

Machine Learning for Signal Processing

Sparse and Overcomplete Representations

Bhiksha Raj
(slides from Sourish Chaudhuri and
Abelino Jimenez)

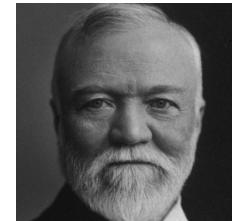
So far

Can we use linear composition to identify
basic units that compose the signal?

So far

$$D \cdot \alpha = X \approx \text{Data}$$

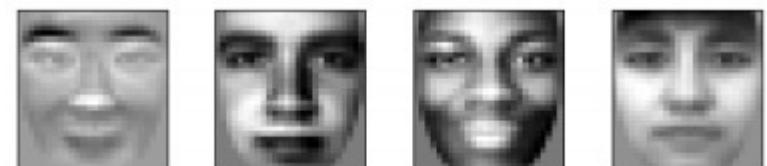
Weights
↓
Basis
↑



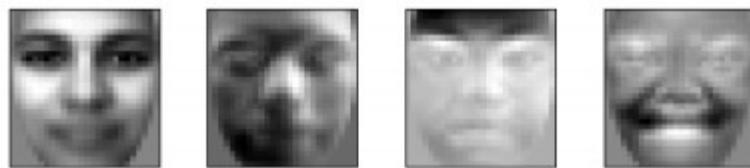
Data Independent



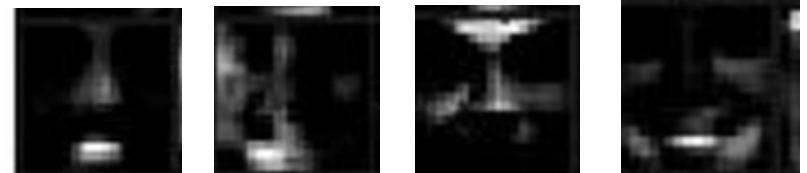
ICA



PCA



NNMF

# basis \leq dim X

Just in case you missed it..

- Remember, #(Basis Vectors) = #unknowns

$$D \cdot \alpha = X$$

Basis Vectors Weights Input data

Standard representations: number of bases \leq dimension of data

A limitation we saw earlier

- Mathematical restrictions on the number of bases have no connection to reality
 - Universe does not respect your mathematical representations of the data
 - In reality: number of building blocks that compose any kind of data is unlimited
- One solution we saw earlier: picking *one* “closest” building block to represent any input
- Today: Learning linear compositional representations without restrictions on the number of basic units

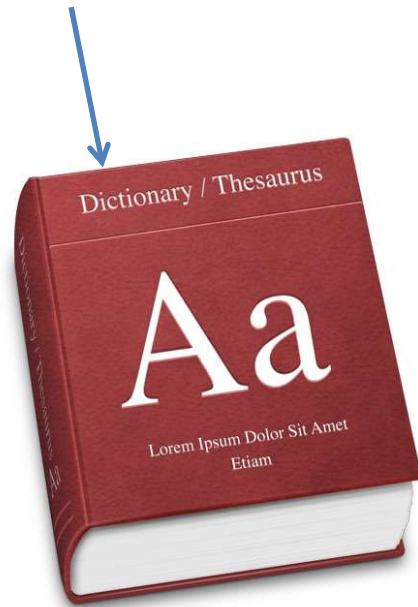
Poll 1

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?

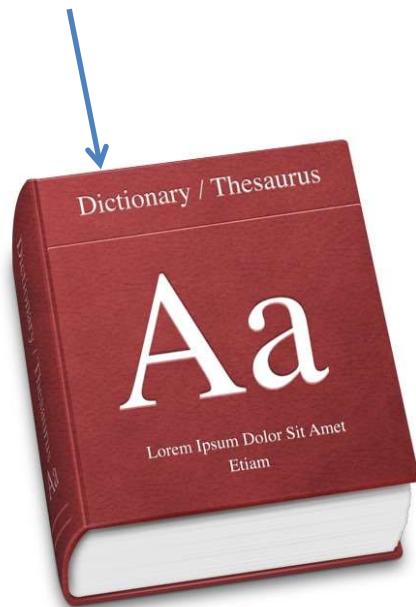
Representing Data

Dictionary (codebook)

