



# *CS545 – Machine Learning for Signal Processing*

# Sparsity, Compressive Sensing and Random Projections

30 October 2023

# Today's lecture

- Sparsity
- Compressive sensing
- Quantization, randomness and high dimensions

# What is sparsity?

- Depends who you ask
- Basic idea: We want most numbers in a collection to be zero
- Too many ways to express that

# A starting point

- Linear equation with multiple solutions:

$$y = \mathbf{a} \cdot \mathbf{x} \Rightarrow 2 = \begin{bmatrix} 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

- Which solution would you pick?

$$\mathbf{x} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

*A sparse answer*

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

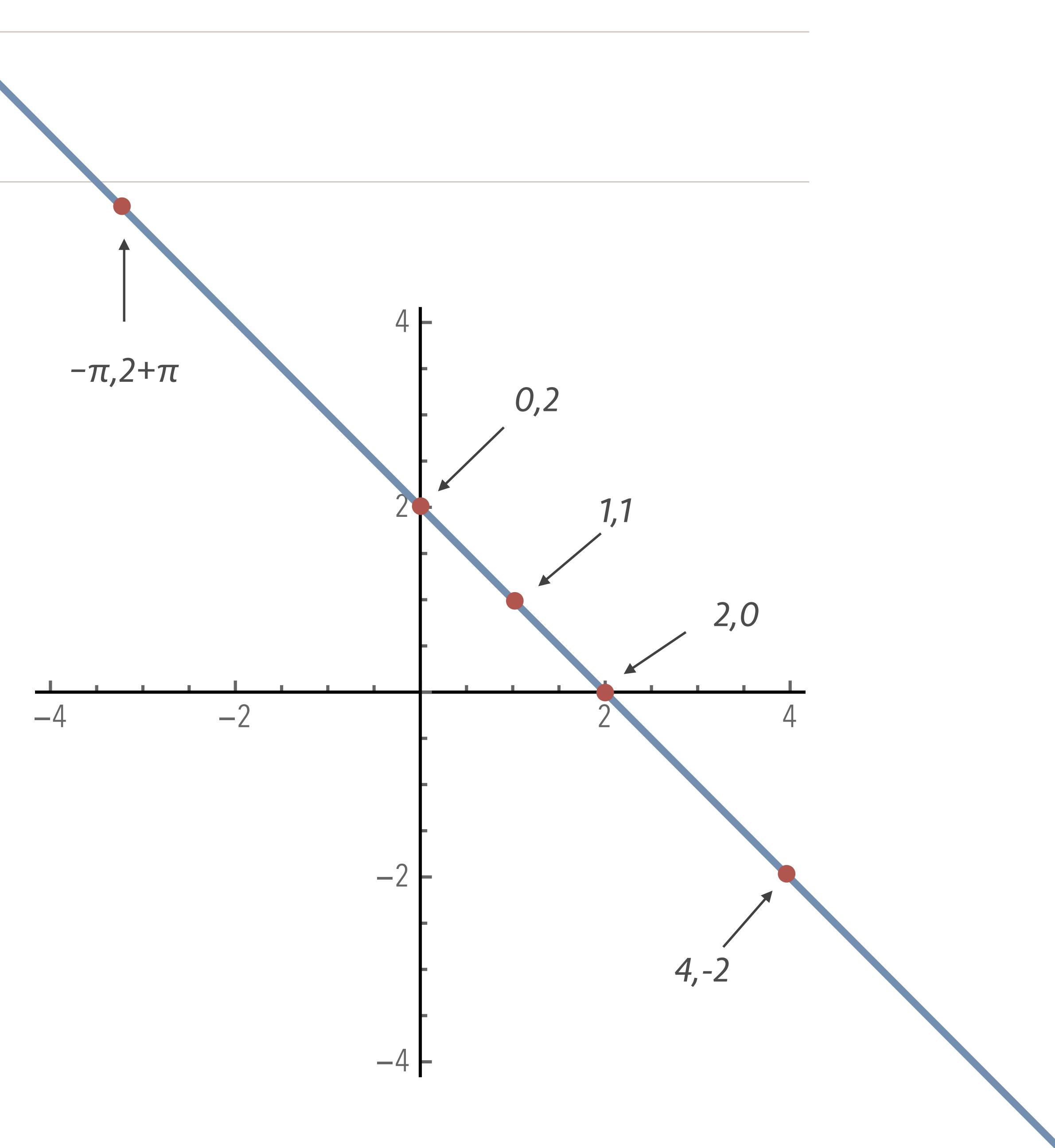
*What most software gives*

$$\mathbf{x} = \begin{bmatrix} 4 \\ -2 \end{bmatrix}$$

*This one is fine too!*

# Infinite solutions

- All solutions lie on a line
- Which one do we pick?
  - Why does most software pick  
 $\mathbf{x} = [1, 1]$ ?
- Does it make a difference?



# The generic answer

- Least squares problem:

$$\begin{aligned} \mathbf{y} &= \mathbf{A} \cdot \mathbf{x} \Rightarrow \mathbf{x} = \mathbf{A}^+ \cdot \mathbf{y} \\ &\Rightarrow \arg \min_{\mathbf{x}} (\|\mathbf{A} \cdot \mathbf{x} - \mathbf{y}\|_2 + \|\mathbf{x}\|_2) \end{aligned}$$

- Find the minimum-norm  $\mathbf{x}$  that minimizes the error
  - But why  $\|\mathbf{x}\|_2$ ? Makes the math easy!

# Least squares, pseudoinverse, and $\ell_2$

- Use Langrangian multipliers:

$$\|\mathbf{x}\|_2^2 + \lambda^\top \cdot (\mathbf{A} \cdot \mathbf{x} - \mathbf{y}) \Rightarrow \hat{\mathbf{x}} = -\frac{1}{2} \mathbf{A}^\top \cdot \lambda$$

- Put back in original equation:

$$\begin{aligned} \mathbf{A} \cdot \hat{\mathbf{x}} &= -\frac{1}{2} \mathbf{A} \cdot \mathbf{A}^\top \cdot \lambda = \mathbf{y} \Rightarrow \lambda = -2 (\mathbf{A} \cdot \mathbf{A}^\top)^{-1} \cdot \mathbf{y} \\ \Rightarrow \hat{\mathbf{x}} &= \mathbf{A}^\top \cdot (\mathbf{A} \cdot \mathbf{A}^\top)^{-1} \cdot \mathbf{y} = \mathbf{A}^+ \cdot \mathbf{y} \end{aligned}$$

# Many more norms

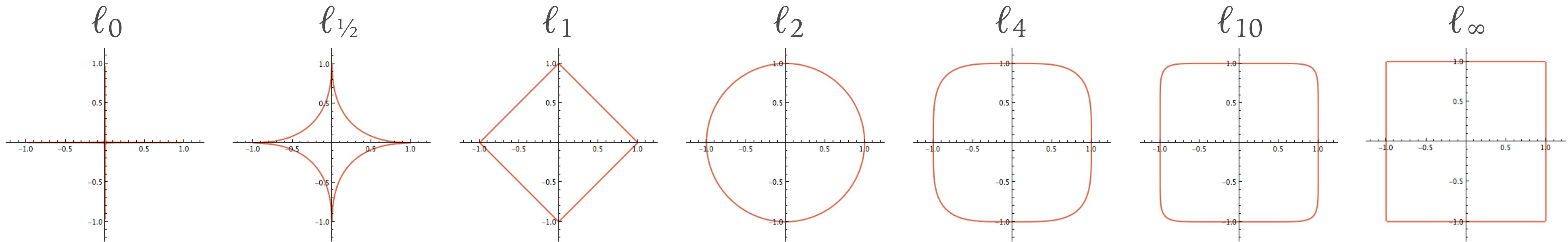
- $p$ -Norm (or  $L_p$  /  $\mathbb{L}_p$  /  $\ell_p$ )

$$\|\mathbf{x}\|_p = \sqrt[p]{\sum |x_i|^p}$$

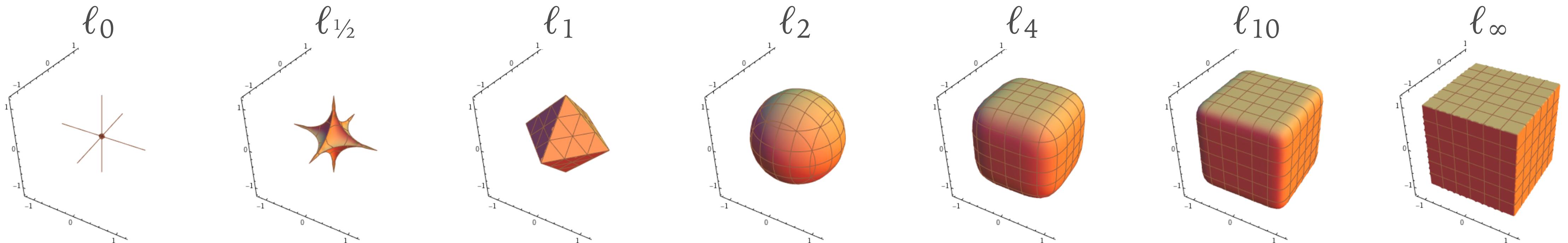
- $\ell_2$  norm is the Euclidean norm
- $\ell_1$  norm is sum of absolute values
- $\ell_0$  norm is the number of non-zero values
- $\ell_\infty$  norm is max of all values

# How do they look?

- Unit norms in 2D

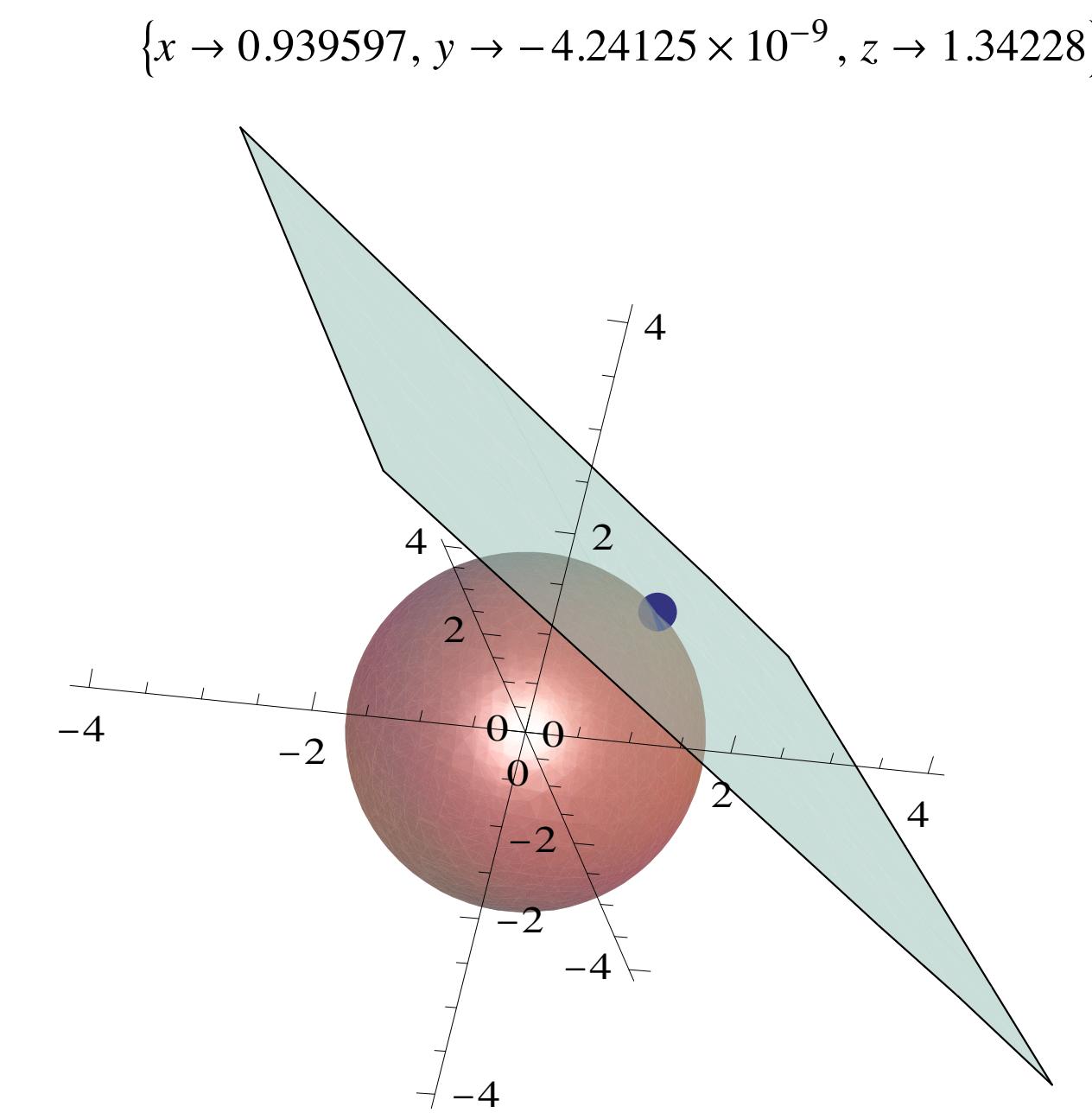
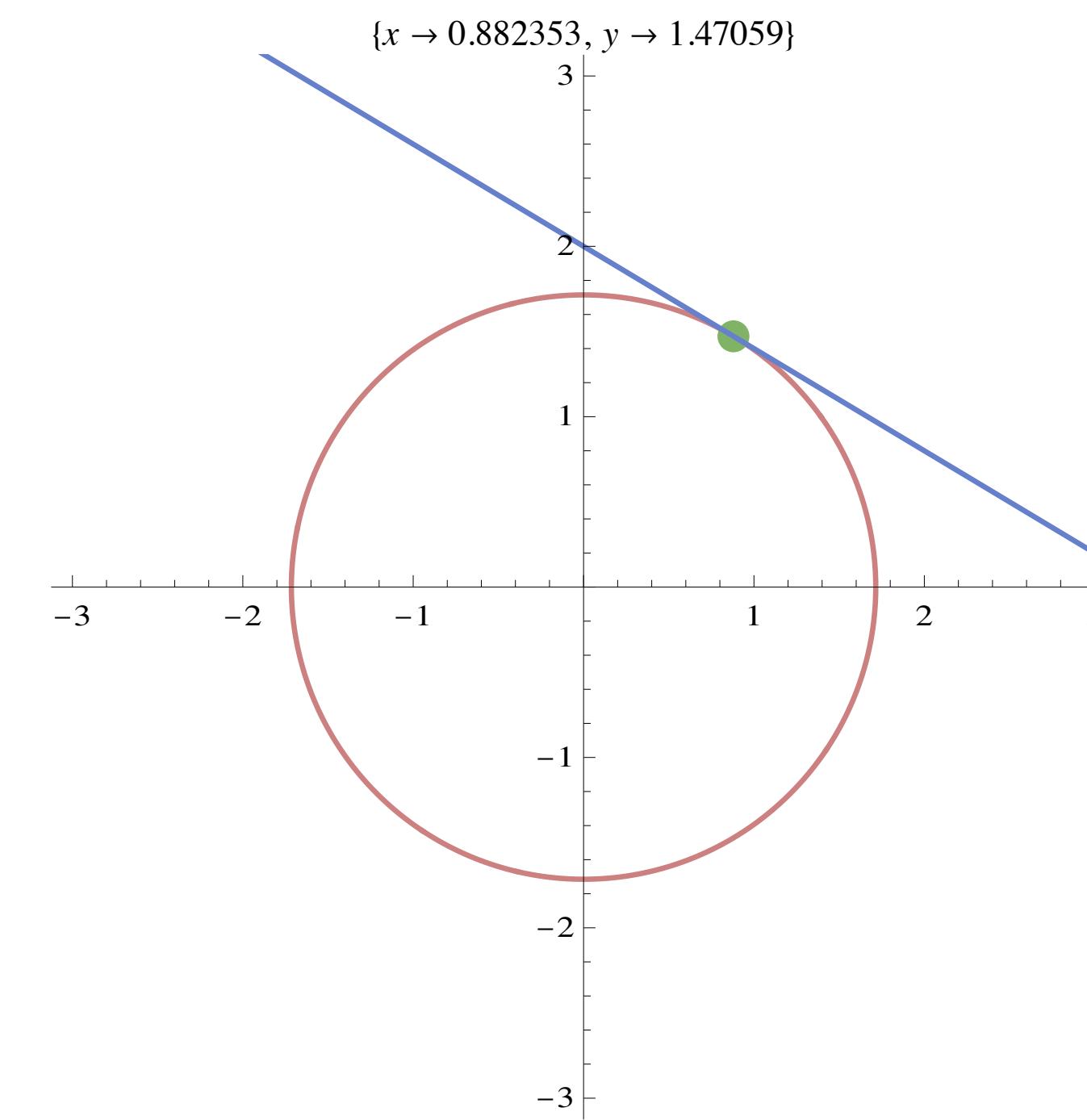


