

An Introduction to Sparse Coding, Sparse Sensing, and Optimization

Speaker: Wei-Lun Chao

Date: Nov. 23, 2011

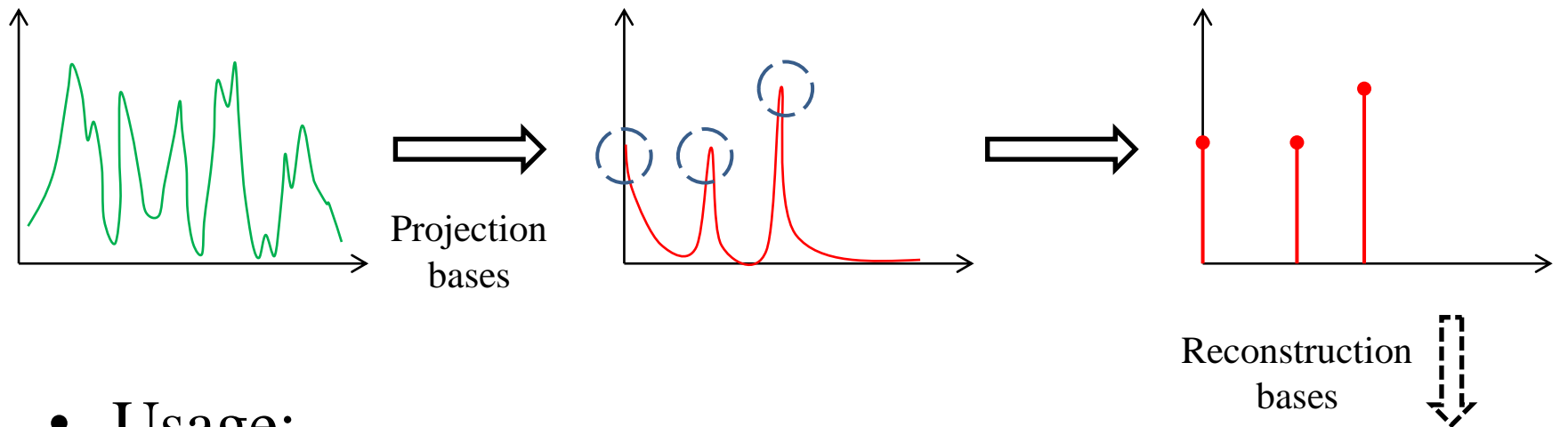
Outline

- Introduction
- The fundamental of optimization
- The idea of sparsity: coding V.S. sensing
- The solution
- The importance of dictionary
- Applications

Introduction

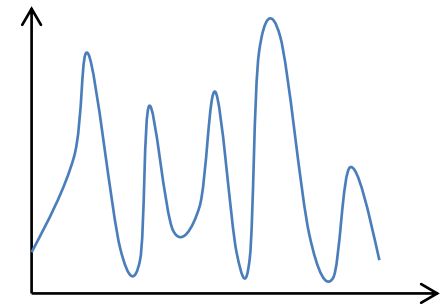
Introduction

- What is sparsity?



- Usage:

- Compression
- Analysis
- Representation
- Fast / sparse sensing



Introduction

- Why do we use **Fourier transform** and its modifications for image and acoustic compression?
 - Differentiability (theoretical)
 - Intrinsic sparsity (data-dependent)
 - Human perception (human-centric)
- Better bases for compression or representation?
 - Wavelets
 - How about data-dependent bases?
 - How about learning?

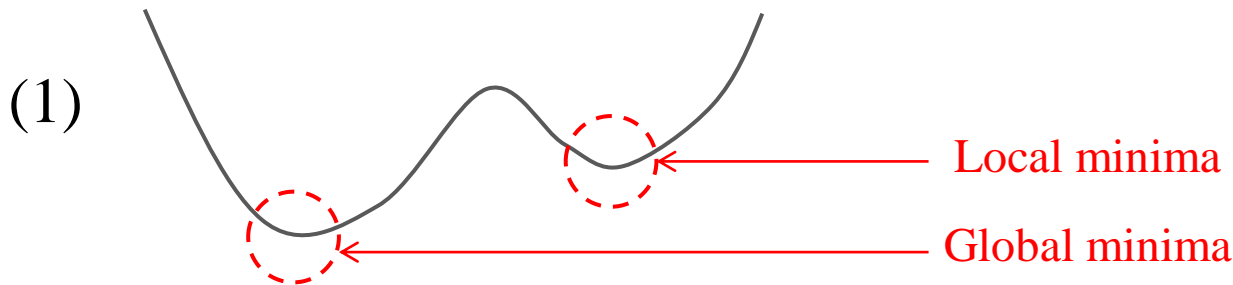
Introduction

- Optimization
 - Frequently faced in algorithm design
 - Used to implement you creative idea
- Issue
 - What kinds of **mathematical form** and its corresponding **optimization algorithms** do guarantee the convergence to **local** or **global optima**?

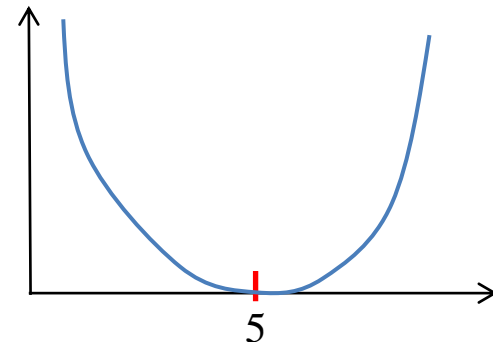
The Fundamental of Optimization

A Warming-up Question

- How do you solve the following problems?



