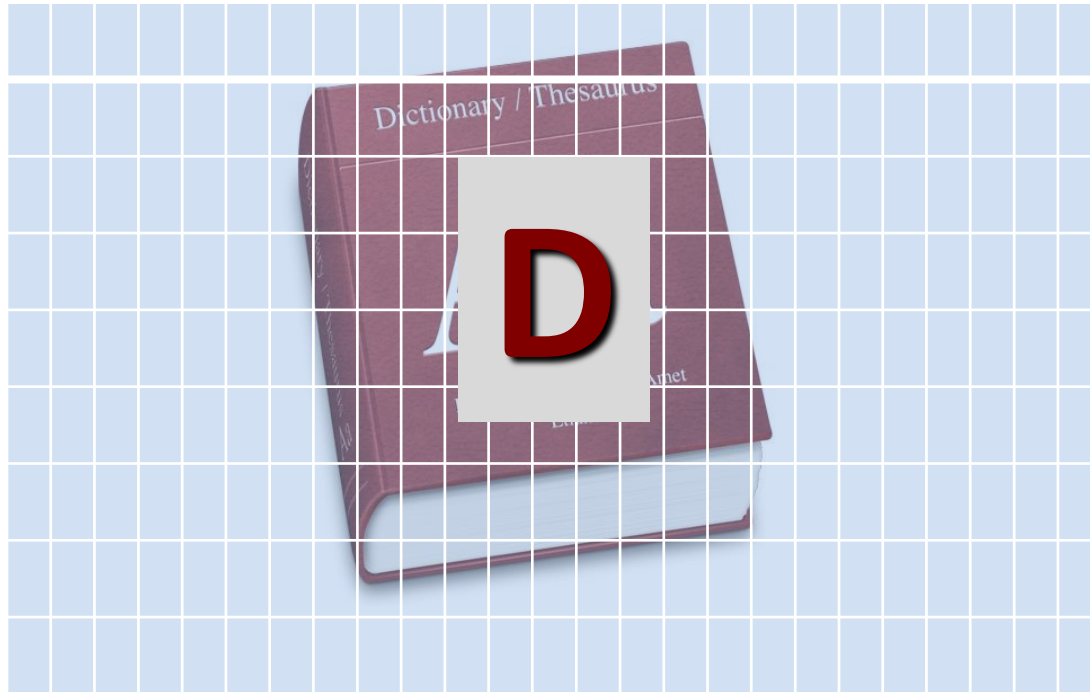
# Problems

- How to obtain the dictionary
  - Which will give us meaningful representations
- How to compute the weights?



# Dictionaries: Compressive Sensing

- Just random vectors!

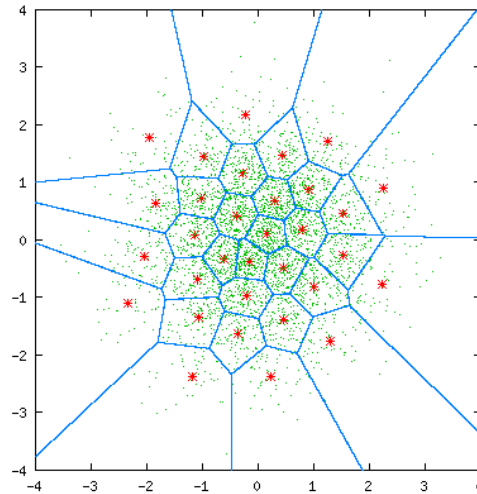


# More Structured ways of Constructing Dictionaries

- Dictionary entries must be structurally “meaningful”
  - Represent true compositional units of data
- Have already encountered two ways of building dictionaries
  - NMF for non-negative data
  - K-means ..



# K-Means for Composing Dictionaries



Train the codebook  
from training data  
using K-means

- Every vector is approximated by the centroid of the cluster it falls into
- Cluster means are “codebook” entries
  - Dictionary entries
  - Also compositional units the compose the data

# K-Means for Dictionaries

Each column is a codeword (centroid) from the codebook

**D**

**$\alpha$**

0  
0  
0  
0  
**1**  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0



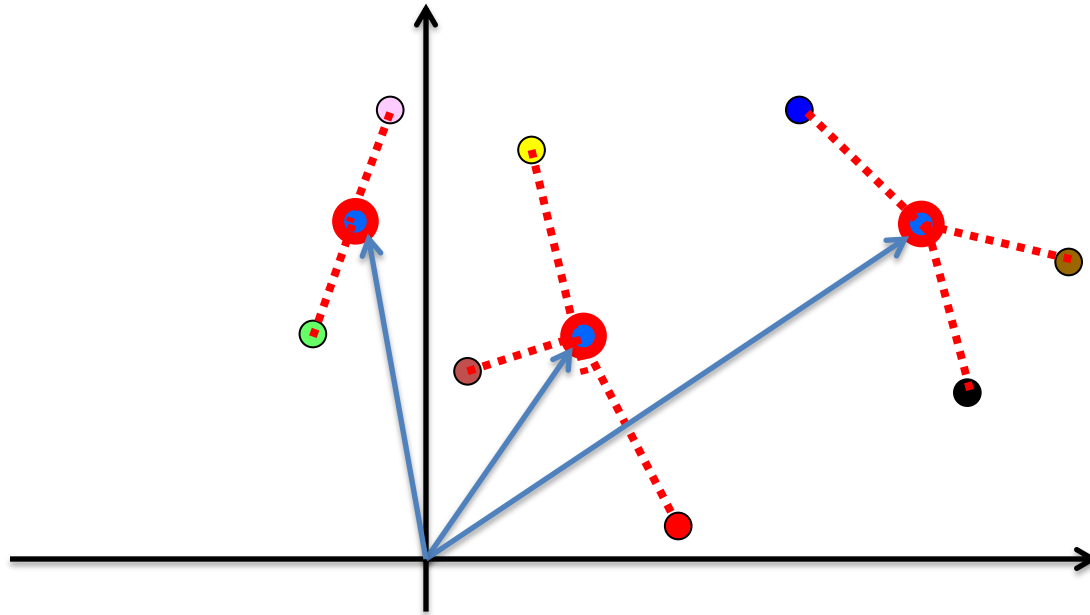
**X**

- $\alpha$  must be 1 sparse
- Only  $\alpha$  entry must 1

$$\|\alpha\|_0 = 1$$

$$\|\alpha\|_1 = 1$$

# K-Means



- Learn Codewords to minimize the total squared length of the training vectors from the closest codeword