- Unit norms in 3D



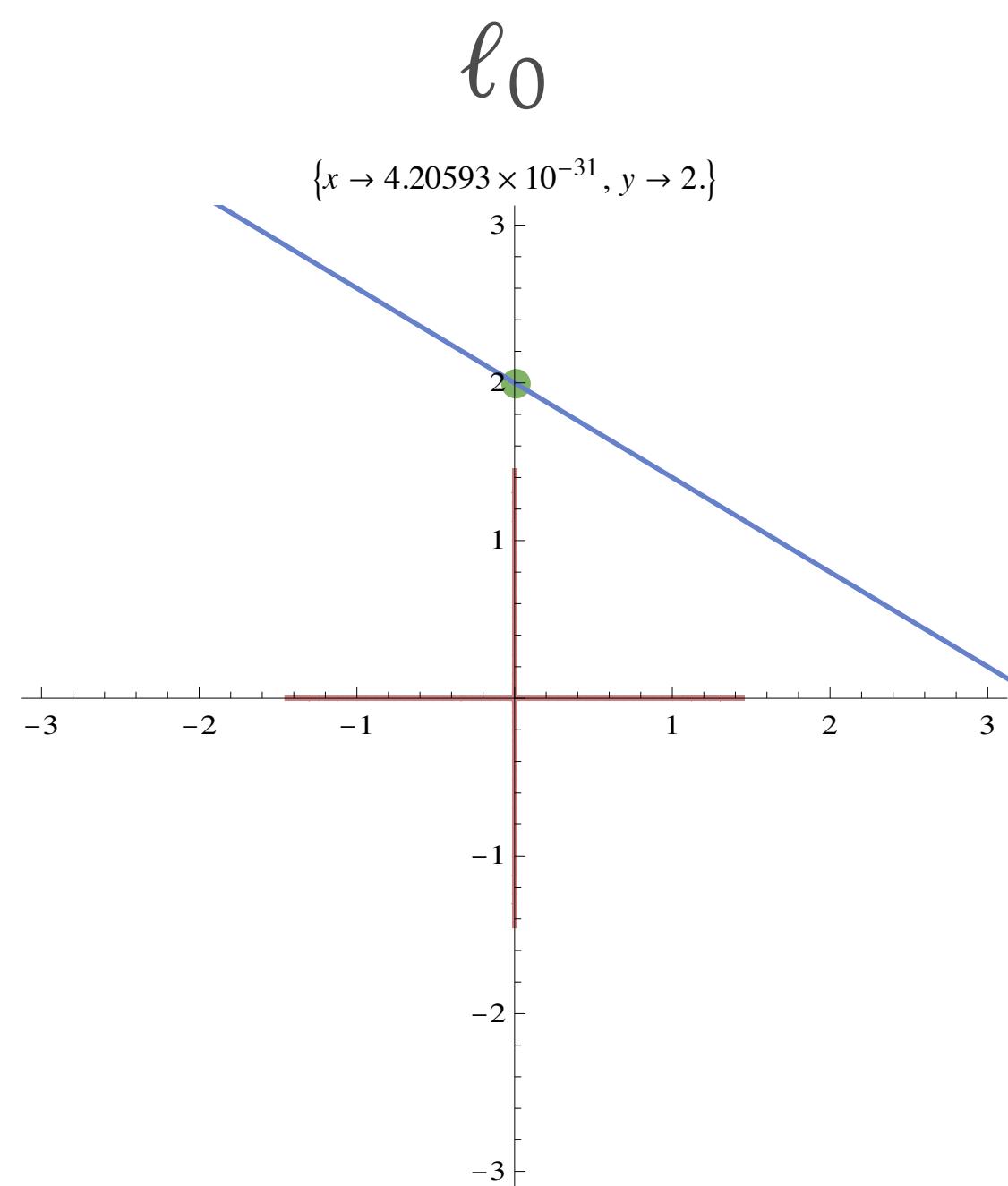
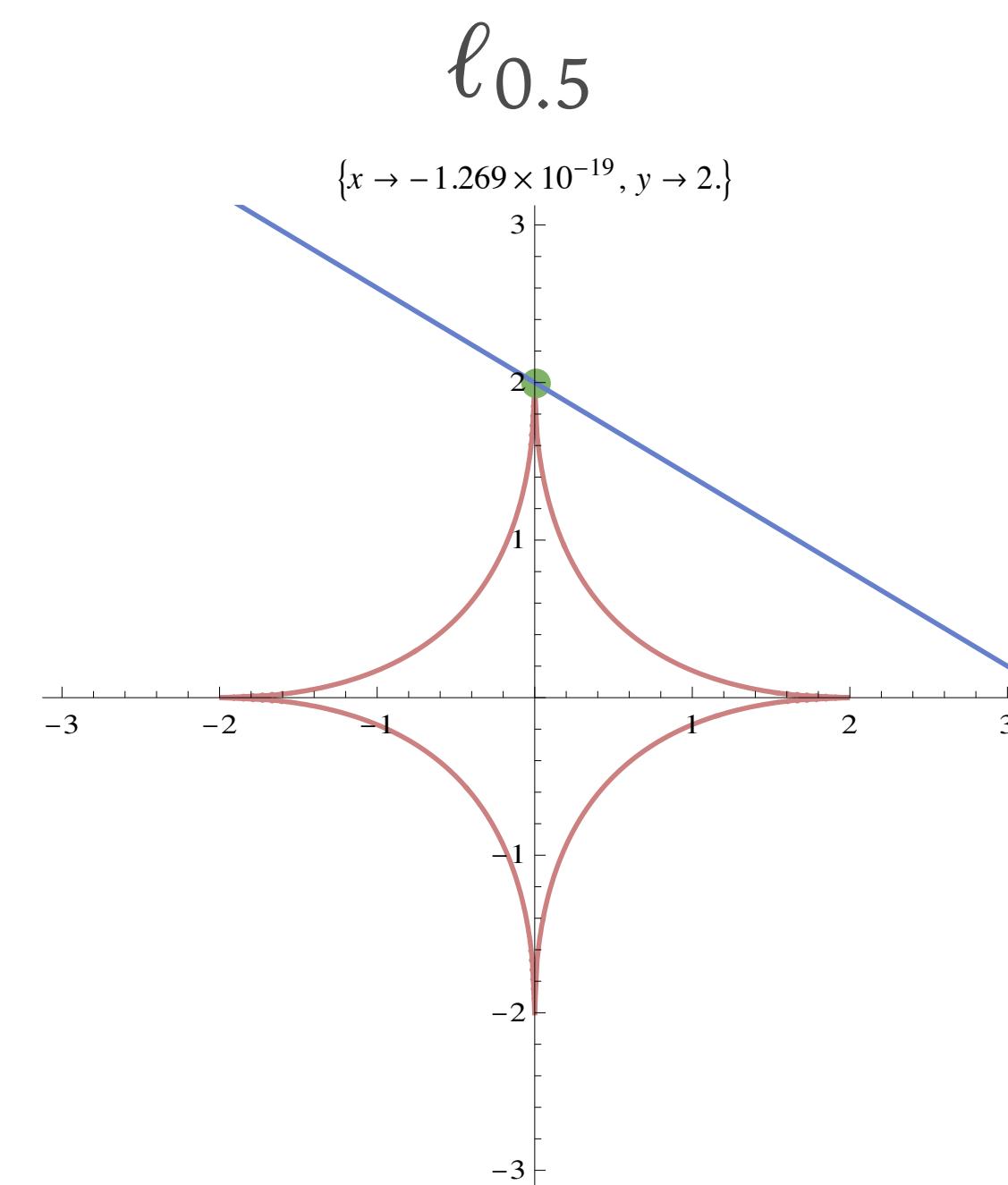
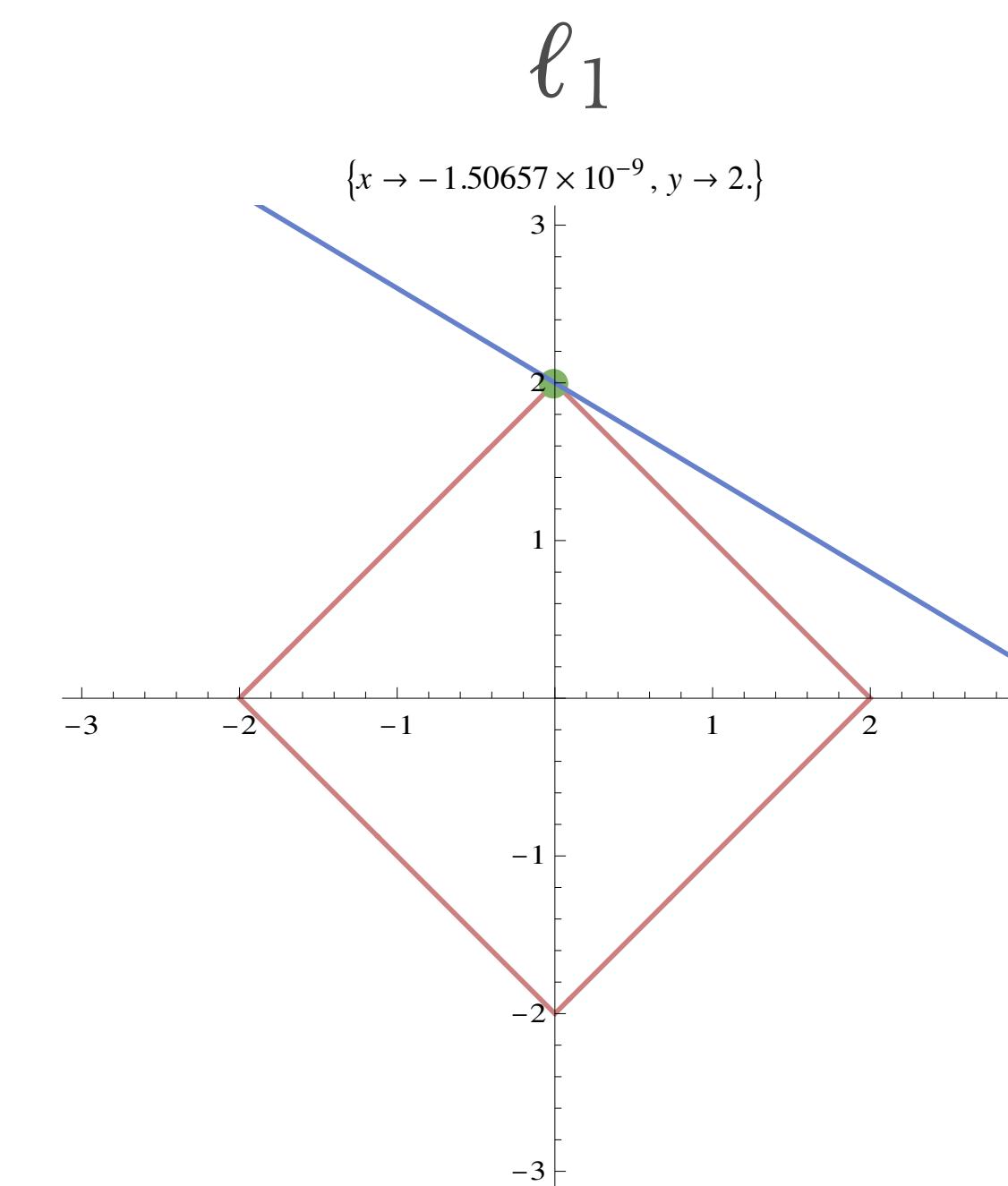
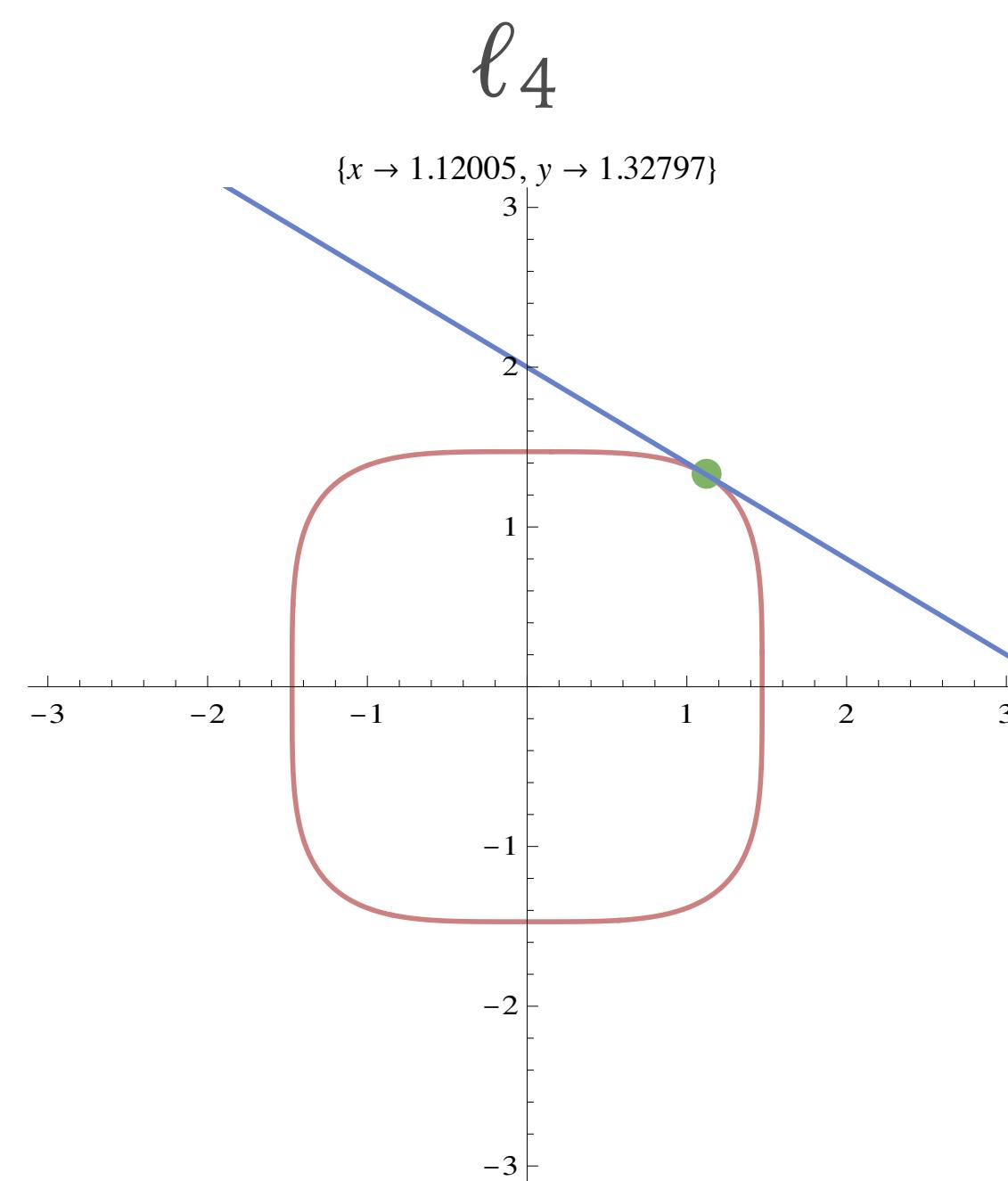
# $\ell_2$ -based pseudoinverse

- All possible solutions lie on a hyperplane
  - Minimum  $\ell_2$  solution will be the point where the smallest possible  $\ell_2$ -ball touches the solutions hyperplane



# Other $\ell_p$ -norm solutions

- With different norms the chosen solution will change
  - Larger  $p$  will produce a “busier” solution
  - Smaller  $p$  will produce *sparser* solution



# Which one to use?

- For sparsity we ideally want minimum  $\ell_0$ 
  - Directly results in smallest number of non-zero values
- But, this is an inconvenient form ...

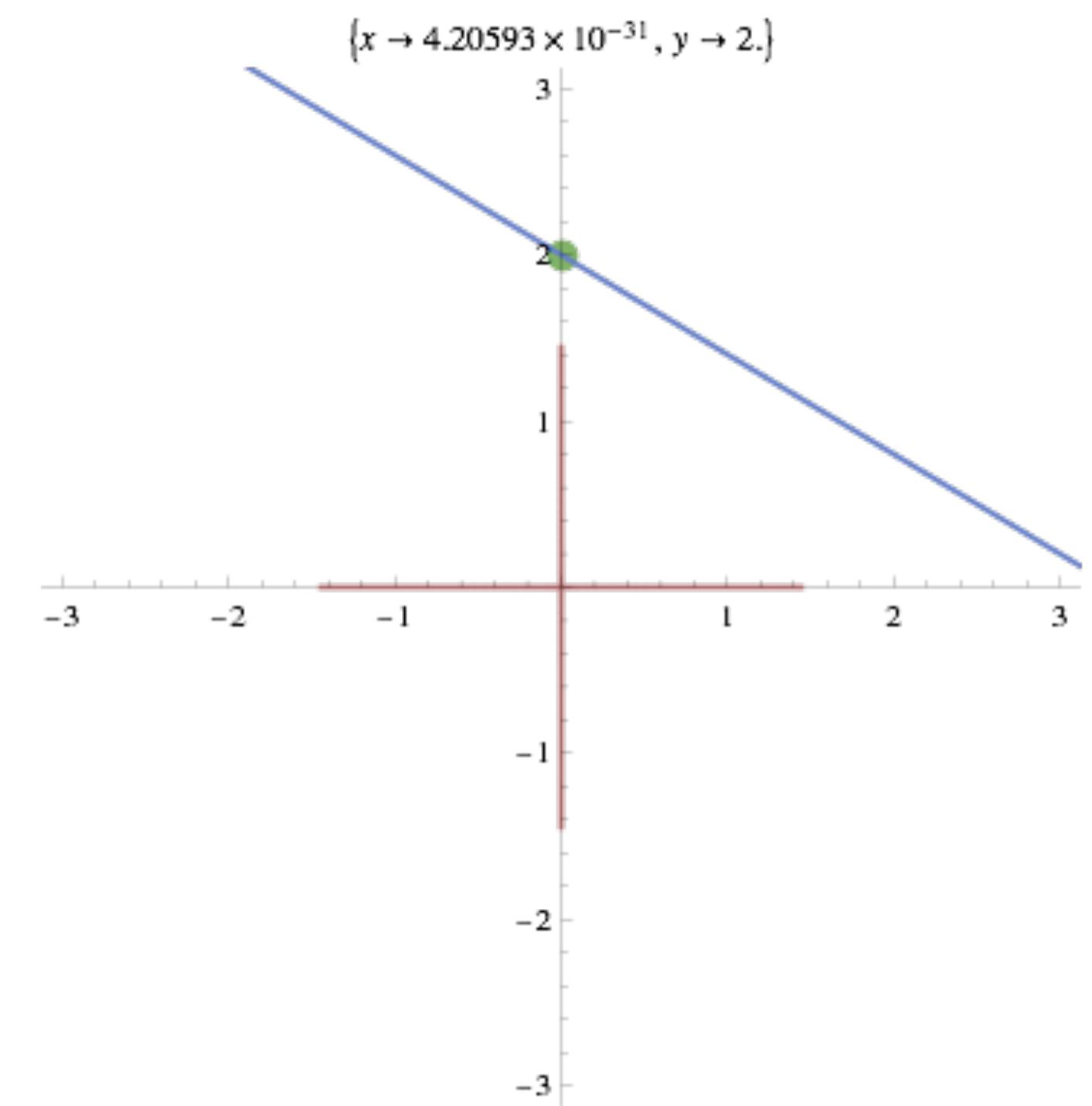
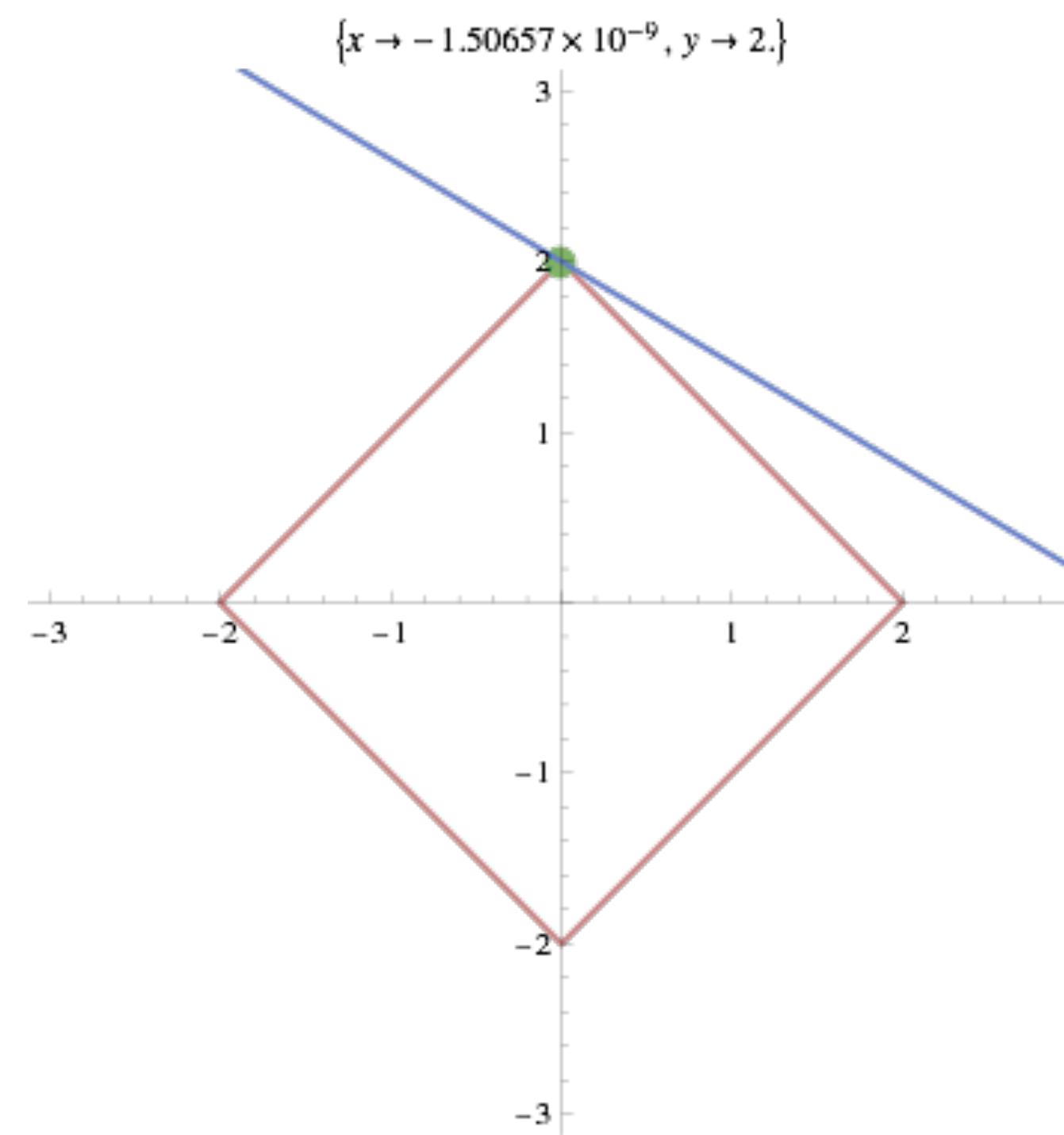
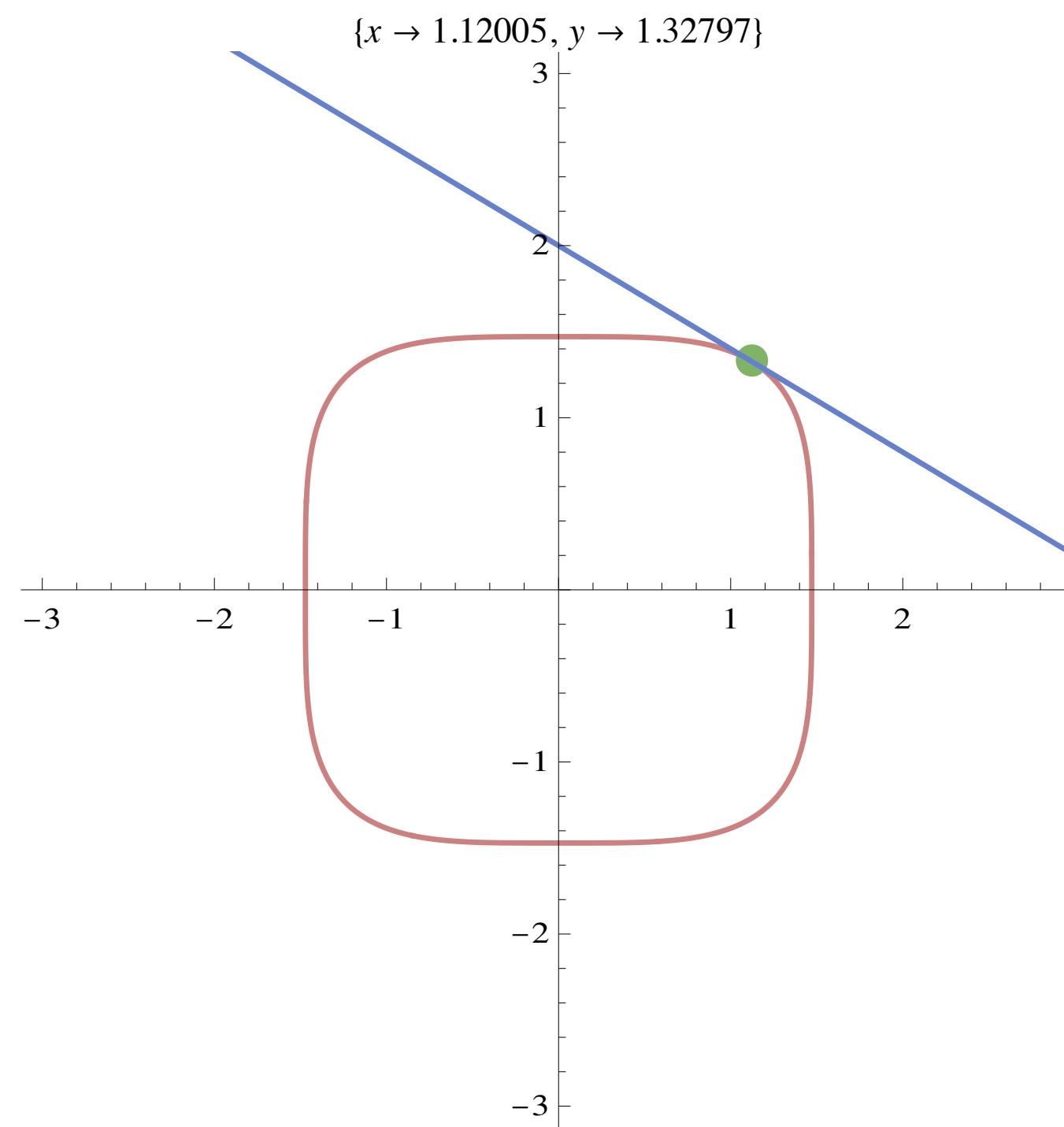
$$\ell_0(\mathbf{x}) = \sum_i [x_i \neq 0] \quad \begin{array}{l} \text{Iverson bracket:} \\ \text{if content evaluates} \\ \text{to true return 1} \\ \text{otherwise return 0} \end{array}$$