(2) $\min_w f(w) = (w-5)^2$ \Longrightarrow (a) Plot



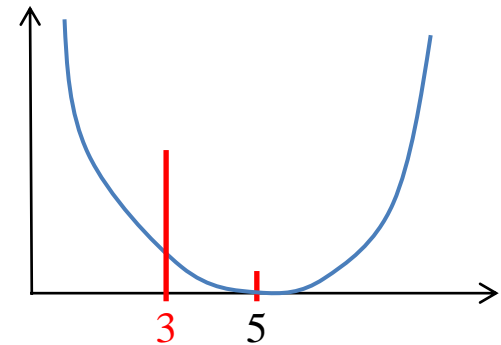
(b) Take derivatives, check = 0

An Advanced Question

- How about the following questions?

(3) $\min_w \sum_{n=1}^N f_n(w)$ \implies (a) Plot? (b) Take derivative = 0?

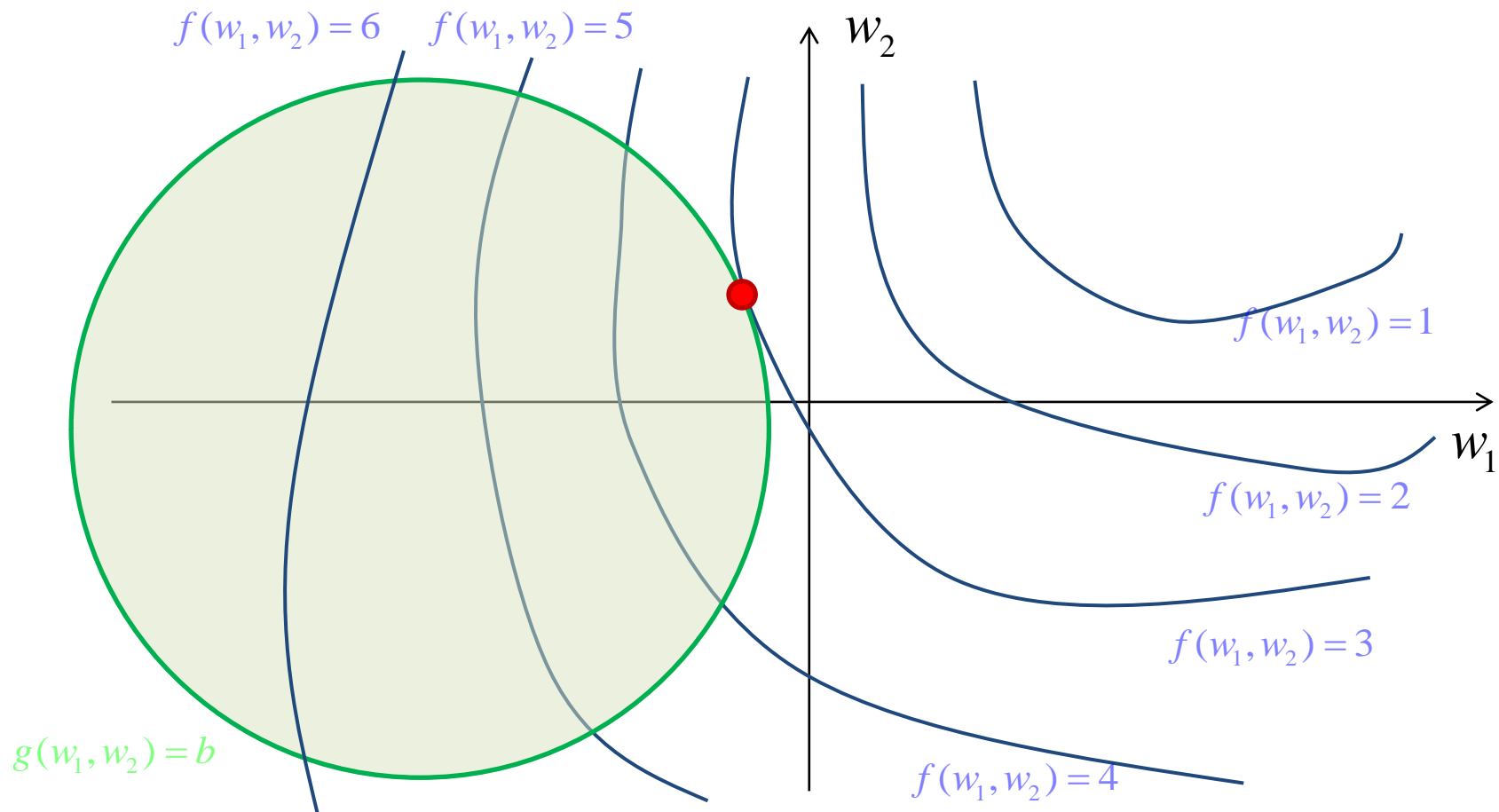
(4) $\min_w f(w) = (w-5)^2$
s.t. $w \leq 3$ \implies Derivative?



(5) $\min_w \sum_{n=1}^N f_n(w)$
s.t. $w_i \leq b_i$ \implies How to do?

Illustration

- 2-D case: $\min_{w_1, w_2} f(w_1, w_2)$, s.t. $g(w_1, w_2) \leq b$



How to Solve?