# Length-unconstrained K-Means for Dictionaries

Each column is a codeword (centroid) from the codebook

**D**

**$\alpha$**

0  
0  
0  
0  
**3**  
0  
0  
0  
0  
0  
0  
0  
0  
0  
0

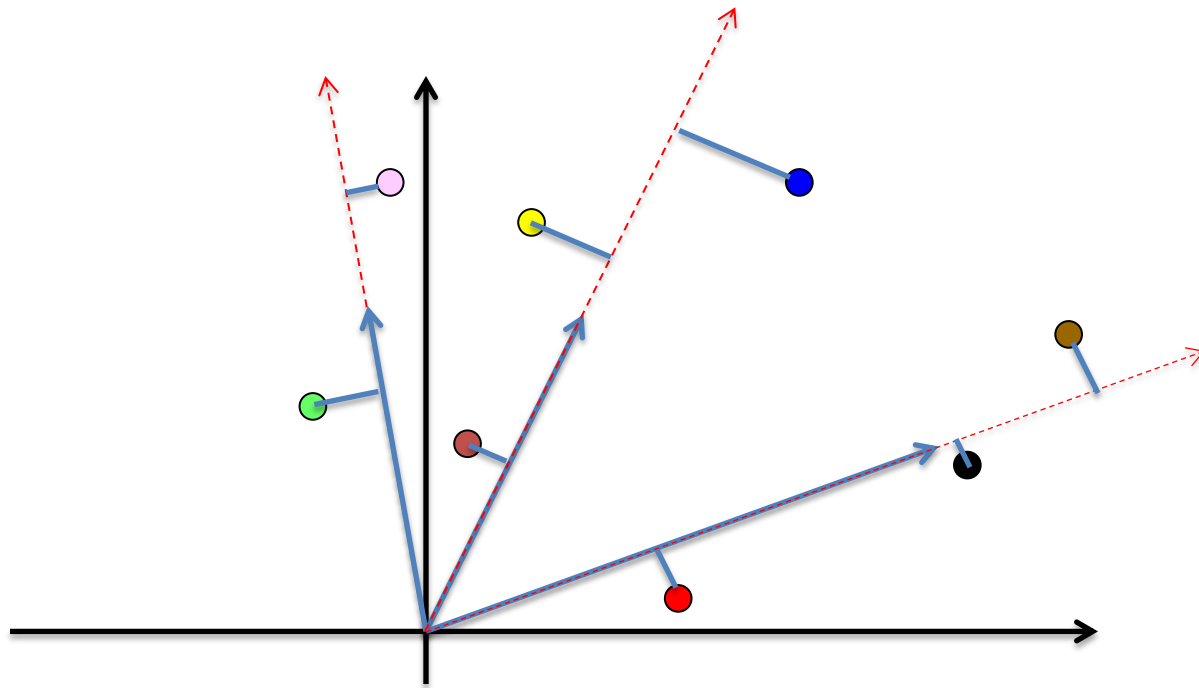


**X**

- $\alpha$  must be 1 sparse
- No restriction on  $\alpha$  value

$$\|\alpha\|_0 = 1$$

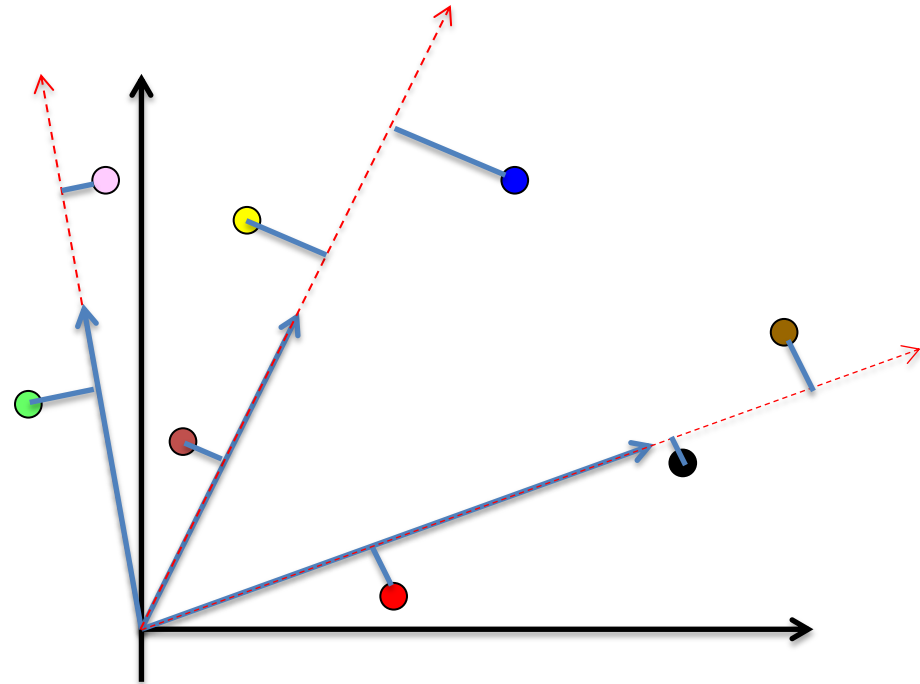
# SVD K-Means



- Learn Codewords to minimize the total squared *projection error* of the training vectors from the closest codeword

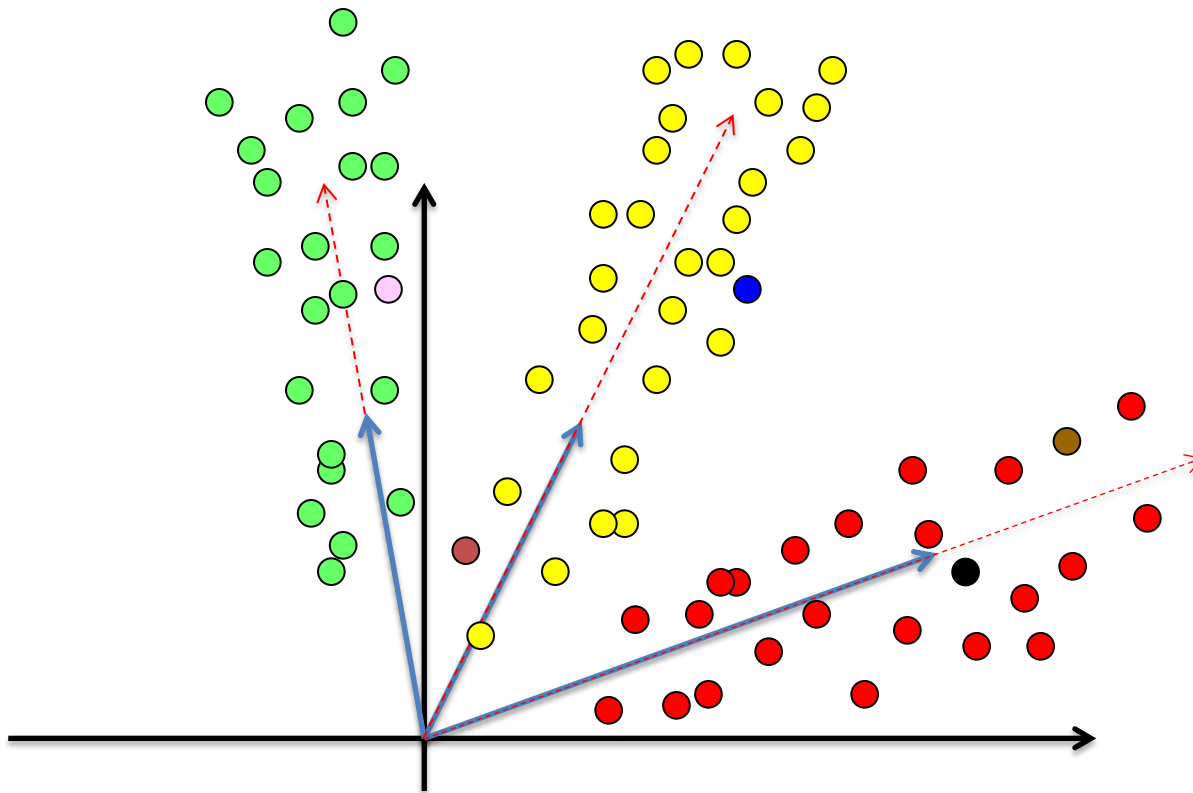
# SVD K-means

1. Initialize a set of centroids randomly
2. For each data point  $x$ , find the projection from the centroid for each cluster
  - $$p_{cluster} = |x^T m_{cluster}|$$
3. Put data point in the cluster of the closest centroid
  - Cluster for which  $p_{cluster}$  is *maximum*
4. When all data points are clustered, recompute centroids