- Discontinuous, no derivative, not convex, etc ...

$$\|\mathbf{x}\|_0 + \lambda^\top \cdot (\mathbf{A} \cdot \mathbf{x} - \mathbf{y}) \longrightarrow \frac{\partial \|\mathbf{x}\|_0}{\mathbf{x}} = ?$$

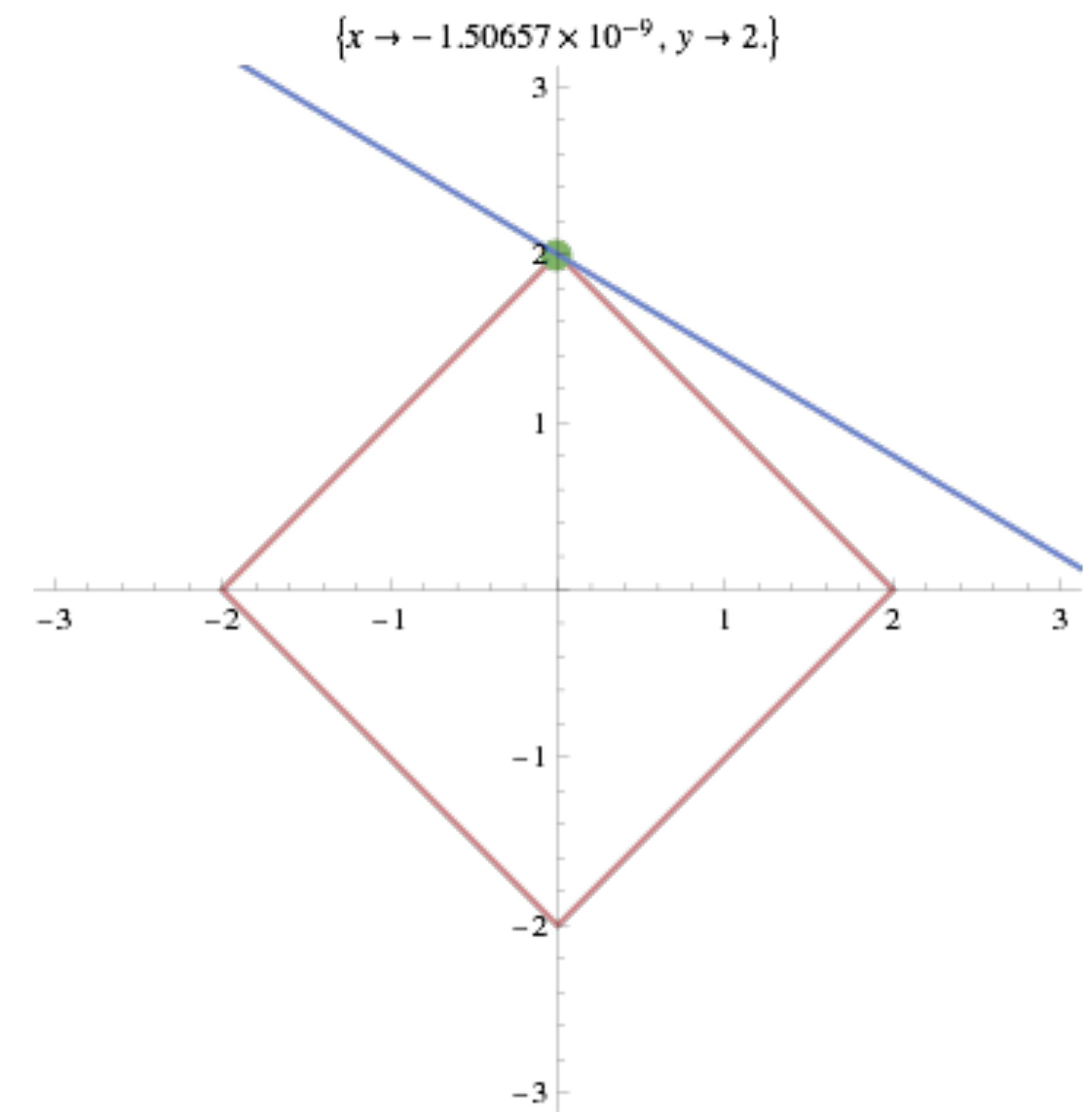
# Let's try something simpler then

- How about using the  $\ell_1$  instead?
  - Seems to produce the same solution



# Why does this work?

- The  $\ell_1$  case is (sort of) convex
  - $\ell_1$  is minimized as we move towards the ideal sparse solution
  - There is one ill-defined scenario, but it is not a big problem
    - Which is it?
- So let's solve that instead



# The problem to solve

- We now have:

$$\arg \min_{\mathbf{x}} (\|\mathbf{A} \cdot \mathbf{x} - \mathbf{y}\|_2 + \|\mathbf{x}\|_1)$$

- Minor glitch: We can't differentiate the absolute values in the  $\ell_1$  norm! Are we stuck again?

$$\sum |x_i| + \lambda^\top \cdot (\mathbf{A} \cdot \mathbf{x} - \mathbf{y})$$

- No, unlike  $\ell_0$  we can use other tools

# Linear programming

- A linear program is defined as:

$$\text{minimize} \quad \mathbf{c}^\top \cdot \mathbf{x}$$

$$\text{subject to} \quad \mathbf{A} \cdot \mathbf{x} \leq \mathbf{y}$$

$$\text{and} \quad x_i \geq 0$$

- A Nobel-prize staple of optimization theory, resource allocation, economics, etc.

# Doesn't exactly match the $\ell_1$ problem

- We would like to change our problem definition to fit the linear programming formulation

*What we have*

$$\text{minimize} \quad \|\mathbf{x}\|_1$$

$$\text{subject to} \quad \mathbf{A} \cdot \mathbf{x} = \mathbf{y}$$

*What we can solve*

$$\text{minimize} \quad \mathbf{c}^\top \cdot \mathbf{x}$$

$$\text{subject to} \quad \mathbf{A} \cdot \mathbf{x} \leq \mathbf{y}$$

$$\text{and} \quad x_i \geq 0$$

# With some shuffling around

- We rewrite the unknown vector  $\mathbf{x}$  as a difference of positive-valued vectors:

$$\mathbf{x} = \mathbf{u} - \mathbf{v}, \quad u_i, v_i \geq 0, \quad \mathbf{z} = \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}$$

- Now our problem can written as:

$$\mathbf{A} \cdot \mathbf{x} = \mathbf{y} \Rightarrow \begin{bmatrix} \mathbf{A}, & -\mathbf{A} \end{bmatrix} \cdot \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix} = \mathbf{y} \Rightarrow \begin{bmatrix} \mathbf{A}, & -\mathbf{A} \end{bmatrix} \cdot \mathbf{z} = \mathbf{y}$$

# Now it is a linear program

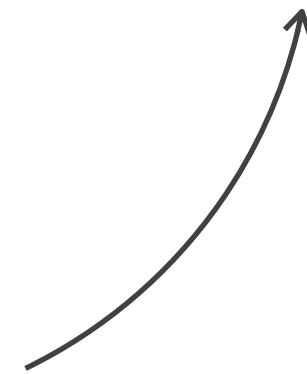
- And we can solve our problem

*Minimum  $\ell_1$  problem*

$$\begin{aligned} \text{minimize} \quad & \|x\|_1 \\ \text{subject to} \quad & A \cdot x = y \end{aligned}$$

*Equivalent linear program*

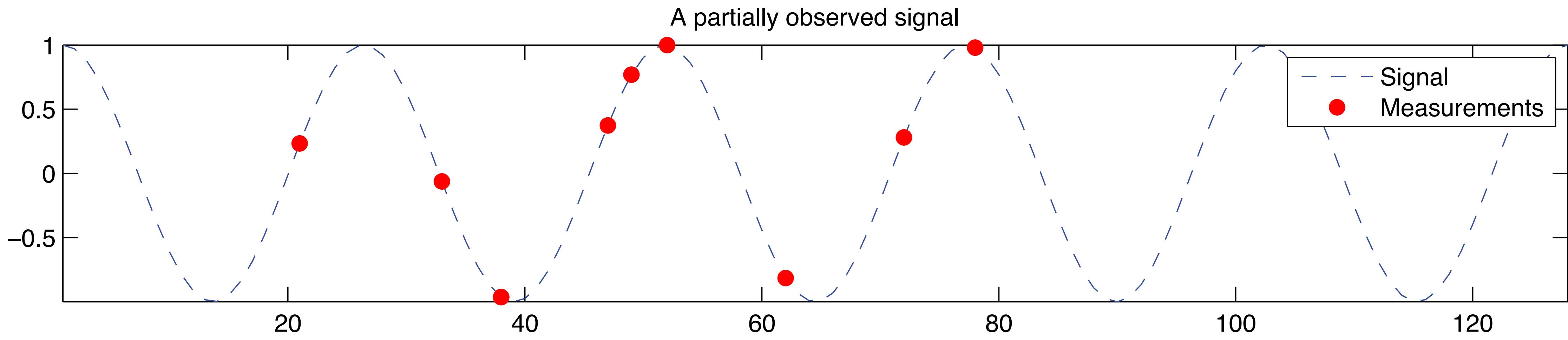
$$\begin{aligned} \text{minimize} \quad & \|z\|_1 = 1^\top \cdot z \\ \text{subject to} \quad & [A, -A] \cdot z = y \\ \text{and} \quad & z_i \geq 0 \end{aligned}$$



*Note that this was an  $\leq$ , but it's easy to change to an  $=$  by adding extra complementary constraints*

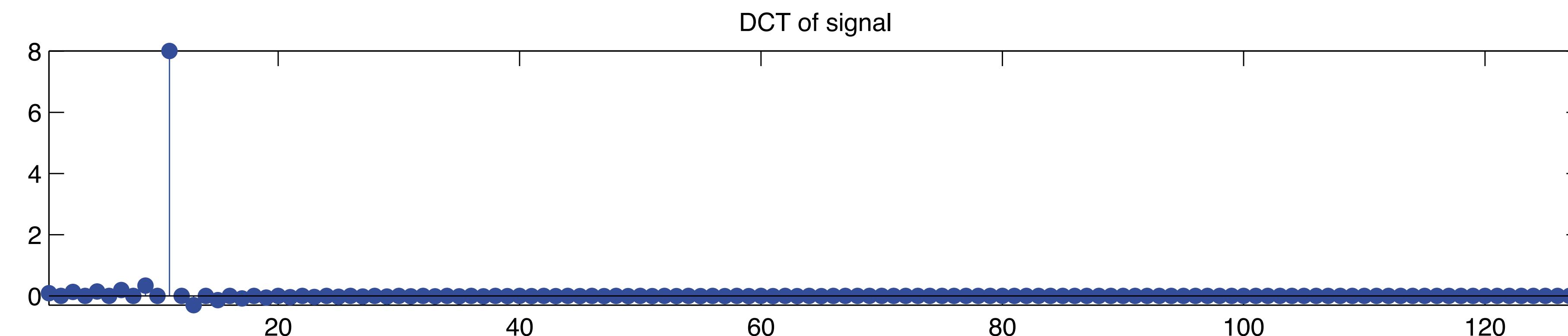
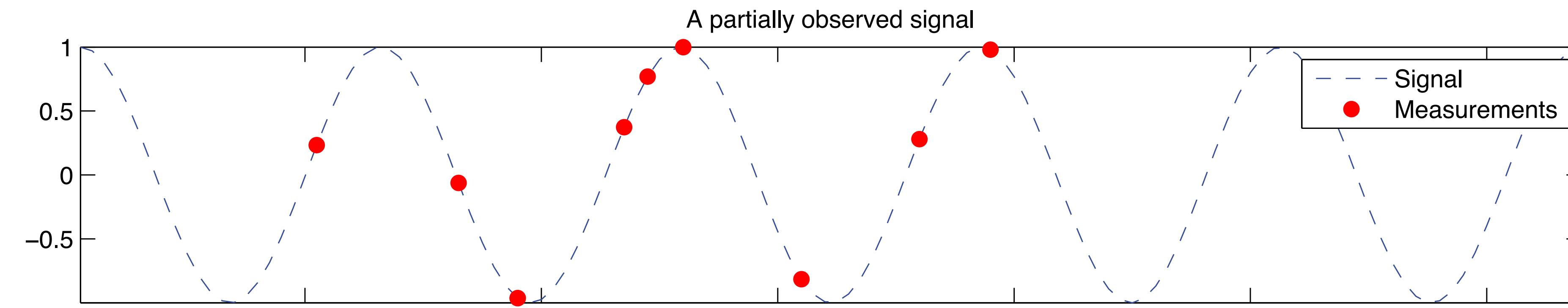
# A simple example

- Suppose you measure an arbitrary sinusoid
  - Can you recover the original signal using only the small number of measurements that we made?



# Key observation

- The original signal is sparse in the frequency domain
  - How about we use that to construct a problem?



# The problem to solve

- Find a set of small coefficients  $\mathbf{x}$  in the DCT domain, and make sure that they explain all our data  $\mathbf{y}$ , i.e.:

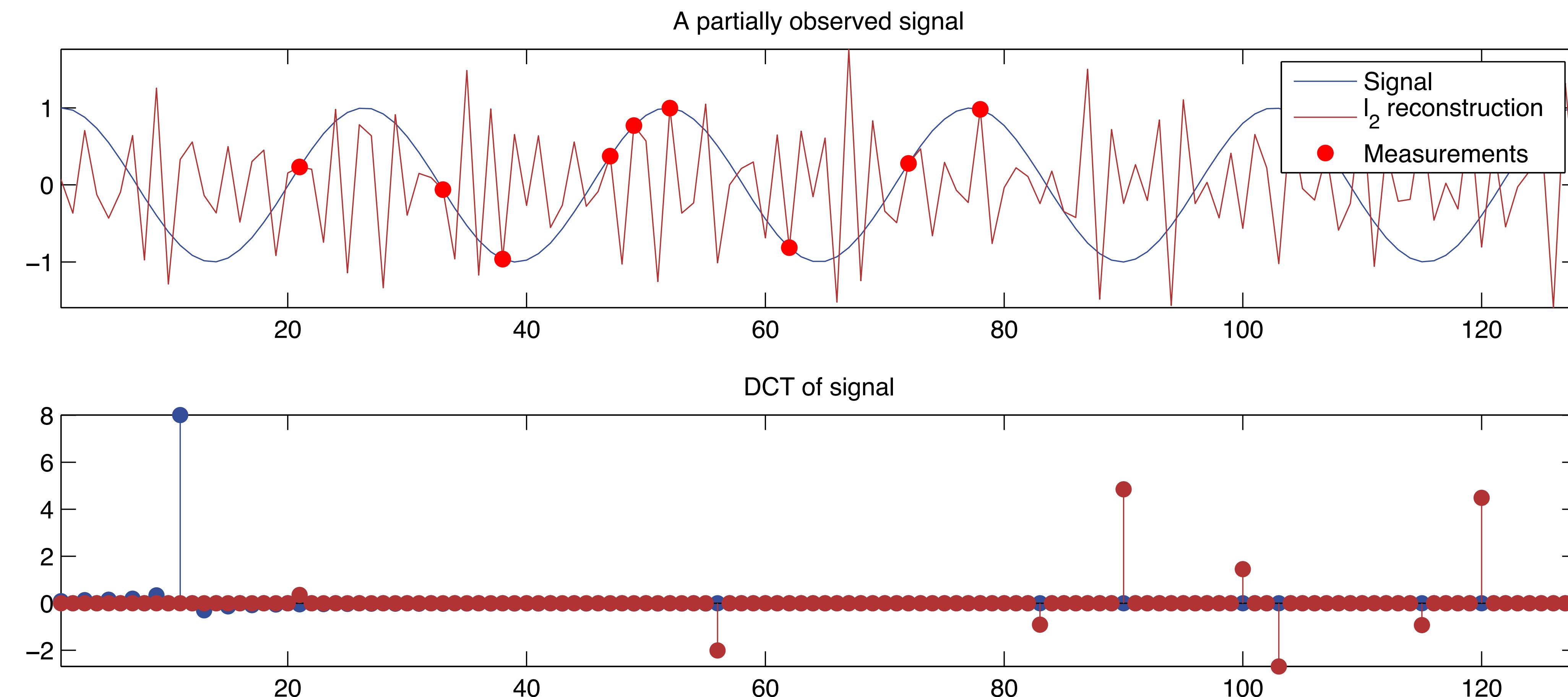
$$\mathbf{A} \cdot \mathbf{C}^{-1} \cdot \mathbf{x} = \mathbf{y}$$

- Where matrix  $\mathbf{A}$  selects only the indices that we observe
- Simple least squares problem using minimum  $\ell_2 \|\mathbf{x}\|$

$$\mathbf{x} = (\mathbf{A} \cdot \mathbf{C}^{-1})^+ \cdot \mathbf{y}$$

# And it doesn't really do much

- But you expected that!