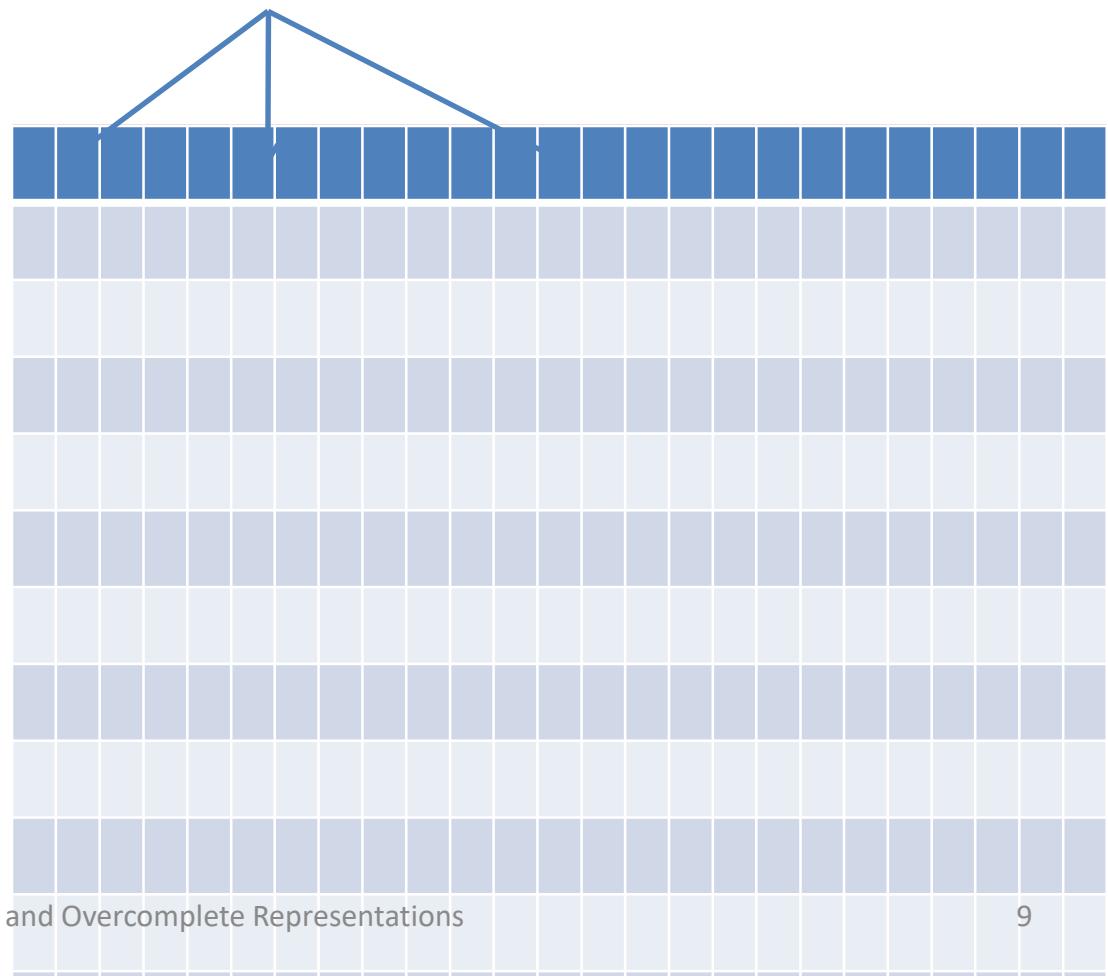


Representing Data

Dictionary

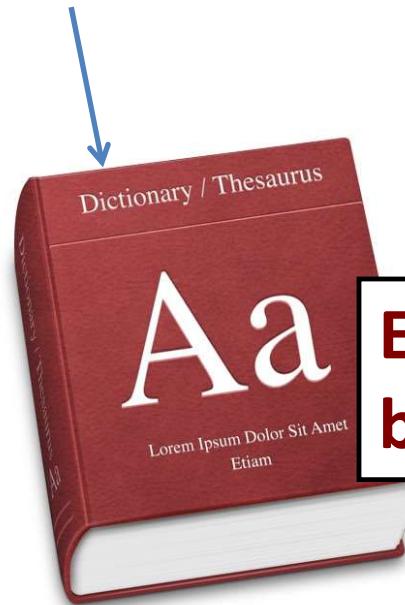


Atoms

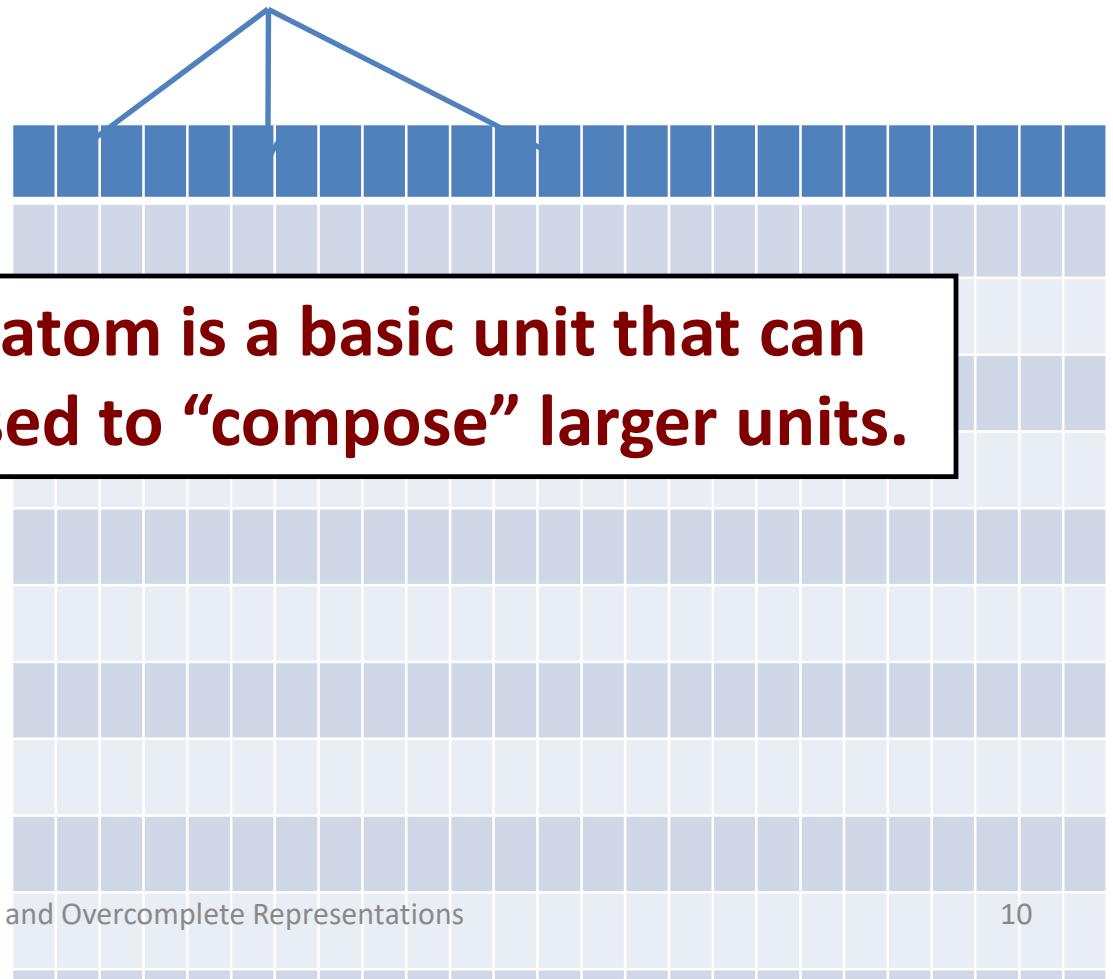


Representing Data

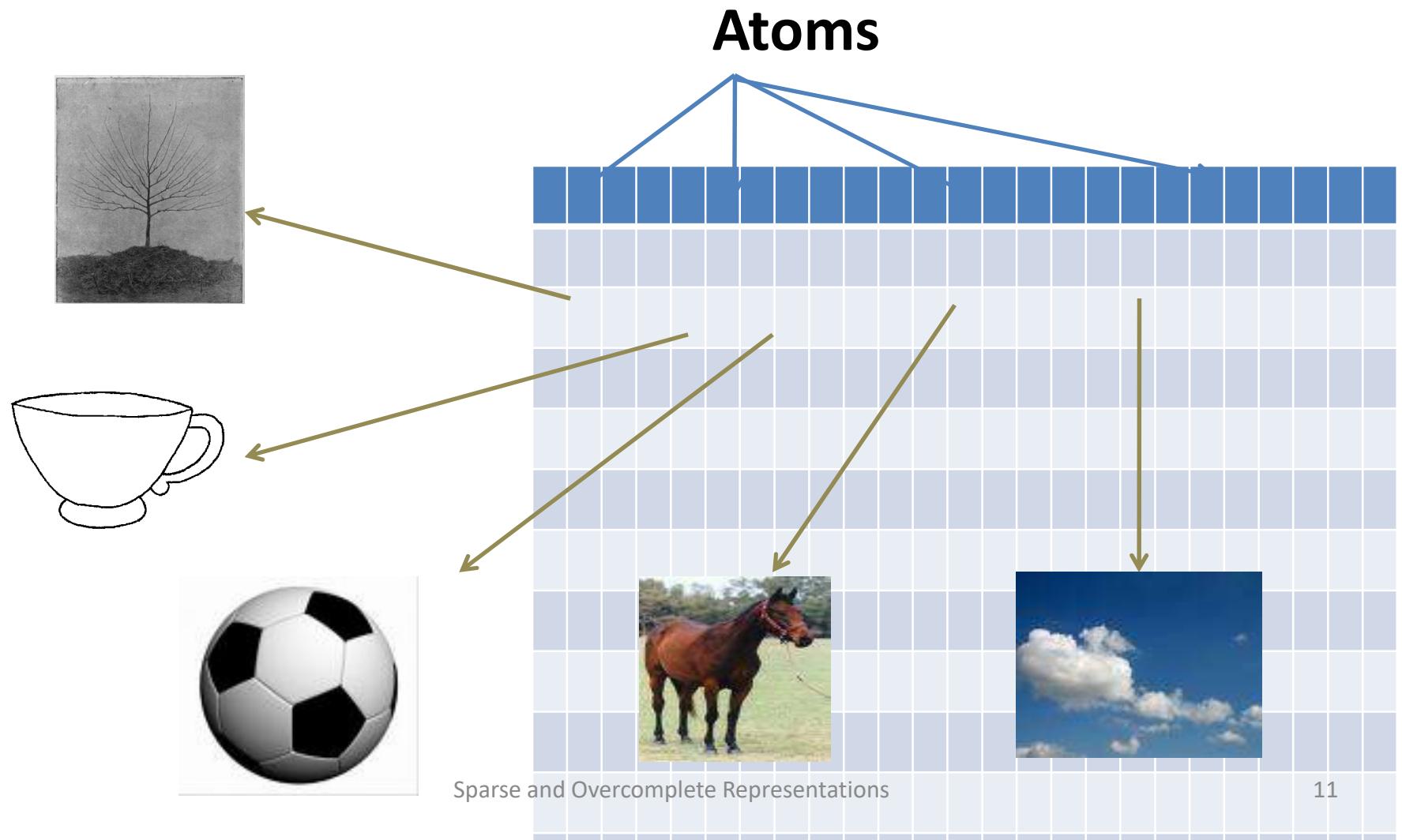
Dictionary



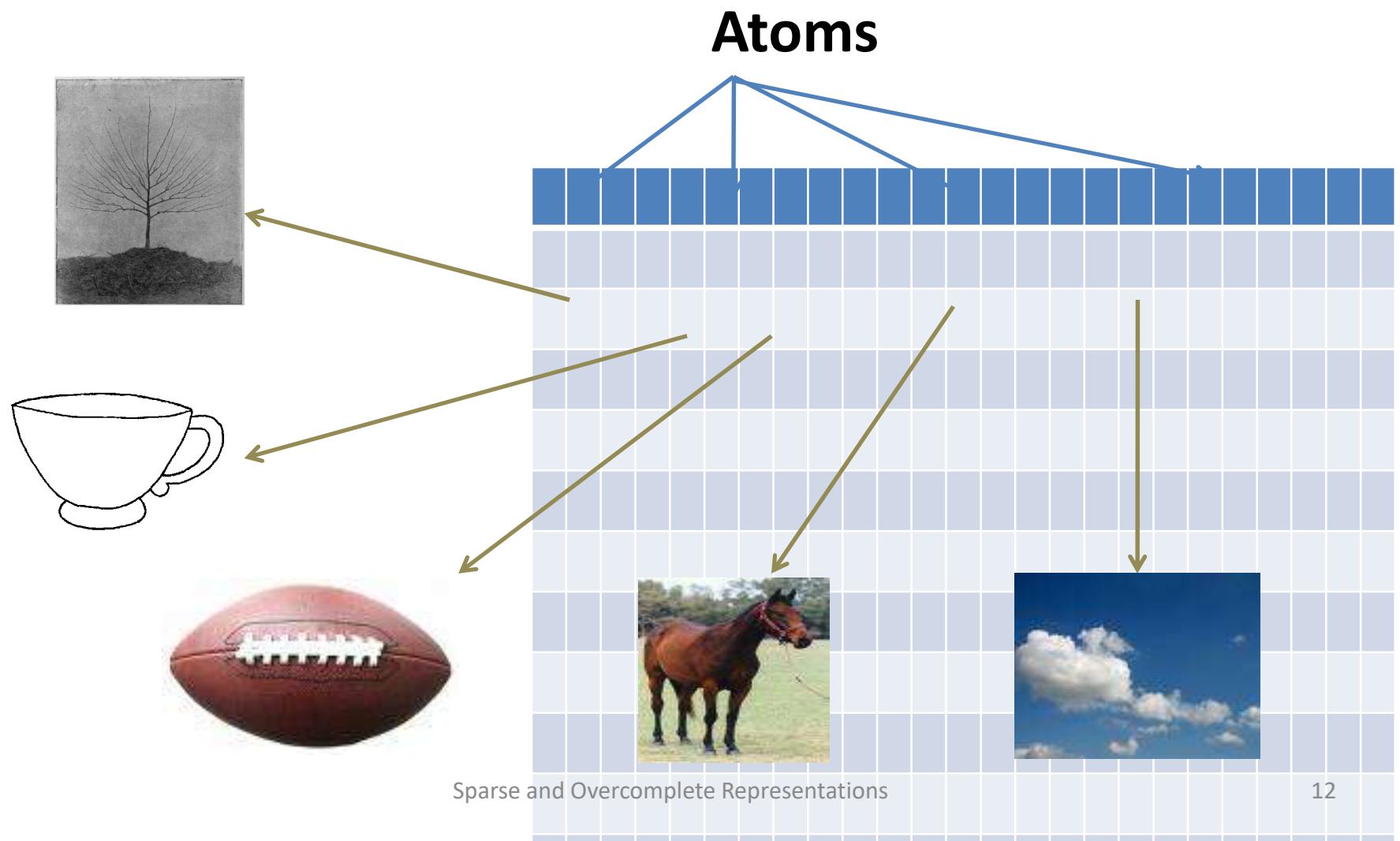
Atoms



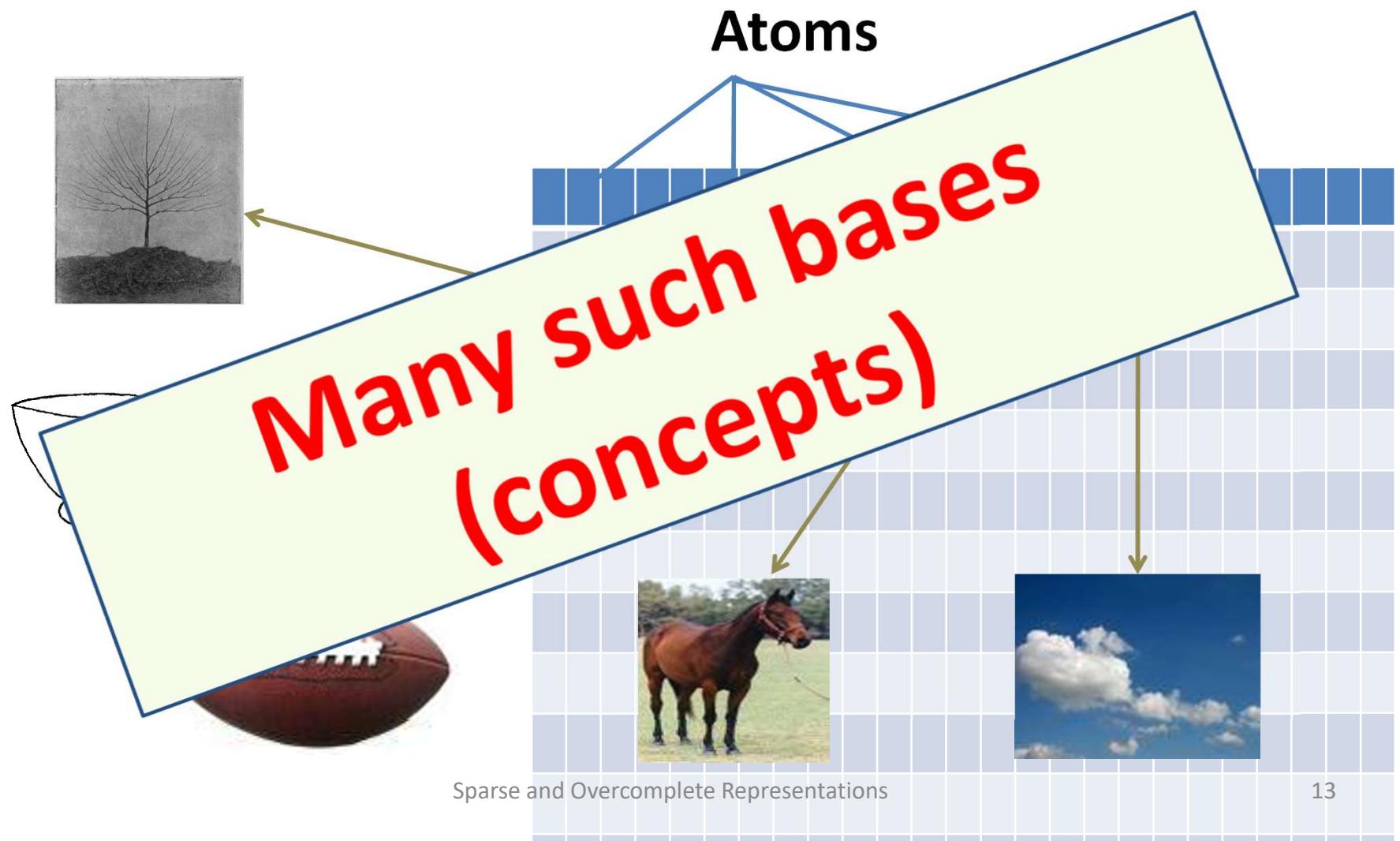
Representing Data



Representing Data



Representing Data



Representing Data



Sparse and Overcomplete Representations

Representing Data

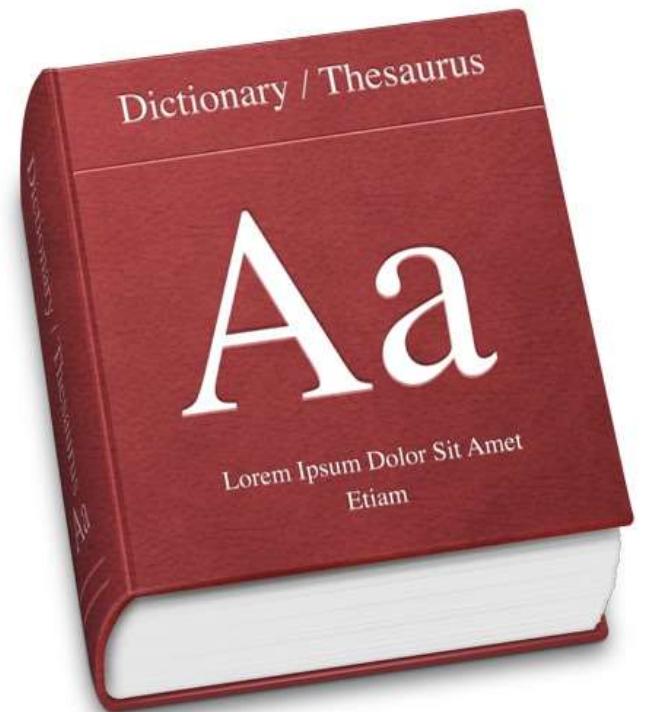


Sparse and Overcomplete Representations

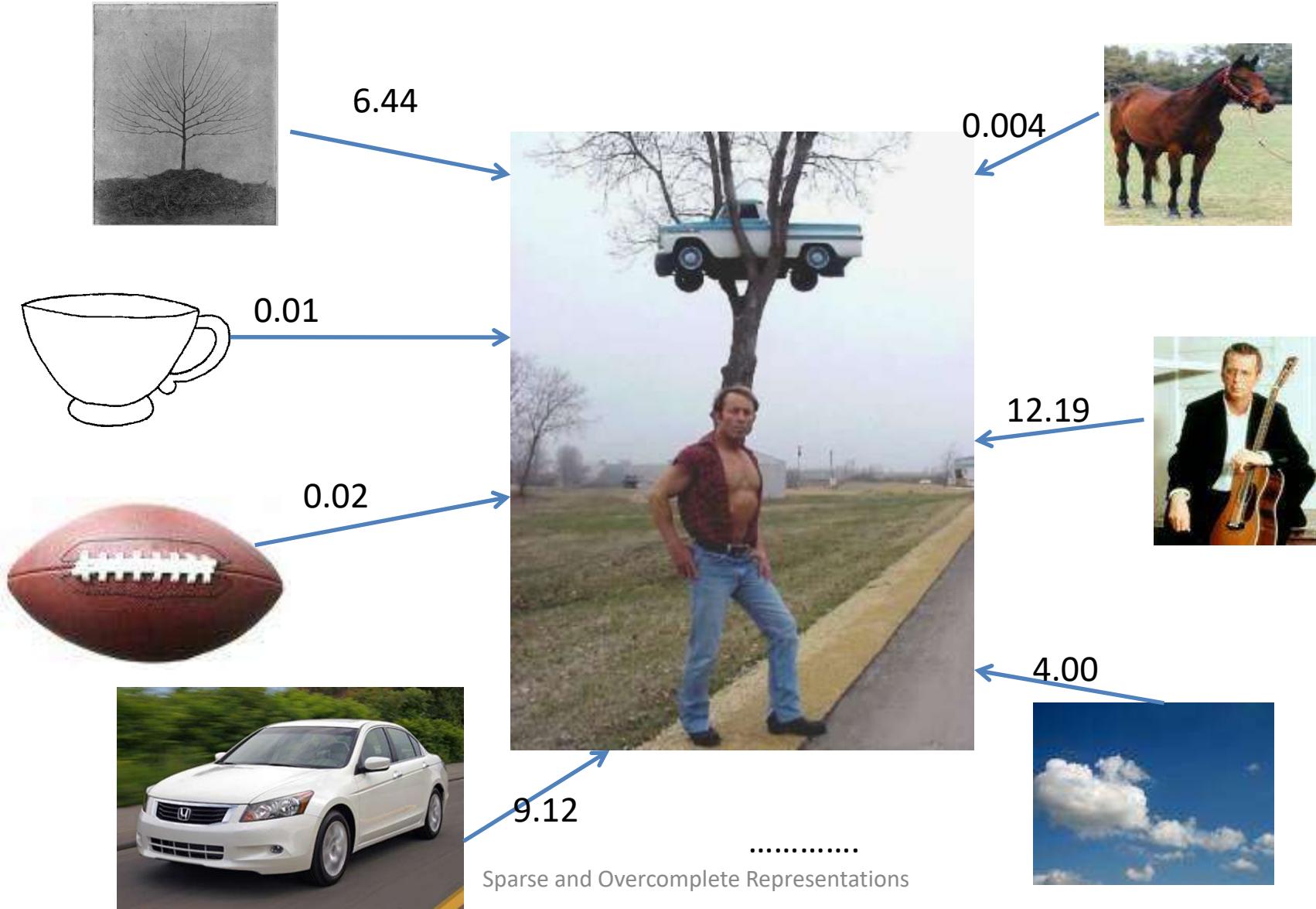
Representing Data



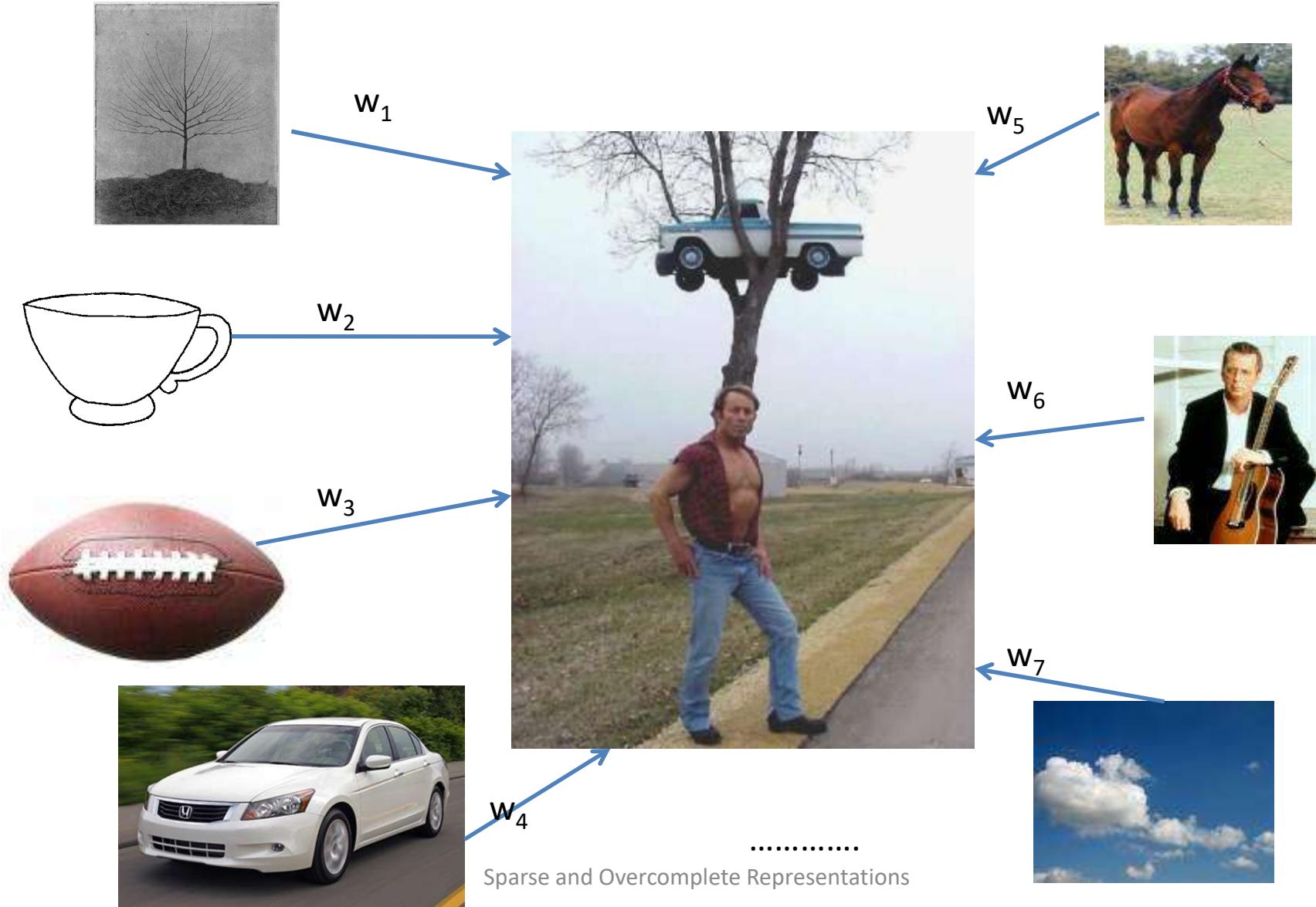
Using concepts that we
know...



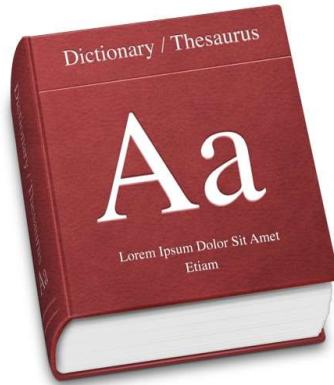
Representing Data



Representing Data



Representing Data

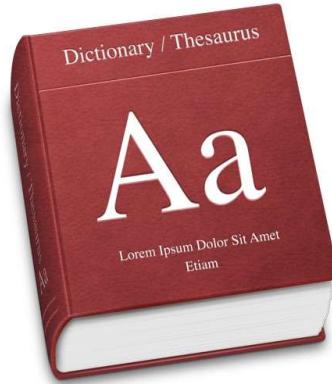


Linear
combination of
elements in the
Dictionary

=

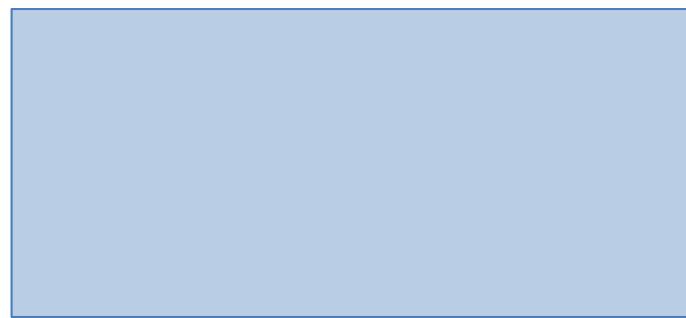


Representing Data



Linear
combination of
elements in the
Dictionary

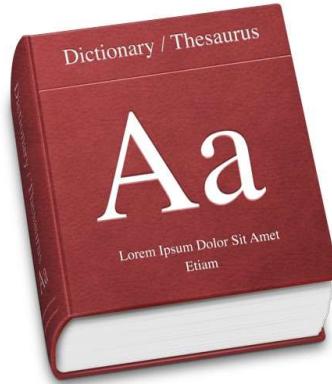
=



=



Representing Data



Linear combination of elements in the Dictionary

=



$$D \alpha = X$$

Quick Linear Algebra Refresher

- Remember, #(Basis Vectors) = #unknowns

$$D \cdot \alpha = X$$

The diagram illustrates the components of the equation $D \cdot \alpha = X$. The term X is labeled "Input data". The term α is labeled "Weights". The term D is labeled "Basis Vectors". Arrows point from each label to its corresponding term in the equation.

Basis Vectors

Weights

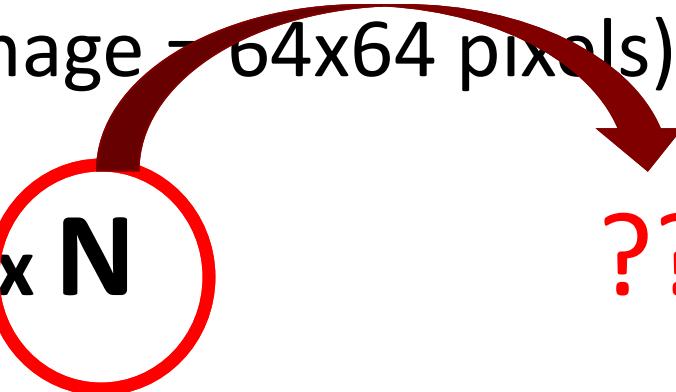
Input data

(from Dictionary)

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)
 - $4096 \times N$

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 64x64 image)

If $N > 4096$ (as it likely is)
we have an **overcomplete** representation

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

➤ $4096 \times N$

VERY LARGE!!!

Overcomplete Representations

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

More generally:

If #(dictionary units) > dimensions of input

- we have an **overcomplete** representation

➤ 4096 x N

VERY LARGE!!!

Quick Linear Algebra Refresher

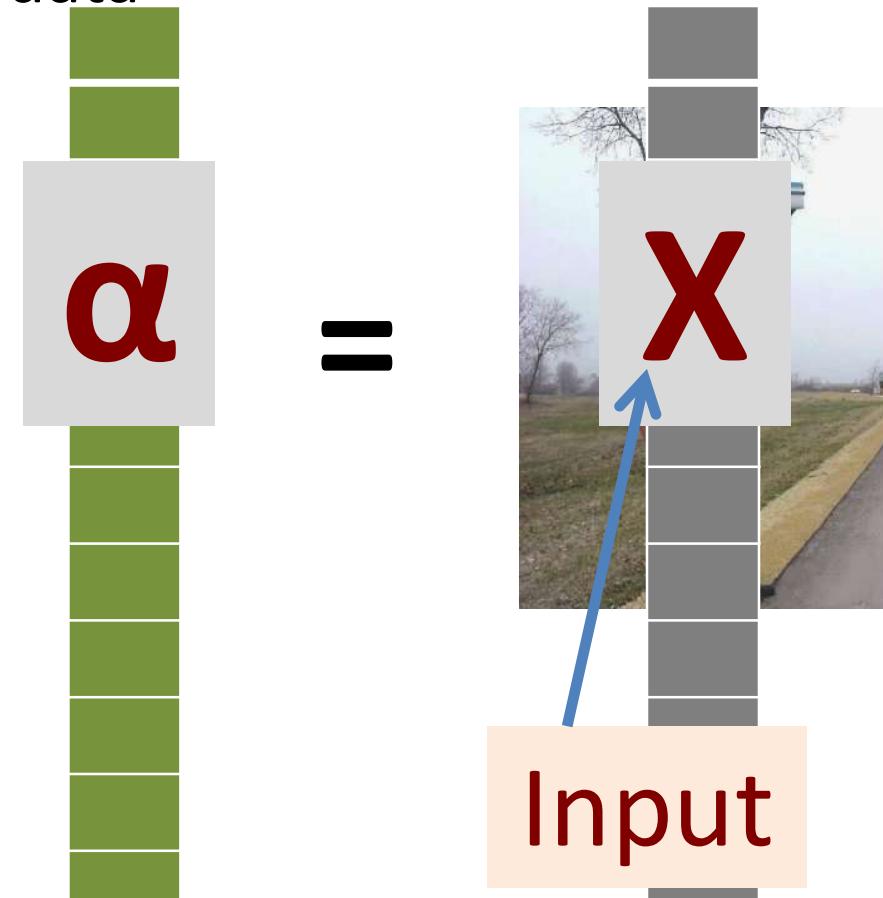
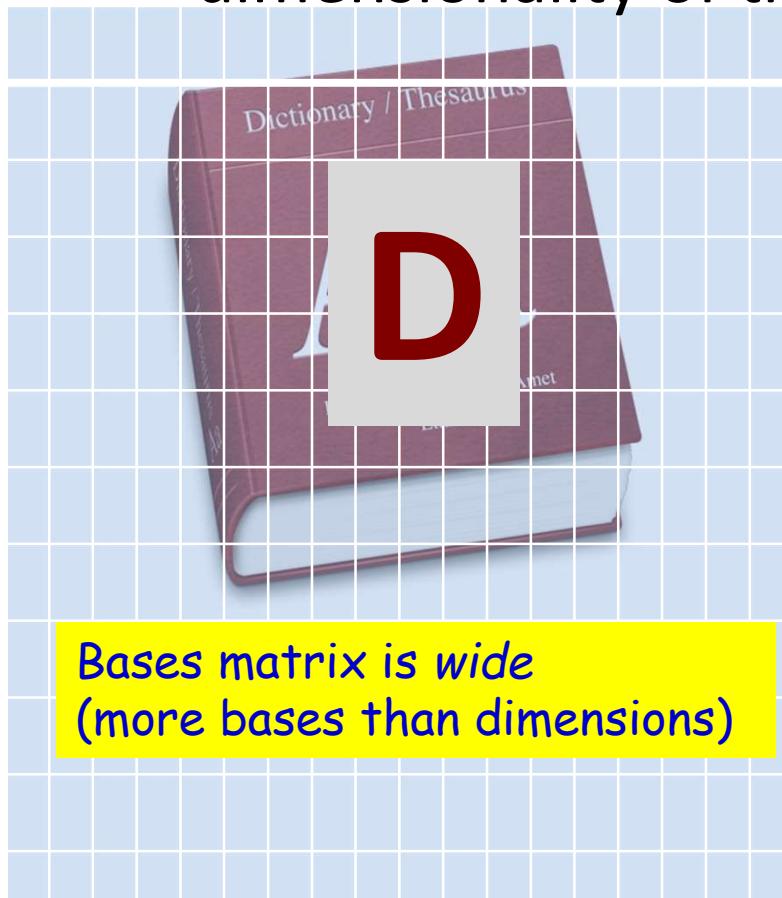
- Remember, #(Basis Vectors) = #unknowns

$$D \cdot \alpha = X$$

The diagram illustrates the components of the equation $D \cdot \alpha = X$. A yellow box labeled "Dictionary Units" contains two arrows: one pointing to the matrix D and another pointing to the vector α . Below the vector α , the word "Weights" is written. To the right of the equation, an arrow points from the word "Input data" to the vector X .