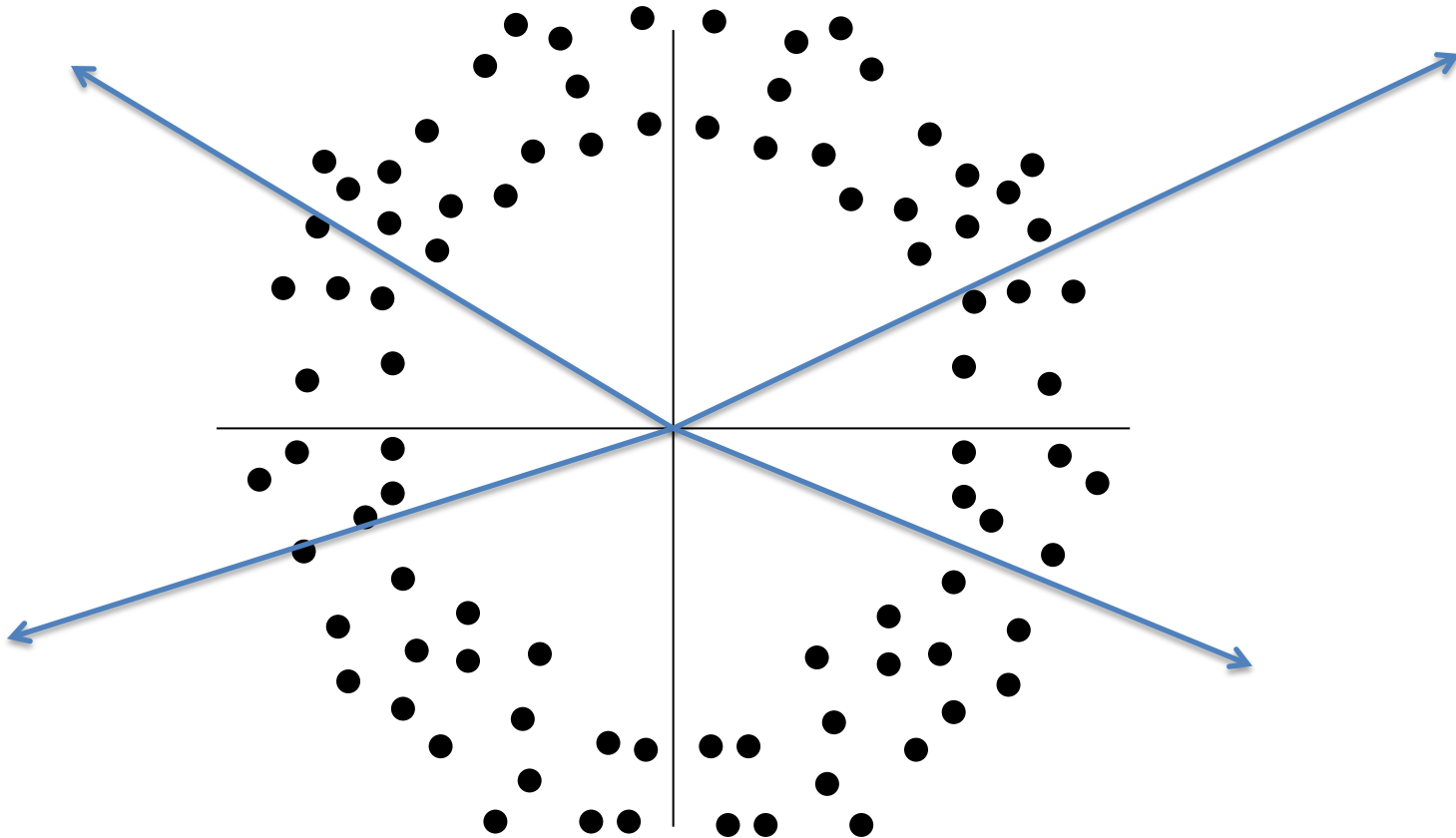
$$m_{cluster} = \text{Principal Eigenvector}(\{x \mid x \in cluster\})$$

# Problem



- Only represents *Radial* patterns

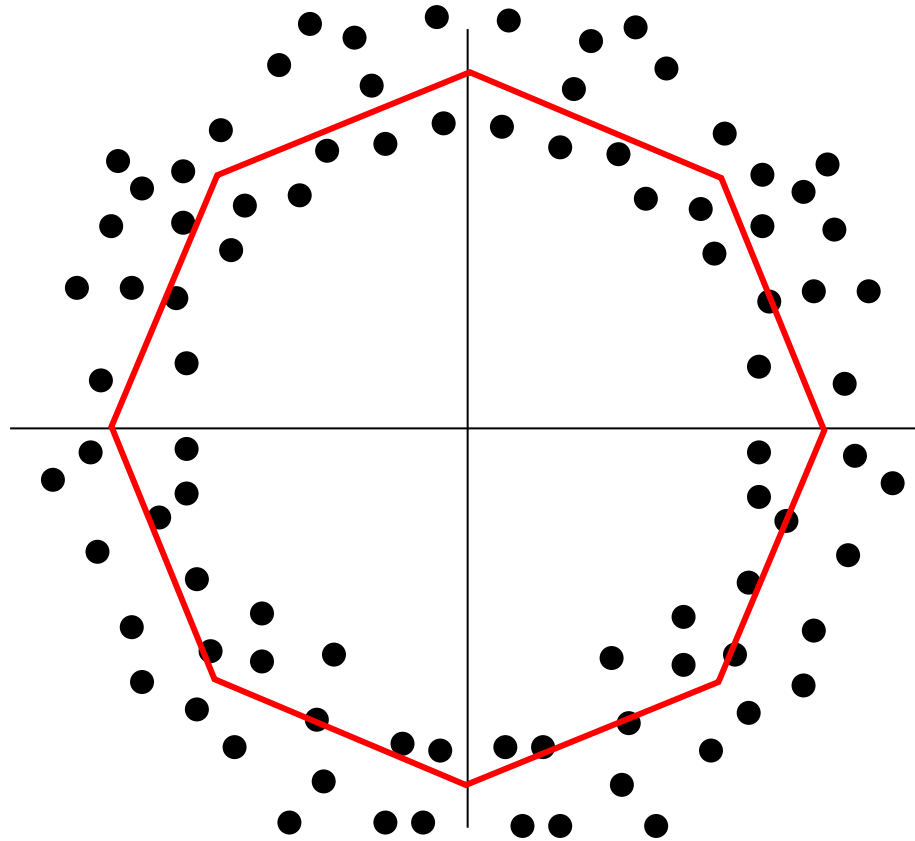
# What about this pattern?



- Dictionary entries that represent radial patterns will not capture this structure
  - 1-sparse representations will not do