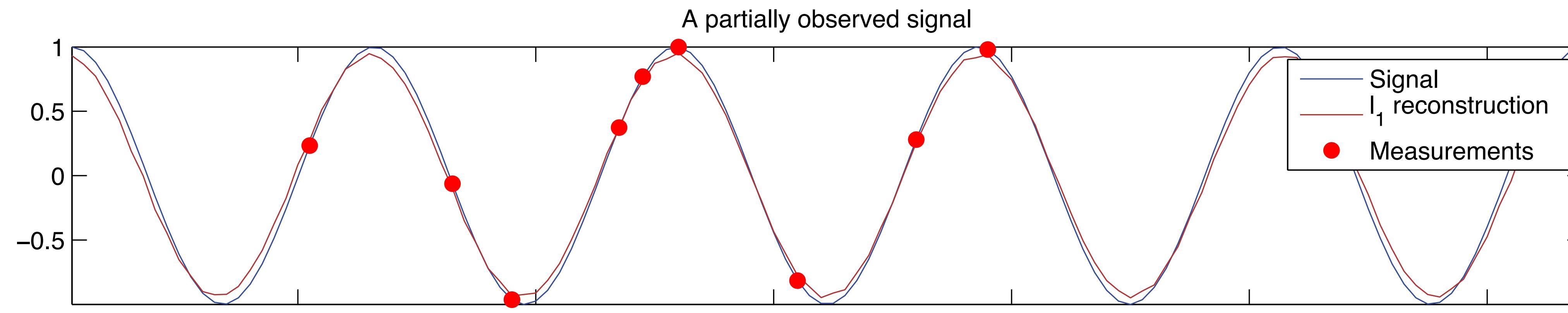


# What happened?

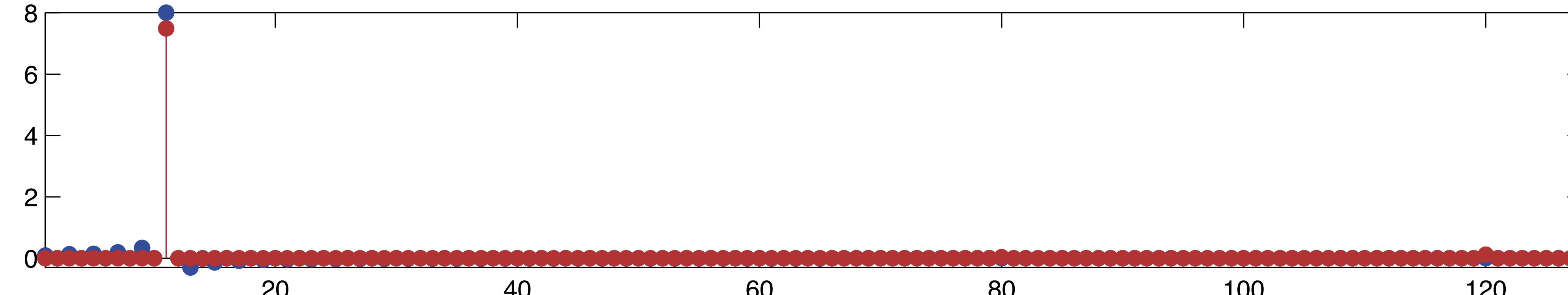
- Minimizing the  $\ell_2$  resulted in adding more frequencies in the signal
  - $\ell_2$  doesn't give sparsity
  - Small coefficients  $\neq \min \ell_2$
- We instead should find a minimal  $\ell_1$  solution
  - Because it actually enforces sparsity

# And the result

- A much better reconstruction



DCT of signal



# A realization

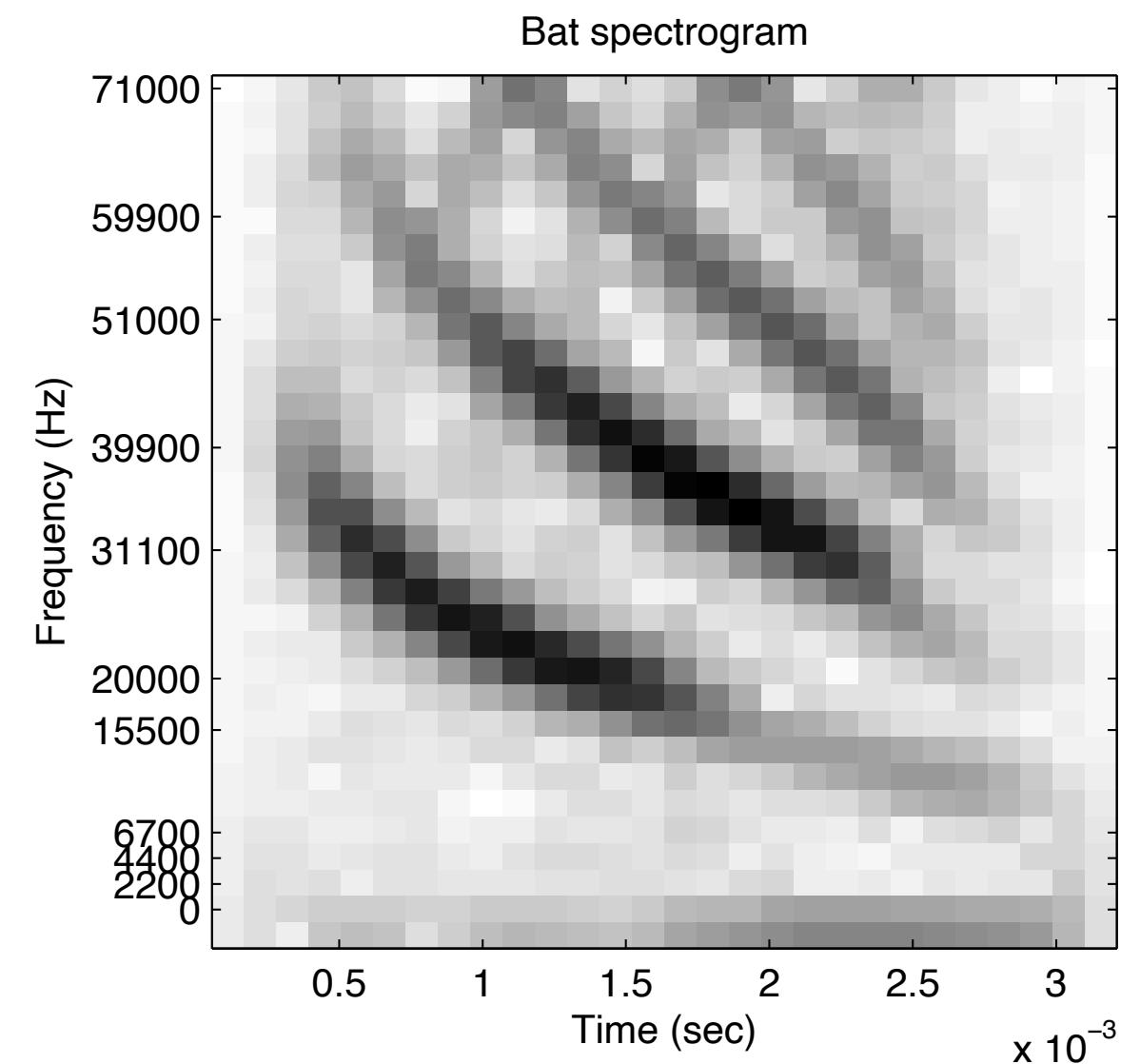
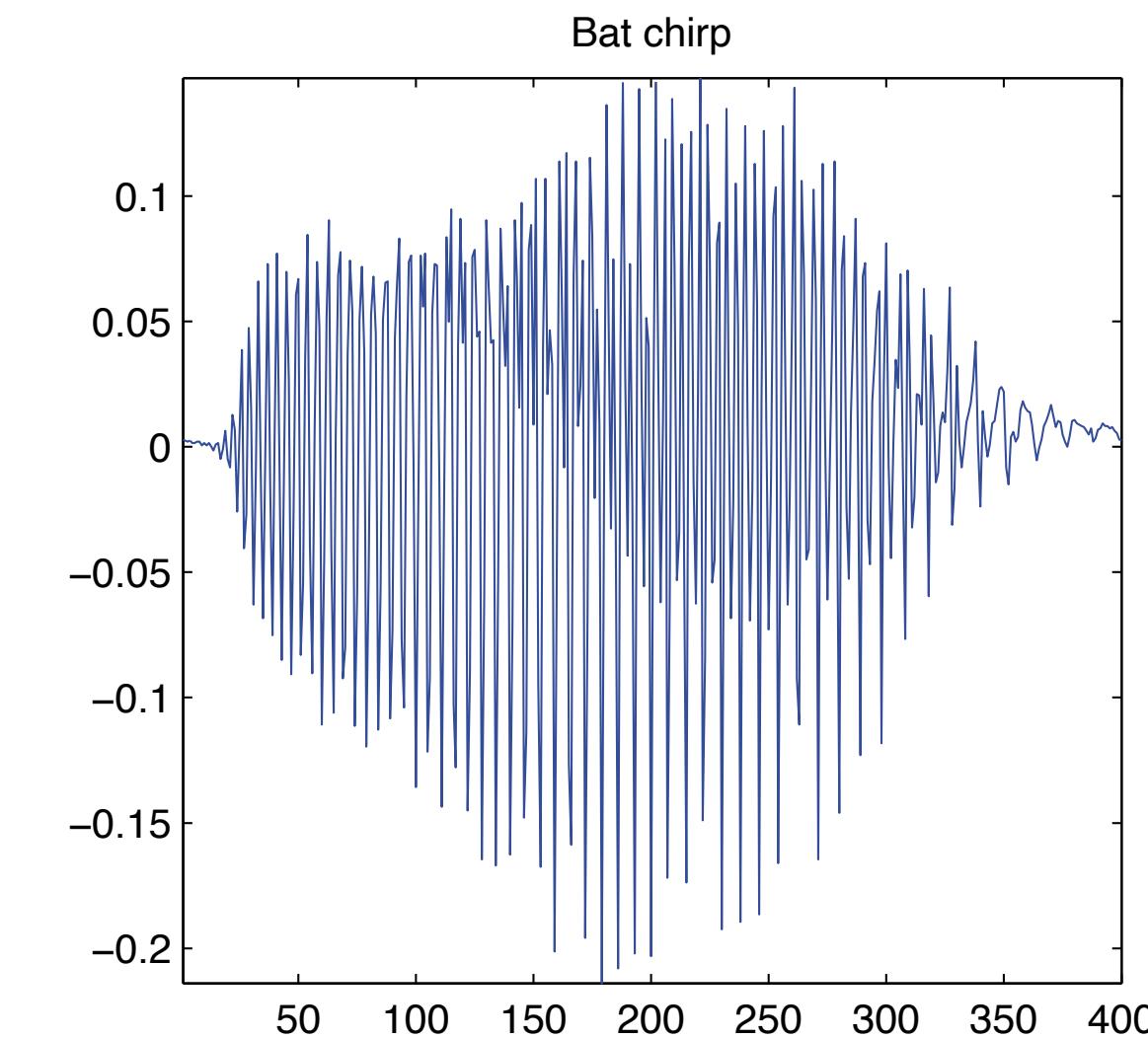
- According to the rules of sampling  
this is impossible!
- What is the magic taking place here?
  - Why is sparsity special?

# The deal with sparsity

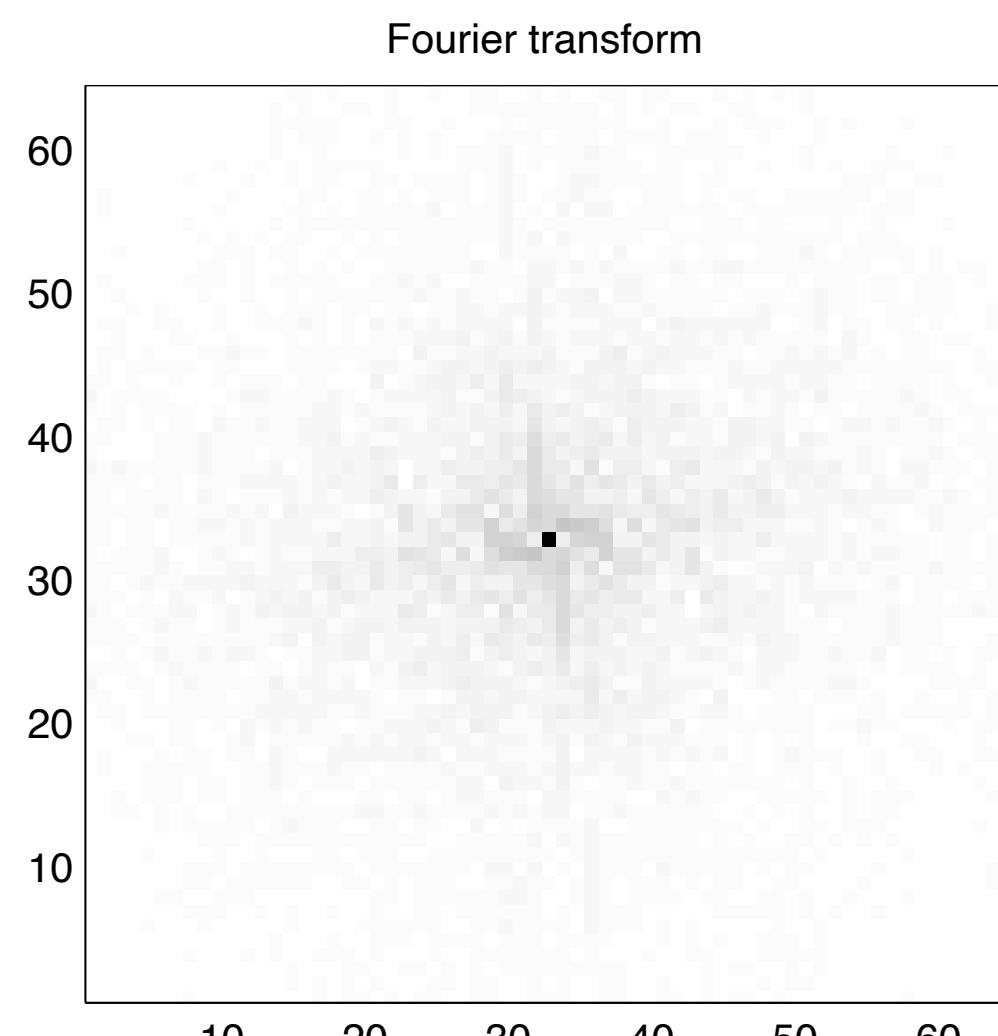
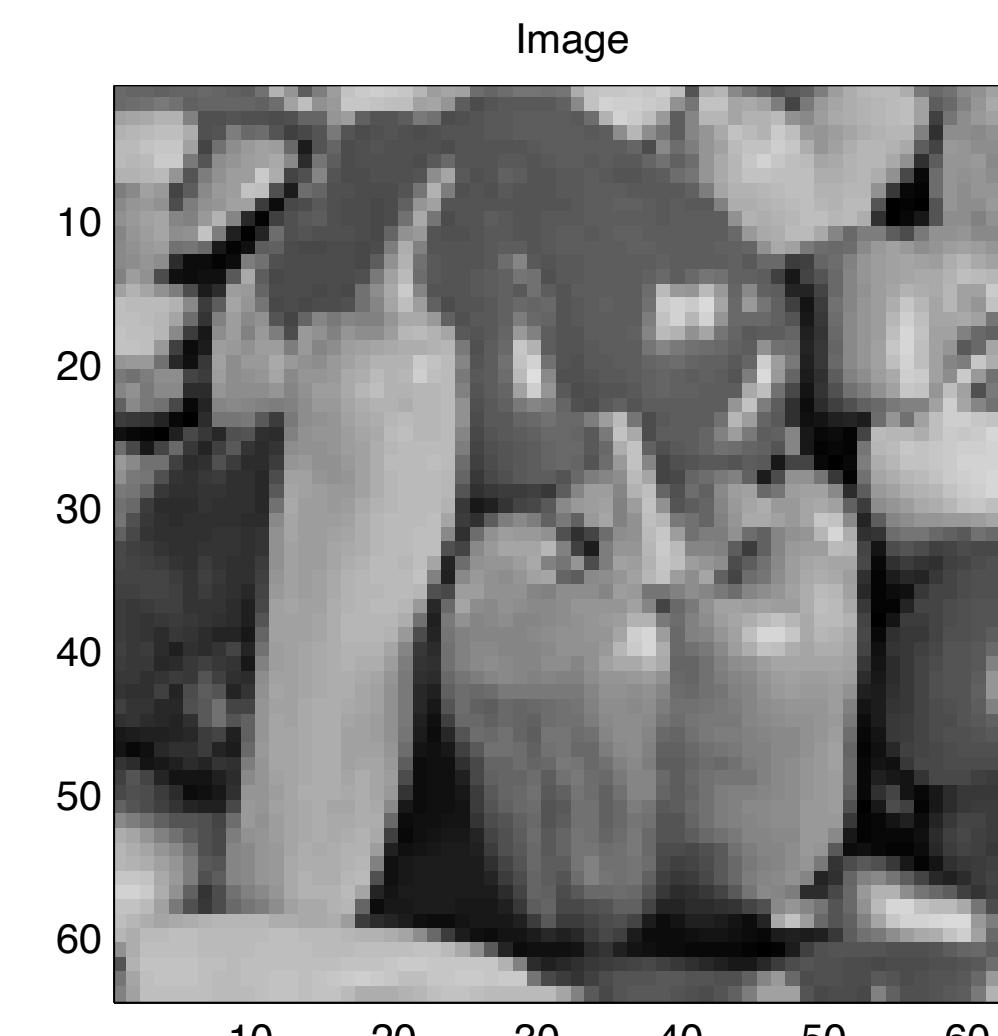
- Sparsity implies structure, and structure is everywhere
- Signals often exhibit sparsity after undergoing the right transformation

# Sparsity in signals

- Most signals are sparse (but in certain domains)
  - e.g. sound spectra
  - Image wavelets
  - ...