- Thanks to.....
 - Lagrange multiplier
 - Linear programming, quadratic programming, and recently, convex optimization
- Standard form:

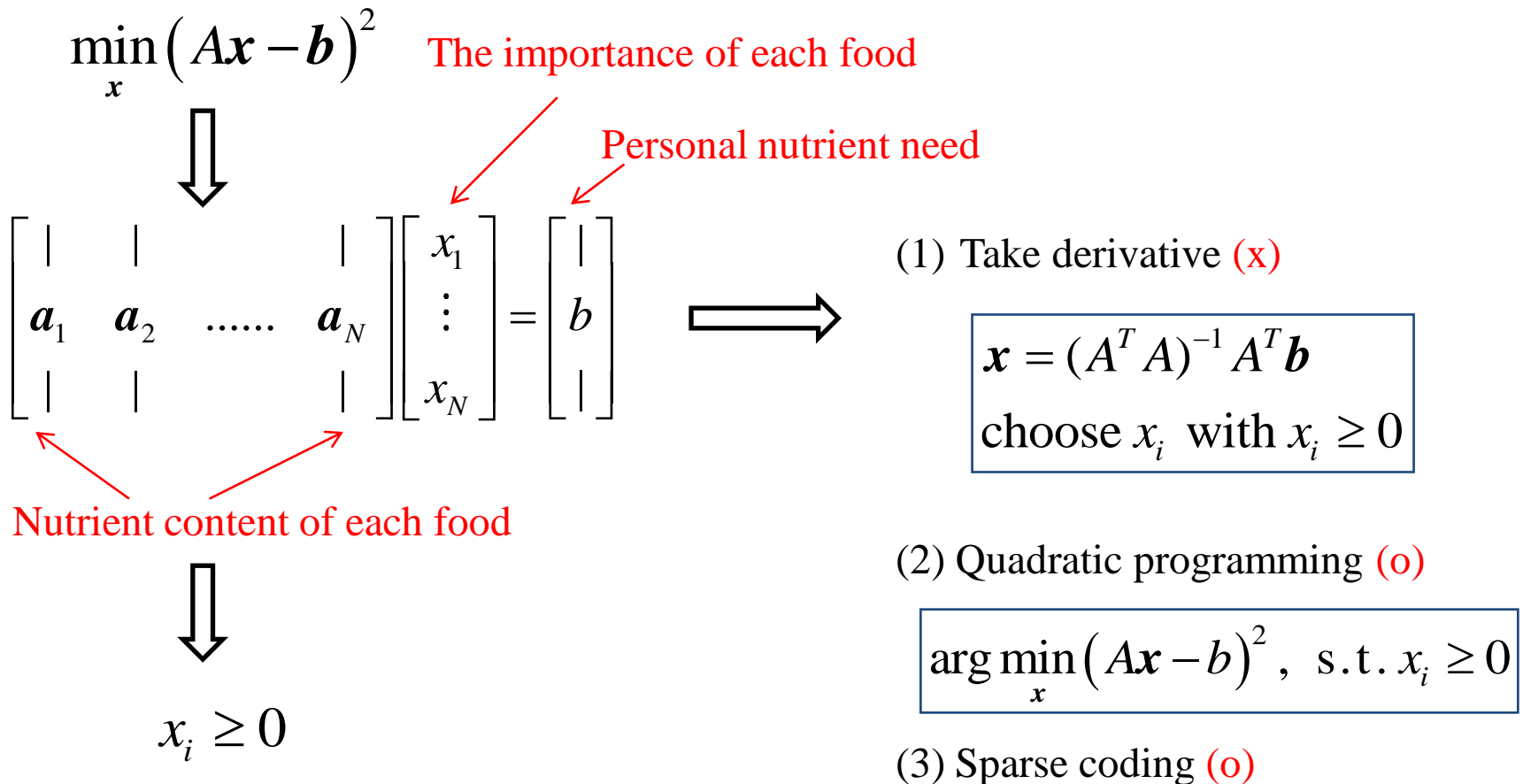
$$\min_{\mathbf{w}} f_0(\mathbf{w})$$

$$\text{s.t. } h_i(\mathbf{w}) \leq b_i, i = 1, \dots, m$$

$$\text{s.t. } g_i(\mathbf{w}) = c_i, i = 1, \dots, n$$

Fallacy

- A quadratic programming problem with constraints



The Idea of Sparsity

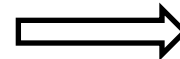
What is Sparsity?

- Think about a problem?

Assume full rank, $N > d$

$$\min_x (Ax - b)^2$$

$$\underbrace{\begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_N \\ | & | & & | \end{bmatrix}}_N \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} | \\ \mathbf{b} \\ | \end{bmatrix} \in R^d$$

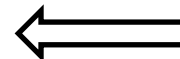


Many x can achieve

$$\min_x (Ax - b)^2 = 0$$



Which do you want?



Choose the x with the least nonzero component

$$\arg \min_x \|\mathbf{x}\|_0, \text{ s.t. } (Ax - \mathbf{b})^2 = 0$$

Why Sparsity?

- The more **concise**, the more better
- In some domain, there naturally exists a **sparse latent vector** that controls the data we saw. (ex. MRI, music)