Dictionary based Representations

- Overcomplete “dictionary”-based representations are linear-composition-based representations with more “atomic building blocks” than the dimensionality of the data



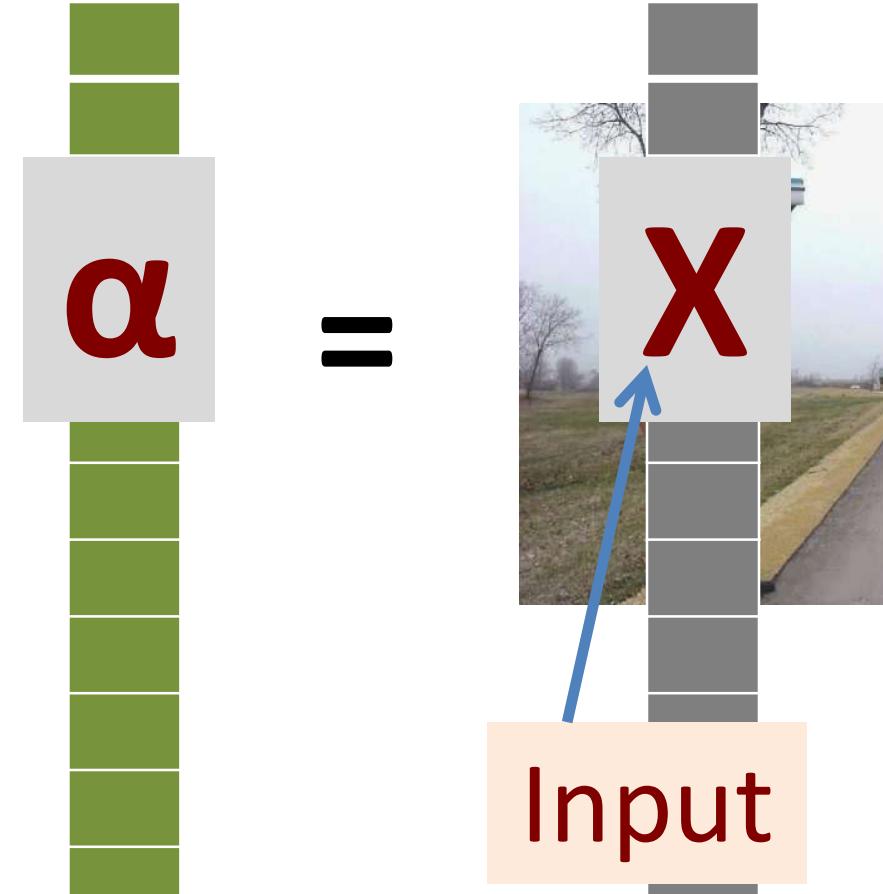
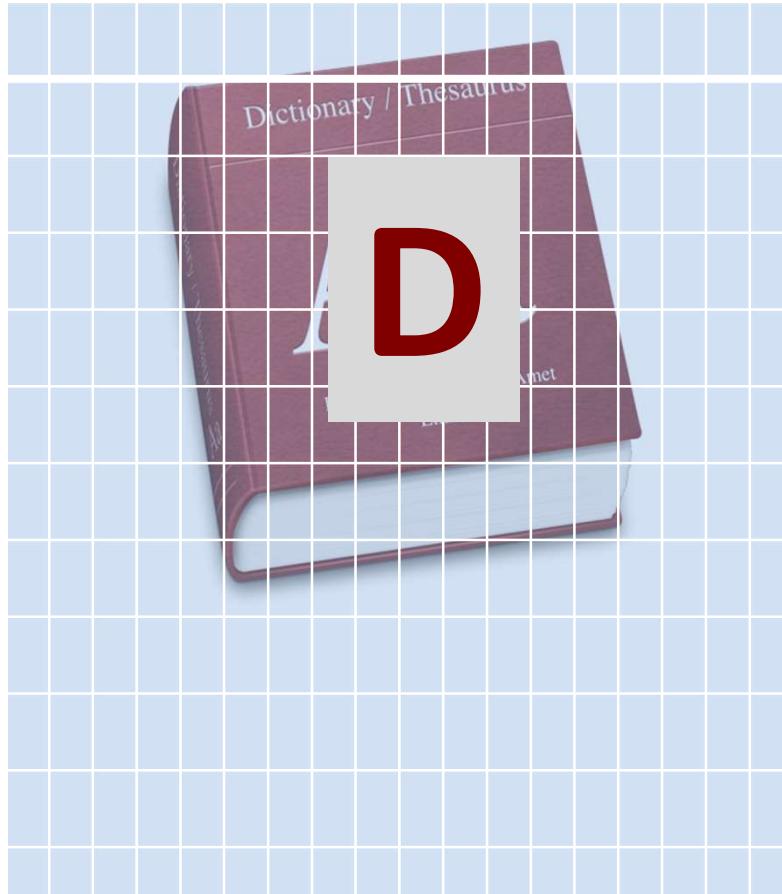
Why Dictionary-based Representations?

- 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 shortly..

Poll 2

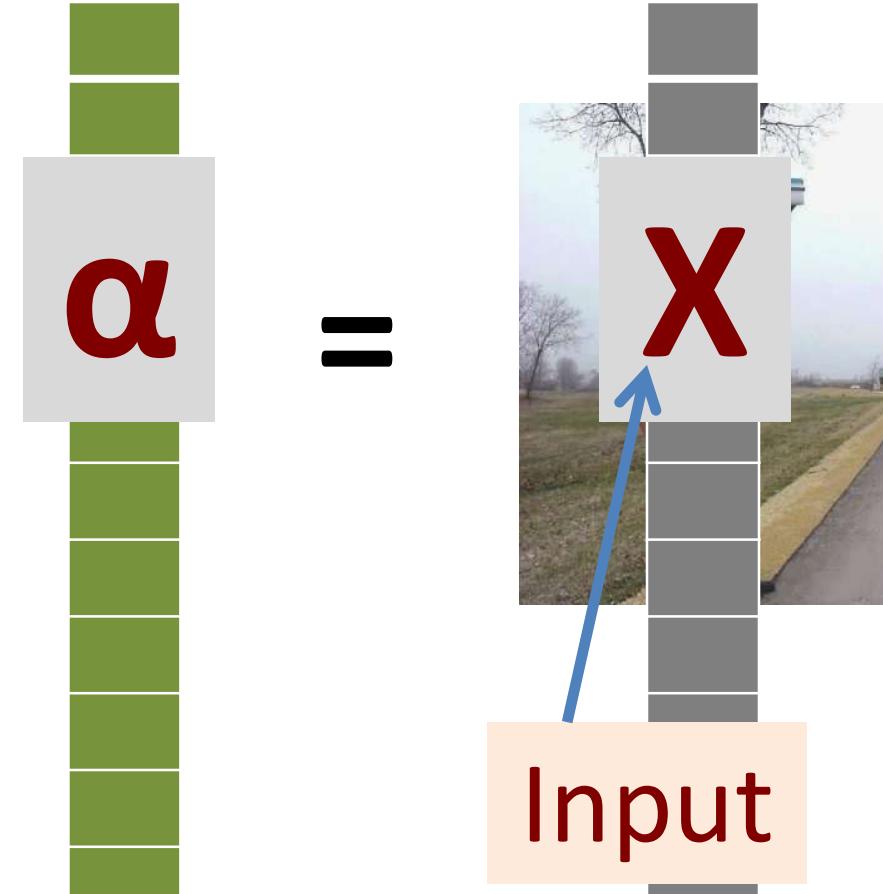
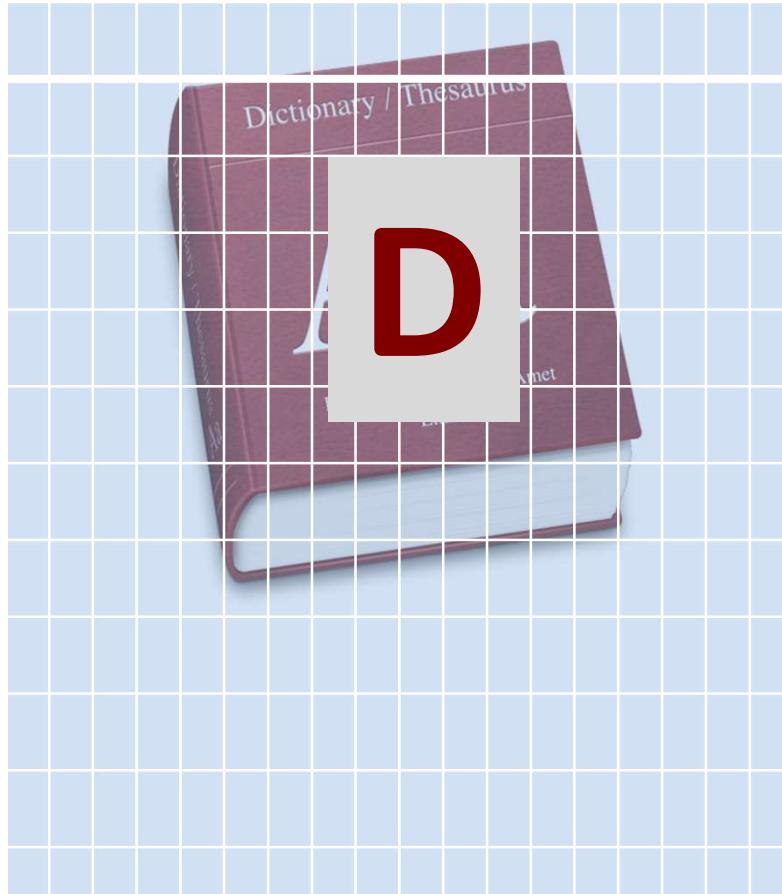
Problems

- 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?



Quick Linear Algebra Refresher

- Remember, #(Basis Vectors) = #unknowns

$$D \cdot \alpha = X$$

The diagram illustrates the components of the equation $D \cdot \alpha = X$. The matrix D is labeled "Dictionary entries". The vector α is labeled "Weights". The vector X is labeled "Input data". Blue arrows point from each label to its corresponding term in the equation.

When can we solve for α ?

Quick Linear Algebra Refresher

$$D \cdot \alpha = X$$

D: full rank

A diagram illustrating a linear equation. On the left, there is a blue rectangle labeled D . To its right is a green rectangle labeled α . To the right of that is a yellow rectangle labeled X . Between the D and the α , and between the α and the X , there is a black equals sign ($=$).

Unique solution

A diagram illustrating a linear equation. On the left, there is a blue rectangle labeled D . To its right is a green rectangle labeled α . To the right of that is a yellow rectangle labeled X . Between the D and the α , and between the α and the X , there is a black equals sign ($=$).

We may have no exact solution

A diagram illustrating a linear equation. On the left, there is a blue rectangle labeled D . To its right is a green rectangle labeled α . To the right of that is a yellow rectangle labeled X . Between the D and the α , and between the α and the X , there is a black equals sign ($=$).

Infinite Solutions

Quick Linear Algebra Refresher

$$D \cdot \alpha = X$$

D: full rank

$$D \quad \alpha = X$$

Unique solution

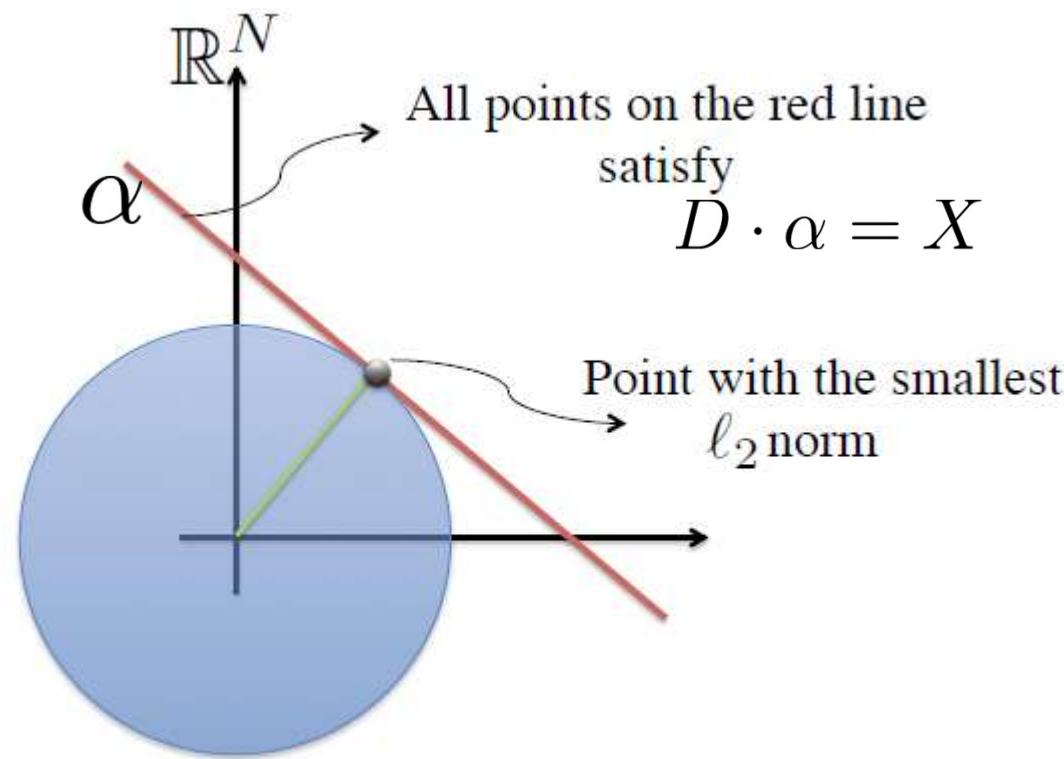
$$D \quad \alpha = X$$

We may have no exact solution

$$D \quad \alpha = X$$

Our Case
Infinite Solutions

Using Pseudo-Inverse?

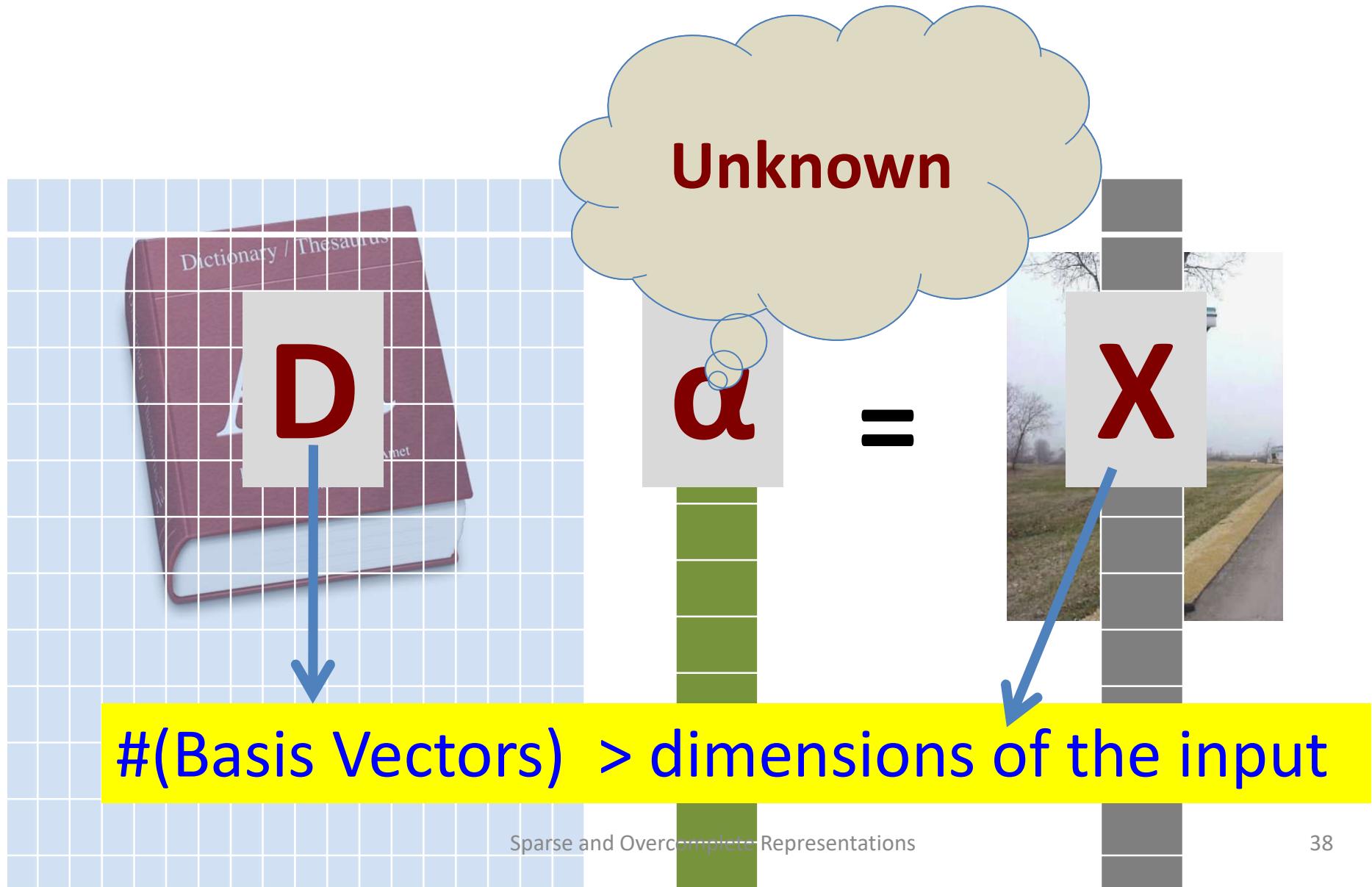


This is equivalent to

$$\text{minimize } \|\alpha\|_2 \text{ subject to } D\alpha = X$$

α will generally be “dense”

Overcomplete Representation



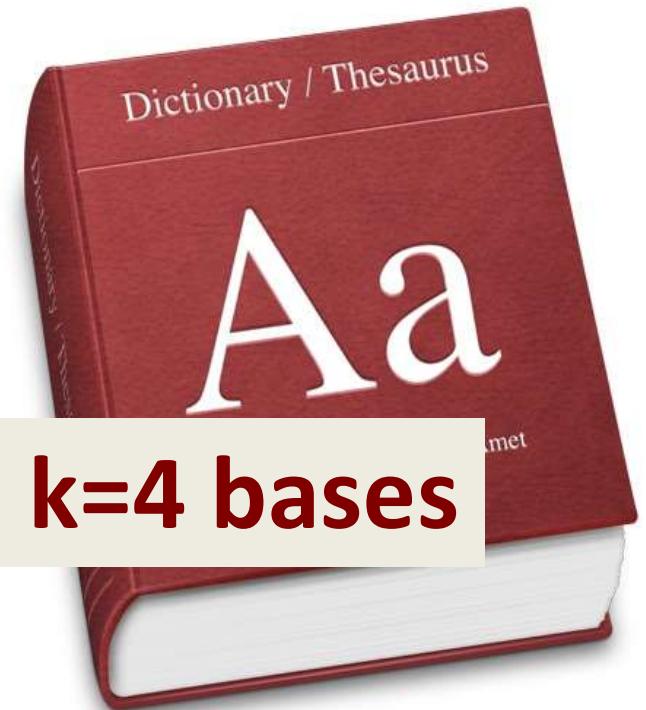
Representing Data



Using bases that we
know...



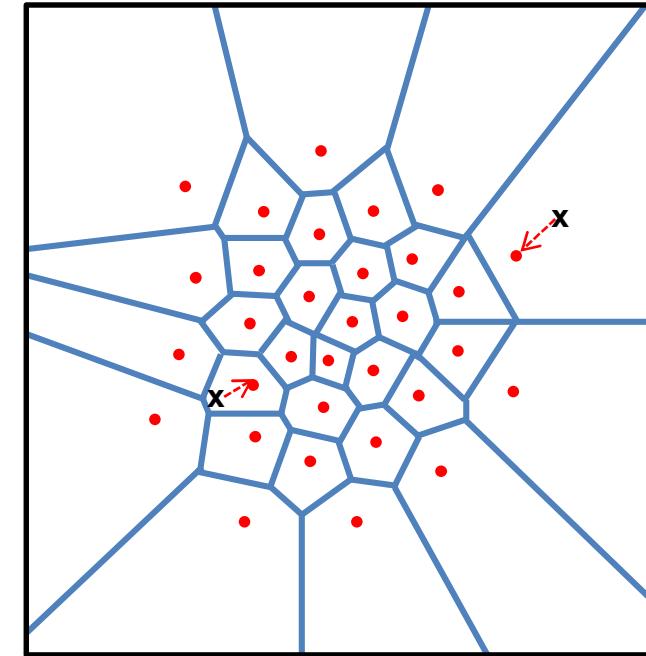
But no more than $k=4$ bases



Alternate view: Recall quantization

$$V = \sum_i w_i d_i$$

$$V = \mathbf{D}\mathbf{w} \quad |\mathbf{w}| = 1 \\ |\mathbf{w}|_0 = 1$$



- d_i are the “representative” vectors of each cluster
- Restriction: only one of the w_i is 1, the rest are 0
 - $\sum_i w_i = 0$
 - \mathbf{w} is unit length and one-sparse
- What if we let *more* than one entry of \mathbf{w} to be non zero?