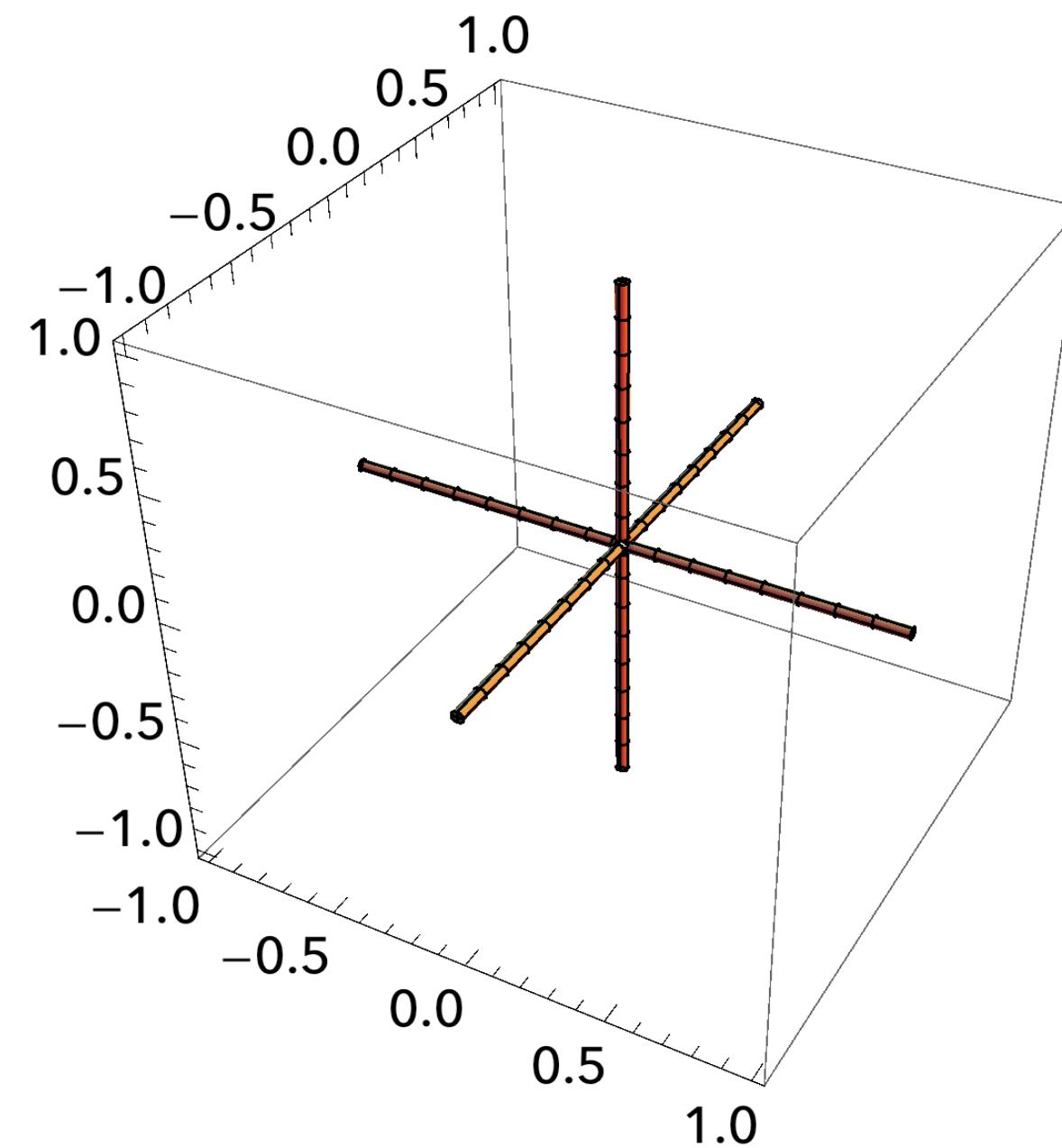
- This is an opportunity to measure signals better
  - or to describe them easier



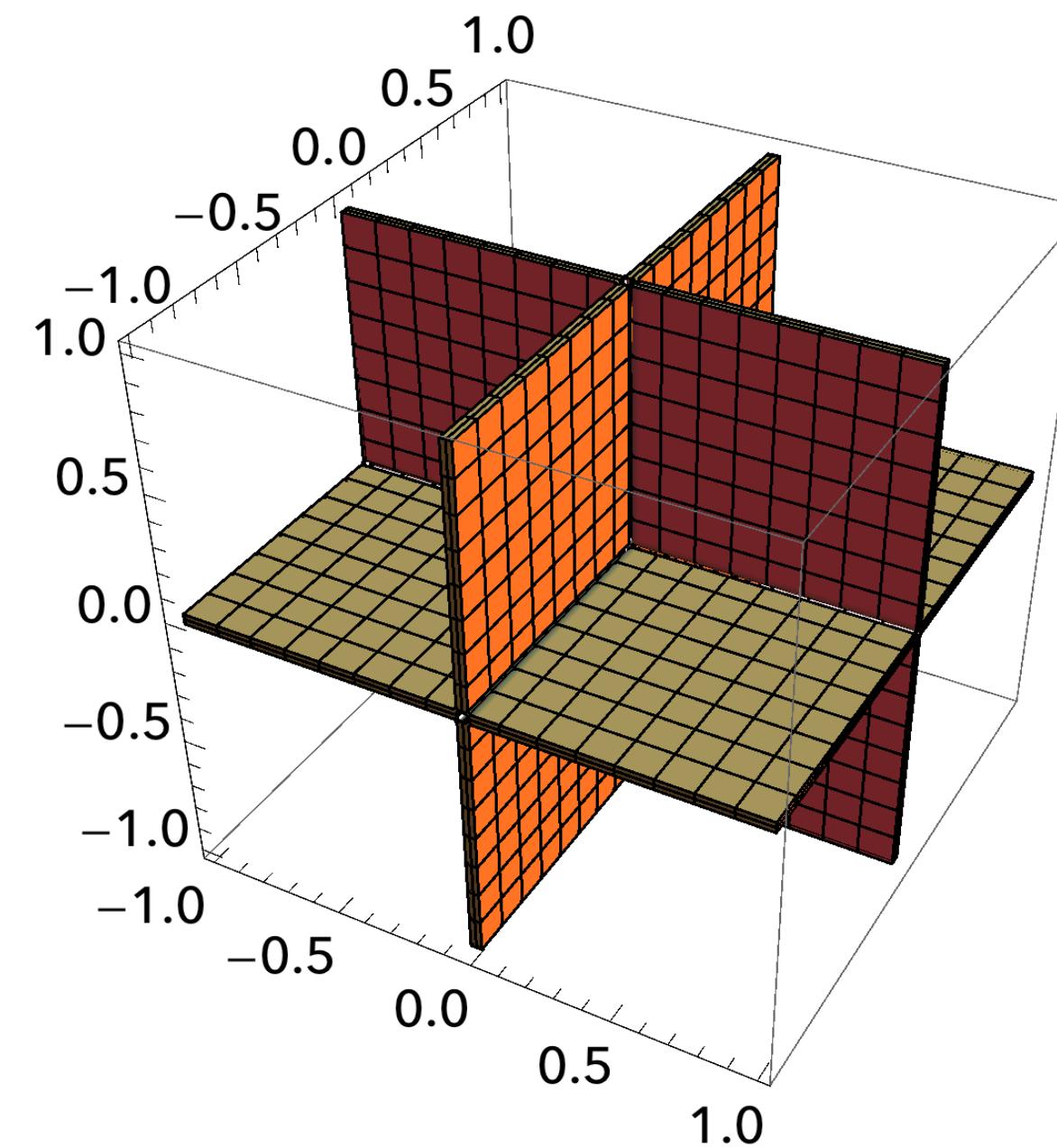
# Vector spaces of sparse signals

- Signals can be sparse in various ways
  - And some are “compressible”

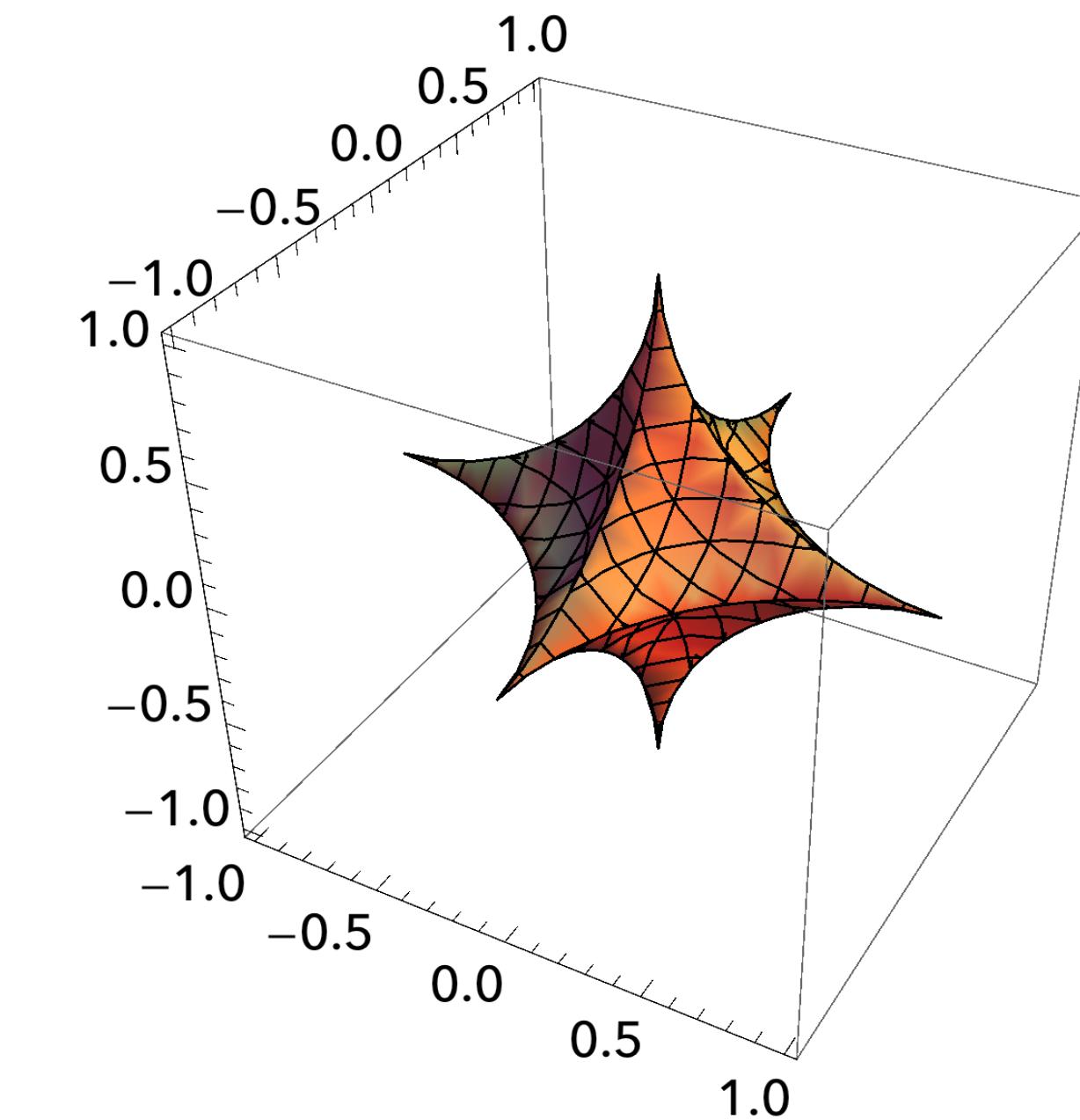
*1-sparse signal space*



*2-sparse signal space*



*Compressible signal space*



# Sparse approximations

- Represent signals using a sparse representation:

$$\mathbf{f} = \sum_k a_k \mathbf{b}_k$$

*Coefficients*      *Bases / Dictionary*

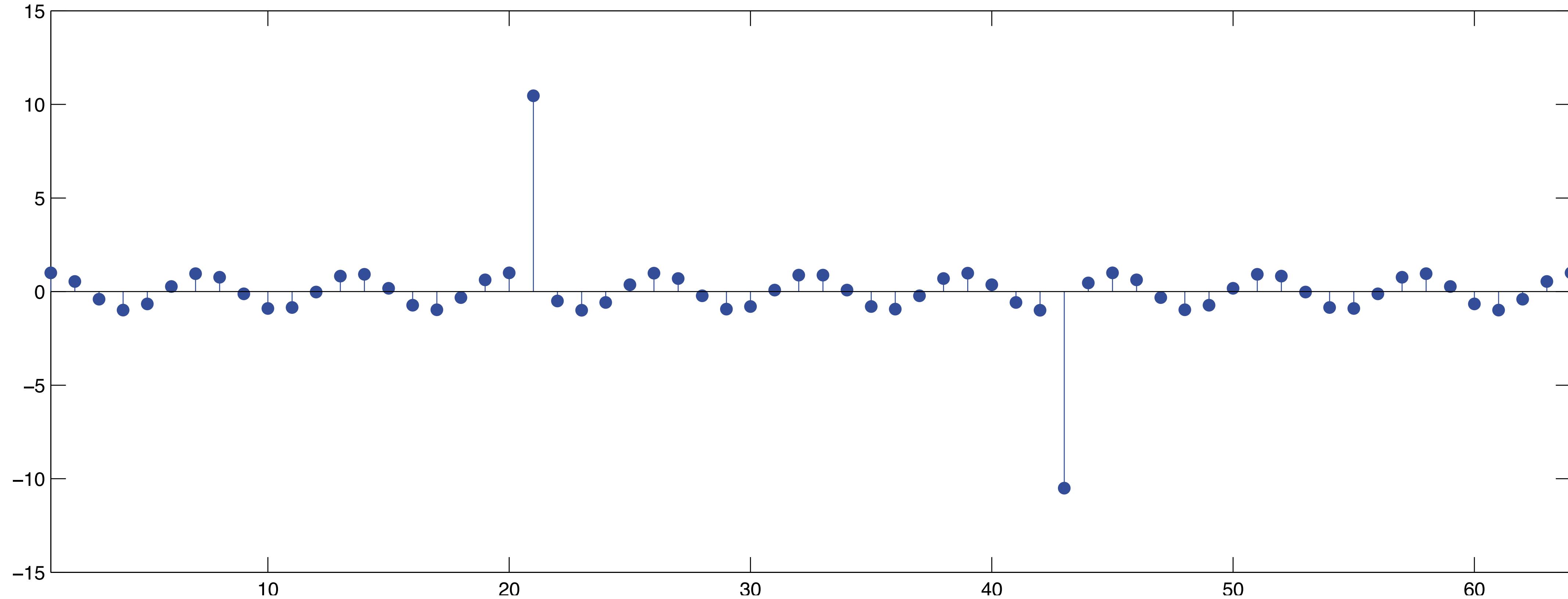
- Two goals:
  - Analysis: Study  $\mathbf{f}$  through structure of  $a$  and  $\mathbf{b}$  (seen that already)
  - Approximation: Reconstruct  $\mathbf{f}$  with a minimal number of terms

# Exposing sparsity via dictionaries

- Can we use dictionaries that produce sparse coefficients?
- Why would they be useful and how would we implement them?

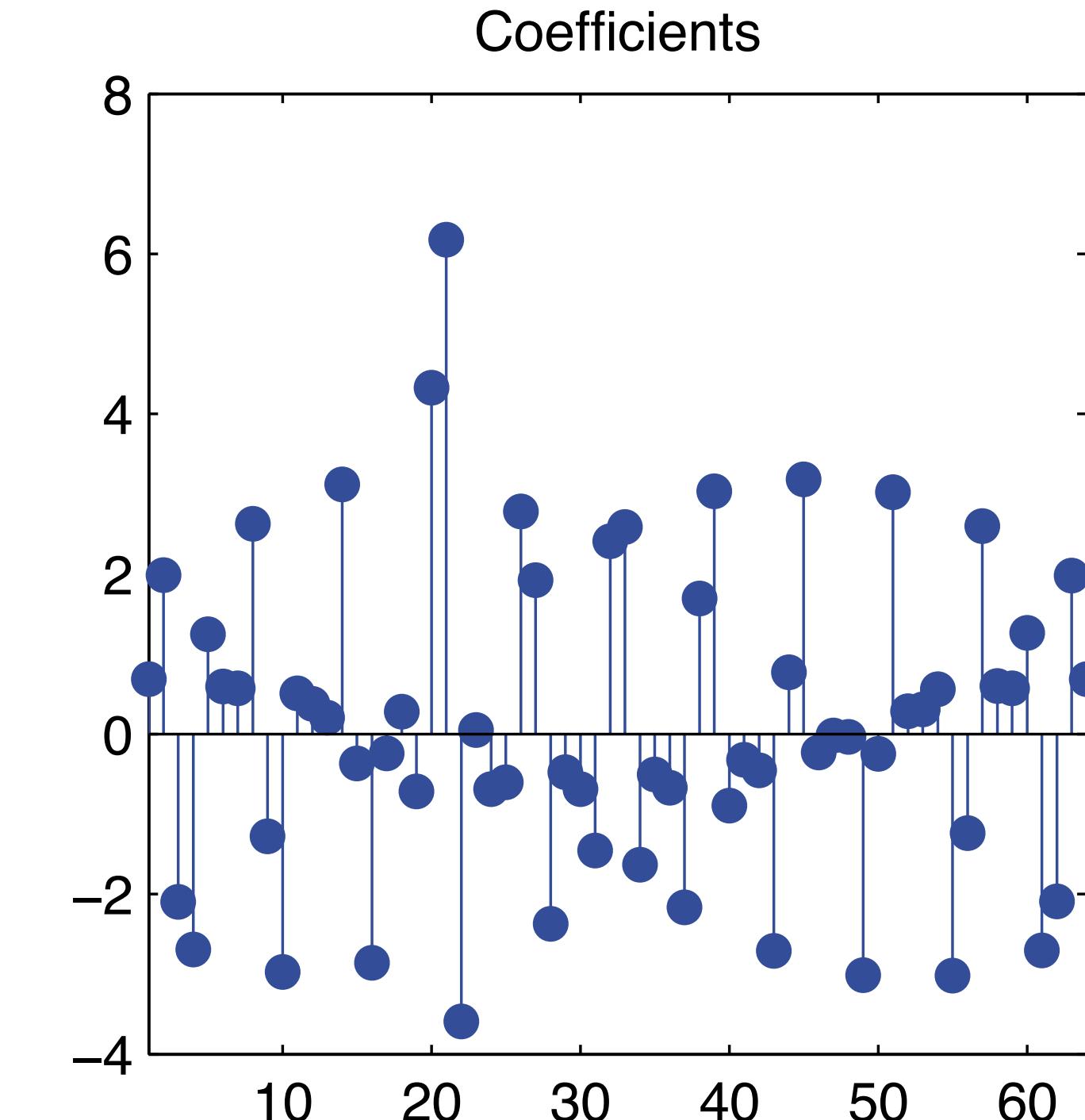
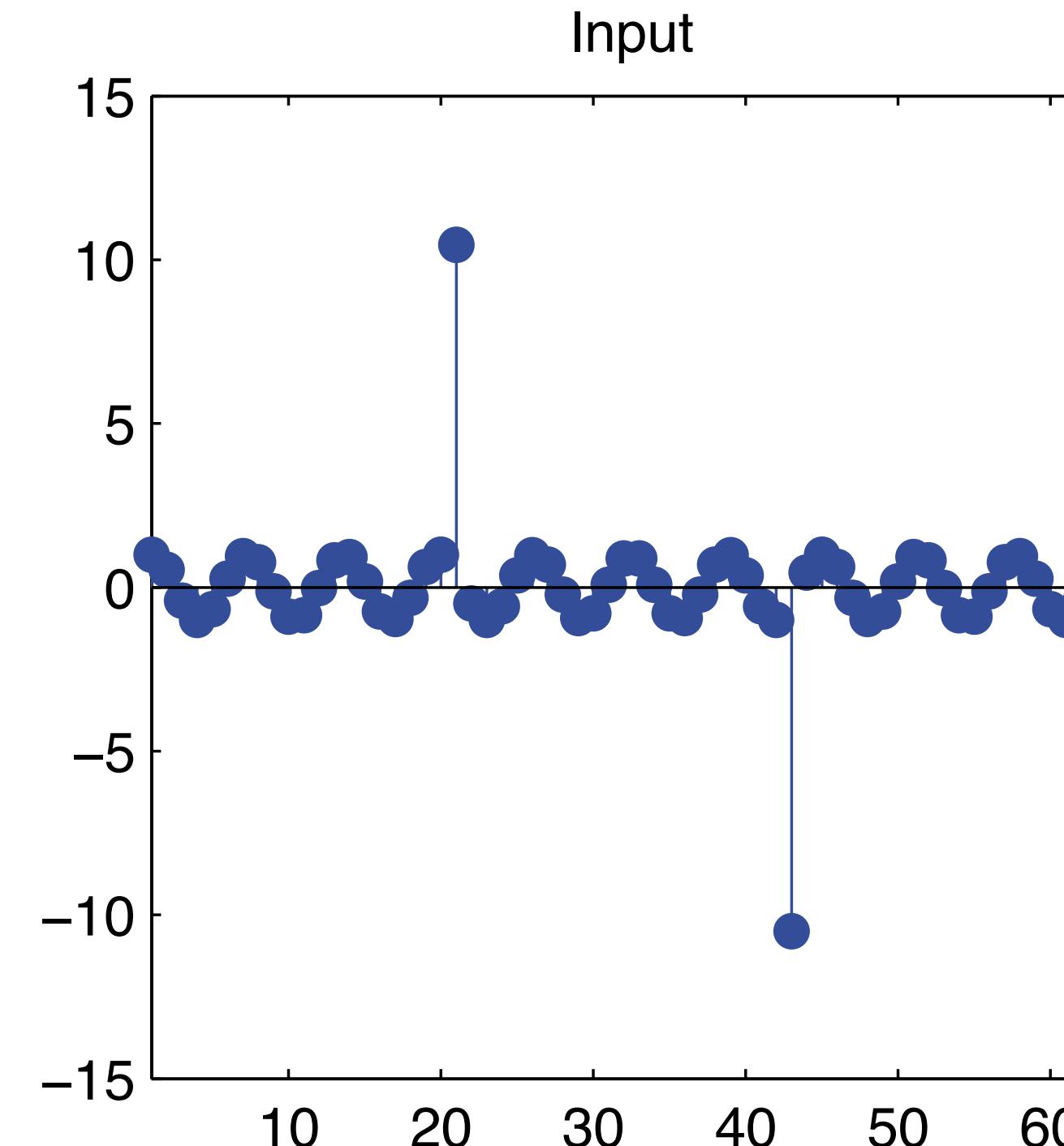
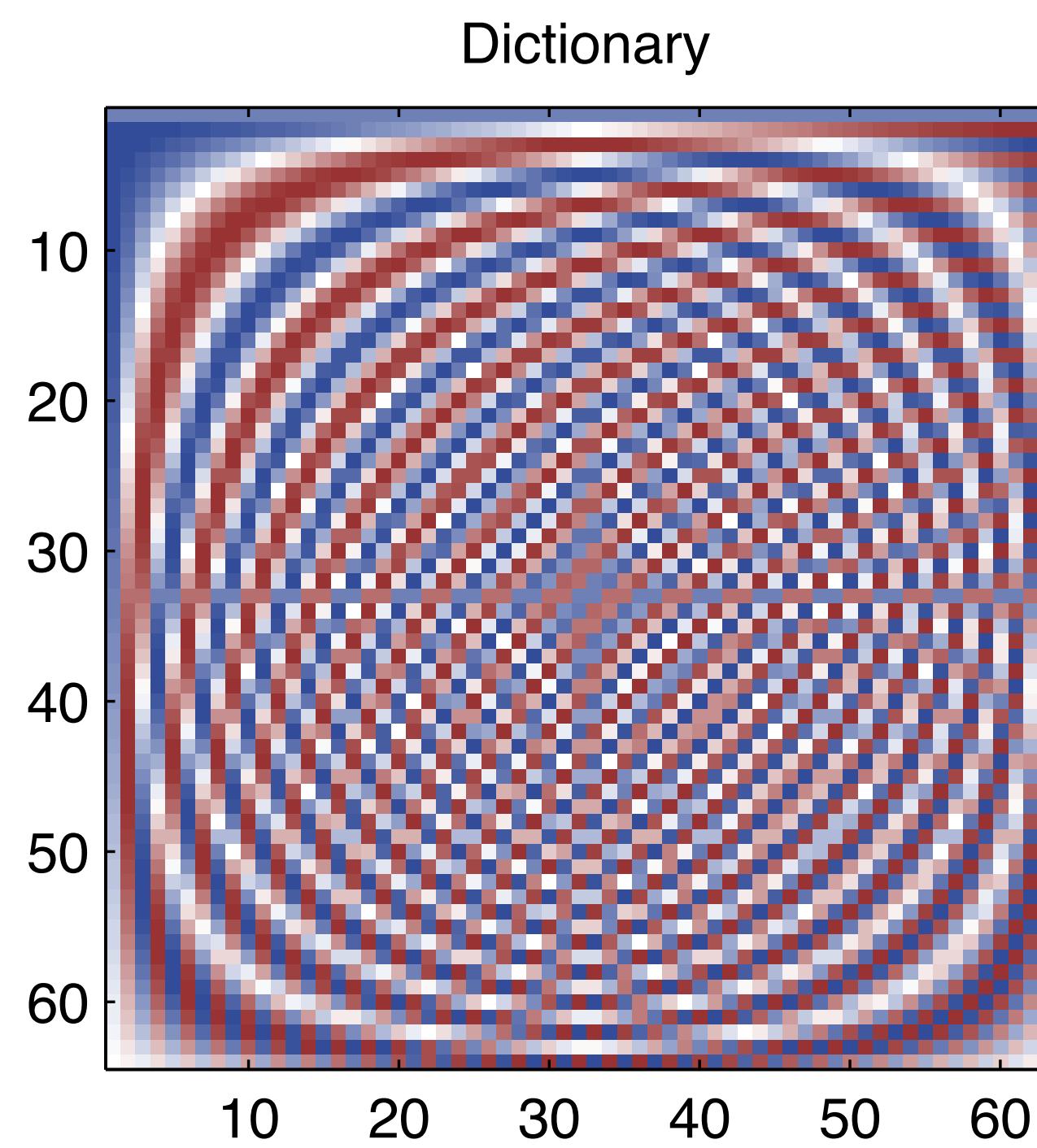
# A simple example