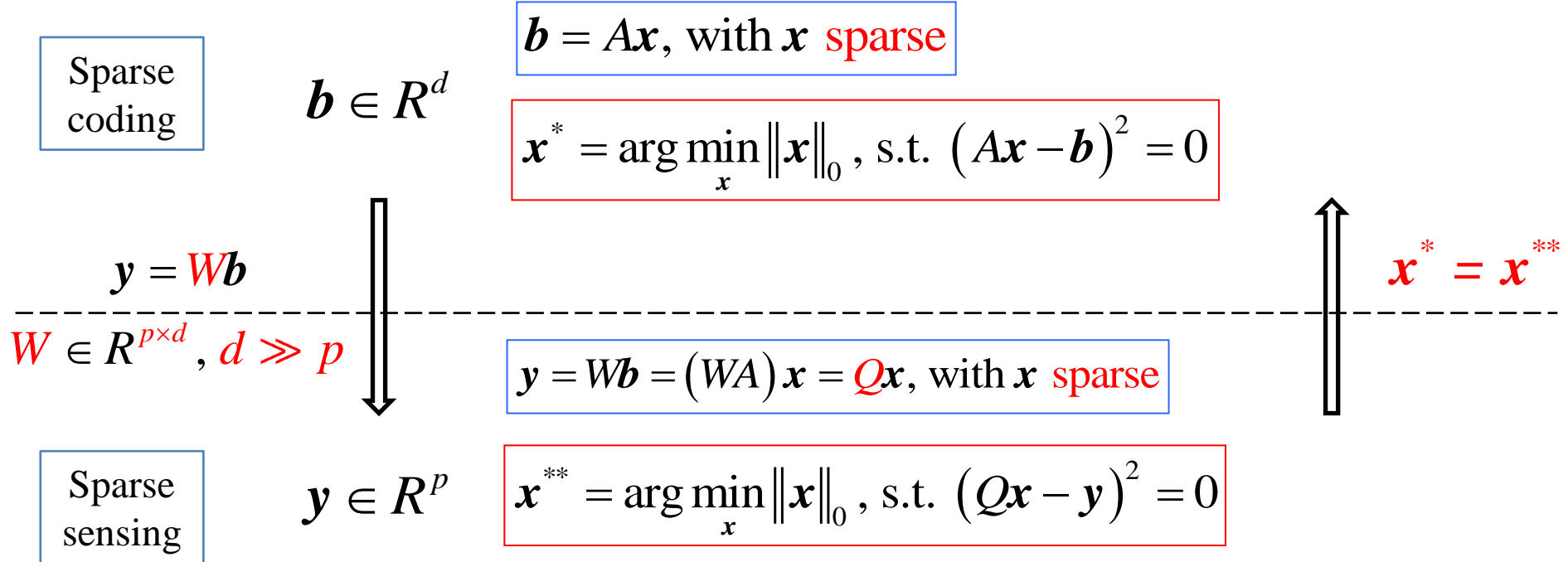
$$\begin{bmatrix} | \\ | \\ \mathbf{b} \\ | \\ | \end{bmatrix} = \begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_d \\ | & | & & | \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ x_i \\ \vdots \\ x_j \\ \vdots \\ 0 \end{bmatrix} + (\text{noise})$$

A ***k*-sparse** domain means that each ***b*** can be constructed by a ***x*** vector with at most ***k*** nonzero element

- In some domain, samples from the **same class** have the sparse property.
- The domain can be **learned**.

Sparse Sensing VS. Sparse Coding

- Assume that: We have $A \in R^{d \times N}$, $N > d$. Now an observation $\mathbf{b} \in R^d$ comes in



Note: p is based on the sparsity of the data (on k)

Sparse Sensing

$$\mathbf{b} \in \mathbb{R}^d \quad \mathbf{b} = \mathbf{A}\mathbf{x}, \text{ with } \mathbf{x} \text{ sparse}$$

$$\mathbf{y} = \mathbf{W}\mathbf{b}$$

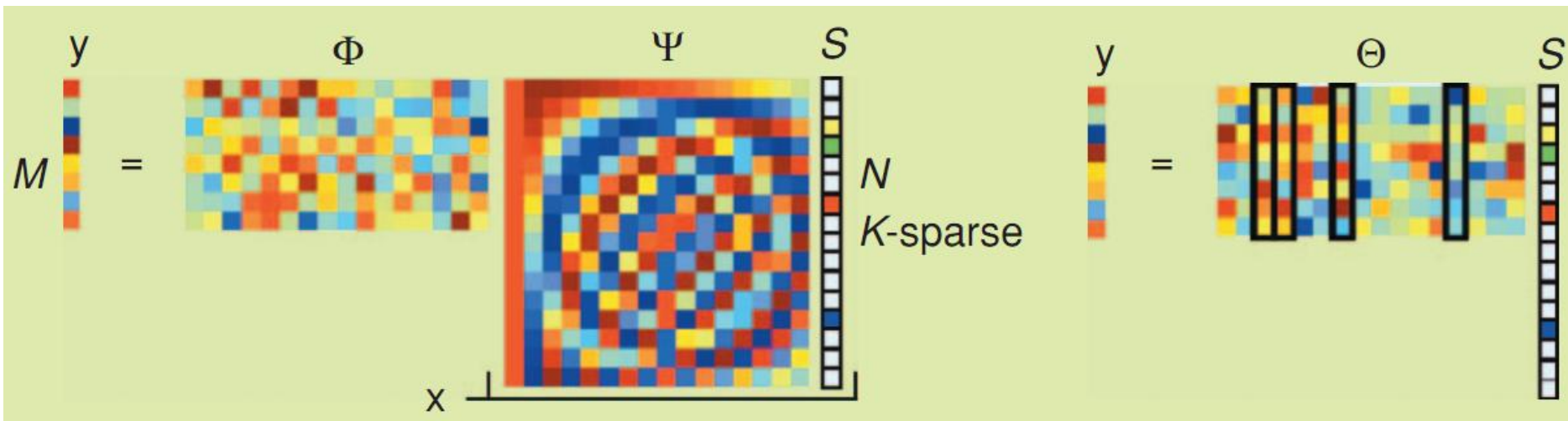
$$\mathbf{W} \in \mathbb{R}^{p \times d}, \quad d \gg p$$



$$\mathbf{y} \in \mathbb{R}^p \quad \mathbf{y} = \mathbf{W}\mathbf{b} = (\mathbf{W}\mathbf{A})\mathbf{x} = \mathbf{Q}\mathbf{x}, \text{ with } \mathbf{x} \text{ sparse}$$



$$\mathbf{x}^* = \mathbf{x}^{**}$$



Sparse Sensing VS. Sparse Coding

- Sparse sensing (compressed sensing):
 - It spends much time or money to get \mathbf{b} , so get \mathbf{y} first then recover \mathbf{b}
- Sparse coding (sparse representation):
 - Believe that there exists the **sparse property** in the data, otherwise sparse representation means nothing.
 - \mathbf{x} is used to be the **feature** of \mathbf{b}
 - \mathbf{x} can be used to **efficiently store** \mathbf{b} and **reconstruct** \mathbf{b}

The Solution

How to Get The Sparse Solution?