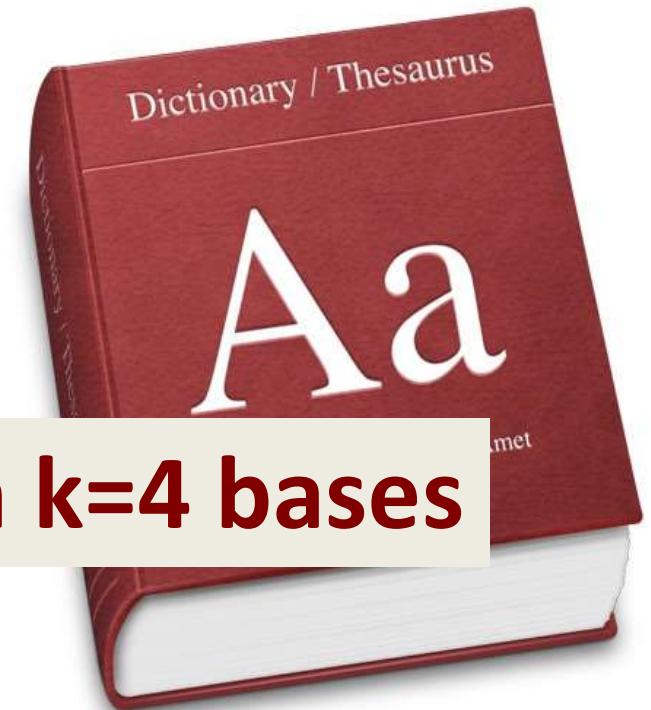
Overcompleteness and Sparsity

- To solve an overcomplete system of the type:

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

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

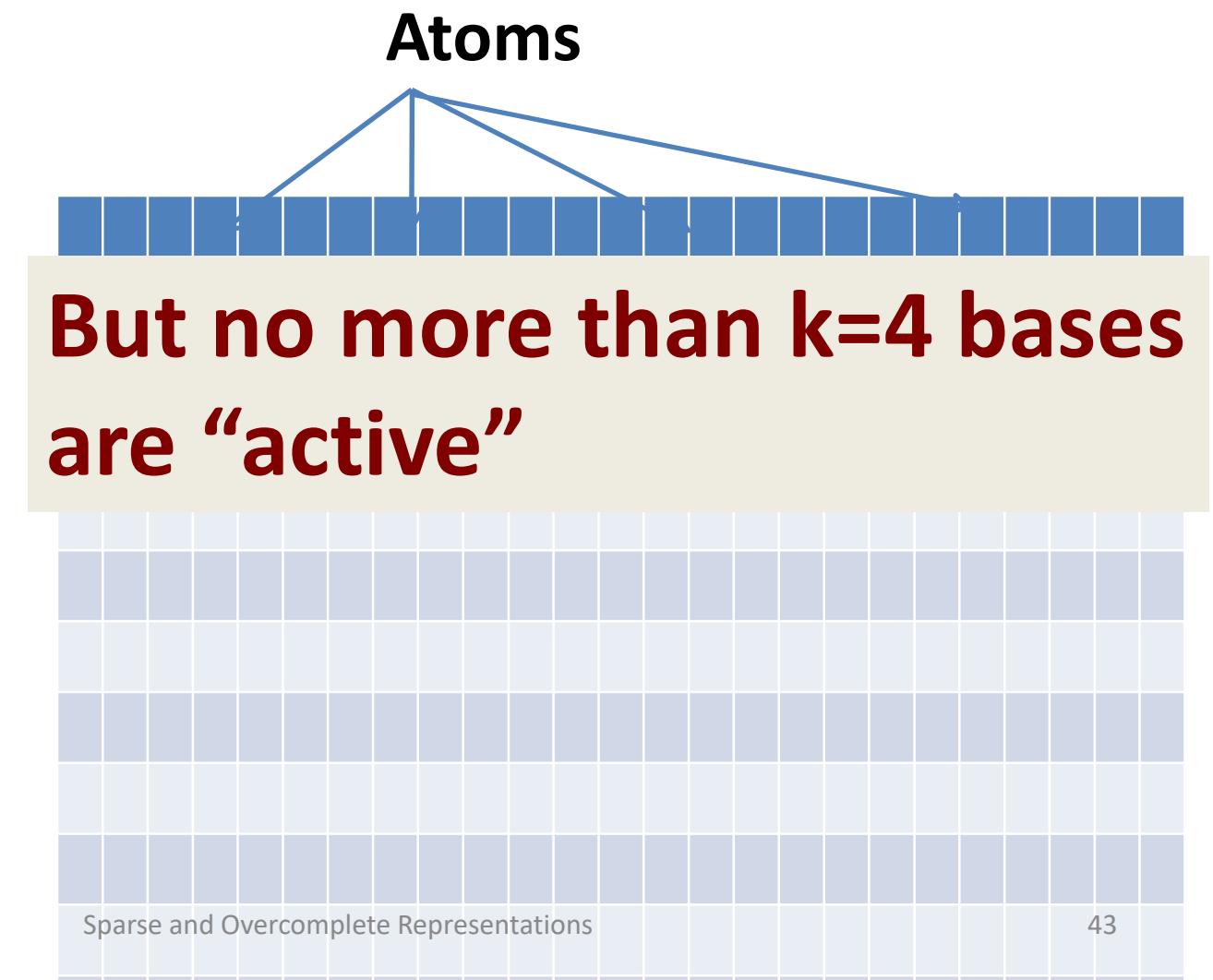
Representing Data



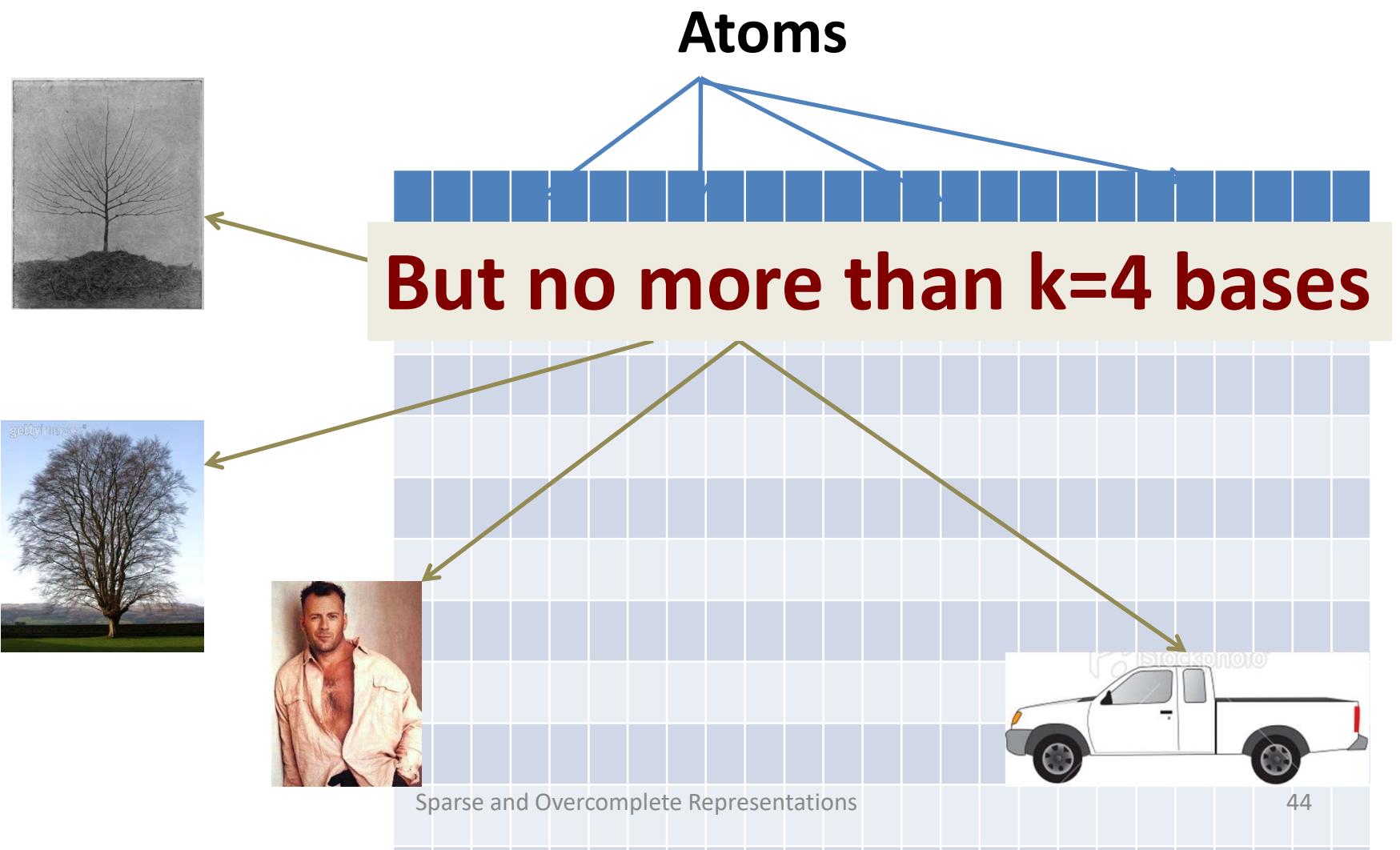
But no more than $k=4$ bases



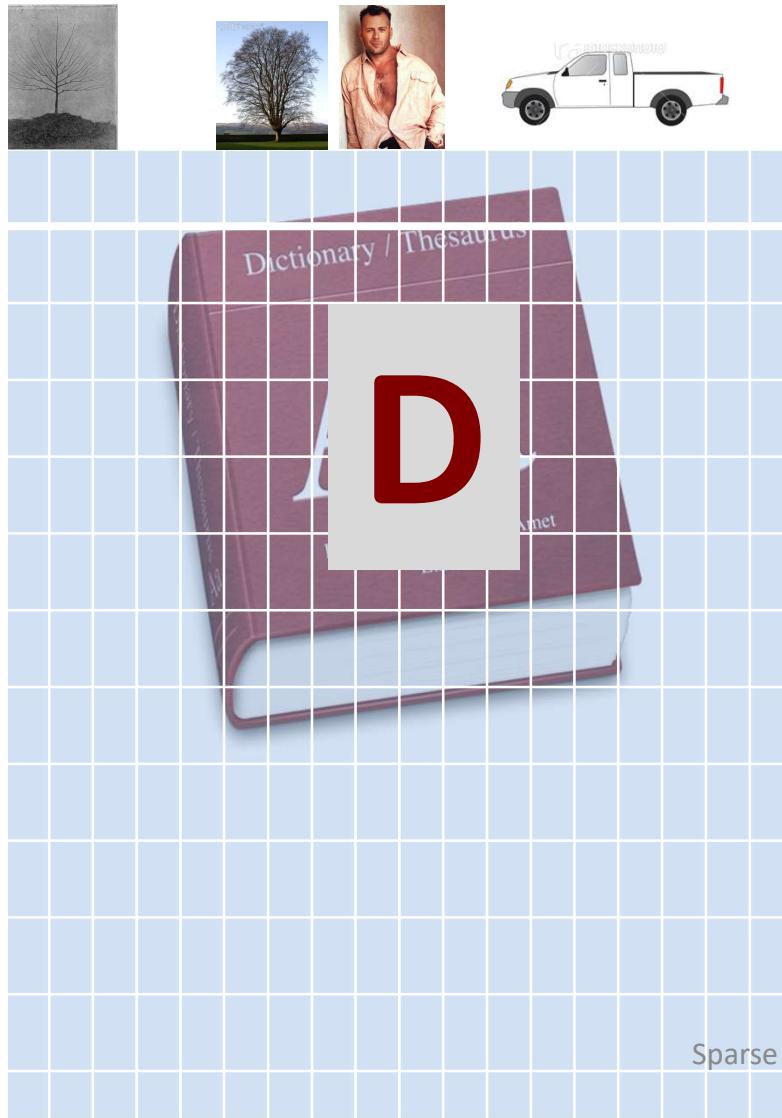
Overcompleteness and Sparsity



Overcompleteness and Sparsity



No more than 4 bases

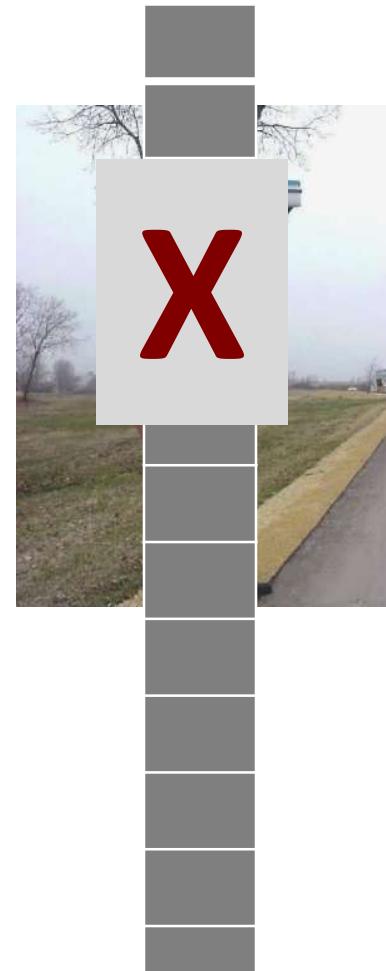


Sparse and Overcomplete Representations

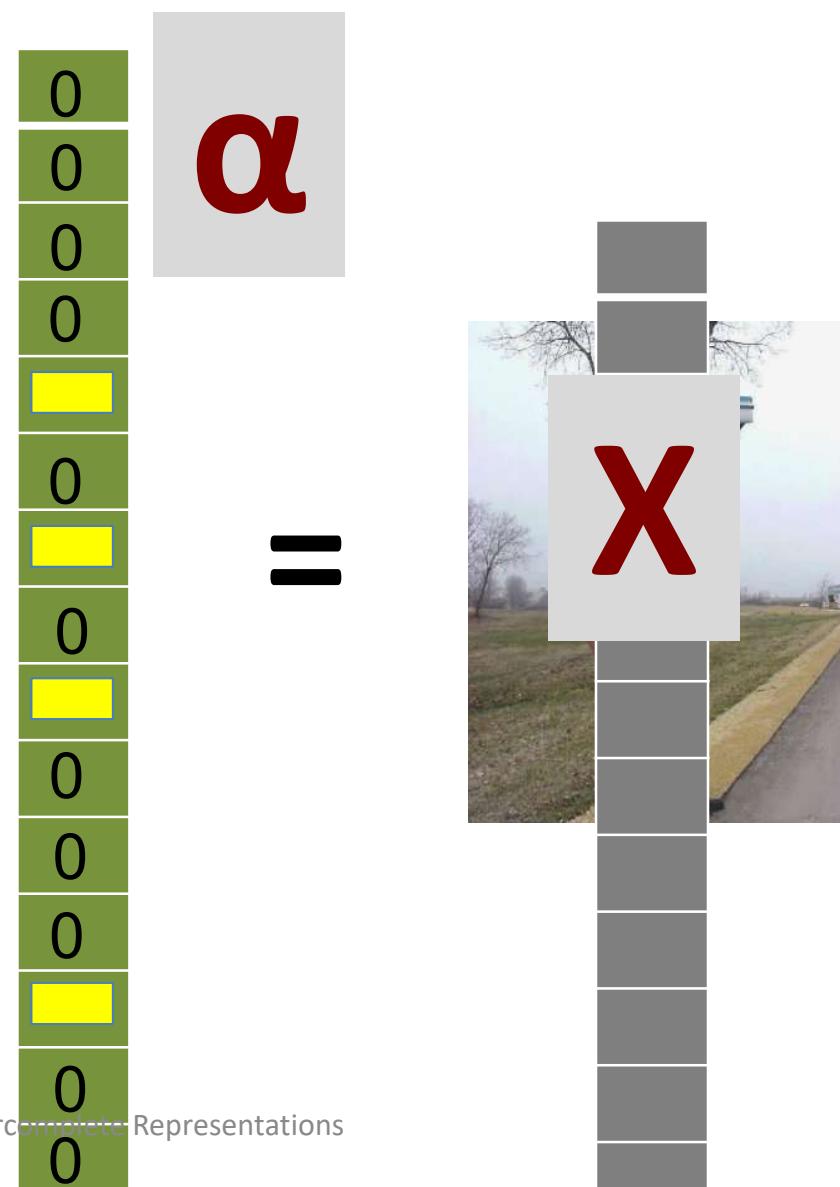
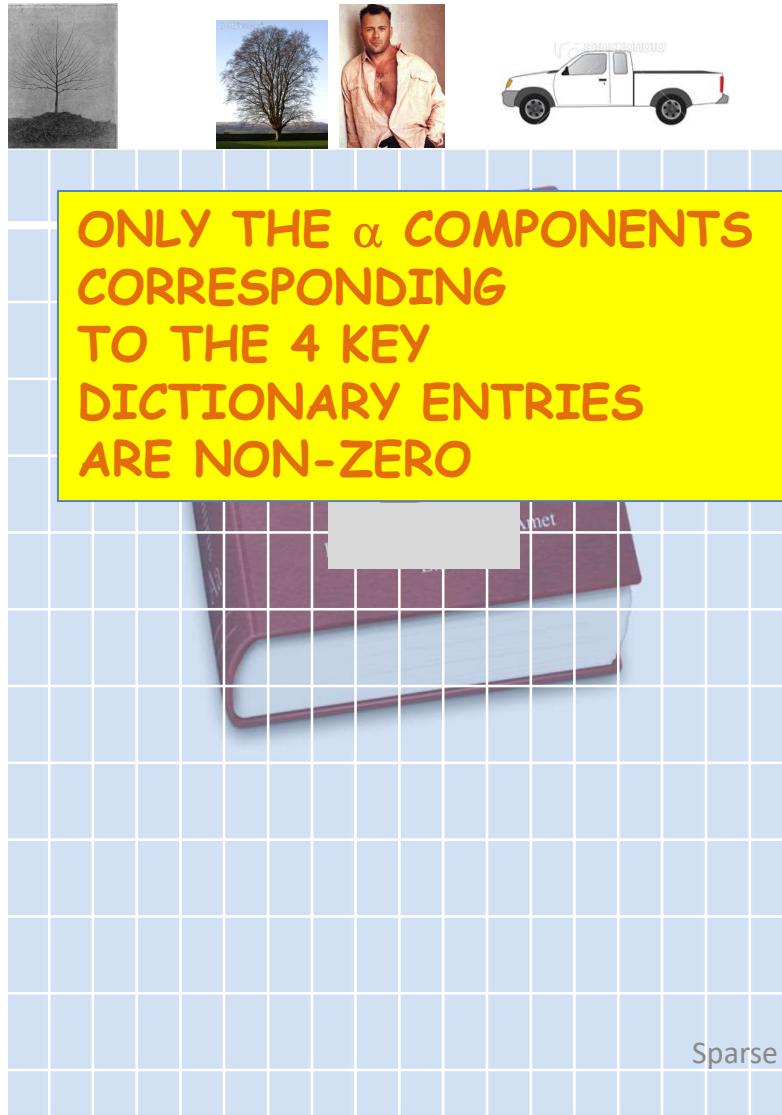
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0
0

$$\alpha$$

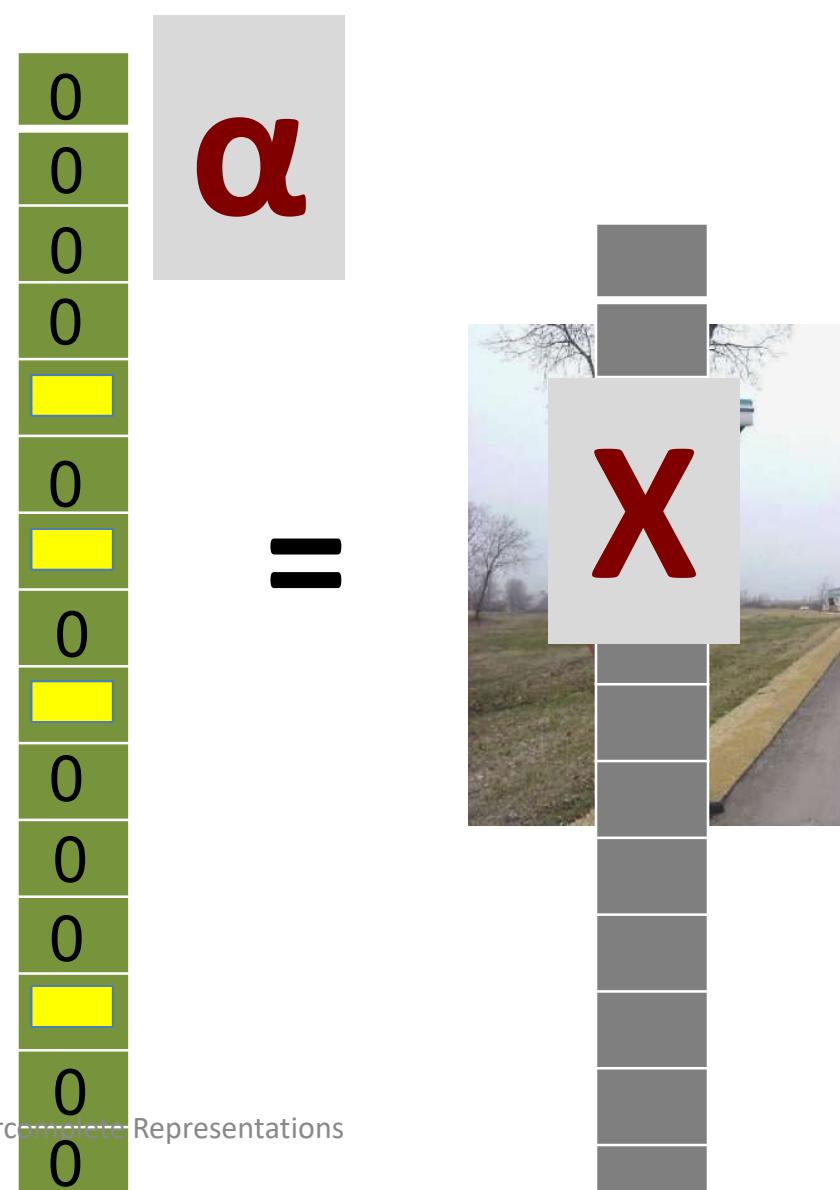
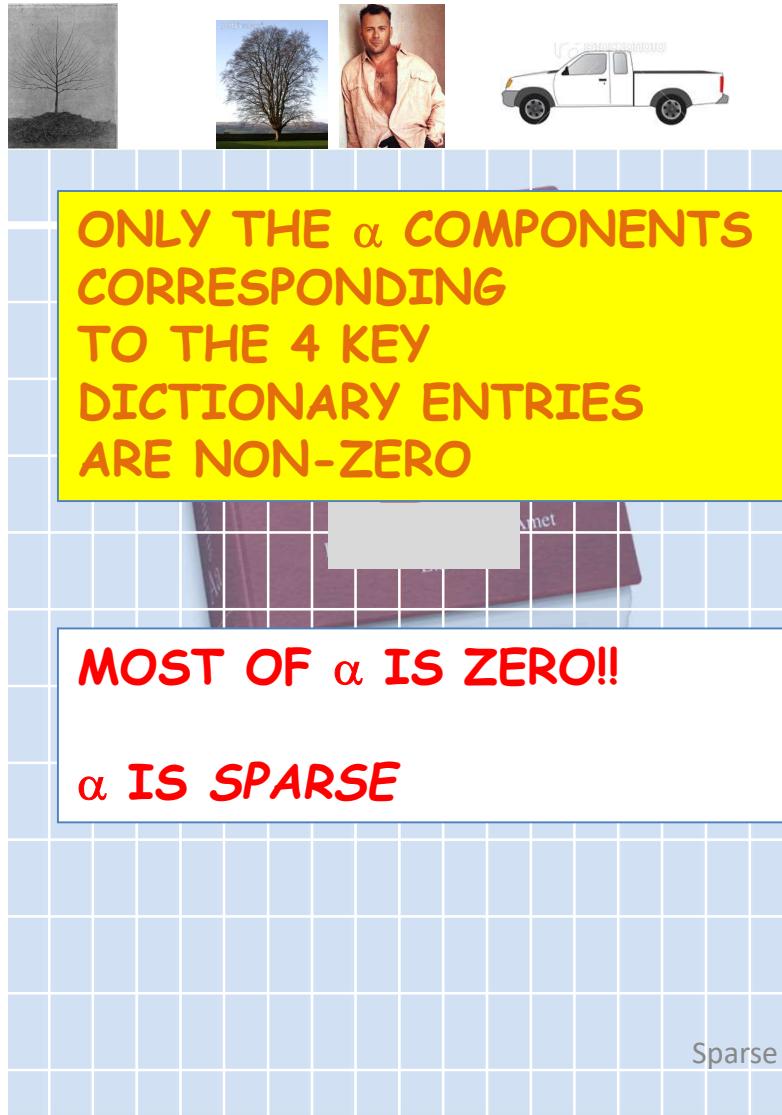
=



No more than 4 bases



No more than 4 bases



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)

The Sparsity Problem

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

The Sparsity Problem

- We want to use as few dictionary entries as possible to do this.