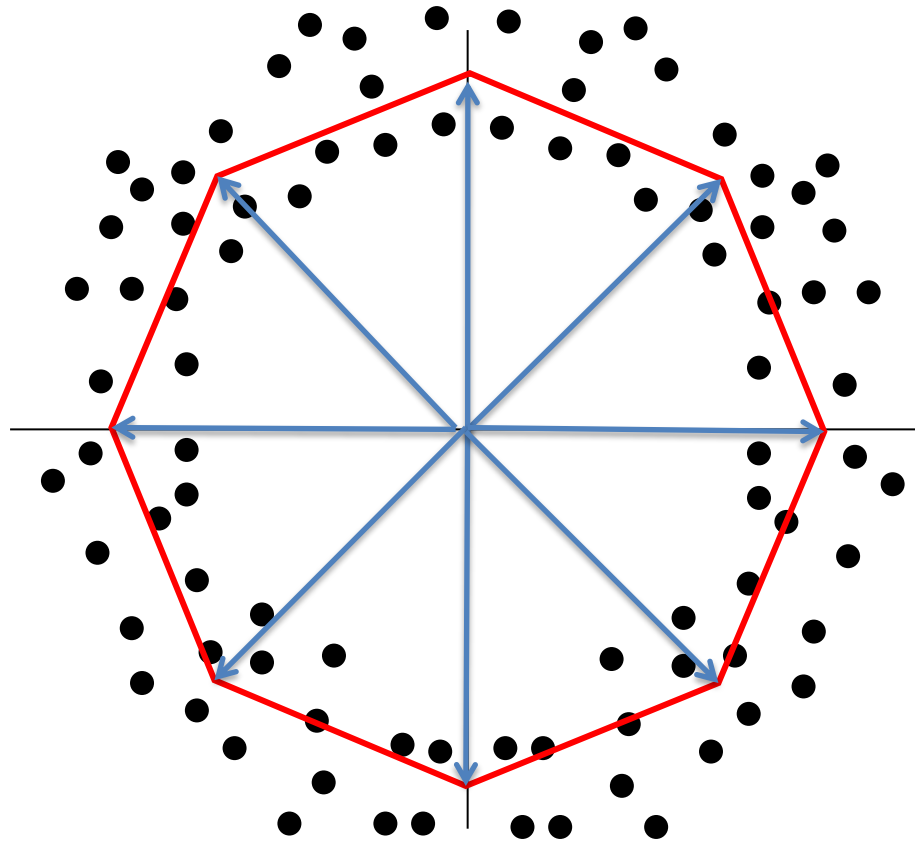


# What about this pattern?



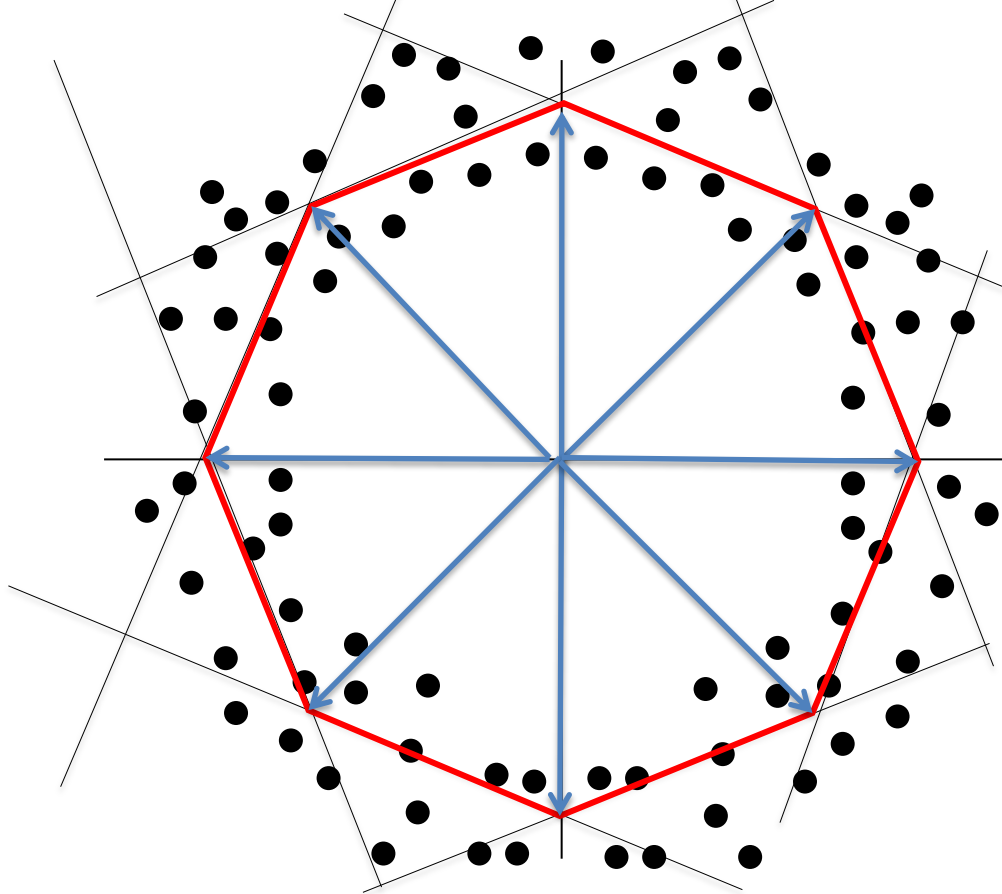
- We need **AFFINE** patterns

# What about this pattern?



- We need **AFFINE** patterns
- *Each vector is modeled by a linear combination of  $K$  (here 2) bases*

# What about this pattern?



Every line is a  
(constrained) combination  
of two bases

**2-sparse**

Constraint:  
Line =  $a \cdot b_1 + (1-a)b_2$

- We need **AFFINE** patterns
- *Each vector is modeled by a linear combination of  $K$  (here 2) bases*

# Codebooks for $K$ sparsity?

Each column is a codeword (centroid) from the codebook

**D**

**$\alpha$**

0  
0  
0  
0  
3  
0  
1  
0  
0  
4  
0  
0  
0  
0  
0



**X**

- $\alpha$  must be  $k$  sparse
- No restriction on  $\alpha$  value

$$\|\alpha\|_0 = k$$

# K SVD

- Initialize Codebook

- For every vector, compute K-sparse alphas

$$D =$$


$$\alpha =$$

0	0	2	0	0
1	0	0	0	0
0	0	0	1	0
0	7	0	0	0
0	0	0	1	0
3	0	0	0	0
0	5	0	0	1
0	0	1	0	2

# K-SVD

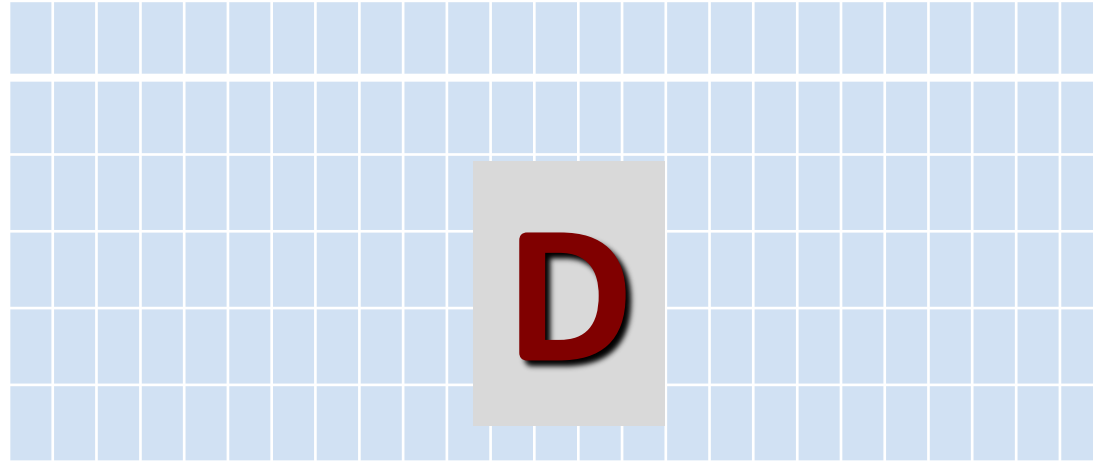
2. For each codeword ( $k$ ):
  - For each vector  $x$ 
    - Subtract the contribution of all other codewords to obtain  $e_k(x)$ 
      - Codeword-specific residual
  - Compute the principal Eigen vector of  $\{e_k(x)\}$
3. Return to step 1

$$D = \begin{matrix} & D_j \quad j \neq 1 \\ \begin{matrix} \text{blue} \\ \text{blue} \\ \text{blue} \end{matrix} & \begin{matrix} \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \\ \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} & \text{yellow} \end{matrix} \end{matrix}$$

$$\alpha = \begin{matrix} & \alpha_j(x) \quad j \neq 1 \\ \begin{matrix} \text{purple} \\ \text{red} \\ \text{red} \\ \text{red} \\ \text{red} \\ \text{red} \\ \text{red} \\ \text{red} \\ \text{red} \end{matrix} & \begin{matrix} 0 & 0 & 2 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 7 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 3 & 0 & 0 & 0 & 0 \\ 0 & 5 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 2 \end{matrix} \end{matrix}$$