- There is no algorithm other than exhaustively searching to solve:

$$\mathbf{x}^* = \arg \min_x \|\mathbf{x}\|_0, \text{ s.t. } (\mathbf{Ax} - \mathbf{b})^2 = 0$$

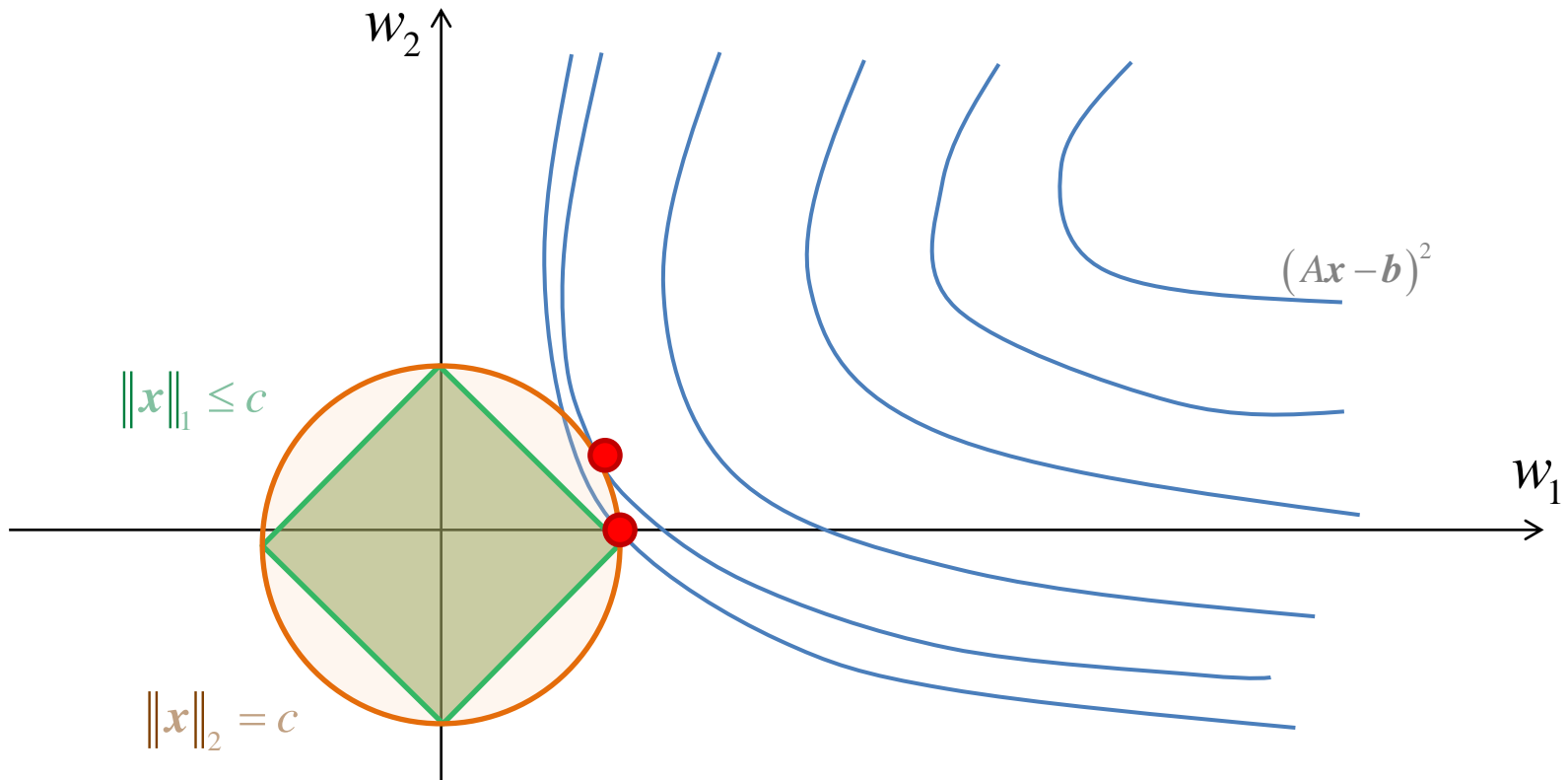
- While in some situations (ex. special form of A), the solution of l_1 minimization approaches the one of l_0 minimization

$$\mathbf{x}^{***} = \arg \min_x \|\mathbf{x}\|_1 = \sum_{n=1}^N |x^{(n)}|, \text{ s.t. } (\mathbf{Ax} - \mathbf{b})^2 = 0$$
$$\mathbf{x}^{***} = \mathbf{x}^*$$

Why l_1 ?

- Question 1: Why l_1 can result in a **sparse** solution?

$$\arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } (A\mathbf{x} - \mathbf{b})^2 = 0 \rightarrow \arg \min_x (A\mathbf{x} - \mathbf{b})^2, \text{ s.t. } \|\mathbf{x}\|_1 \leq c$$



Why l_1 ?