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

The Sparsity Problem

- We want to use as few dictionary entries as possible to do this.

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

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

The Sparsity Problem

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

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

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

Poll 3

The Sparsity Problem

- We want to use as few dictionary entries as possible to do this.

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

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,
norm(ResidualInputSignal) < threshold

Matching Pursuit

Algorithm Matching Pursuit

Input: Signal: $f(t)$.

Output: List of coefficients: (a_n, g_{γ_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: $\|Rf_n\| < threshold$)

From 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,

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

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

Basis Pursuit

- Remember,

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

In the general case, this is intractable

Basis Pursuit

- Remember,

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

In the general case, this is intractable

Requires combinatorial optimization

Basis Pursuit

- Replace the intractable expression by an expression that is solvable

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

Basis Pursuit

- Replace the intractable expression by an expression that is solvable

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

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

$$\begin{array}{ll} \min_{\underline{\alpha}} & \|\underline{\alpha}\|_1 \\ \text{s.t.} & \underline{X} = \mathbf{D}\underline{\alpha} \end{array}$$

Objective

Constraint

Basis Pursuit

- We can formulate the optimization term as:

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

The diagram shows the optimization equation for Basis Pursuit. It consists of a rectangular frame containing the mathematical expression. Two blue arrows point from the text "Constraint" and "Objective" at the bottom to the terms $\|\underline{X} - \mathbf{D}\underline{\alpha}\|^2$ and $\lambda \|\underline{\alpha}\|_1$ respectively within the frame.

Constraint Objective

Basis Pursuit

- We can formulate the optimization term as:

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

λ 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](#)

<http://www-stat.stanford.edu/~tibs/ftp/lasso.ps>


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

λ is a penalty term on the non-zero elements and promotes sparsity

Basis Pursuit

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

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

- $\|\underline{\alpha}\|_1$ is not differentiable at $\alpha_j = 0$
- Gradient of $\|\underline{\alpha}\|_1$ 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

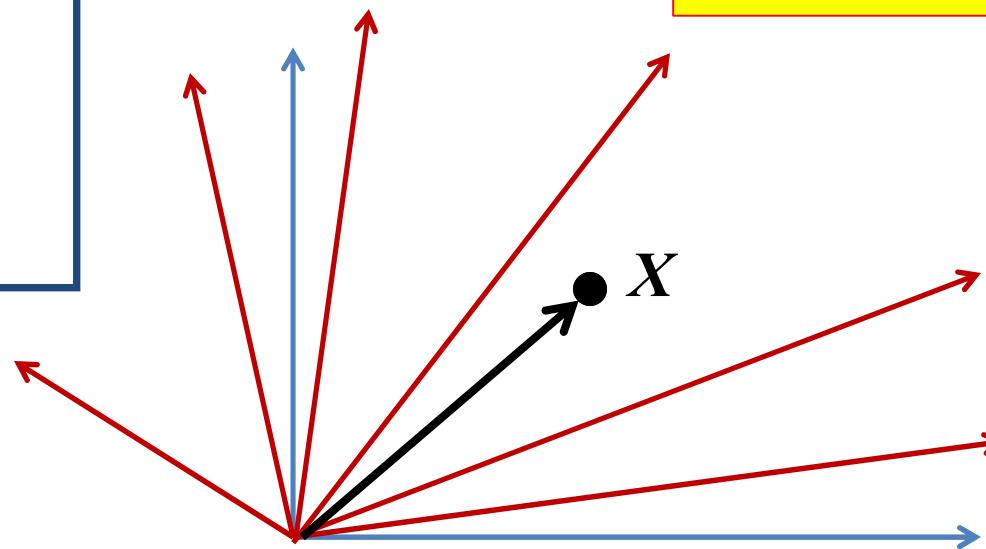
- There are efficient ways to solve the LASSO formulation.
 - http://web.stanford.edu/~hastie/glmnet_matlab/
- Simplest solution: Coordinate descent algorithms
 - On webpage..

L_1 vs L_0

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

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

Overcomplete set
of 6 "bases"

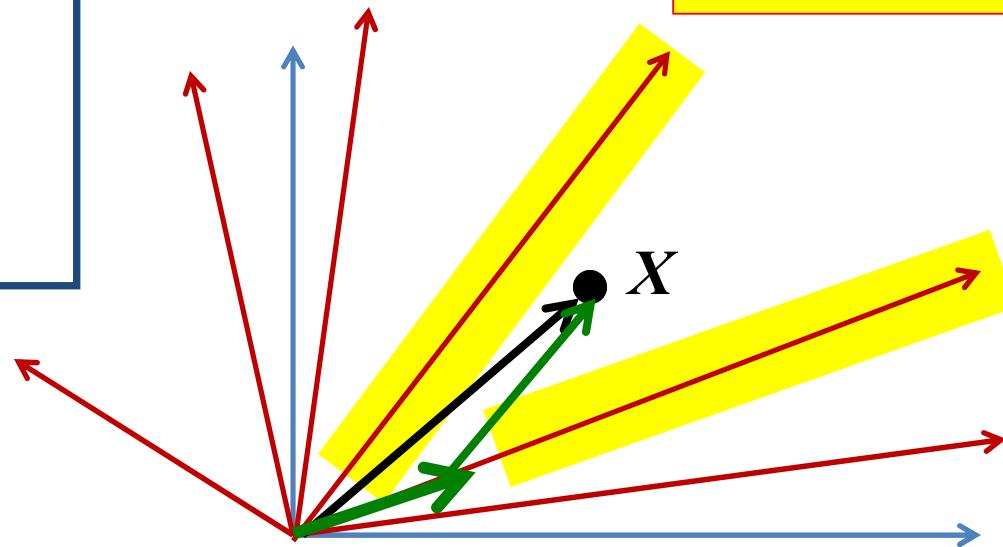


- **L_0 minimization**
 - Two-sparse solution
 - ANY pair of bases can explain X with 0 error

L_1 vs L_0

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

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 (remember the equations)	Soft thresholding
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

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

$$\begin{aligned} & \underset{\underline{\alpha}}{\text{Min}} \|\underline{X} - \mathbf{D}\underline{\alpha}\|_2^2 \\ & \text{s.t. } \|\underline{\alpha}\|_0 < C \end{aligned}$$

- L_1 minimization
- L_1 constrained optimization

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

$$\begin{aligned} & \underset{\underline{\alpha}}{\text{Min}} \|\underline{X} - \mathbf{D}\underline{\alpha}\|_2^2 \\ & \text{s.t. } \|\underline{\alpha}\|_1 < C \end{aligned}$$

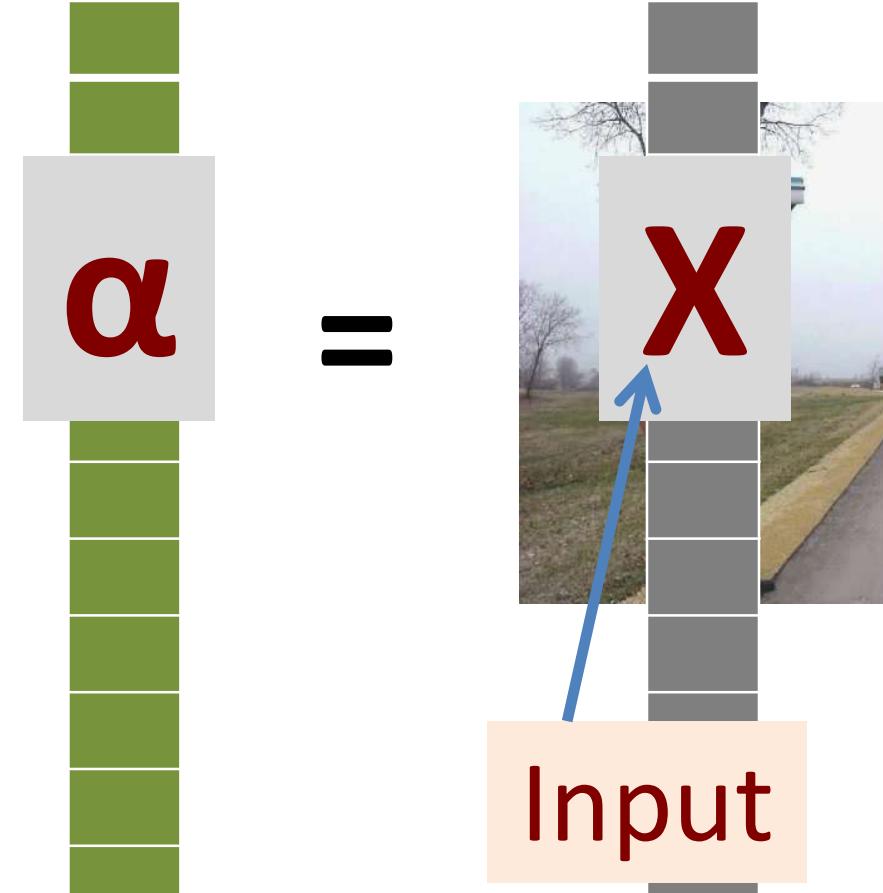
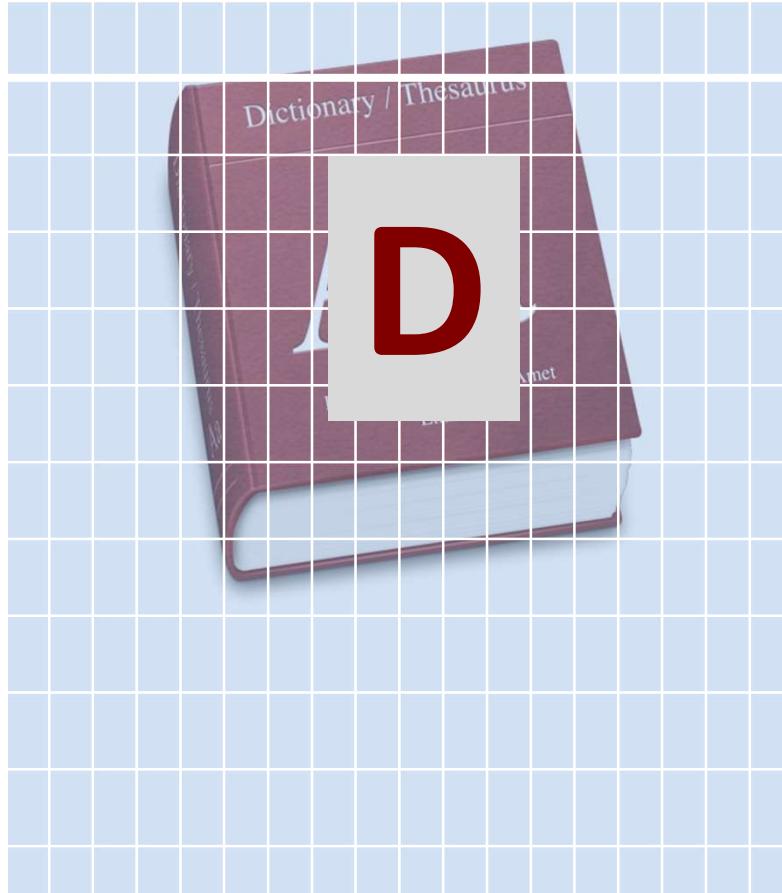
Many Other Methods..

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

Poll 4

Problems

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

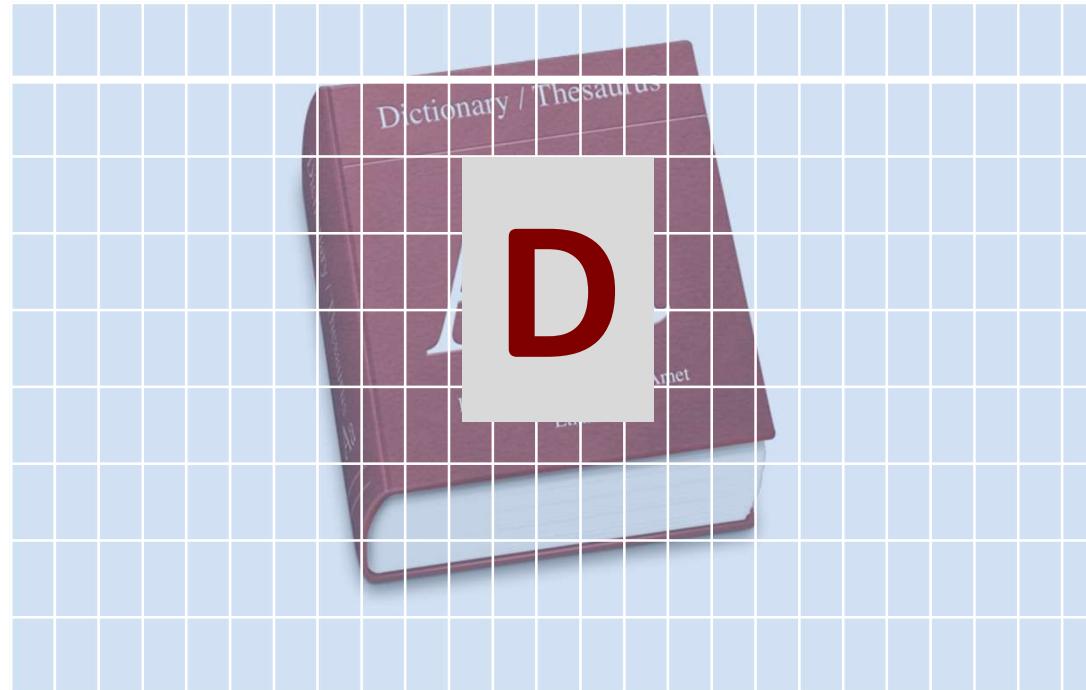


Trivial Solution

- $D = \text{Training data}$
- Impractical in most situations
 - Popular approach: sample random vectors from training data

Dictionaries: Compressive Sensing

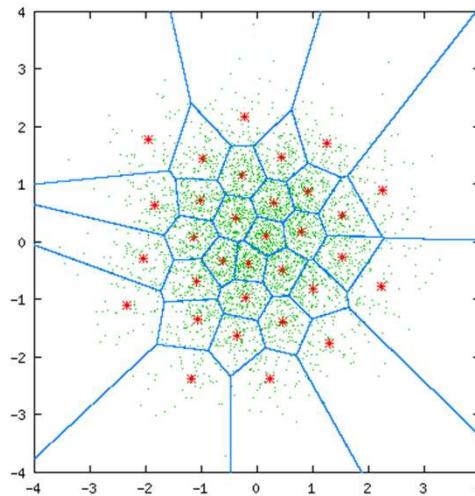
- Just random vectors!



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 ..

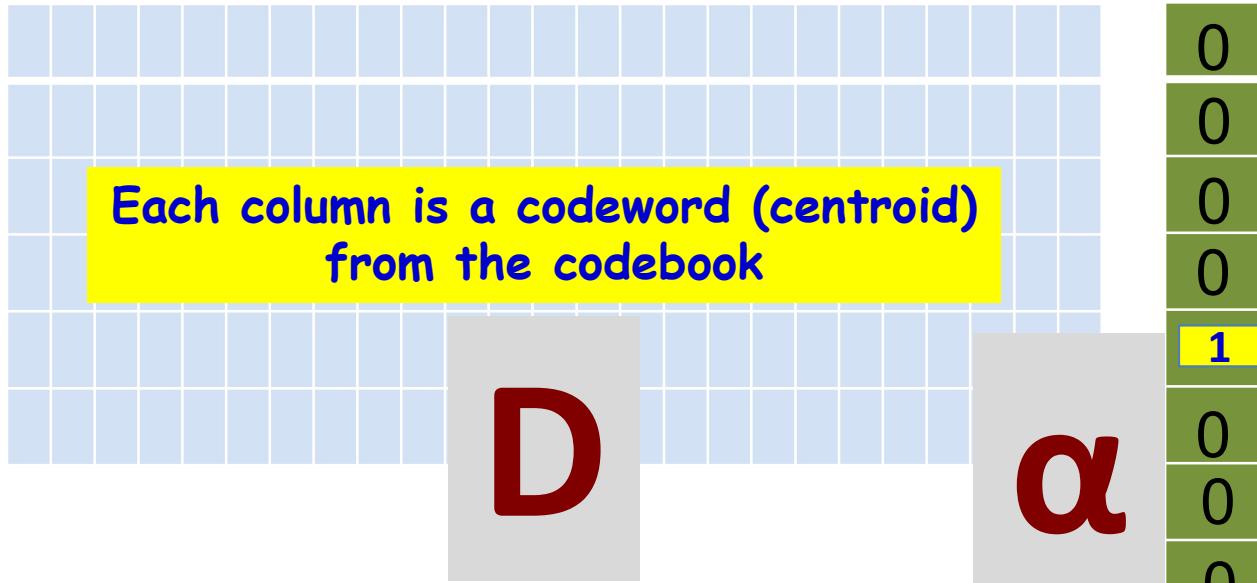
K-Means for Composing Dictionaries



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



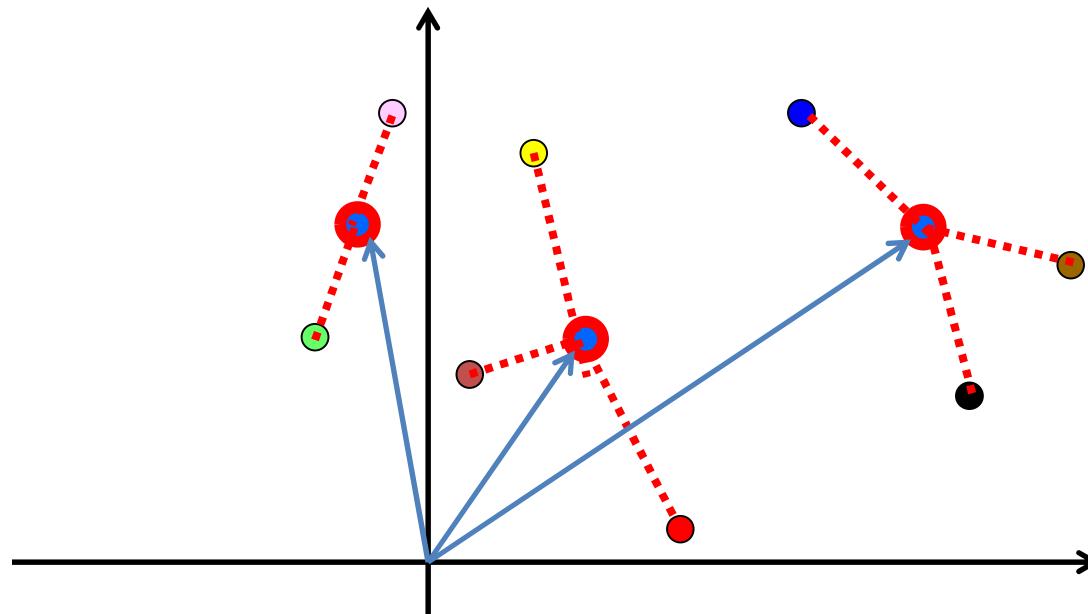
X

- α must be 1 sparse
- Only α entry must 1

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

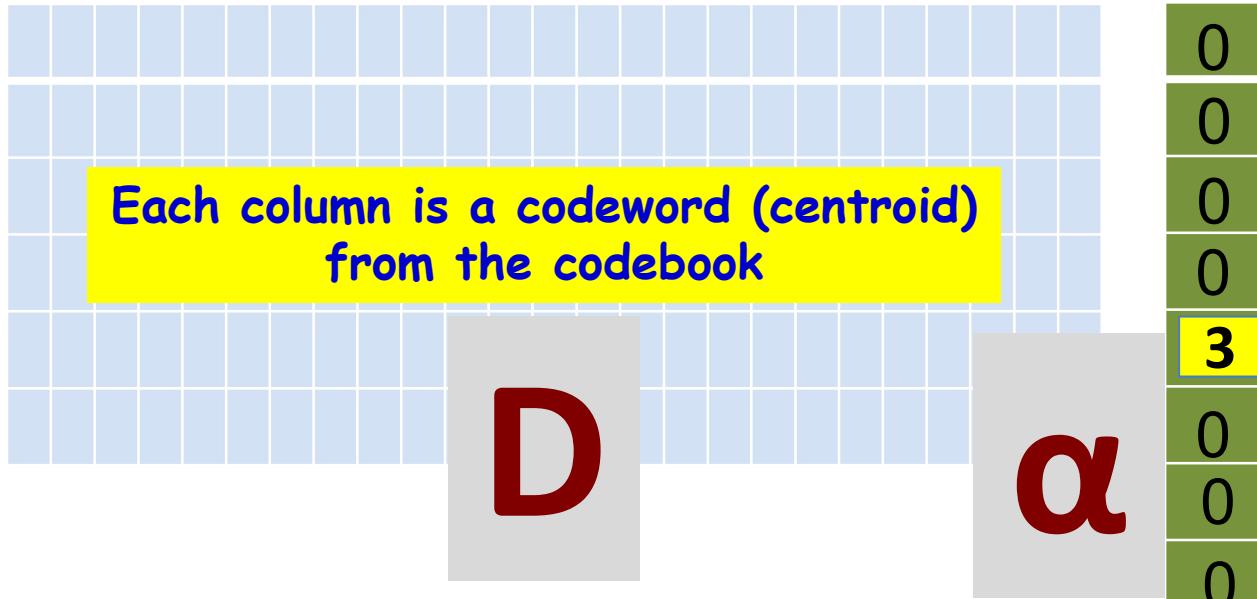
$$\|\alpha\|_1 = 1$$

K-Means



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

Length-unconstrained K-Means for Dictionaries



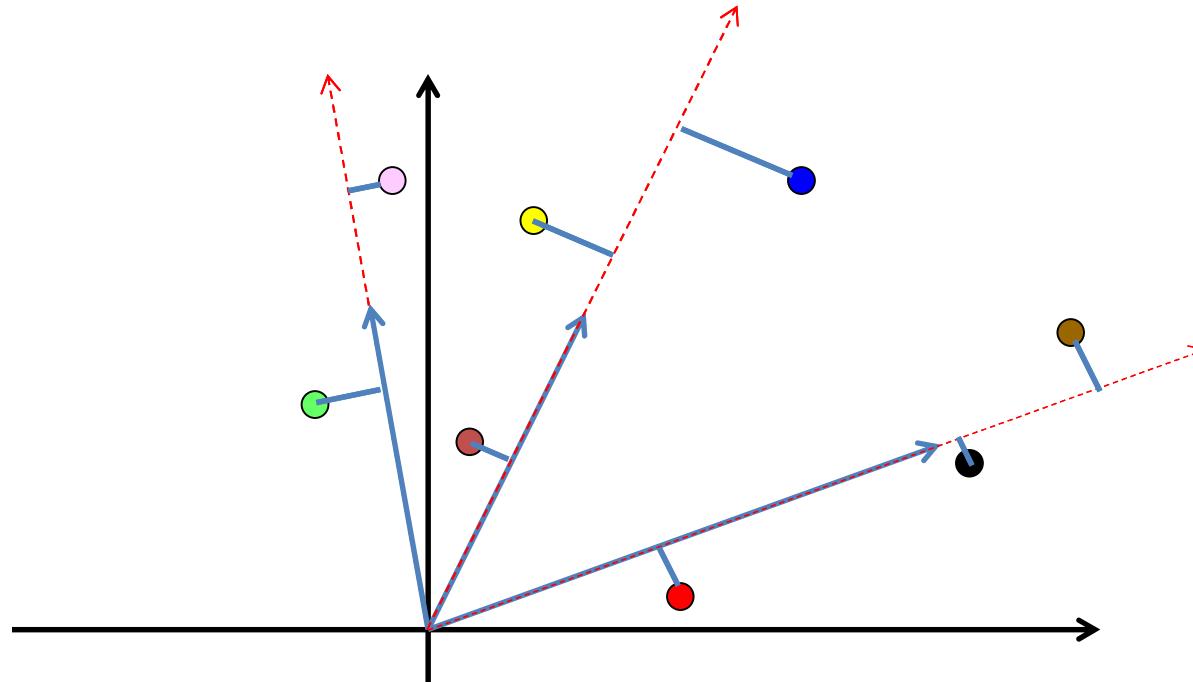
- α must be 1 sparse
- No restriction on α value