$$e_k(x) = x - \sum_{j \neq k} \alpha_j D_j$$

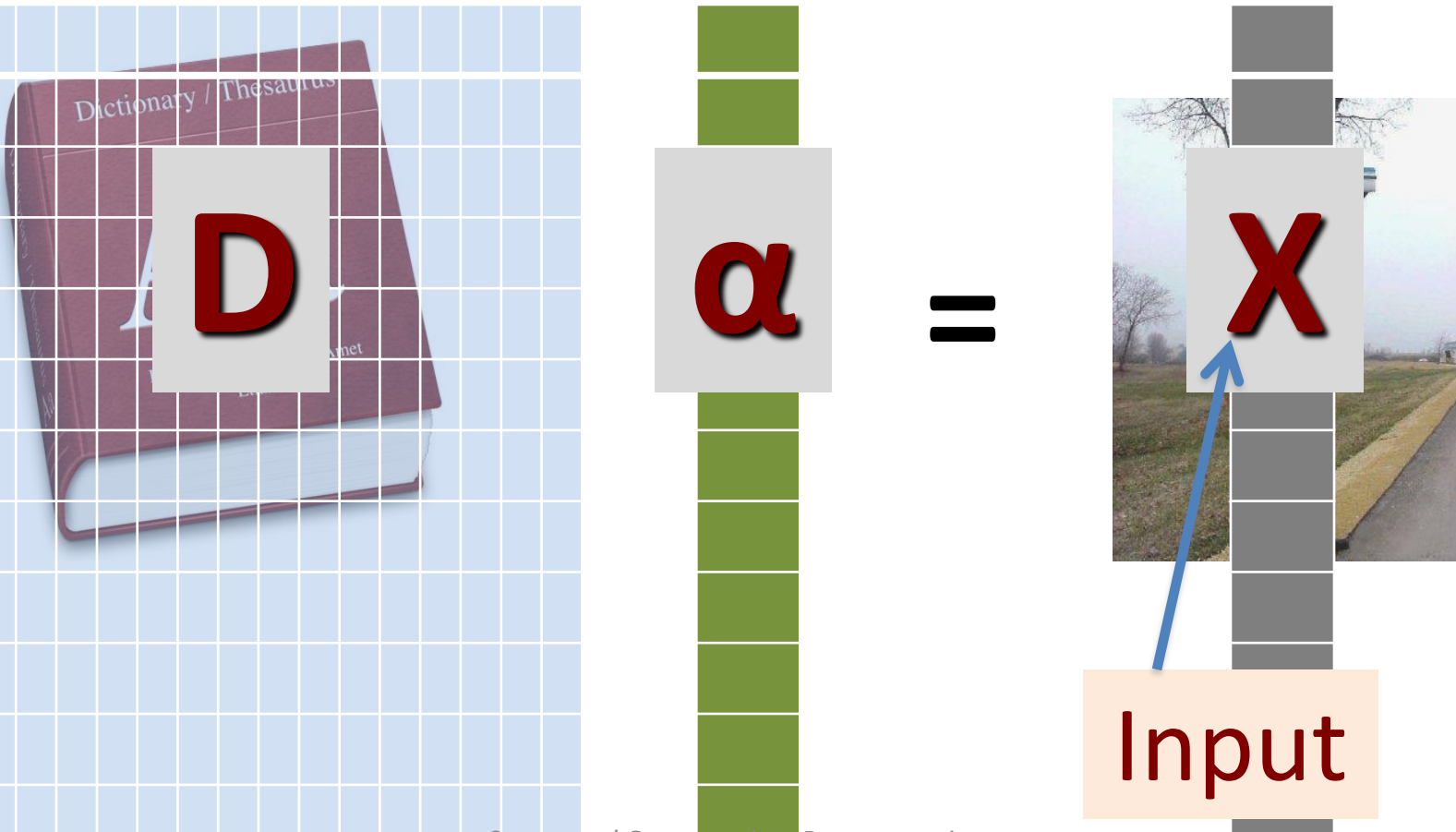
# K-SVD



- Termination of each iteration: Updated dictionary
- Conclusion: A dictionary where any data vector can be composed of at most  $K$  dictionary entries
  - More generally, sparse composition

# Problems

- How to obtain the dictionary
  - Which will give us meaningful representations
- How to compute the weights?





# Applications of Sparse Representations

- Many many applications
  - Signal representation
  - Statistical modelling
  - ..
  - We've seen one: Compressive sensing
- Another popular use
  - **Denoising**

# Denoising

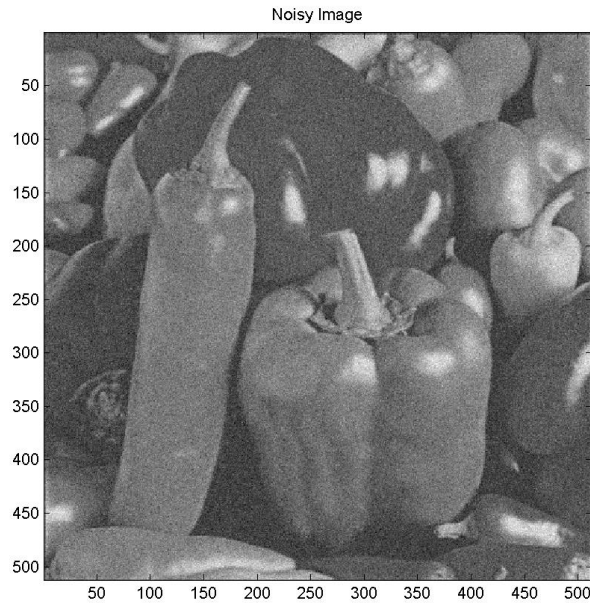
- As the name suggests, remove noise!

# Denoising

- As the name suggests, remove noise!
- We will look at image denoising as an example

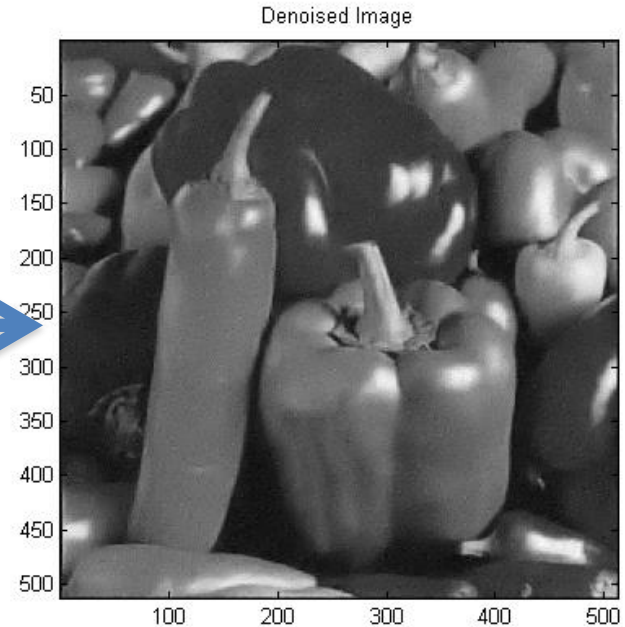
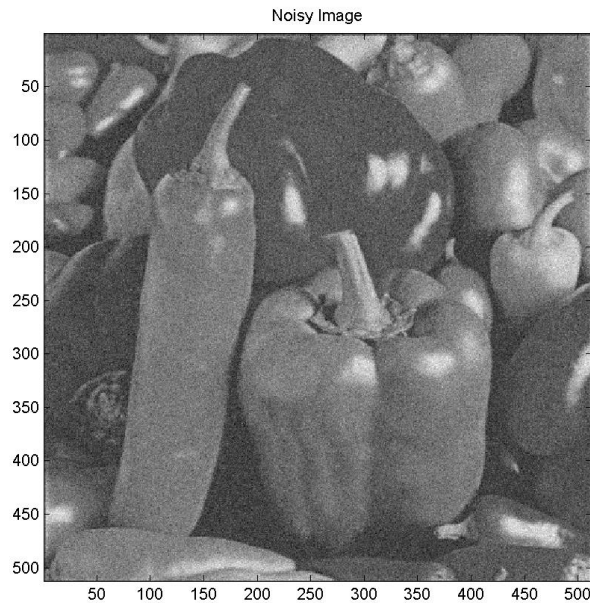
# Image Denoising