- How would you describe this signal?
  - A sinusoid with two spikes



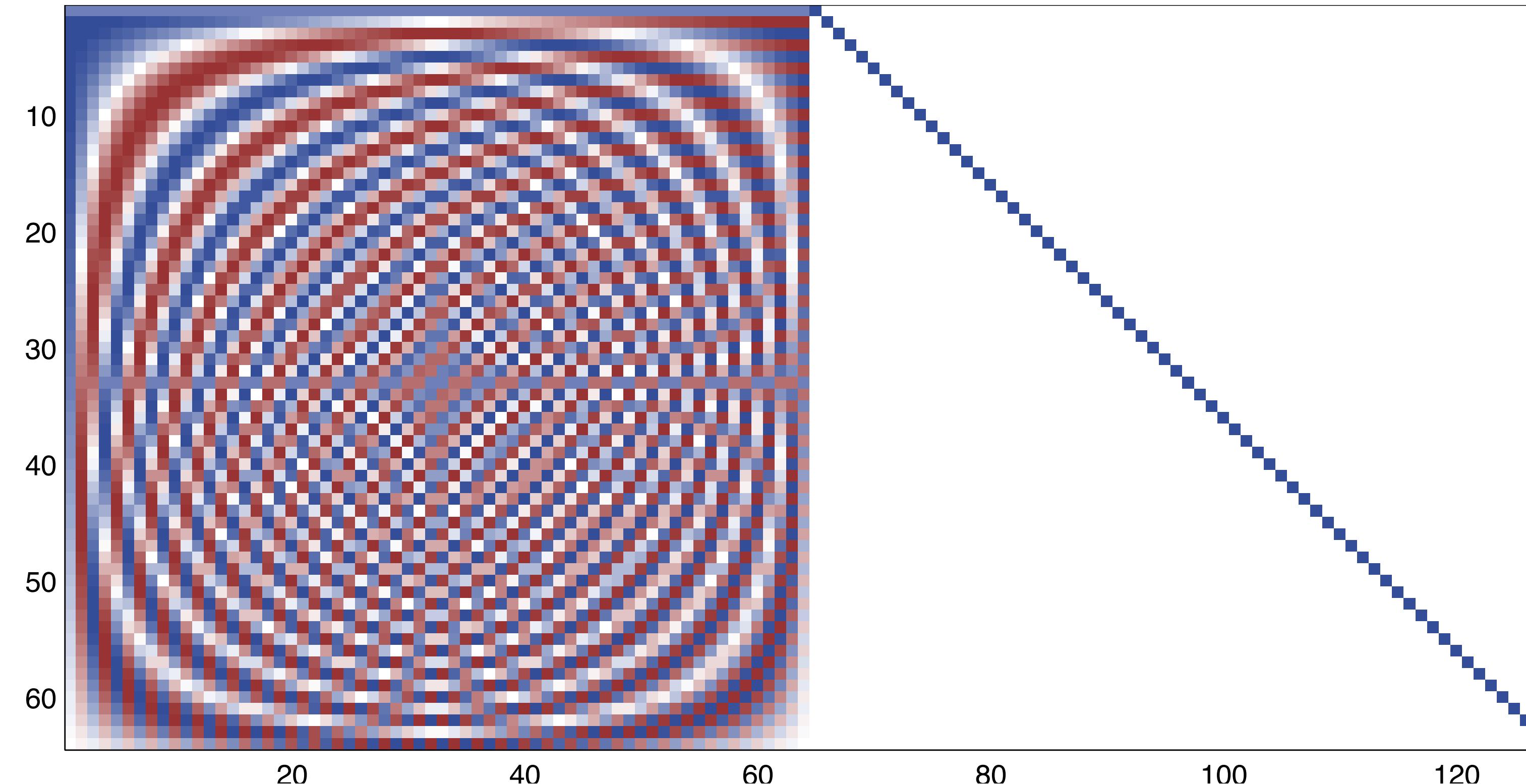
# Using a generic dictionary

- Analyzed via the DCT
  - Resulting coefficients are not sparse
    - Multiple sines are used to approximate the spikes



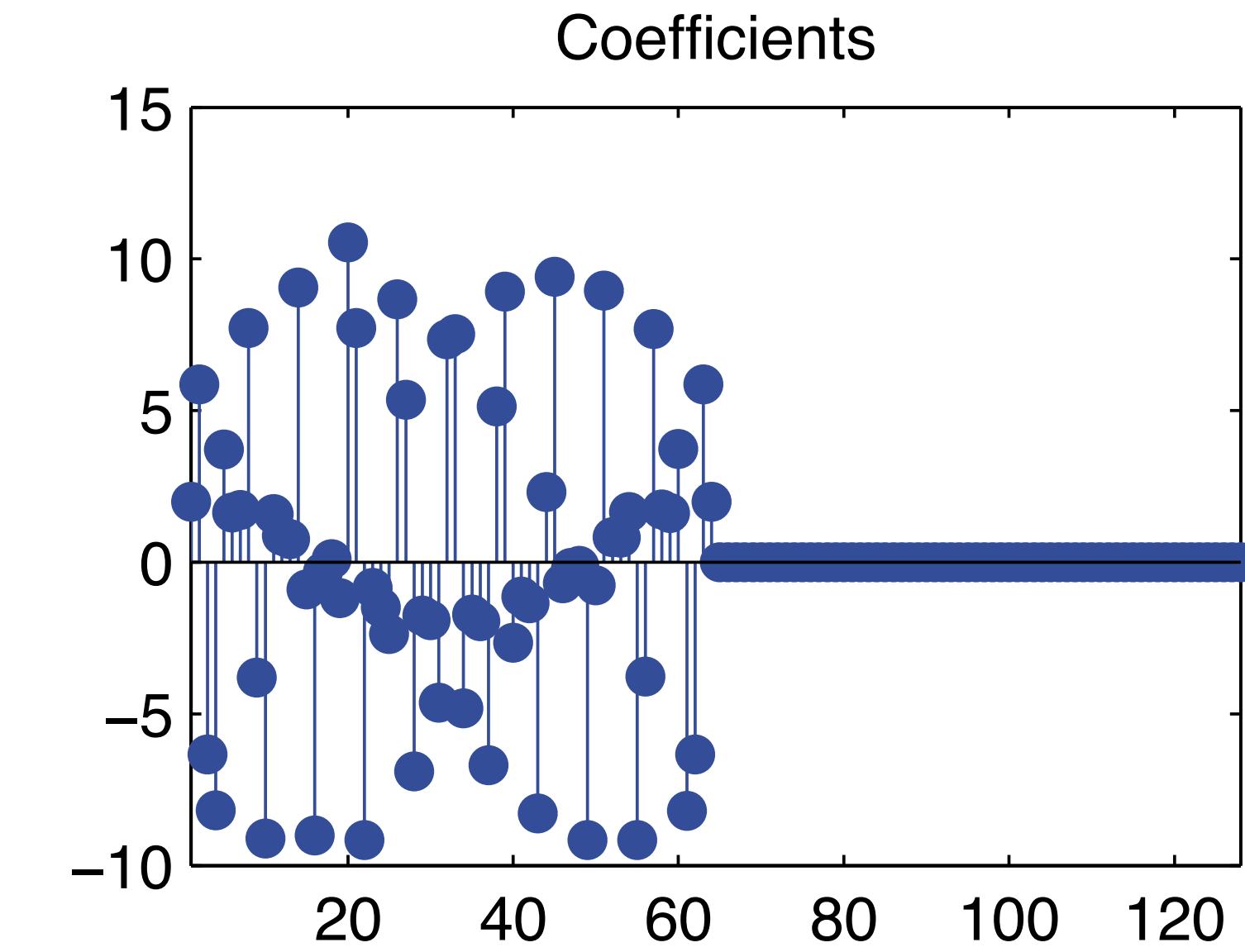
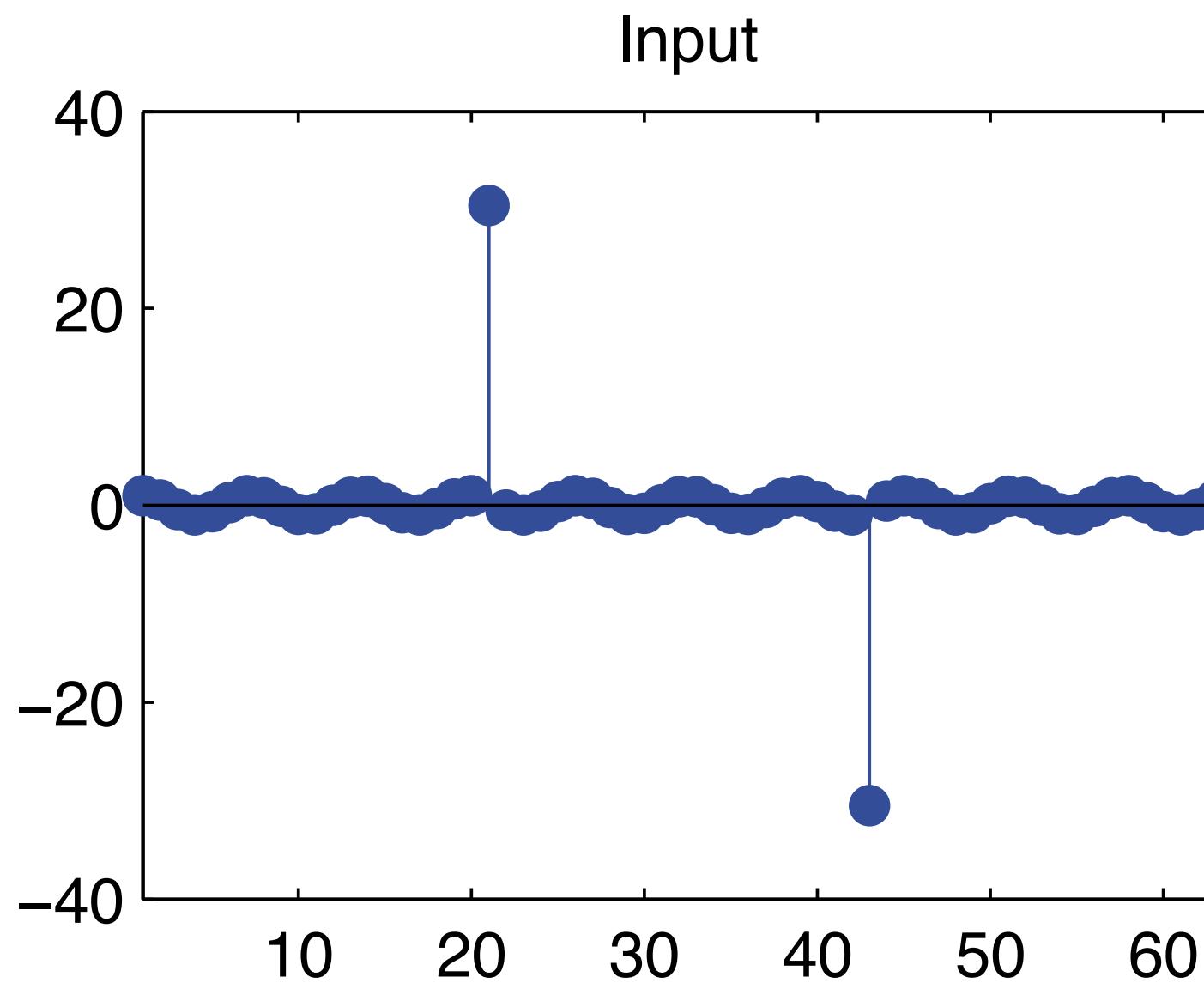
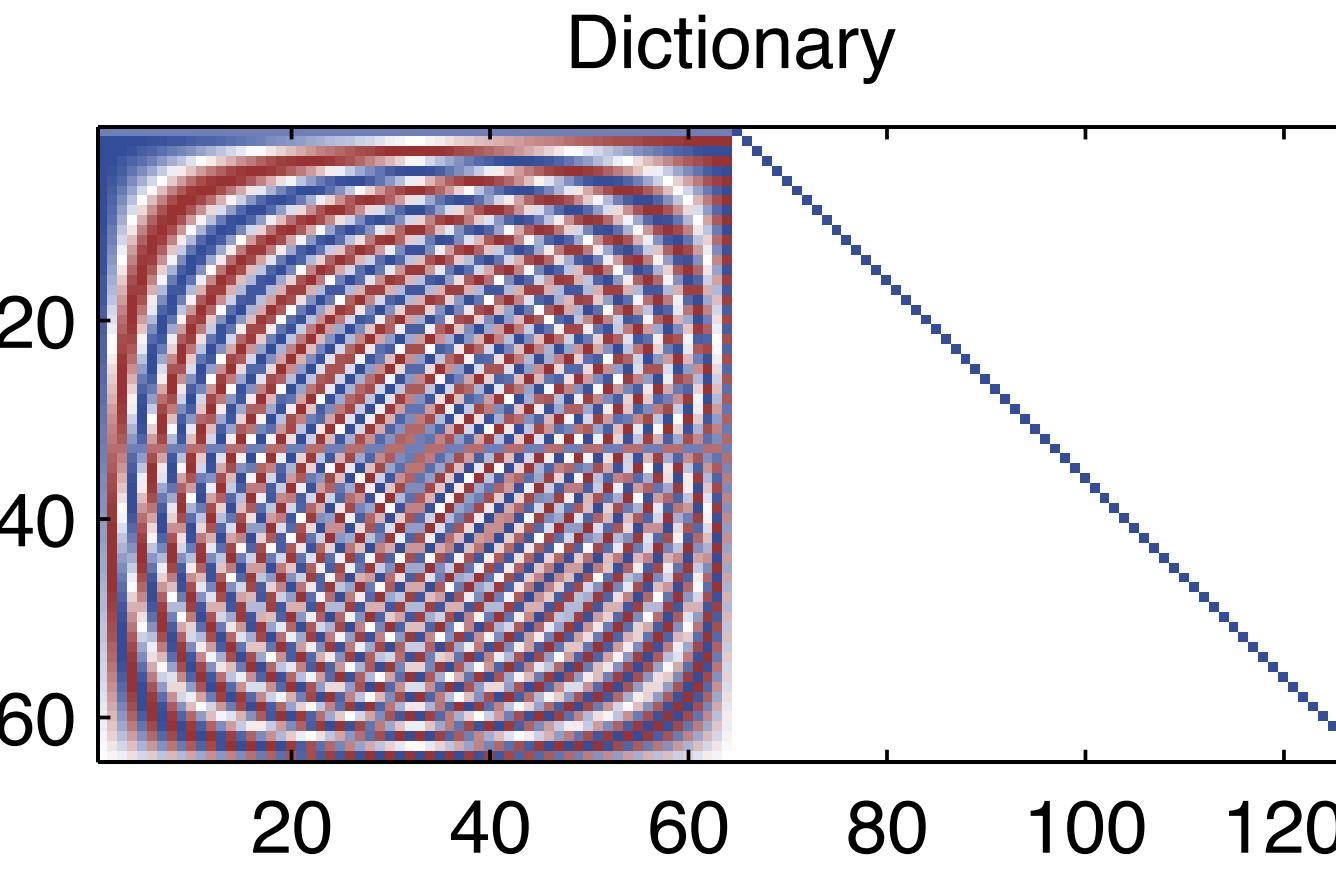
# A “better” dictionary

- Use both sinusoids and spikes!
  - Now we won’t use as many sines to represent the spikes