- Question 2: Why the sparse solution achieved by l_1 minimization approaches the one of l_0 minimization?
 - This is a matter of Mathematics
 - No matter how, sparse representation based on l_1 minimization has been widely used for pattern recognition.
 - In addition, if one doesn't care about using the sparse solution for **representation (feature)**, it seems **OK** if these two solutions are not the same.

$$\mathbf{b} = A\mathbf{x}^{***}$$

$$\mathbf{b} = A\mathbf{x}^*$$

Noise

- Sometimes, the data is observed with **noise**

$$\begin{bmatrix} | \\ | \\ | \end{bmatrix} \mathbf{b} = \begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \dots & \mathbf{a}_d \\ | & | & & | \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ x_i^* \\ \vdots \\ x_i^* \\ \vdots \\ 0 \end{bmatrix} \xrightarrow{l_0(l_1) \text{ minimization}} \mathbf{b} = \mathbf{b} + \text{noise} \xrightarrow{\quad} \mathbf{x} \approx \mathbf{x} ???$$

- The answer seems to be **negative**

$$\mathbf{b} = \mathbf{A}\mathbf{x}^*, \mathbf{y}^* = \arg \min_y \|\mathbf{y}\|_1, \text{ s.t. } (\mathbf{A}\mathbf{y} - \text{noise})^2 = 0$$

$\mathbf{y}^* + \mathbf{x}^* \neq \mathbf{x}^*$, and \mathbf{y}^* is usually **not sparse**

$\mathbf{x} = \arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } (\mathbf{A}\mathbf{x} - \mathbf{b})^2 = 0$ is neither equal to $\mathbf{y}^* + \mathbf{x}^*$ nor to \mathbf{x}^*

possibly not sparse

Noise

- Several ways to overcome this:

$$\begin{aligned} \arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } (\mathbf{Ax} - \mathbf{b})^2 = 0 &\rightarrow \arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } \|\mathbf{Ax} - \mathbf{b}\|_2 \leq c \\ &\rightarrow \arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } \|\mathbf{Ax} - \mathbf{b}\|_1 \leq c \end{aligned}$$

$$\arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } (\mathbf{Ax} - \mathbf{b})^2 = 0 \rightarrow \arg \min_z \|\mathbf{z}\|_1, \text{ s.t. } ([\mathbf{A} \mid \mathbf{I}]\mathbf{z} - \mathbf{b})^2 = 0, \text{ where } \mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{t} \end{bmatrix}$$

- What is the difference between:

$$\|\mathbf{Ax} - \mathbf{b}\|_2 \leq c \text{ and } \|\mathbf{Ax} - \mathbf{b}\|_1 \leq c$$

Equivalent form

- You may also see several forms for the problem:

$$\begin{aligned} \arg \min_x \|\mathbf{x}\|_1, \text{ s.t. } \|\mathbf{Ax} - \mathbf{b}\|_1 \leq \mathbf{c} &\rightarrow \arg \min_x \|\mathbf{x}\|_1 + \lambda \|\mathbf{Ax} - \mathbf{b}\|_1 \\ &\rightarrow \arg \min_x \|\mathbf{Ax} - \mathbf{b}\|_1, \text{ s.t. } \|\mathbf{x}\|_1 \leq \mathbf{d} \end{aligned}$$

- These equivalent forms are derived from Lagrange multiplier
- There have been several publications aiming at how solving the l_1 minimization problem.

The Importance of Dictionary

Dictionary generation

- If the preceding sections, we generally assume that the (**over-complete**) bases A is existed and known
- However in practice, we usually need to build it:
 - Wavelet + Fourier + Haar +
 - Learning based on data
- How to learn?

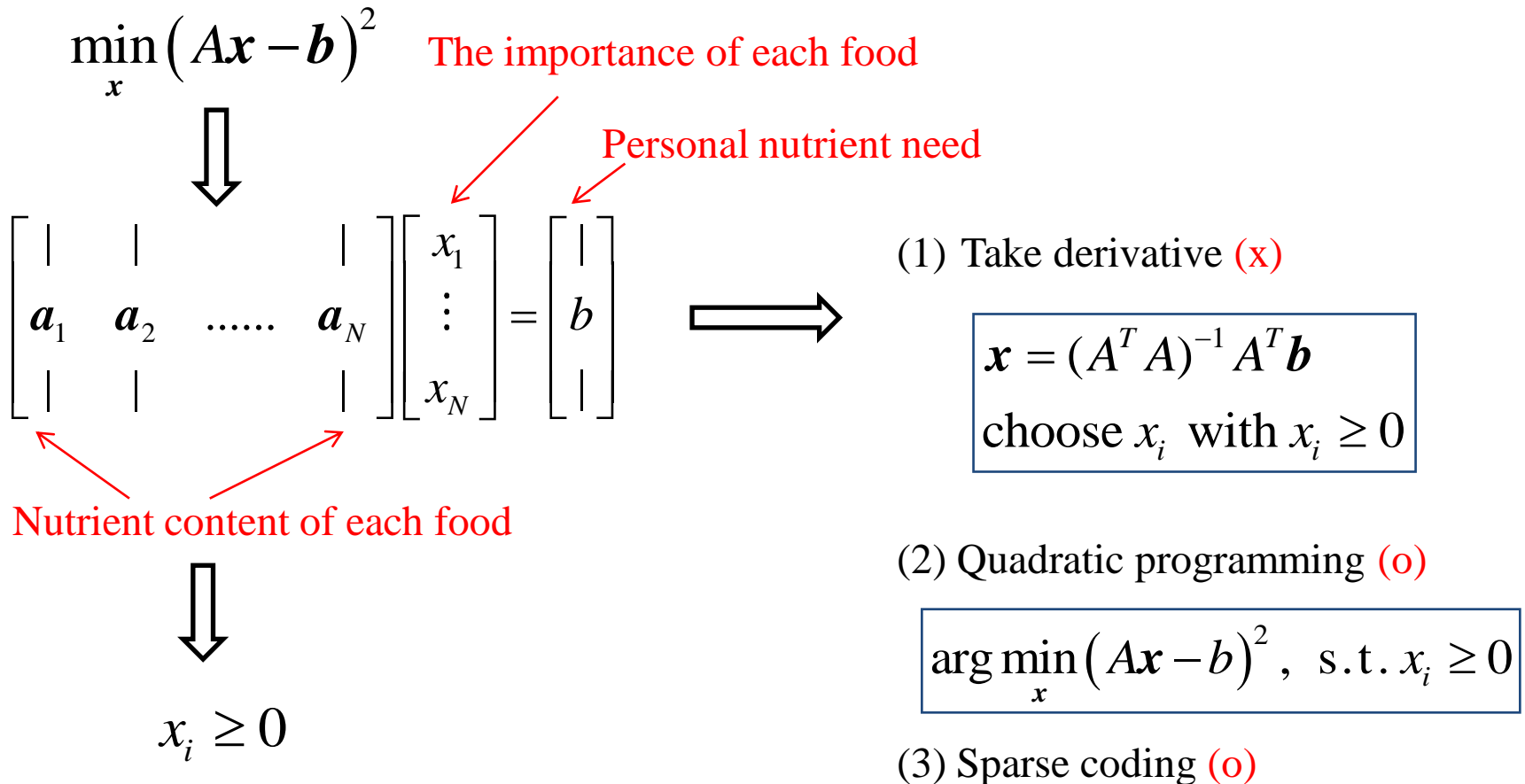
$$\text{Given a training set } \{\mathbf{b}^{(i)} \in R^d\}_{i=1}^N, \text{ form } B \text{ as } B = [\mathbf{b}^{(1)} \ \mathbf{b}^{(2)} \ \dots \ \mathbf{b}^{(N)}]$$
$$(A^*, X^*) = \arg \min_{A, X} \|B - AX\|_F^2 + \lambda \|X\|_1, \text{ where } X = [\mathbf{x}^{(1)} \ \mathbf{x}^{(2)} \ \dots \ \mathbf{x}^{(N)}]$$