- Here's what we want



# Image Denoising

- Here's what we want



# Image Denoising

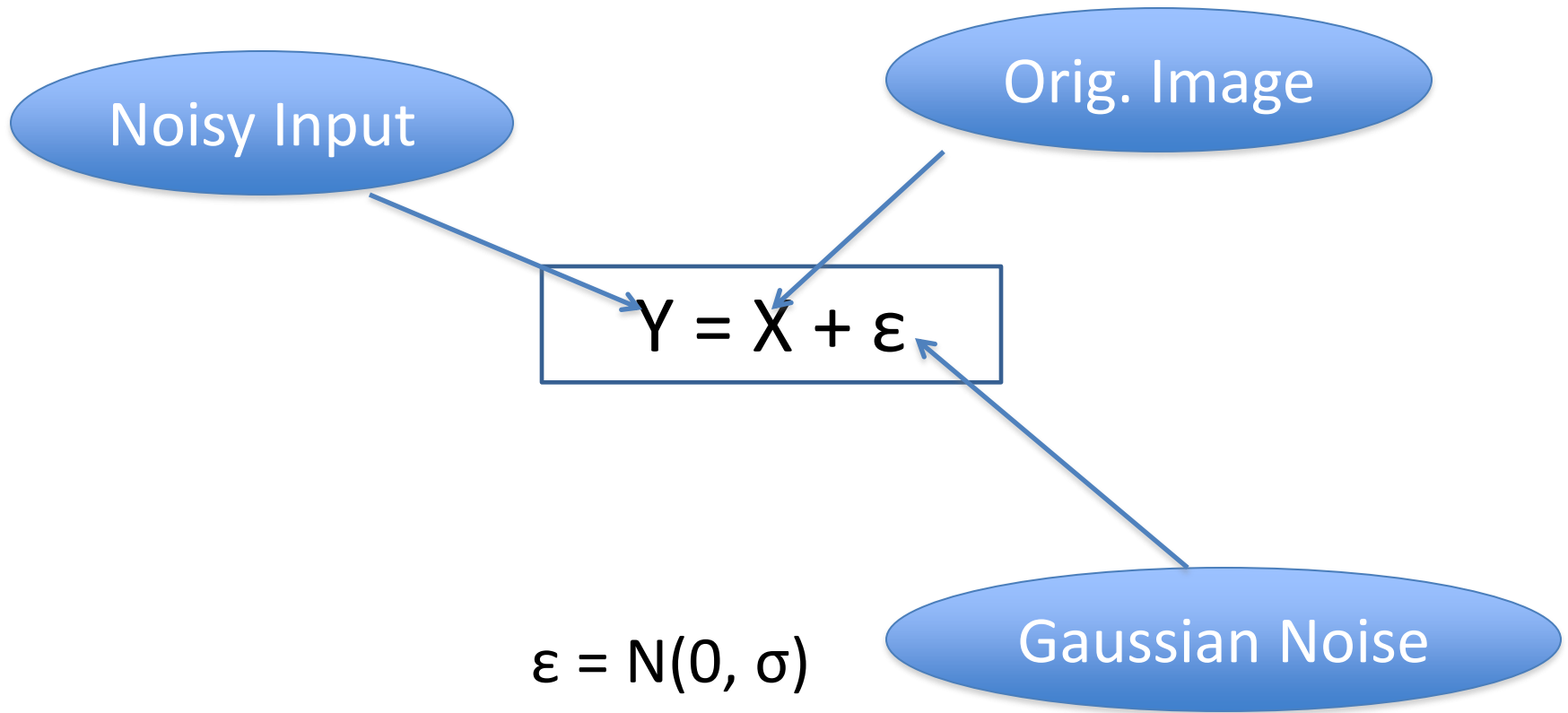
- Here's what we want



# The Image Denoising Problem

- Given an image
- Remove Gaussian additive noise from it

# Image Denoising





# Image Denoising

- Remove the noise from  $\mathbf{Y}$ , to obtain  $\mathbf{X}$  as best as possible.

# Image Denoising

- Remove the noise from  $\mathbf{Y}$ , to obtain  $\mathbf{X}$  as best as possible
- Using sparse representations over learned dictionaries

# Image Denoising

- Remove the noise from  $\mathbf{Y}$ , to obtain  $\mathbf{X}$  as best as possible
- Using sparse representations over learned dictionaries
- We will *learn* the dictionaries

# Image Denoising

- Remove the noise from  $\mathbf{Y}$ , to obtain  $\mathbf{X}$  as best as possible
- Using sparse representations over learned dictionaries
- We will *learn* the dictionaries
- **What data will we use?** *The corrupted image itself!*

# Image Denoising

- We use the data to be denoised to learn the dictionary.
- Training and denoising become an iterated process.
- We use image patches of size  $\sqrt{n} \times \sqrt{n}$  pixels (i.e. if the image is 64x64, patches are 8x8)

# Image Denoising

- The data dictionary  $D$ 
  - Size =  $n \times k$  ( $k > n$ )
  - This is known and fixed, to start with
  - Every image patch can be sparsely represented using  $D$

# Image Denoising

- Recall our equations from before.
- We want to find  $\underline{\alpha}$  so as to minimize the value of the equation below:

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_0 \}$$

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

# Image Denoising

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

- In the above,  $X$  is a patch.



# Image Denoising

$$\underset{\underline{\alpha}}{\text{Min}} \{ \|\underline{X} - \mathbf{D}\underline{\alpha}\|^2 + \lambda \|\underline{\alpha}\|_1 \}$$

- In the above,  $X$  is a patch.
- If the larger image is fully expressed by the every patch in it, how can we go from patches to the image?

# Image Denoising

$$\begin{aligned} \underset{\underline{\alpha}_{ij}, X}{Min} \{ & \mu \|\underline{X} - Y\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 \\ & + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \} \end{aligned}$$

# Image Denoising

$$\underset{\alpha_{ij}, X}{\text{Min}} \left\{ \mu \|\underline{X} - Y\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\}$$

$(X - Y)$  is the error between the input and denoised image.  $\mu$  is a penalty on the error.

# Image Denoising

$$\underset{\underline{\alpha}_{ij}, X}{Min} \left\{ \mu \left\| \underline{X} - Y \right\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\}$$

Error bounding in each patch

- $R_{ij}$  selects the  $(ij)^{\text{th}}$  patch
- Terms in summation = no. of patches

# Image Denoising

$$\underset{\underline{\alpha}_{ij}, X}{\text{Min}} \left\{ \mu \left\| \underline{X} - Y \right\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 \right.$$

$$\left. + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\}$$



$\lambda$  forces sparsity

# Image Denoising

- But, we don't "*know*" our dictionary  $D$ .
- We want to estimate  $D$  as well.

# Image Denoising

- But, we don't "*know*" our dictionary  $\mathbf{D}$ .
- We want to estimate  $\mathbf{D}$  as well.

$$\underset{D, \alpha_{ij}, X}{\text{Min}} \left\{ \mu \left\| \underline{X} - Y \right\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\}$$

We can use the previous equation itself!!!

# Image Denoising

$$\begin{aligned} \underline{\underset{D, \alpha_{ij}, X}{Min}} \quad & \left\{ \mu \left\| \underline{X} - Y \right\|_2^2 + \sum_{ij} \left\| \underline{R_{ij} X} - \mathbf{D} \underline{\alpha_{ij}} \right\|_2^2 \right. \\ & \left. + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha_{ij}} \right\|_0 \right\} \end{aligned}$$

How do we estimate all 3 at once?



# Image Denoising

$$\begin{aligned} \underline{\underset{D, \alpha_{ij}, X}{Min}} \quad & \left\{ \mu \left\| \underline{X} - Y \right\|_2^2 + \sum_{ij} \left\| \underline{R_{ij} X} - \mathbf{D} \underline{\alpha_{ij}} \right\|_2^2 \right. \\ & \left. + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha_{ij}} \right\|_0 \right\} \end{aligned}$$

How do we estimate all 3 at once?

**We cannot estimate them at the same time!**

# Image Denoising

$$\begin{aligned} \underline{\underset{D, \alpha_{ij}, X}{Min}} \{ & \mu \|\underline{X} - Y\|_2^2 + \sum_{ij} \left\| \underline{R_{ij} X} - \mathbf{D} \underline{\alpha_{ij}} \right\|_2^2 \\ & + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha_{ij}} \right\|_0 \} \end{aligned}$$

How do we estimate all 3 at once?

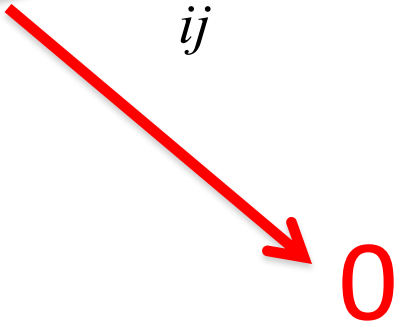
Fix 2, and find the optimal 3<sup>rd</sup>.

# Image Denoising

$$\underbrace{\text{Min}}_{D, \alpha_{ij}, X} \left\{ \mu \left\| \underline{X} - Y \right\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\}$$

Initialize  $X = Y$

# Image Denoising

$$\underset{\underline{\alpha}_{ij}}{\text{Min}} \left\{ \mu \left\| \underline{X} - \underline{Y} \right\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} \underline{X} - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\}$$


Initialize  $X = Y$ , initialize  $D$

You know how to solve the remaining portion for  $\alpha$  – MP, BP!