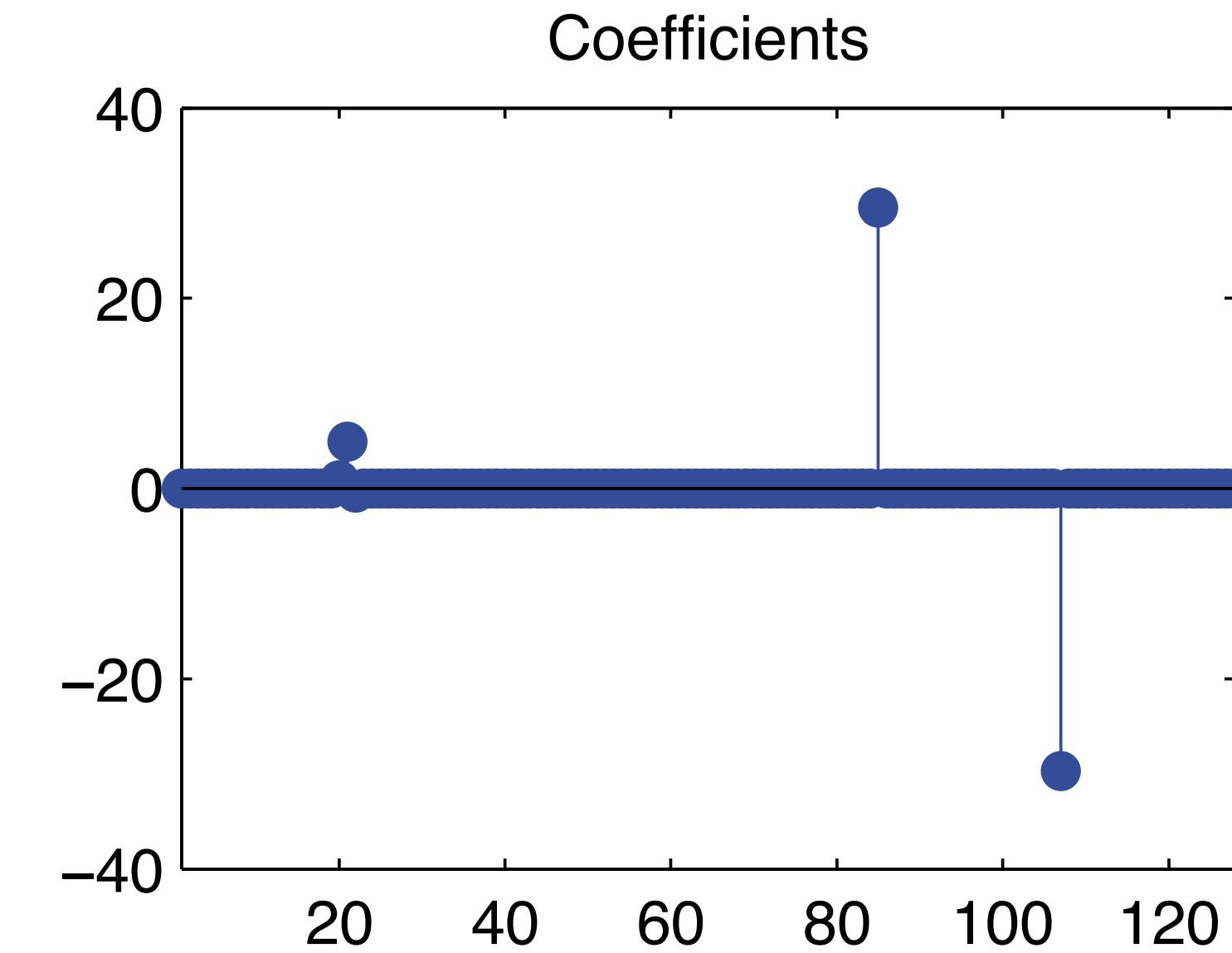
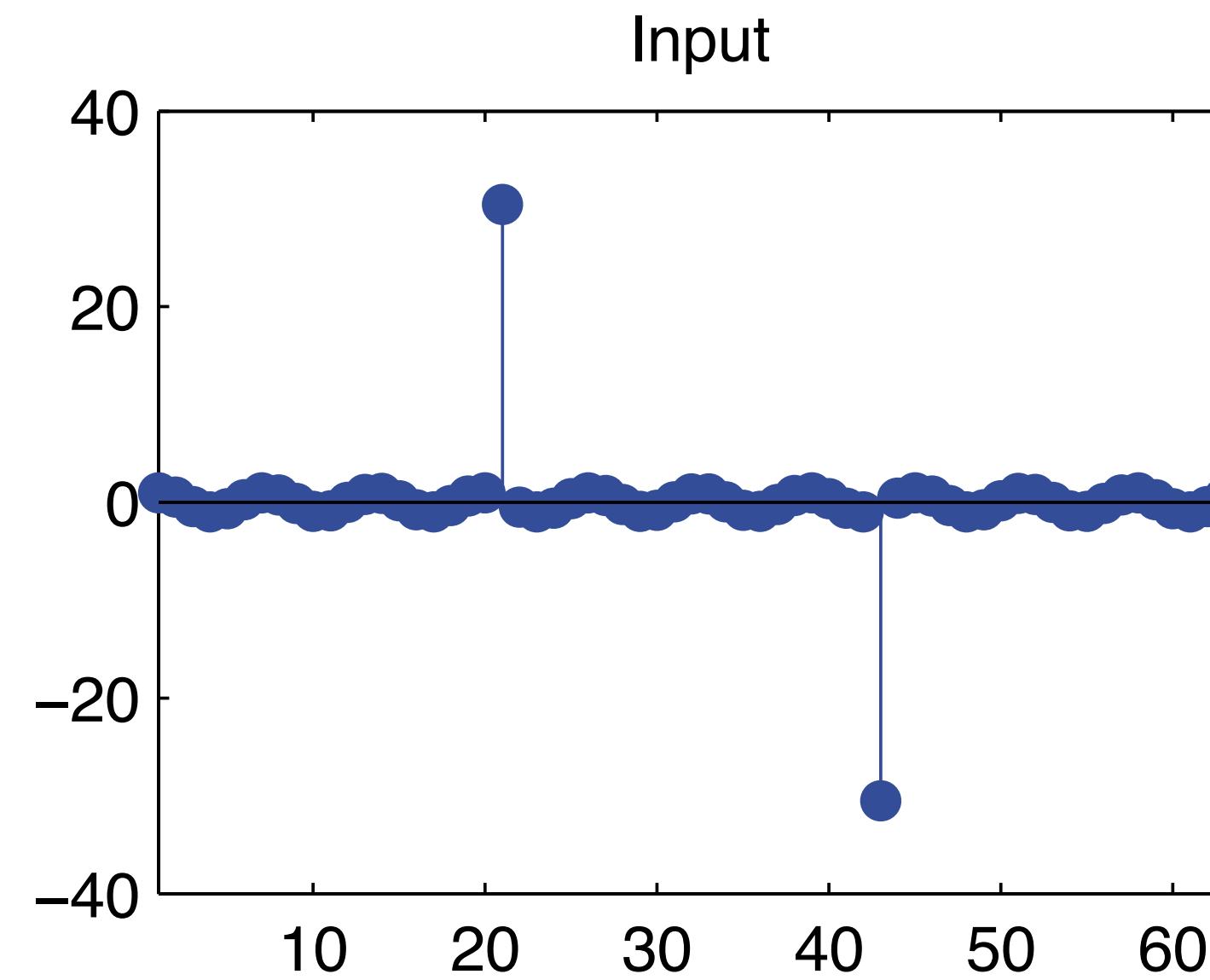
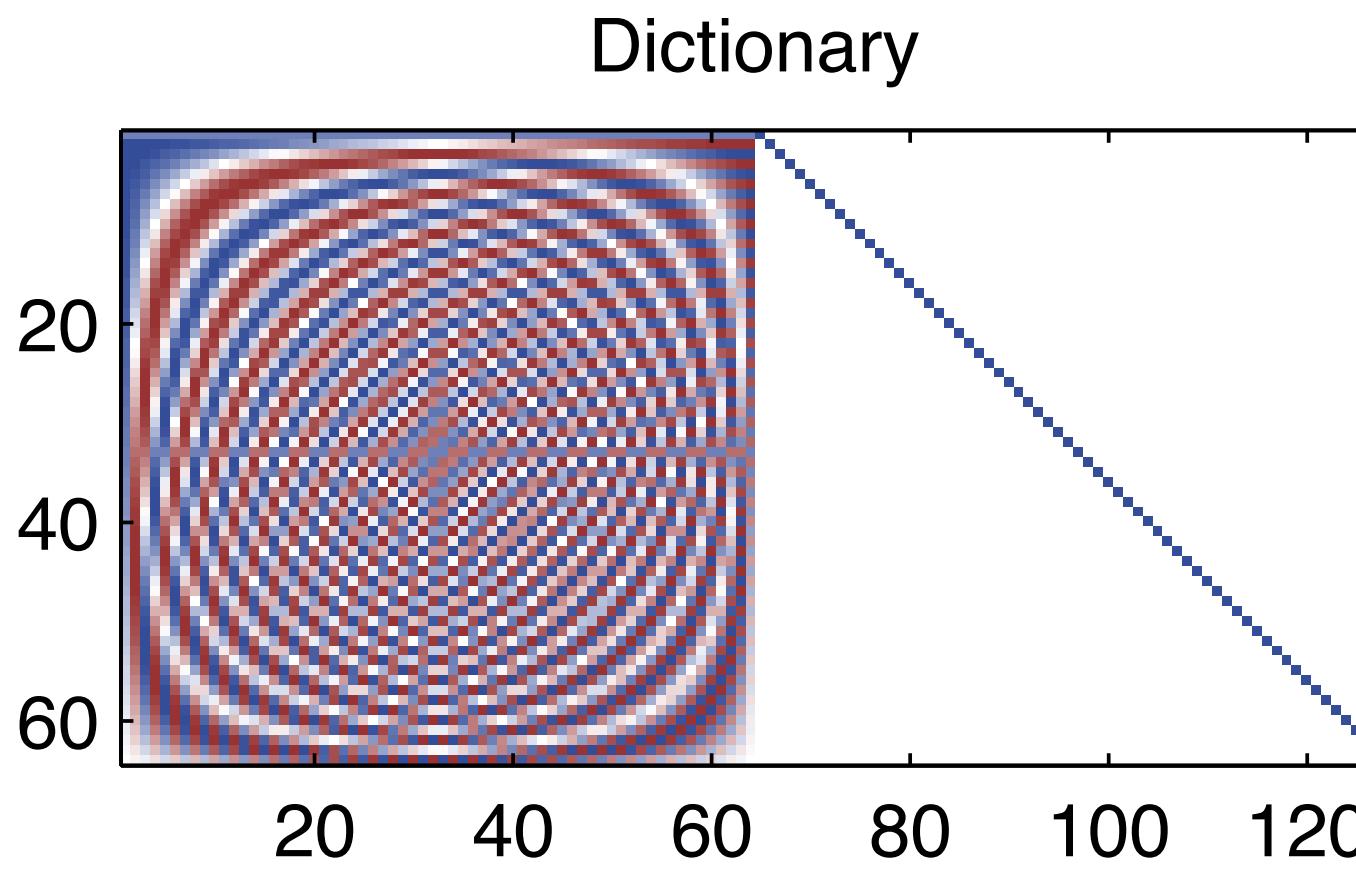
# Applying the dictionary

- Using the straightforward decomposition won't help
  - Spike elements are not utilized
  - Minimal  $\ell_2$  cost penalizes the bases describing the loud spikes



# Doing it the right way

- This time we ask for minimum  $\ell_1$  coefficients
  - And we get a perfect description of the input!



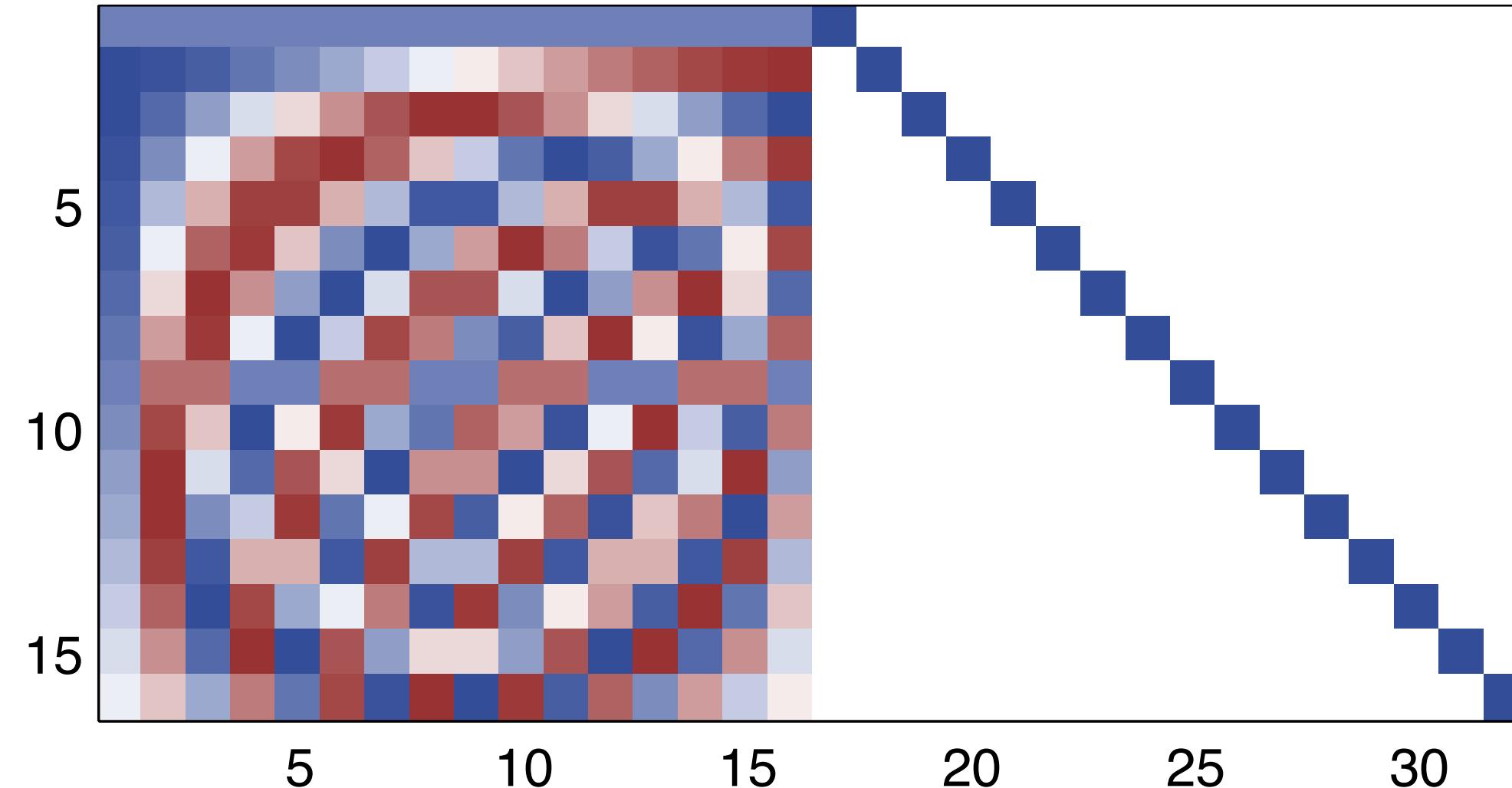
# Overcomplete dictionaries

- Use dictionaries that contain “everything”!
  - Use compact descriptions of elements
- Some problems
  - Large size/computations
  - Lack of fast algorithms (e.g. FFT)
  - Problems with *coherence*
    - This is a big one

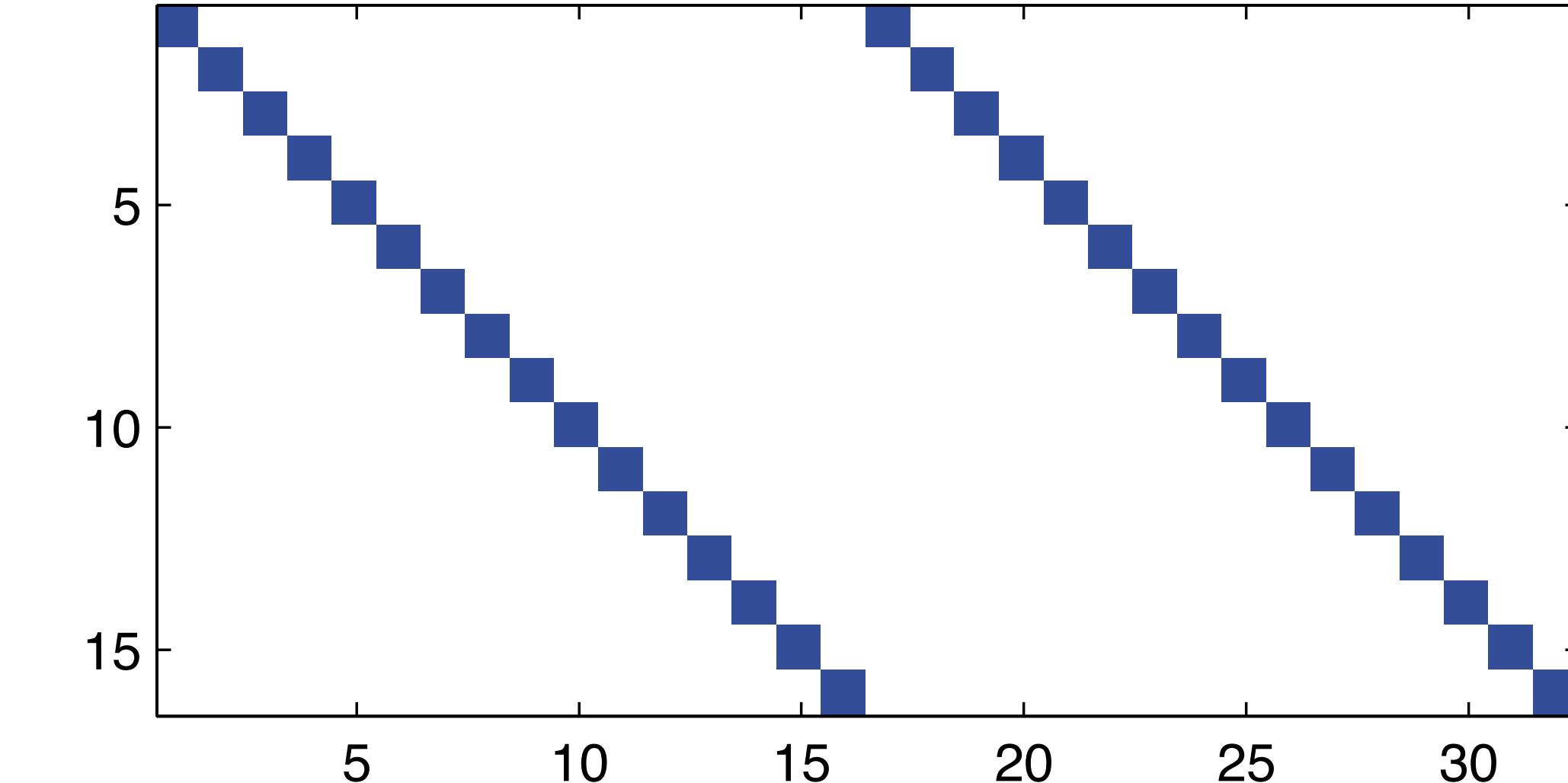
# Dictionary coherence

- Make sure that the dictionary elements don't result in ambiguous coefficients
  - A simple measure of coherence:  $\mu = \max_{i,j} |\mathbf{d}_i^\top \cdot \mathbf{d}_j|$

Ok

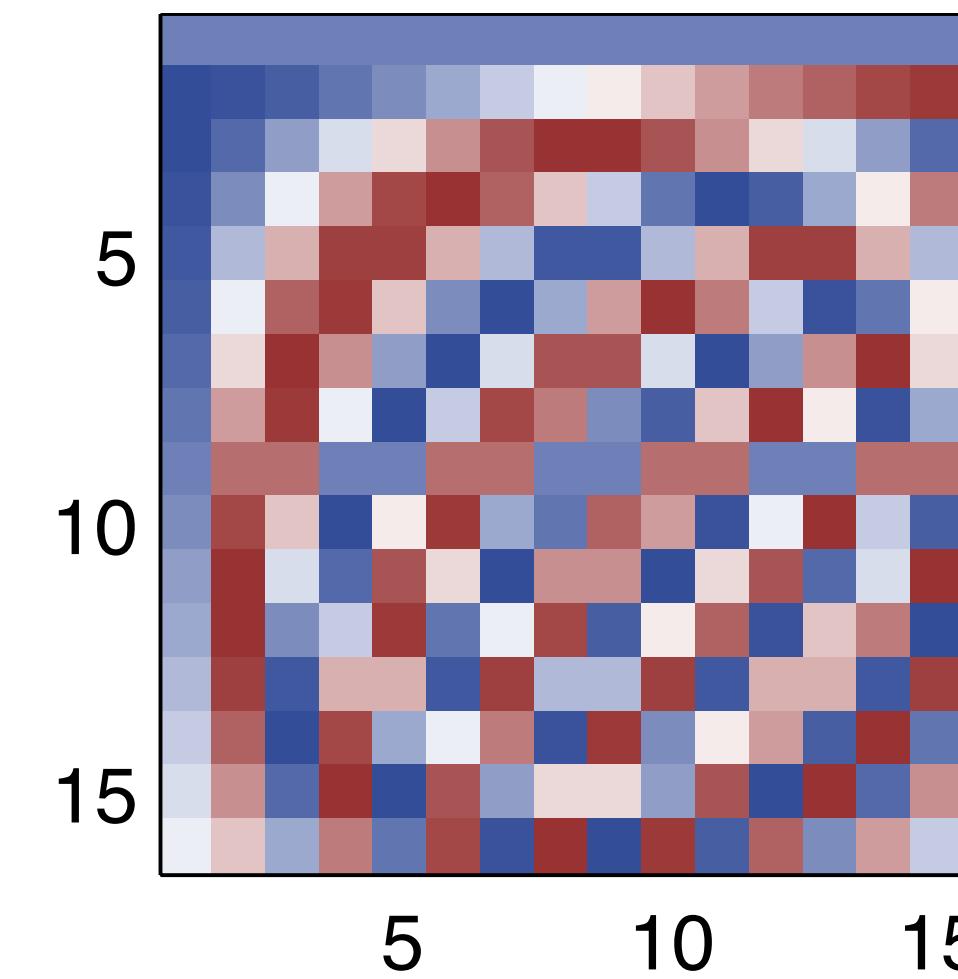


Bad!