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

Sparse and Overcomplete Representations



SVD K-Means



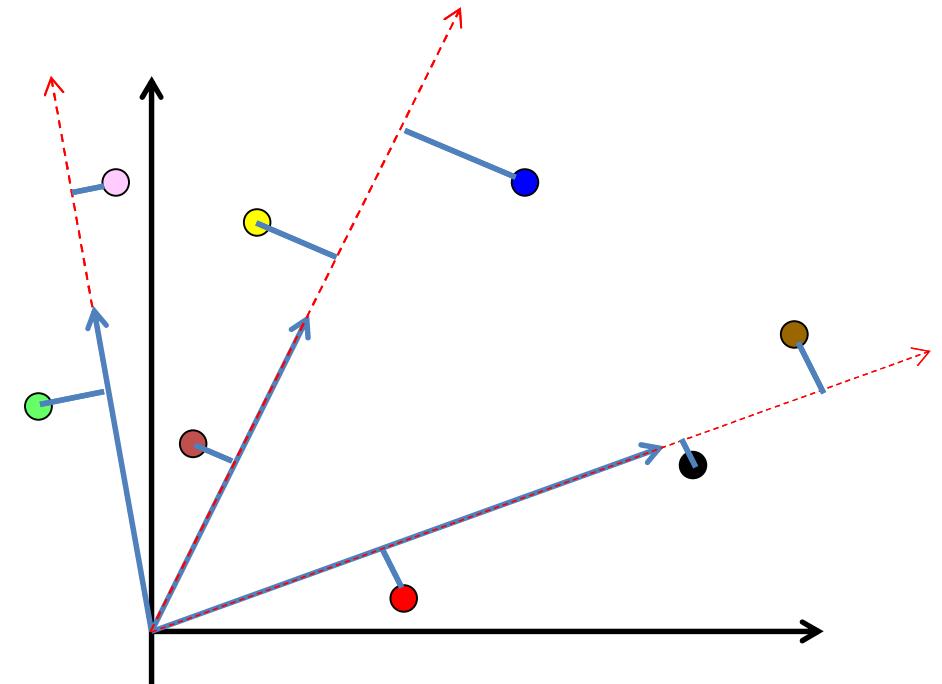
- Learn Codewords to minimize the total squared *projection error* of the training vectors from the closest codeword

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

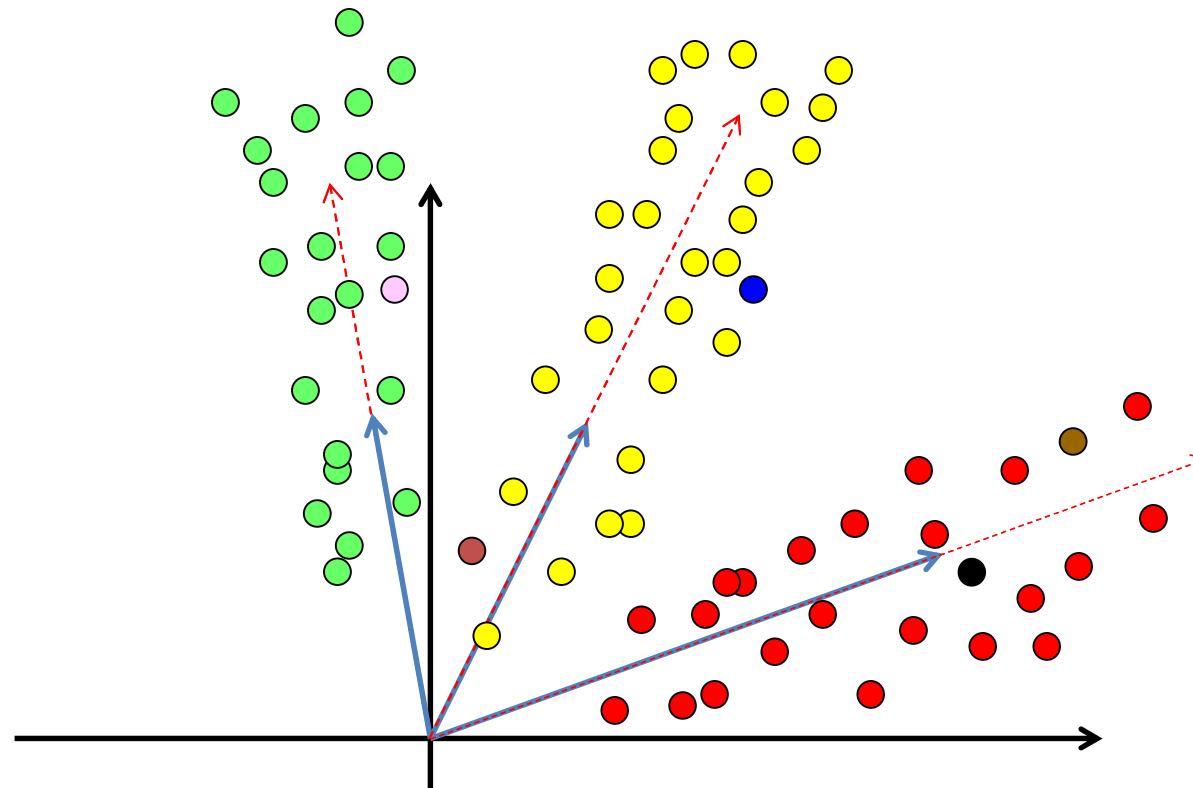
- $$p_{cluster} = |x^T m_{cluster}|$$

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



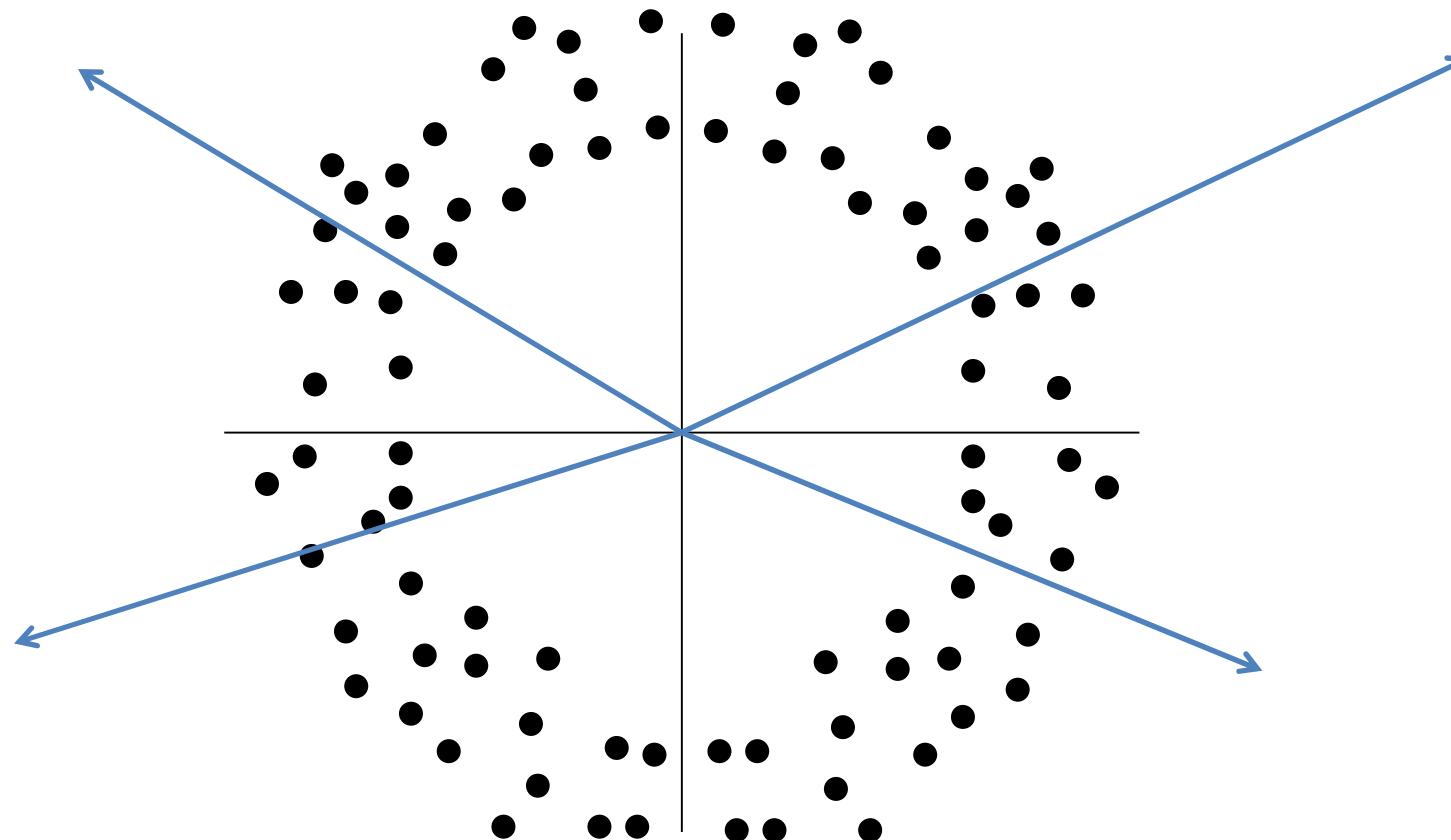
$$m_{cluster} = \text{Principal Eigenvector}(\{x \mid x \in \text{cluster}\})$$

Problem



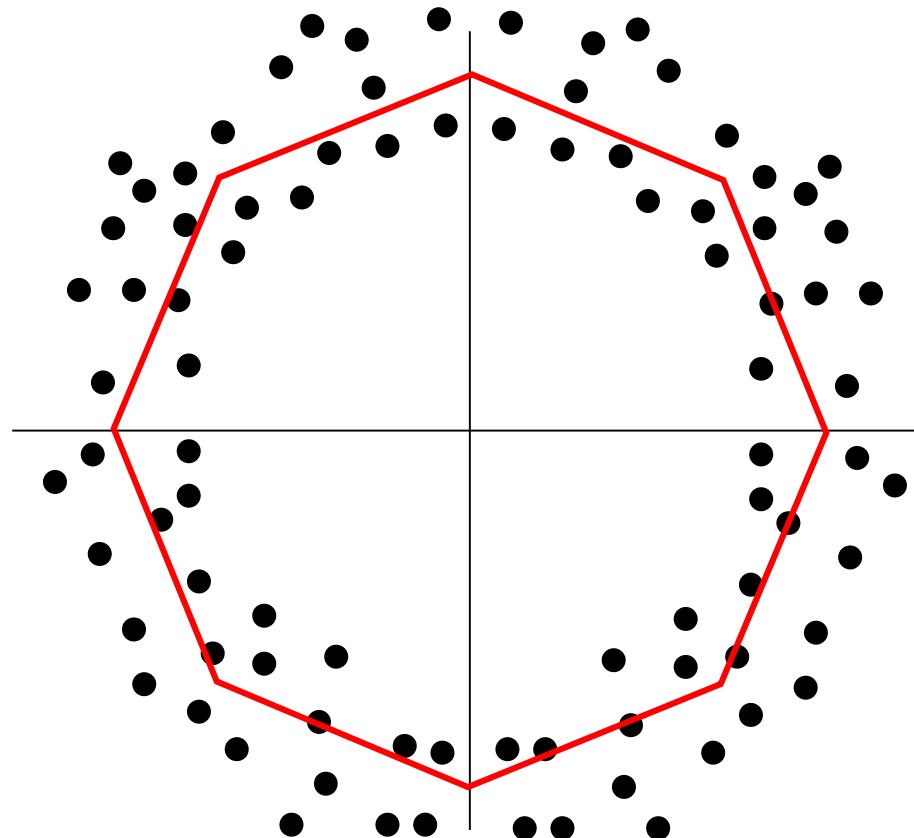
- Only represents *Radial* patterns

What about this pattern?



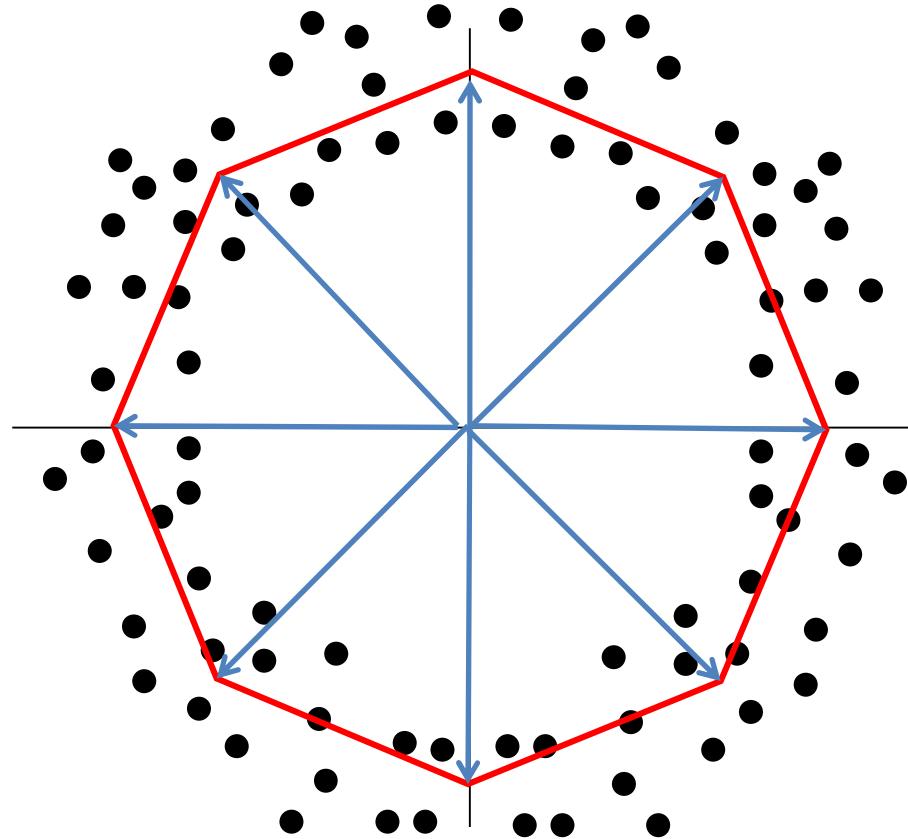
- Dictionary entries that represent radial patterns will not capture this structure
 - 1-sparse representations will not do

What about this pattern?



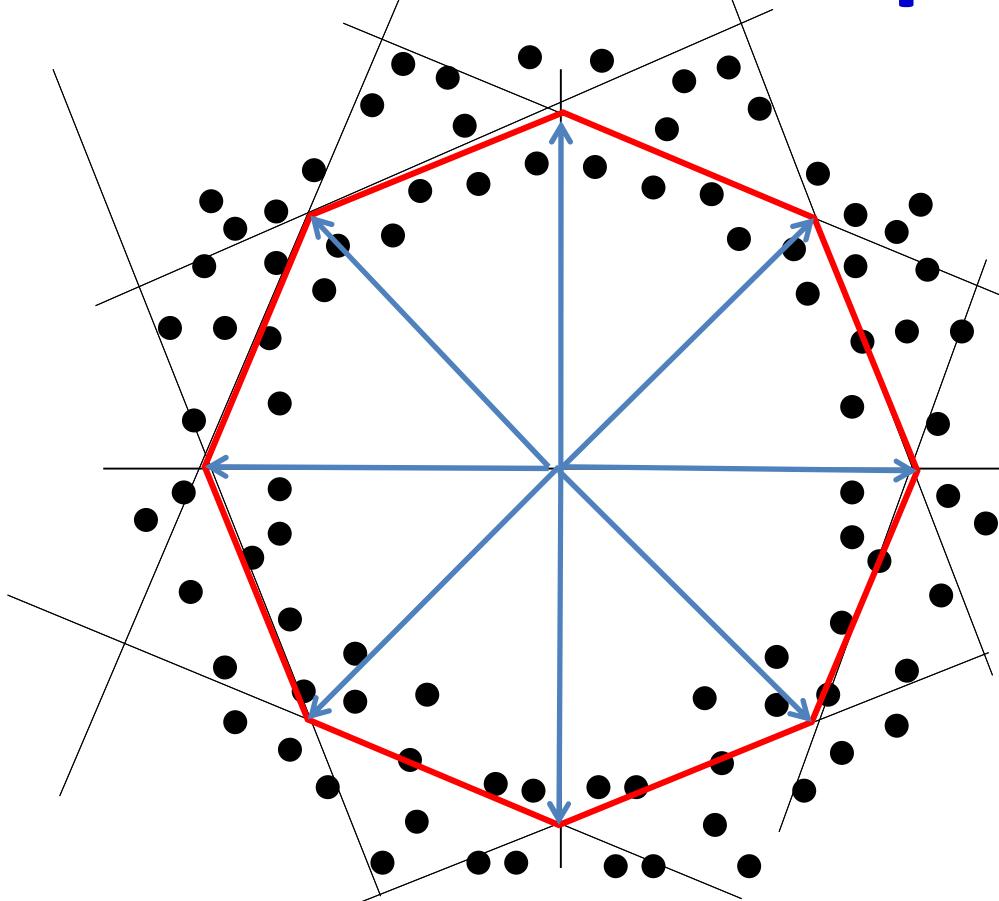
- We need **AFFINE** patterns

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?



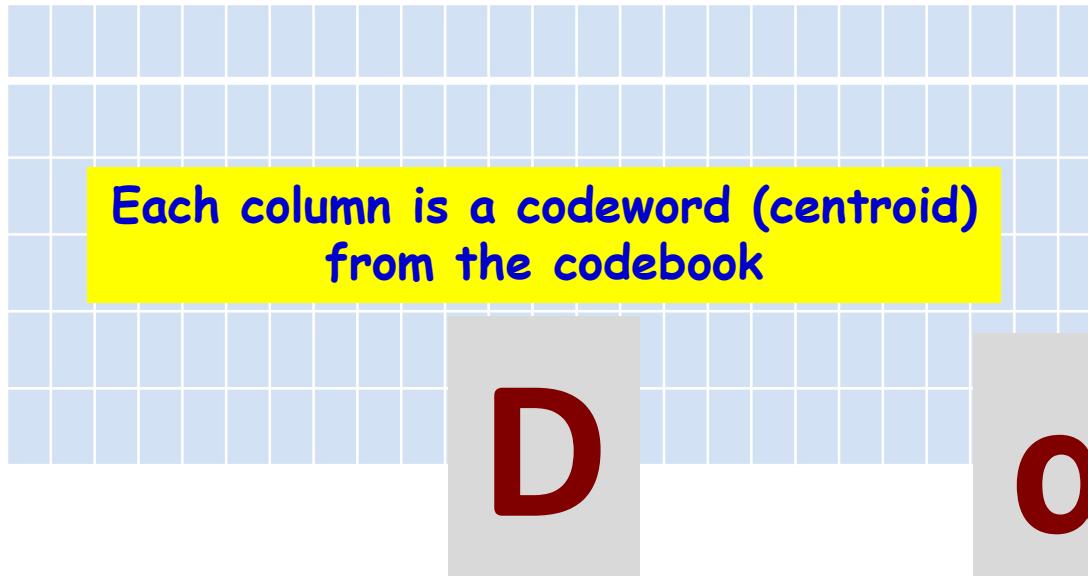
Every line is a
(constrained) combination
of two bases

2-sparse

Constraint:
 $\text{Line} = a.b_1 + (1-a)b_2$

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

Codebooks for K sparsity?