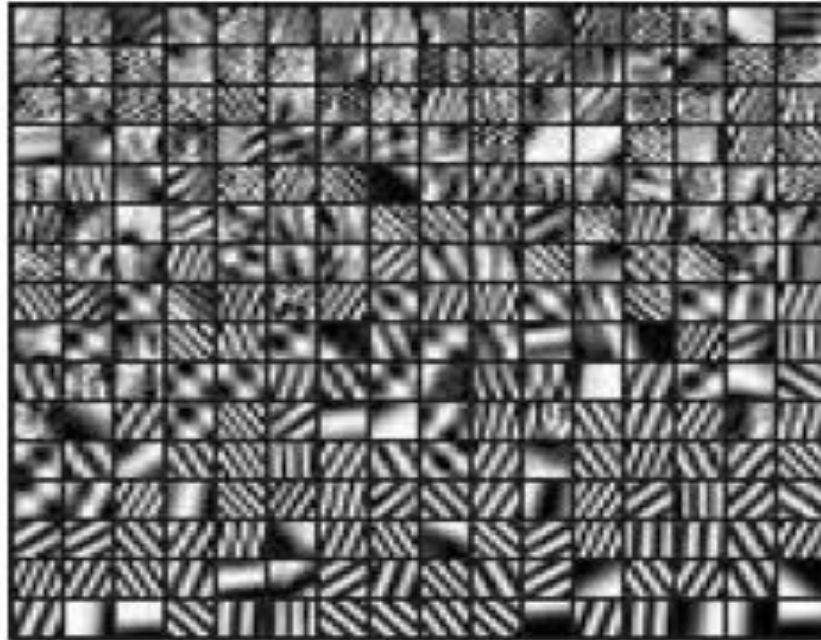
# Image Denoising

- Now, update the dictionary  $D$ .
- Update  $D$  one column at a time, following the [K-SVD algorithm](#)
- K-SVD maintains the sparsity structure

# Image Denoising

- Now, update the dictionary  $D$ .
- Update  $D$  one column at a time, following the [K-SVD algorithm](#)
- K-SVD maintains the sparsity structure
- Iteratively update  $\alpha$  and  $D$

# Image Denoising



## Learned Dictionary for Face Image denoising

From: M. Elad and M. Aharon, *Image denoising via learned dictionaries and sparse representation*, CVPR, 2006.

# Image Denoising

$$\underset{X}{\text{Min}} \left\{ \mu \|\underline{X} - Y\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 \right. \\ \left. + \sum_{ij} \lambda_{ij} \left\| \underline{\alpha}_{ij} \right\|_0 \right\} \rightarrow \text{Const. wrt } X$$

We know  $\mathbf{D}$  and  $\alpha$

The quadratic term above has a closed-form solution



# Image Denoising

$$\underset{X}{\text{Min}} \left\{ \mu \|\underline{X} - Y\|_2^2 + \sum_{ij} \left\| \underline{R}_{ij} X - \mathbf{D} \underline{\alpha}_{ij} \right\|_2^2 \right. \\ \left. + \sum_{ij} \lambda_{ij} \|\underline{\alpha}_{ij}\|_0 \right\} \rightarrow \text{Const. wrt } X$$

We know  $D$  and  $\alpha$

$$X = (\mu I + \sum_{ij} R_{ij}^T R_{ij})^{-1} (\mu Y + \sum_{ij} R_{ij}^T D \alpha_{ij})$$

# Image Denoising

- Summarizing... We wanted to obtain 3 things

# Image Denoising

- Summarizing... We wanted to obtain 3 things
  - Weights  $\alpha$
  - Dictionary  $D$
  - Denoised Image  $X$

# Image Denoising

- Summarizing... We wanted to obtain 3 things
  - Weights  $\alpha$  – Your favorite pursuit algorithm
  - Dictionary  $\mathbf{D}$  – Using K-SVD
  - Denoised Image  $\mathbf{X}$

# Image Denoising

- Summarizing... We wanted to obtain 3 things
  - Weights  $\alpha$  – Your favorite pursuit algorithm
  - Dictionary  $\mathbf{D}$  – Using K-SVD
  - Denoised Image  $\mathbf{X}$

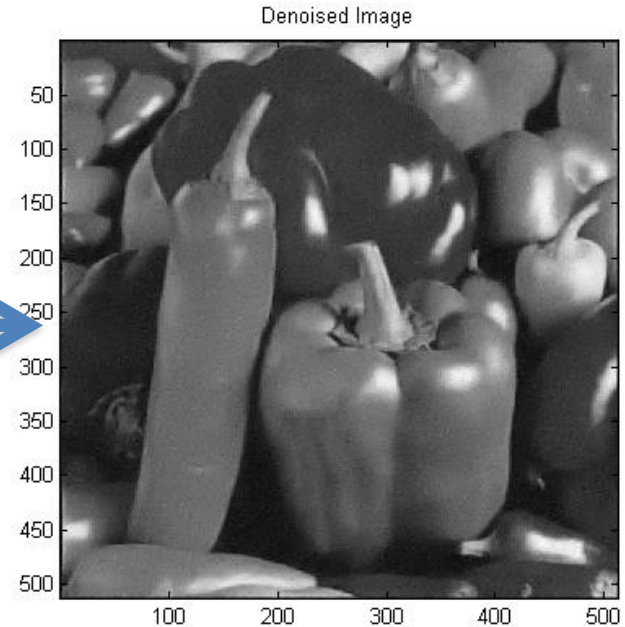
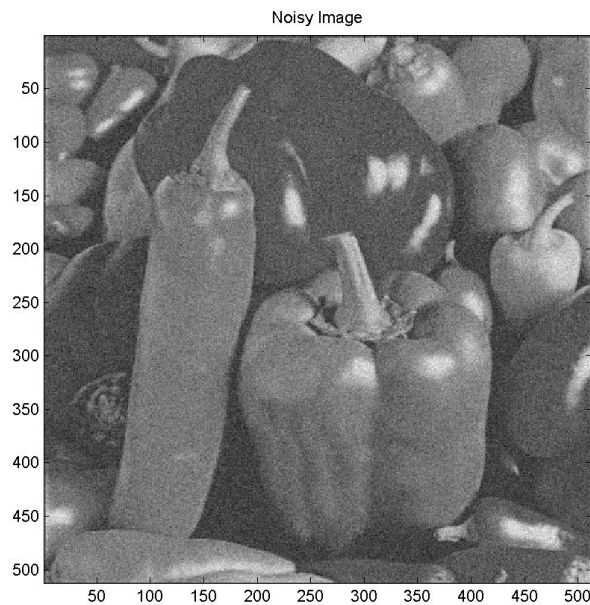
Iterating

# Image Denoising

- Summarizing... We wanted to obtain 3 things
  - Weights  $\alpha$
  - Dictionary  $\mathbf{D}$
  - Denoised Image  $\mathbf{X}$ - Closed form solution