- May result in **over-fitting**

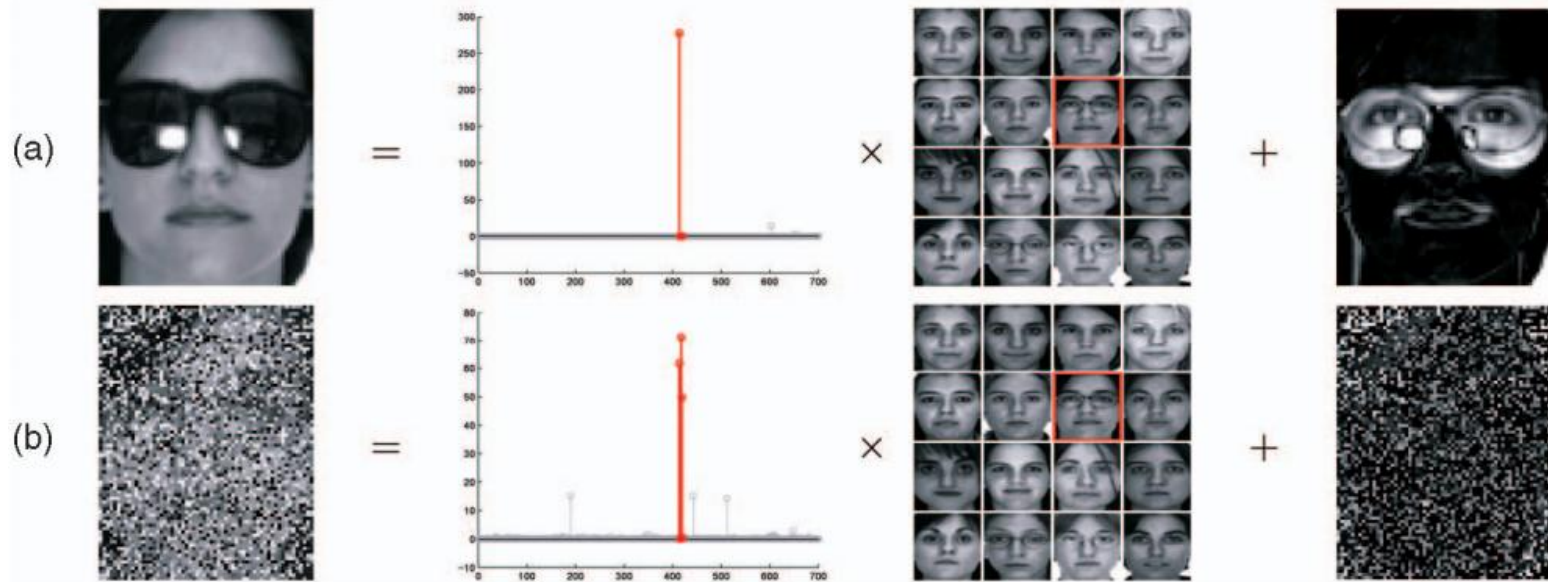
Applications

Back to the problem we have

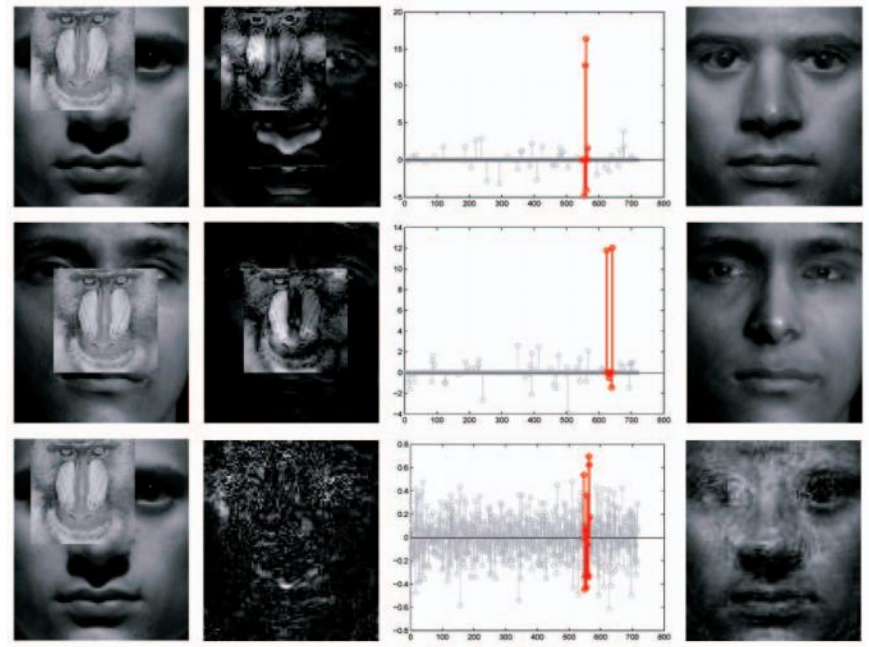
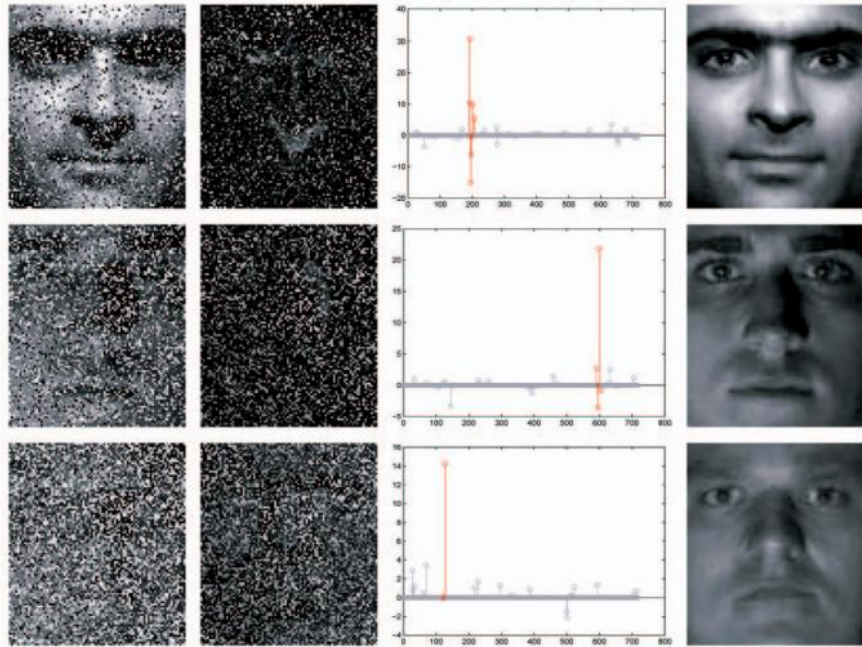
- A quadratic programming problem with constraints



Face Recognition (1)



Face Recognition (2)



An important issue

- When using sparse representation as a way of feature extraction, you may wonder, even if there exists the sparsity property in the data, does sparse feature really come up with **better results**? Does it contain any **semantic meaning**?
- Successful areas:
 - Face recognition
 - Digit recognition
 - Object recognition (with **careful design**):
Ex. K-means → Sparse representation

De-noising

Learn a **patch dictionary**.
For each patch, compute the **sparse representation**
then use it to **reconstruct**
the patch.

$$\mathbf{x}^* = \arg \min_x \|\mathbf{x}\|_1 + \lambda \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_1$$