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



X

- α must be k sparse
- No restriction on α value

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

Formalizing

Given training data

$$\{X_1, X_2, \dots, X_T\}$$

We want to find a dictionary D , such that

$$D\alpha_i = X_i$$

With α_i sparse

Formalizing

Two objectives:

- Approximation
- Sparsity in coefficients

$$\begin{aligned} \|D\alpha_i - X_i\| \\ \|\alpha_i\|_1 \end{aligned}$$

$$\min_{D, \alpha_i} \sum_{i=1}^T \|X_i - D\alpha_i\|^2 + \|\alpha_i\|_1$$

NON-Convex!!!

An iterative method

- Given D , estimate α_i to get sparse solution
 - We can use any method

$$\min_{\alpha_i} \sum_{i=1}^T \|X_i - D\alpha_i\|^2 + \|\alpha_i\|_1$$

- Given α_i , estimate D

$$\min_D \sum_{i=1}^T \|X_i - D\alpha_i\|^2$$

Difficult!

K SVD

- Initialize Codebook

$D =$

1. For every vector, compute K-sparse alphas

- Using any pursuit algorithm

$\alpha =$

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 that used k
 - 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

$D =$

D _j j != 1									

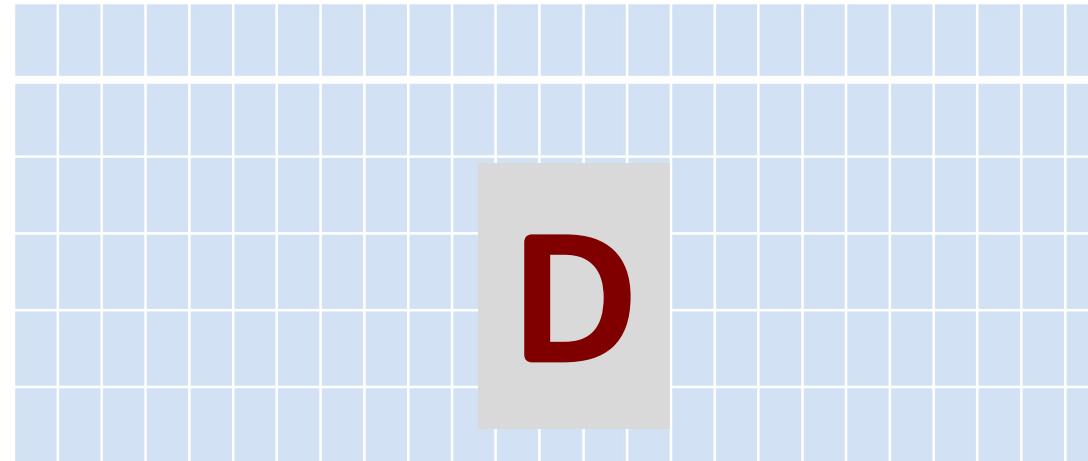
$\alpha =$

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

$\alpha_j(x) \quad j != 1$

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

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

K-SVD algorithm (skip)

Initialization : Set the random normalized dictionary matrix $\mathbf{D}^{(0)} \in \mathbb{R}^{n \times K}$. Set $J = 1$.

Repeat until convergence,

Sparse Coding Stage: Use any pursuit algorithm to compute \mathbf{x}_i for $i = 1, 2, \dots, N$

$$\min_{\mathbf{x}} \{ \|\mathbf{y}_i - \mathbf{D}\mathbf{x}\|_2^2 \} \quad \text{subject to} \quad \|\mathbf{x}\|_0 \leq T_0.$$

Codebook Update Stage: For $k = 1, 2, \dots, K$

- Define the group of examples that use \mathbf{d}_k ,
 $\omega_k = \{i | 1 \leq i \leq N, \mathbf{x}_i(k) \neq 0\}$.
- Compute

$$\mathbf{E}_k = \mathbf{Y} - \sum_{j \neq k} \mathbf{d}_j \mathbf{x}^j,$$

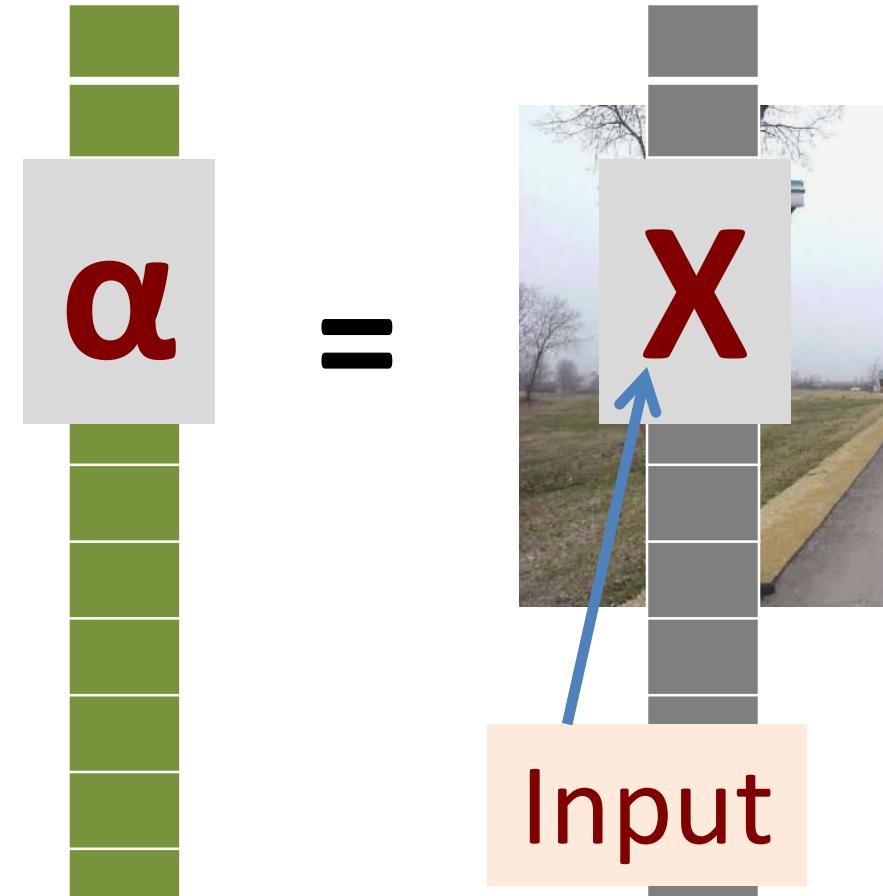
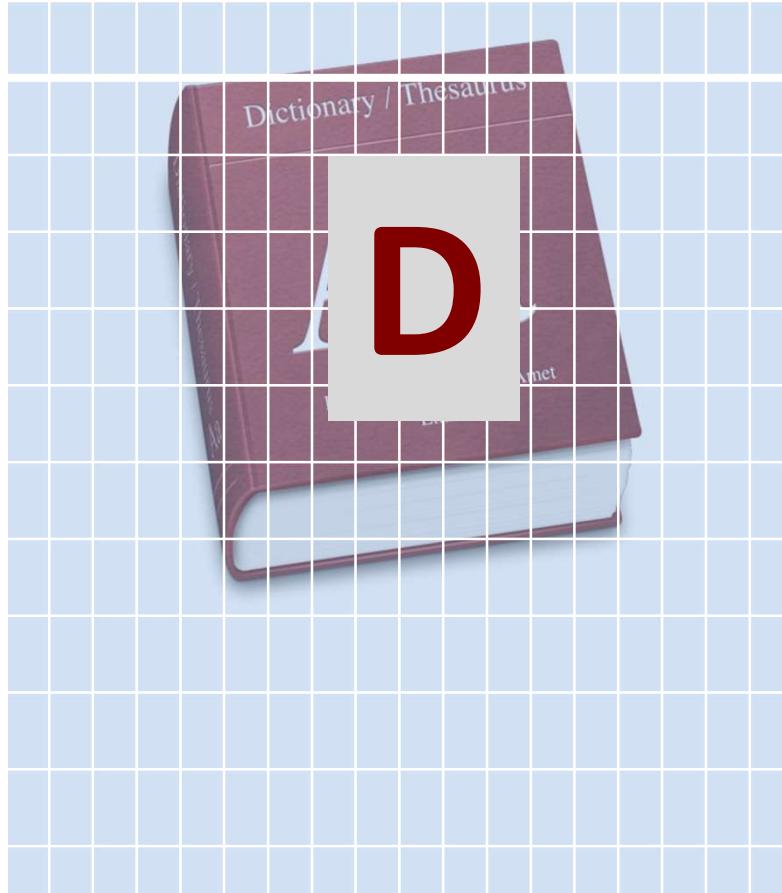
- Restrict \mathbf{E}_k by choosing only the columns corresponding to those elements that initially used \mathbf{d}_k in their representation, and obtain \mathbf{E}_k^R .
- Apply SVD decomposition $\mathbf{E}_k^R = \mathbf{U}\mathbf{\Delta}\mathbf{V}^T$. Update:
 $\mathbf{d}_k = \mathbf{u}_1, \mathbf{x}_R^k = \mathbf{\Delta}(1, 1) \cdot \mathbf{v}_1$

Set $J = J + 1$.

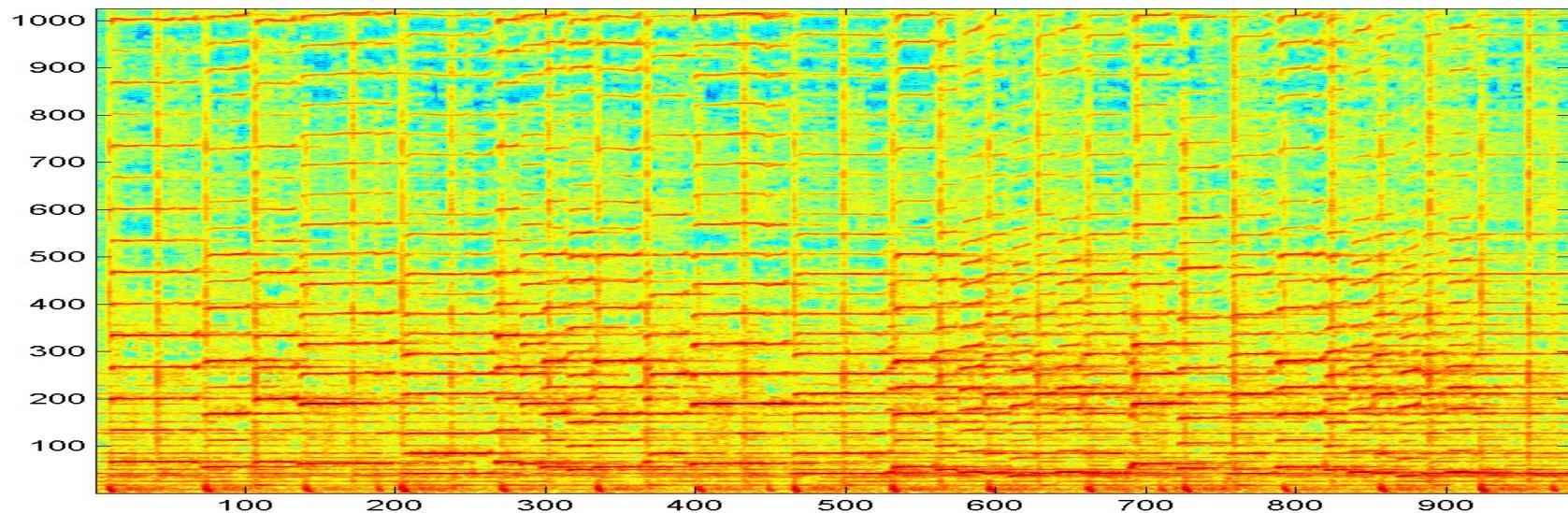
Sparse and Overcomplete Representations

Problems

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

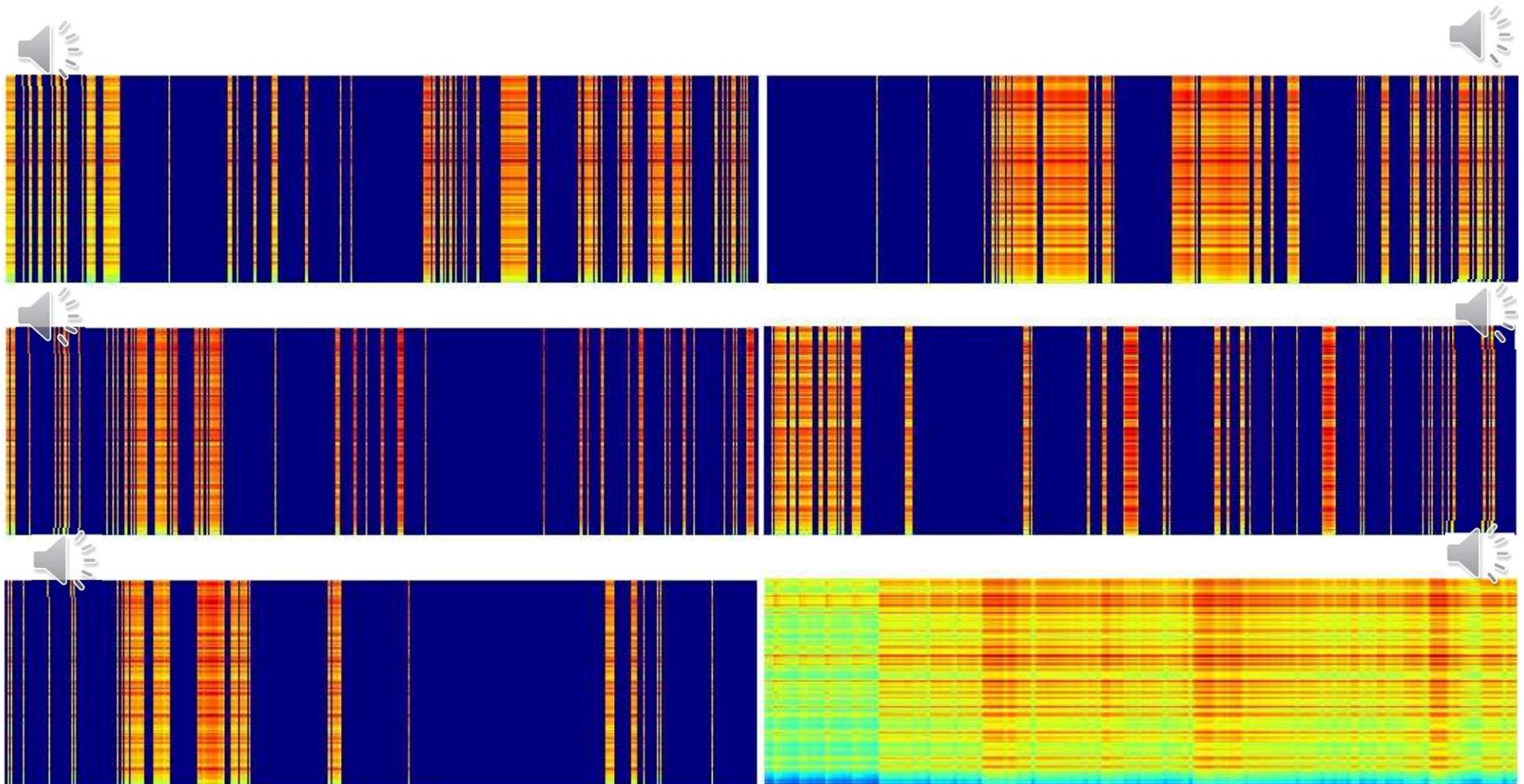


So how does that work



- In case you forgot this music...
- 975 vectors (1025 dimensions)
- $N=12, K=5$

K-SVD bases



Applications of Sparse Representations

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

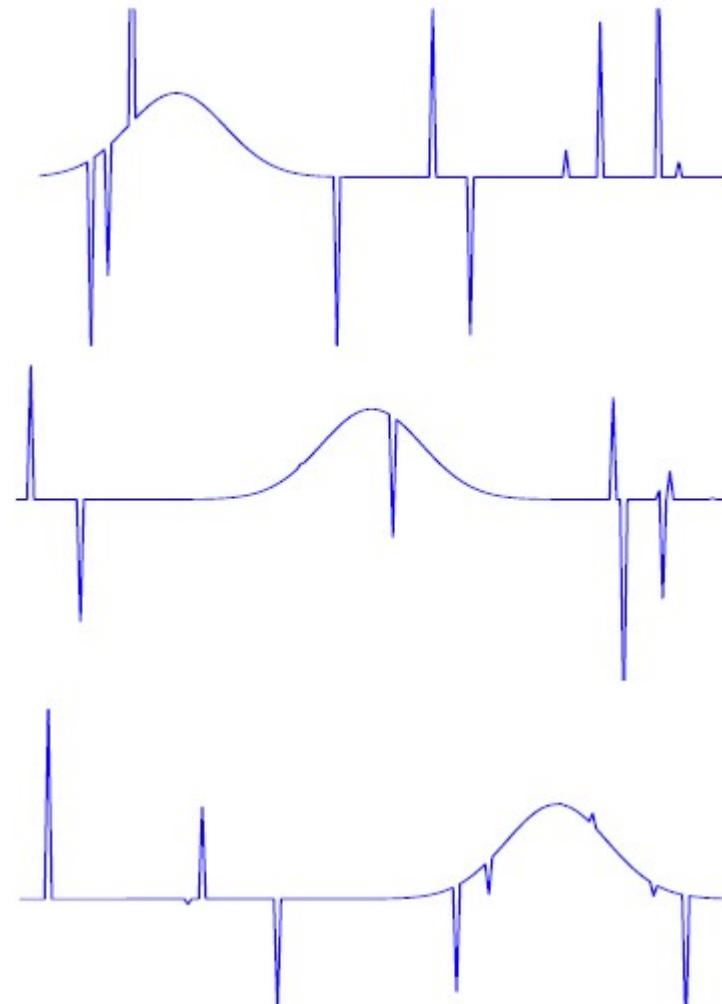
Denoising

- As the name suggests, remove noise!

Denoising

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

A toy example



A toy example

$$D = [I \ G]$$

I Identity matrix
 G Translation of a Gaussian pulse

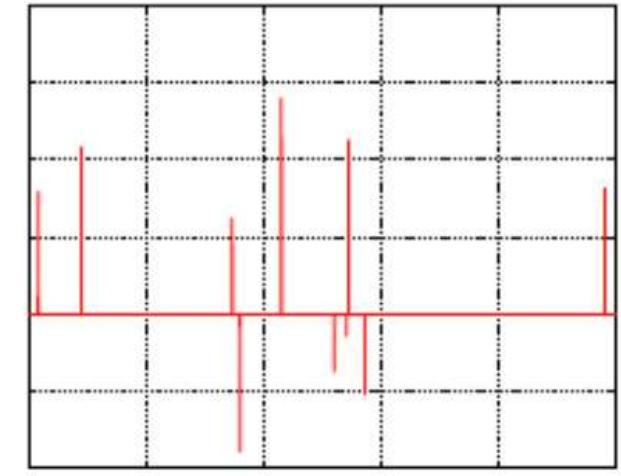
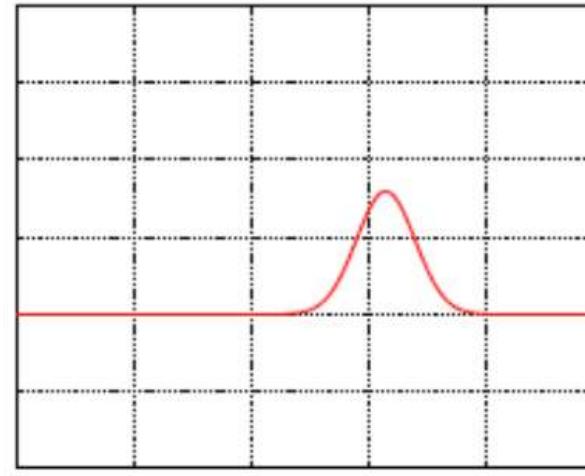
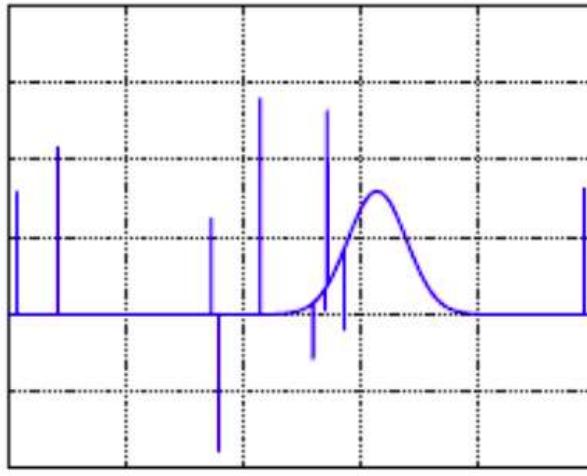