# Image Denoising

- Here's what we want



# Image Denoising

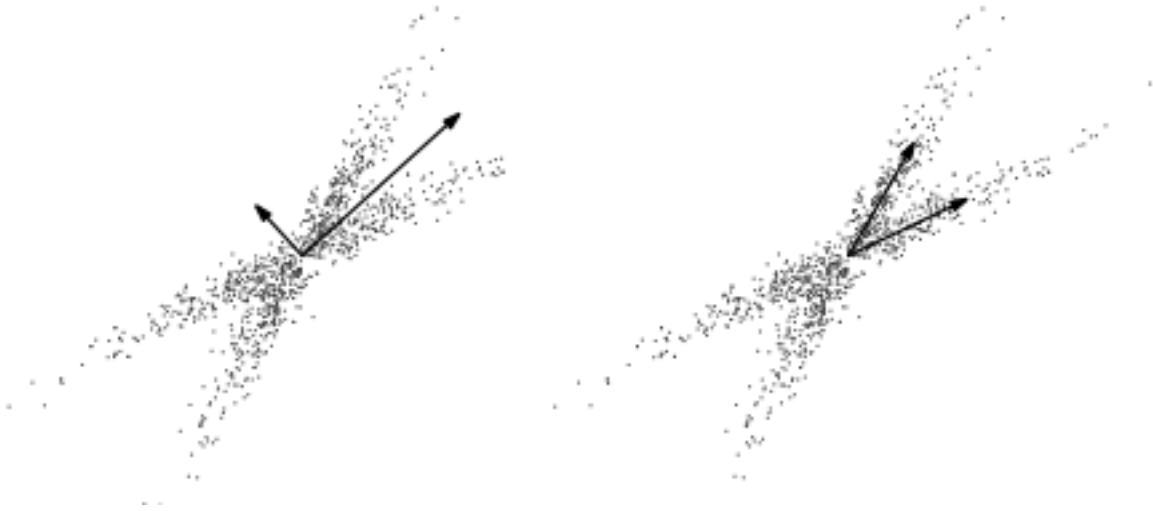
- Here's what we want





# Comparing to Other Techniques

Non-Gaussian data



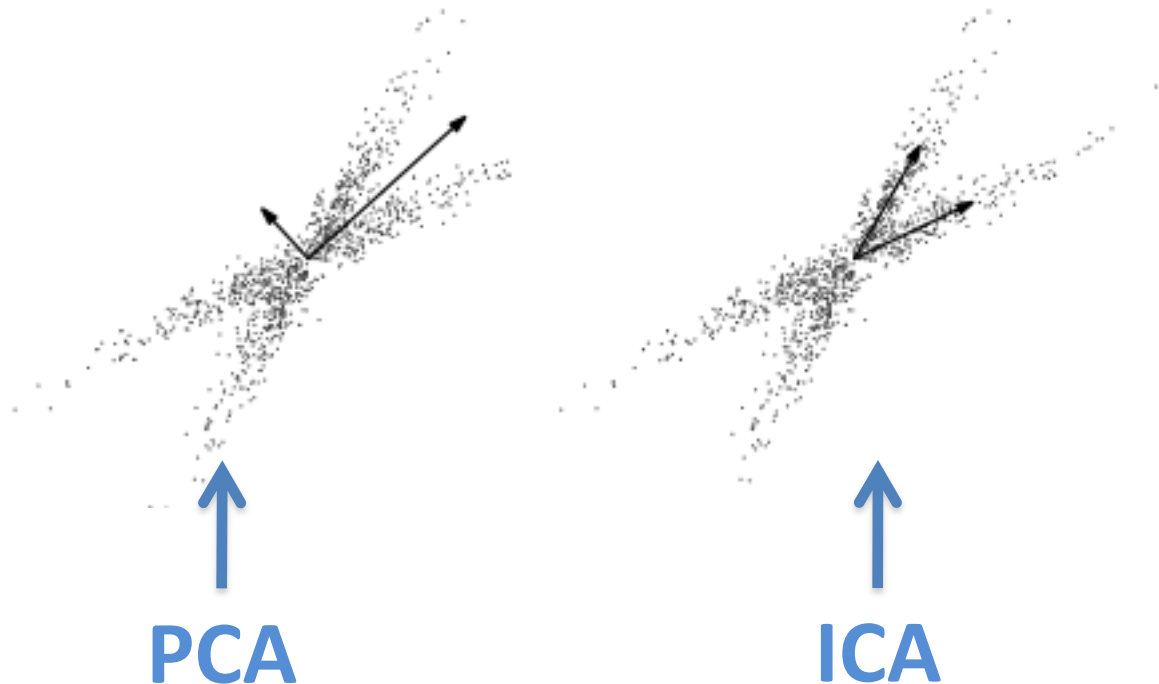
PCA of ICA

Which is which?

Images from Lewicki and Sejnowski, *Learning Overcomplete Representations*, 2000.

# Comparing to Other Techniques

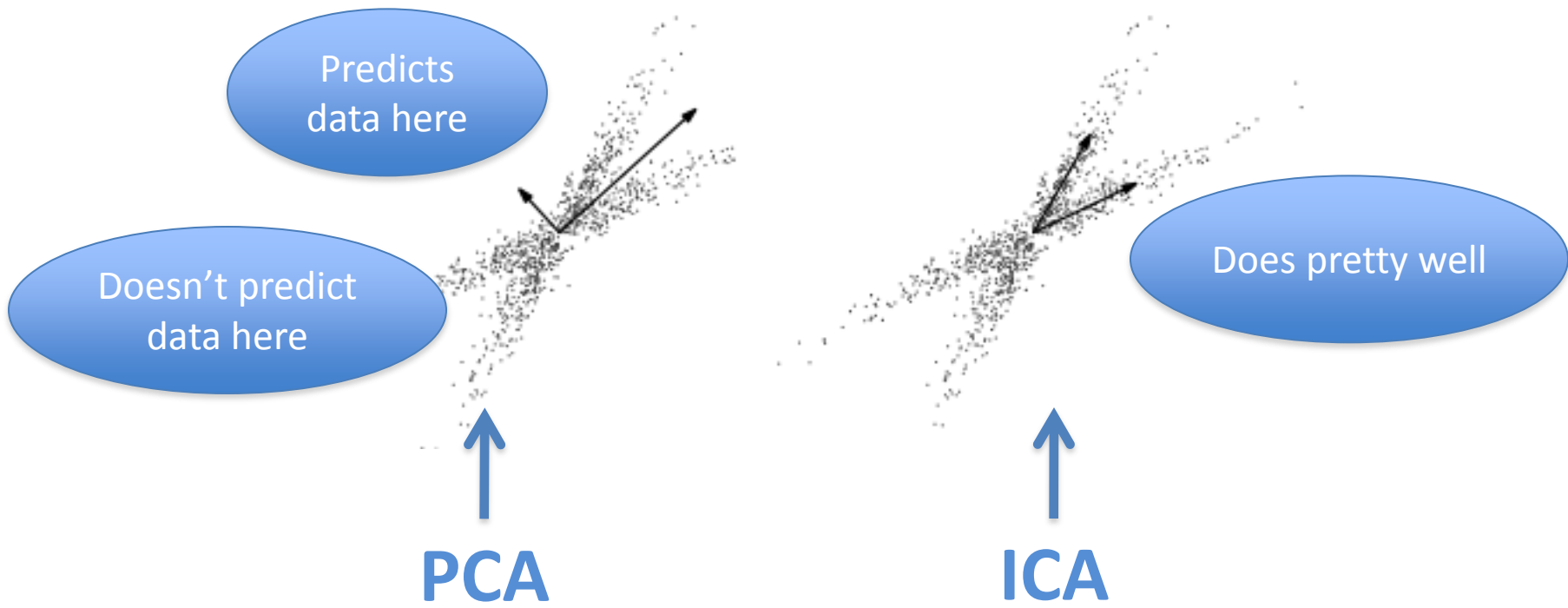
Non-Gaussian data



Images from Lewicki and Sejnowski, *Learning Overcomplete Representations*, 2000.

# Comparing to Other Techniques

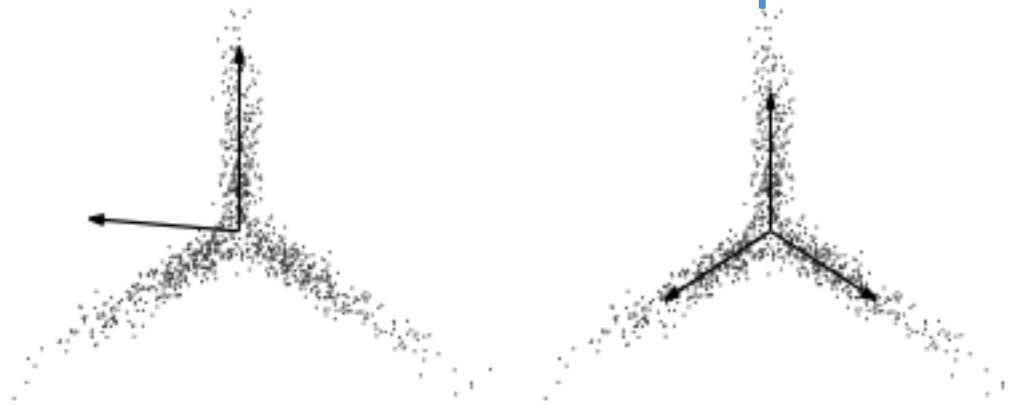
## Non-Gaussian data



Images from Lewicki and Sejnowski, *Learning Overcomplete Representations*, 2000.

# Comparing to Other Techniques

Data still in 2-D space



**ICA**

**Overcomplete**

Doesn't capture the underlying representation,  
which Overcomplete representations can do...

# Summary

- Overcomplete representations can be more powerful than component analysis techniques.
- Dictionary can be learned from data.
- Relative advantages and disadvantages of the pursuit algorithms.