$$\mathbf{b} = \mathbf{A}\mathbf{x}^*$$



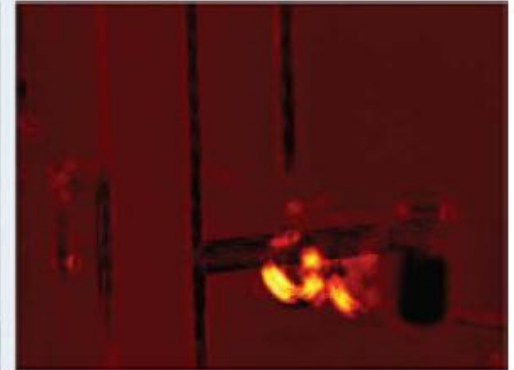
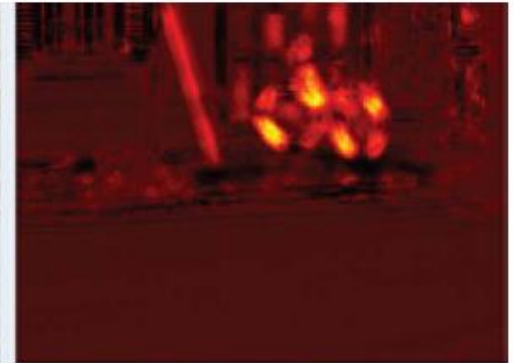
Detection based on reconstruction

Learn a **patch dictionary for a specific object**. For each patch in the image, compute the **sparse representation** and use it to **reconstruct** the image. Check the **error** for each patch, and identify those with **small error** as detected object.

$$\mathbf{x}^* = \arg \min_x \|\mathbf{x}\|_1 + \lambda \|\mathbf{Ax} - \mathbf{b}\|_1$$

$$\mathbf{b} = \mathbf{Ax}^*$$

$$\text{check } \|\mathbf{b} - \mathbf{b}\|_2^2$$



Maybe not **over-complete**

Other cases: Foreground-background detection, pedestrian detection,

Conclusion

What you should know

- What is the form of standard optimization?
- What is sparsity?
- What is sparse coding and sparse sensing?
- What kind of optimization method to solve it?
- Try to use it !!

Thank you for listening