Image Denoising

- Here's what we want

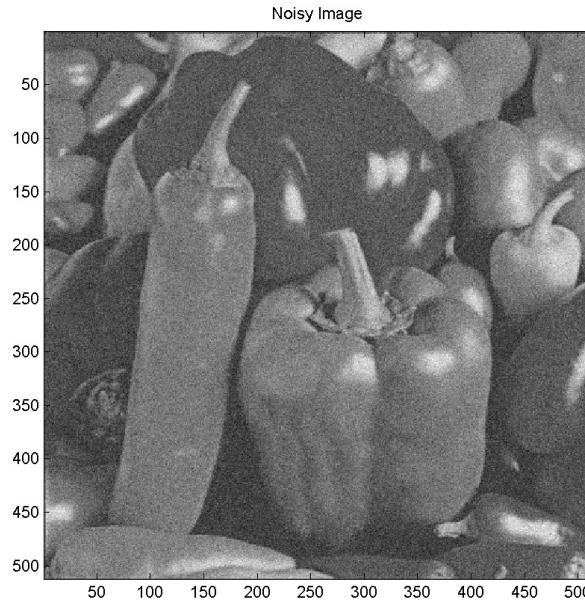


Image Denoising

- Here's what we want

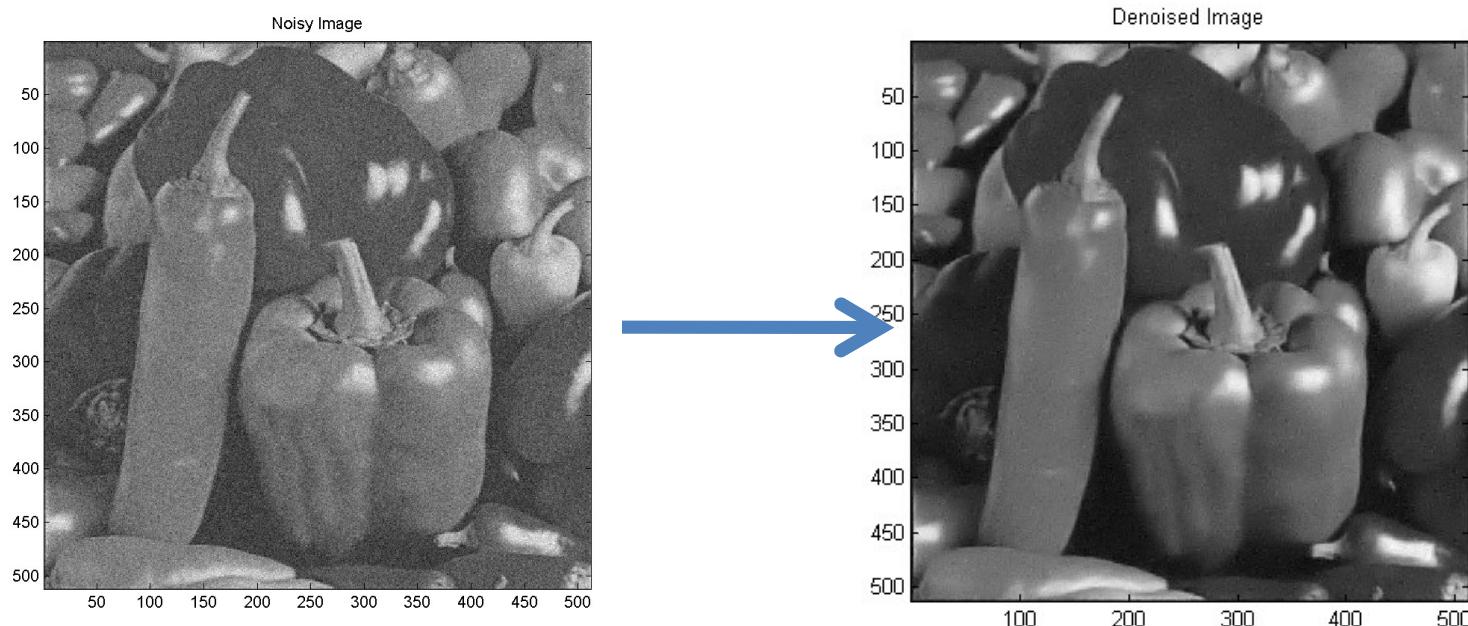


Image Denoising

- Here's what we want



The Image Denoising Problem

- Given an image
- Remove Gaussian additive noise from it

Image Denoising

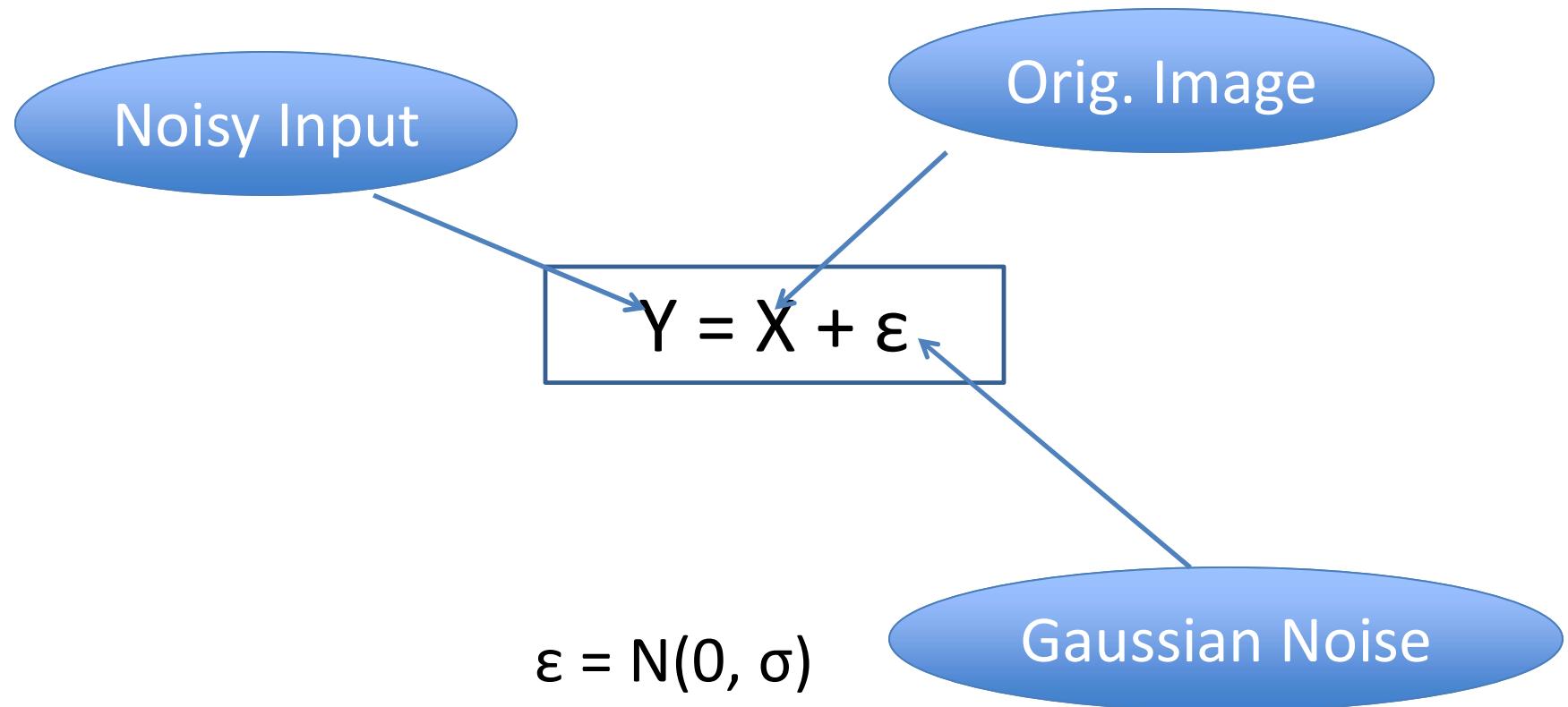


Image Denoising

- 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:

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

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

Image Denoising

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

- In the above, X is a patch.

Image Denoising

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

- 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

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

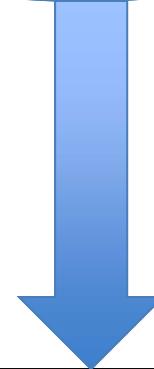
Image Denoising

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

$(\underline{X} - Y)$ is the error between the input and denoised image. μ is a penalty on the error.

Image Denoising

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



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

Image Denoising

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



λ 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 D.
- We want to estimate D as well.

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

We can use the previous equation itself!!!

Image Denoising

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

How do we estimate all 3 at once?

Image Denoising

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

How do we estimate all 3 at once?

We cannot estimate them at the same time!

Image Denoising

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

How do we estimate all 3 at once?

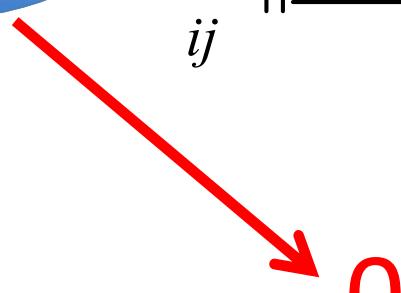
Fix 2, and find the optimal 3rd.

Image Denoising

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

Initialize $\mathbf{X} = \mathbf{Y}$

Image Denoising

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


Initialize $\underline{X} = Y$, initialize D

You know how to solve the remaining portion for α – MP, BP!

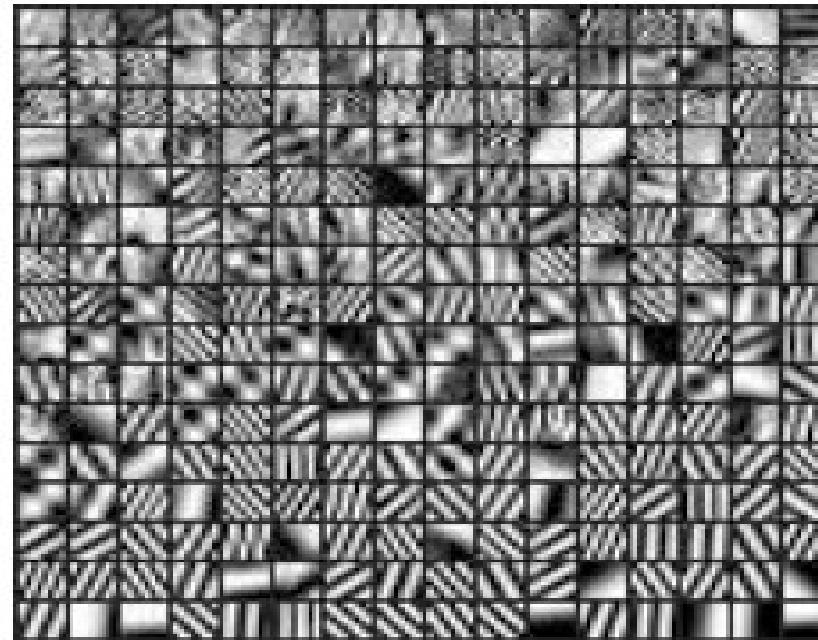
Image Denoising

- 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 α 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

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

} → Const. wrt X

We know D and α

The quadratic term above has a closed-form solution

Image Denoising

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

} → Const. wrt X

We know D and α

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

Image Denoising

- Summarizing... We wanted to obtain 3 things

Image Denoising

- Summarizing... We wanted to obtain 3 things
 - Weights α
 - Dictionary D
 - Denoised Image X

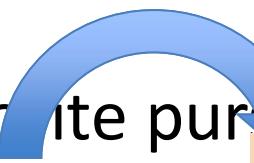
Image Denoising

- Summarizing... We wanted to obtain 3 things
 - Weights α – Your favorite pursuit algorithm
 - Dictionary D – Using K-SVD
 - Denoised Image X

Image Denoising

- Summarizing... We wanted to obtain 3 things

- Weights α – Your favorite pursuit algorithm
- Dictionary D – Using K-SVD
- Denoised Image X



Iterating

Image Denoising

- Summarizing... We wanted to obtain 3 things
 - Weights α
 - Dictionary D
 - Denoised Image X - Closed form solution

Image Denoising

- Here's what we want

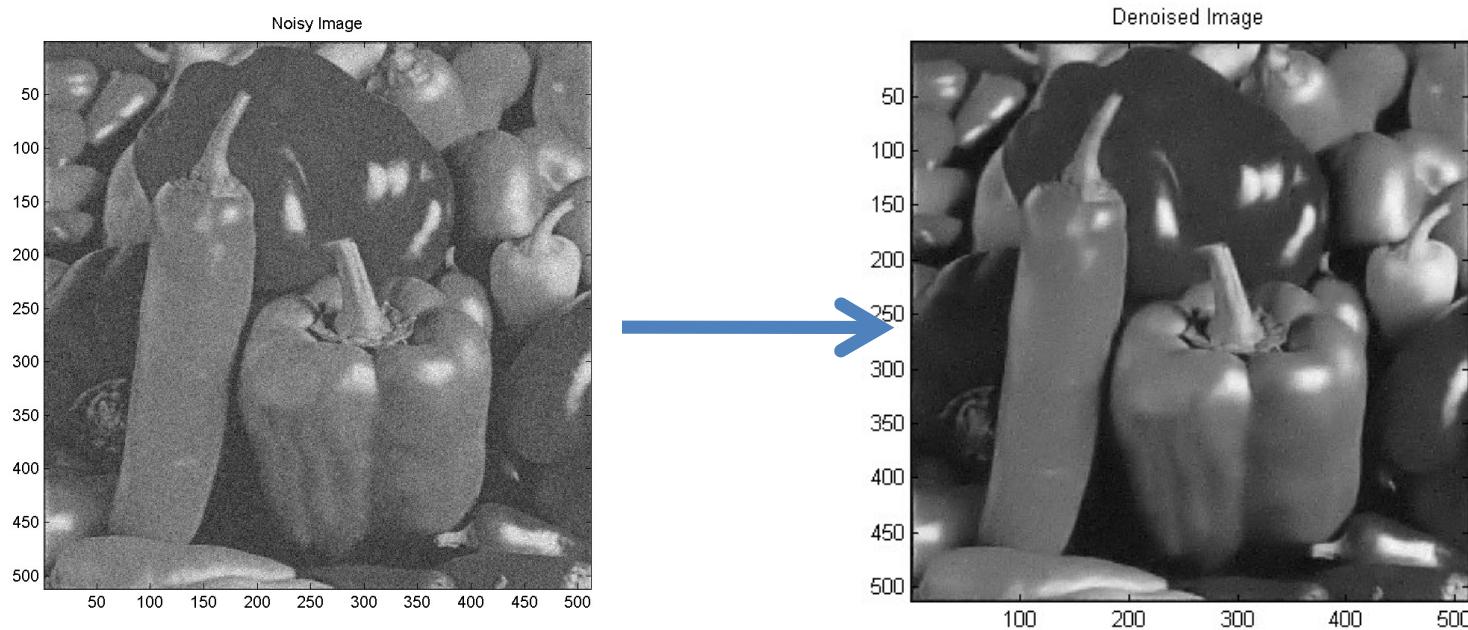


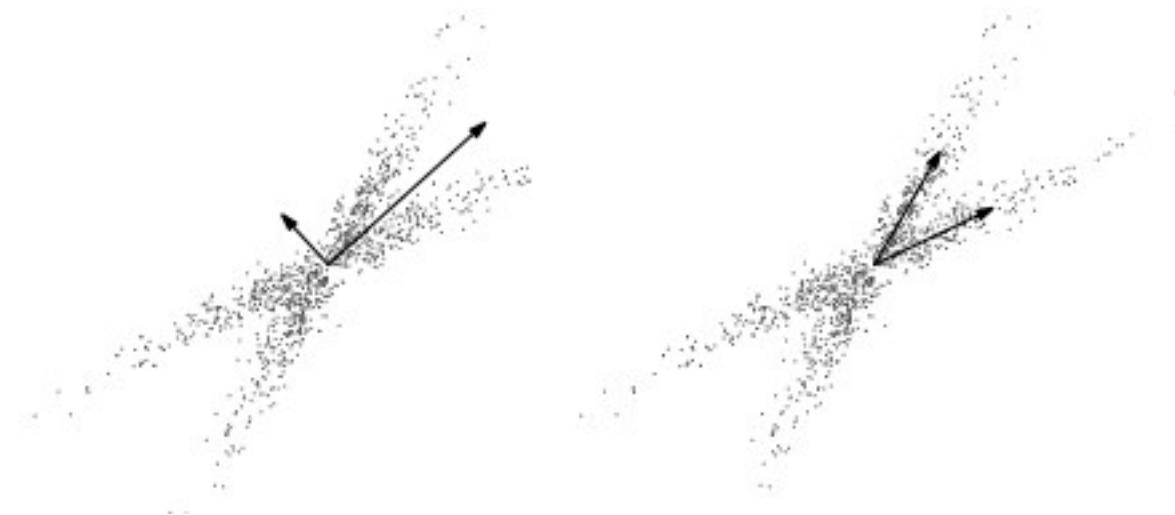
Image Denoising

- Here's what we want



Comparing to Other Techniques

Non-Gaussian data



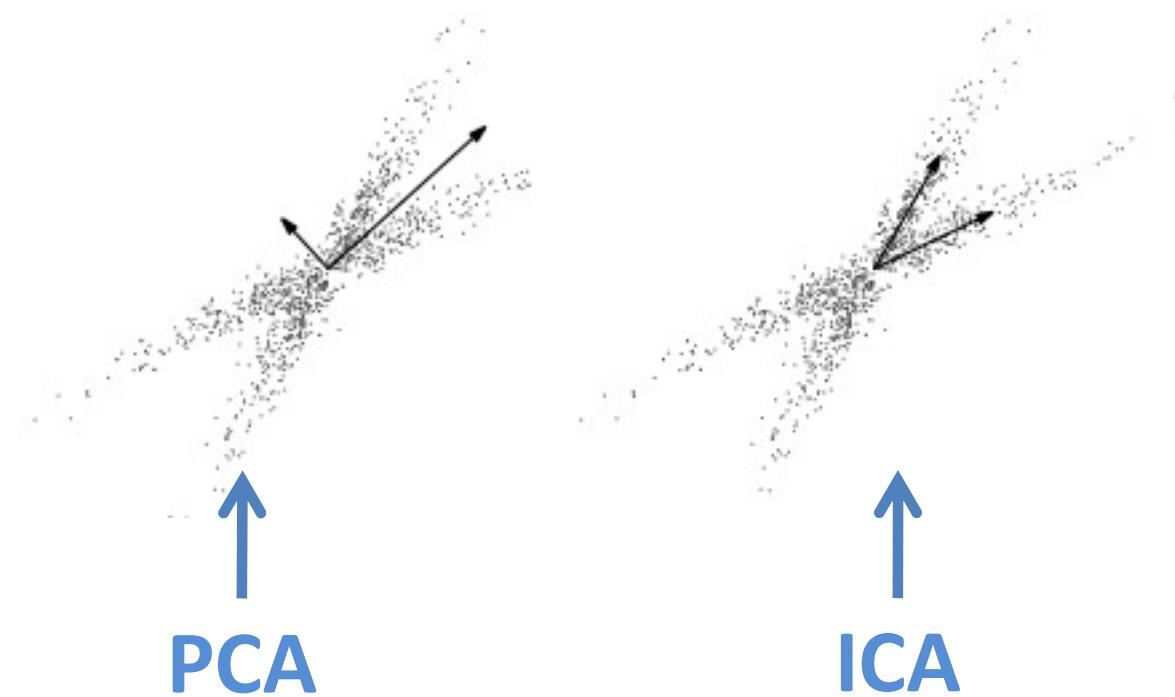
PCA of ICA

Which is which?

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

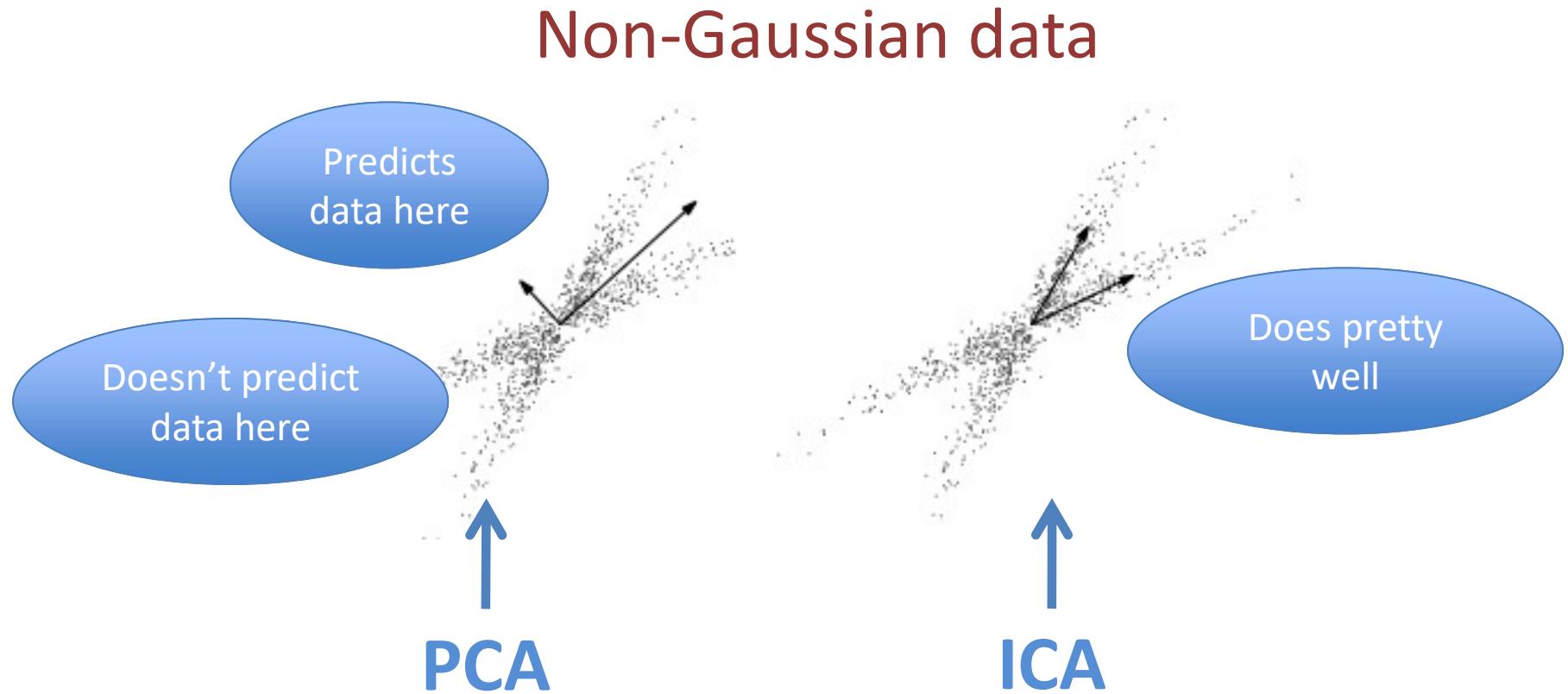
Comparing to Other Techniques

Non-Gaussian data



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

Comparing to Other Techniques



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.