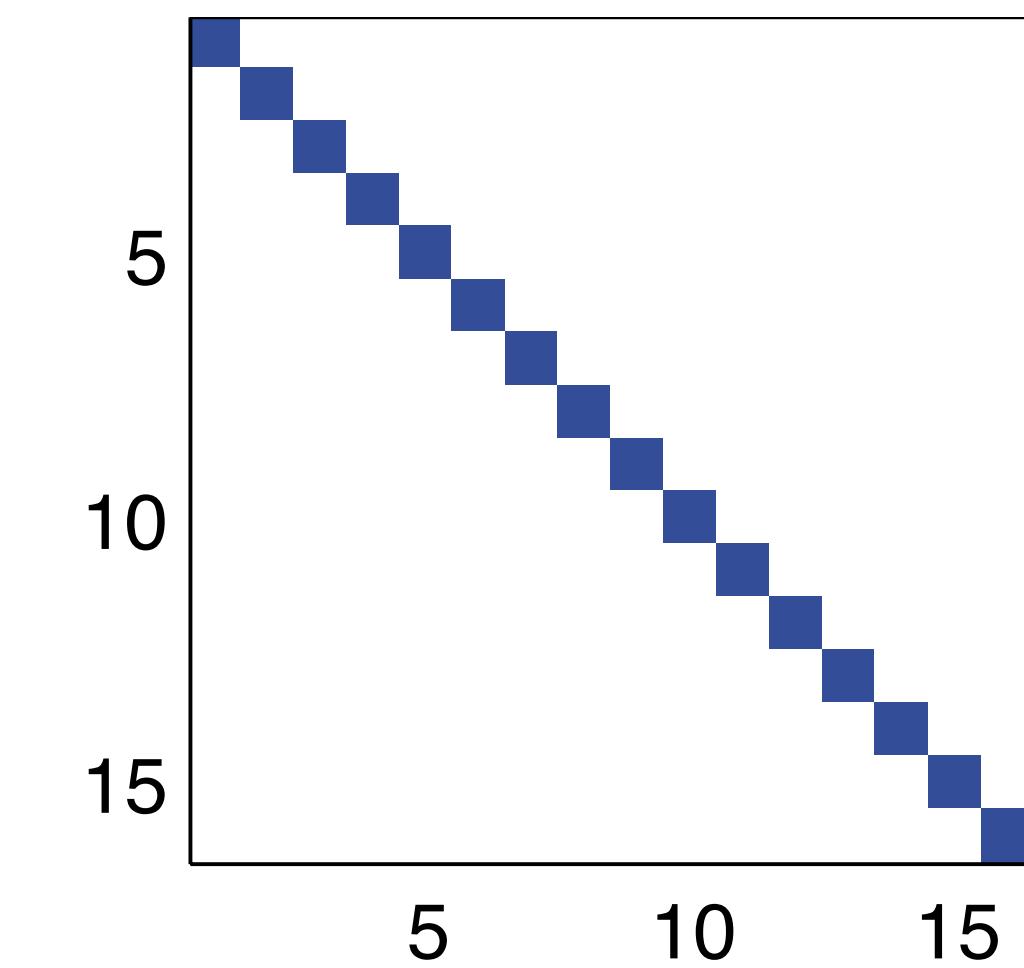


# Examples of incoherent dictionaries

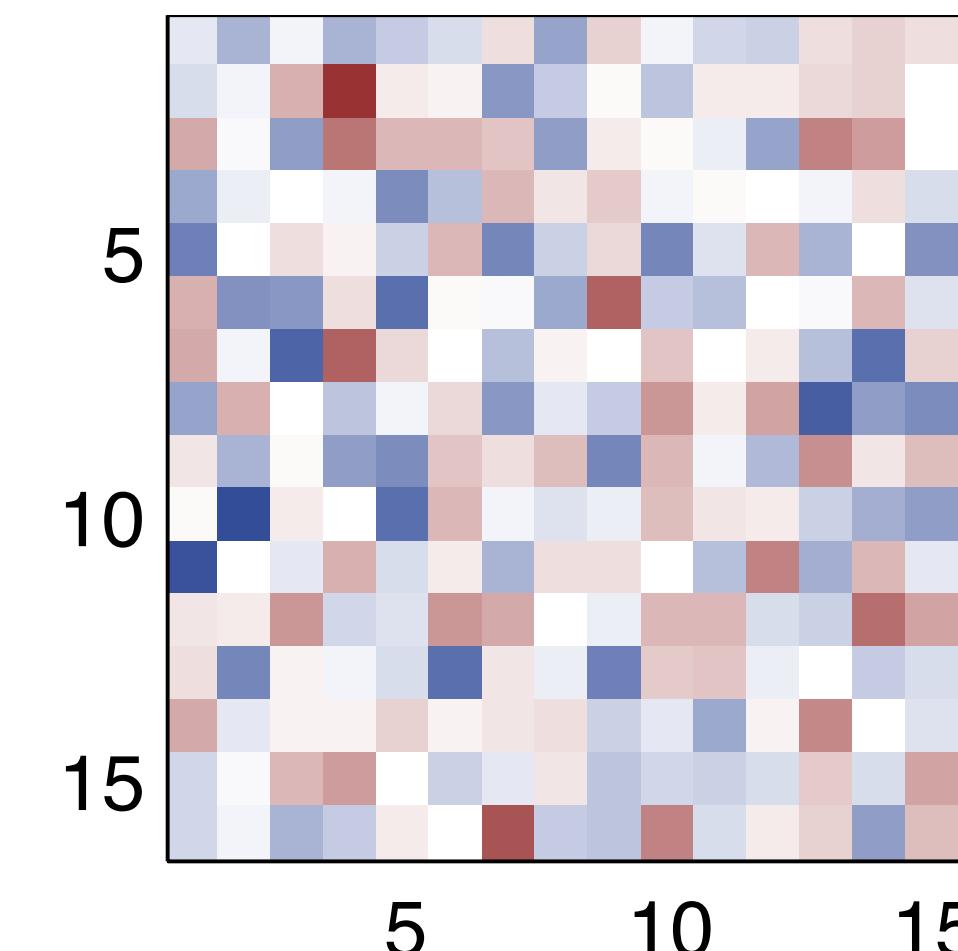
Fourier-like bases



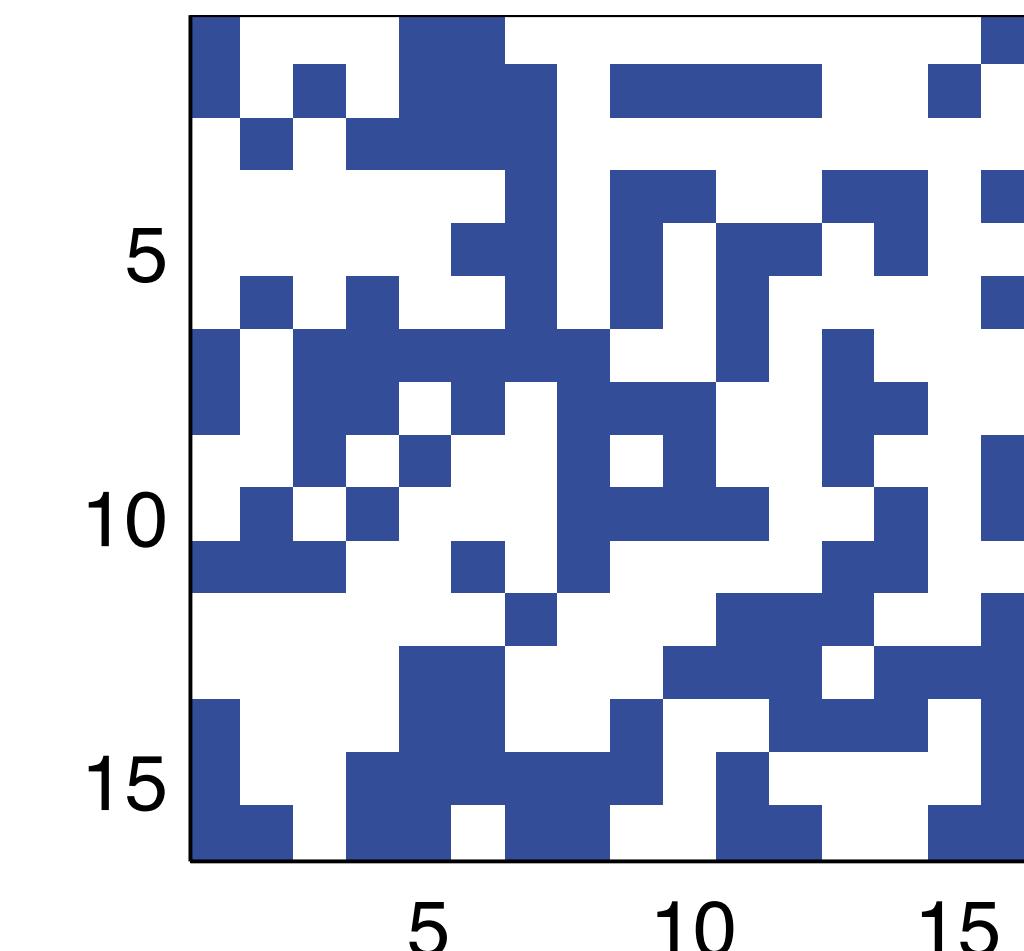
Impulse bases



Gaussian noise bases



Binary noise bases

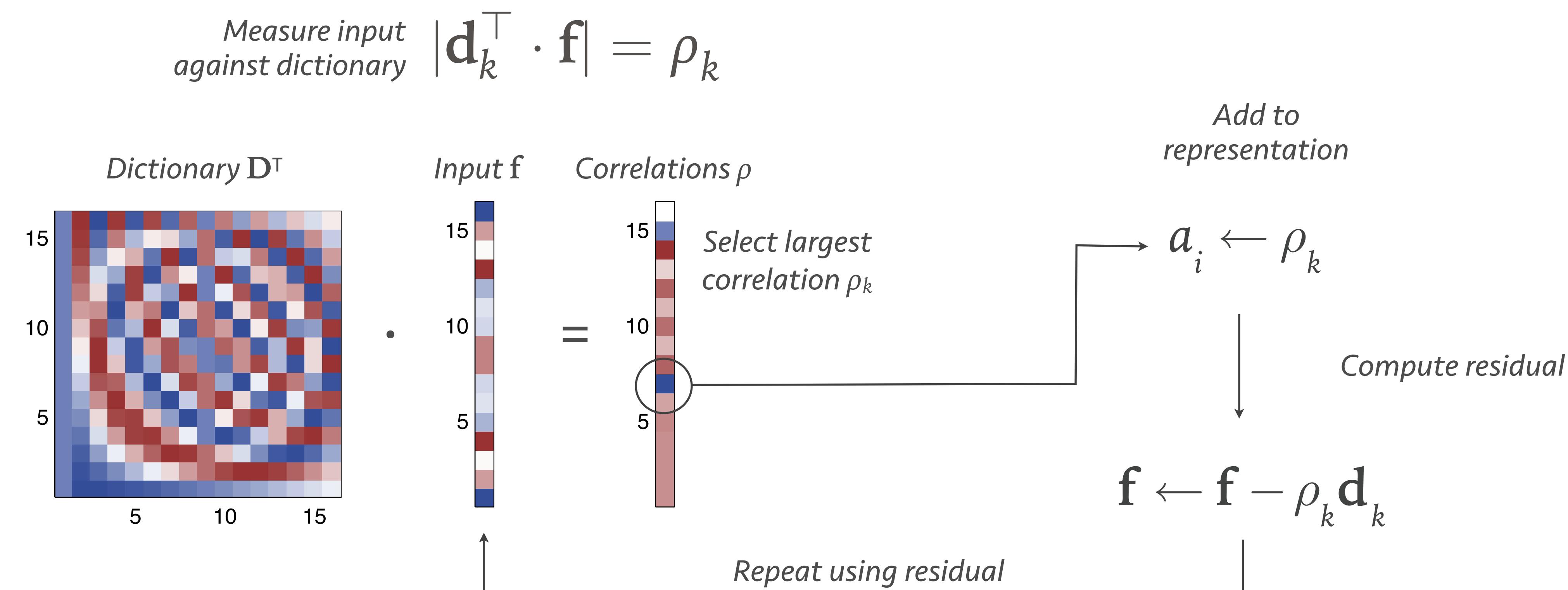


# Getting greedy

- The linear programming approach will fail now
  - It is slow for large dictionaries
  - It looks for exact equality, not approximation
- We can instead use a greedy approach to resolving sparse approximations

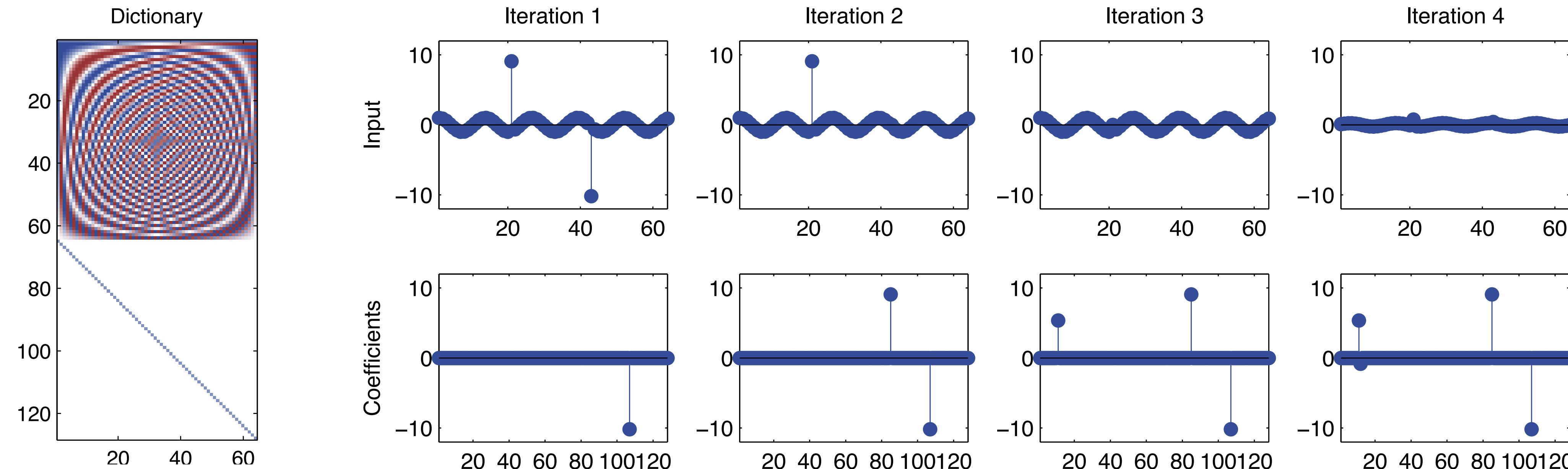
# Matching pursuit (MP)

- Family of many approaches based on successive fits
  - Each new fit explains what the previous ones couldn't



# On a familiar example

- Each iteration knocks off an element
  - by 4th iteration there's nothing significant left to represent
  - Faster! For  $N = 1024$ , MP: 0.005 sec, Linear Programming: 63 sec



# Revisiting sampling

- In traditional signal acquisition we sample uniformly
  - e.g. constant sample rates in audio, CCD grids in camera
- Foundation: Nyquist/Shannon sampling theory
  - Sample at twice the highest frequency
  - Projects to a complete basis that spans all the signal space

# A redundancy in the loop

- Take a picture
  - Using dense sampling → lots of data
- Transform to a sparse domain and quantize
  - i.e. MPEG/JPEG compression → fewer data
- Process, transmit, view, etc.

# Compressive sensing

- Why sample and then compress?
  - Do both at once!
- Sample fewer samples and use signal sparsity
  - Helps in finding a unique and plausible sparse signal

# The compressive sensing pipeline

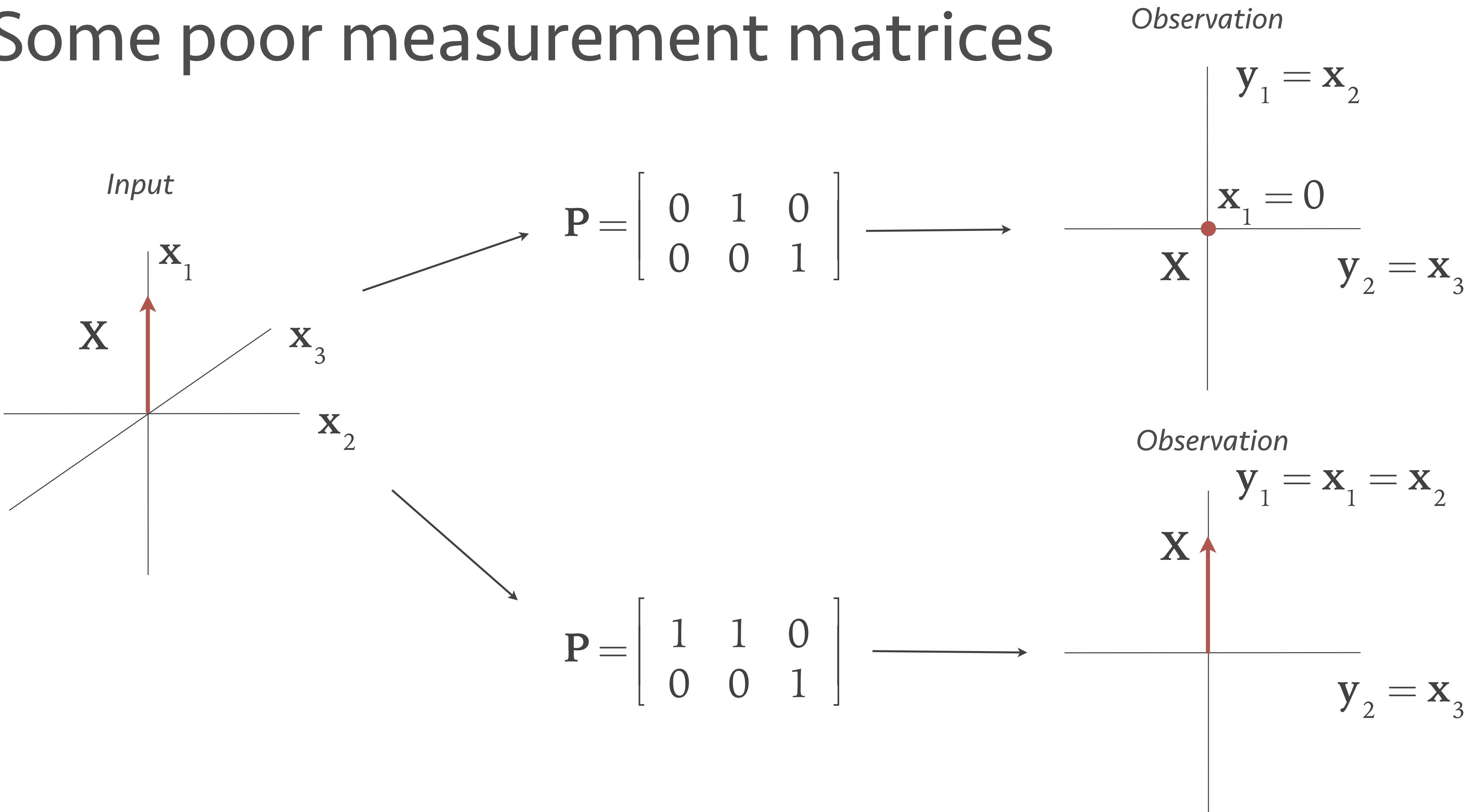
- Acquire signal using underdetermined measurements
  - e.g. linear combinations of a few samples:  $\mathbf{y}_i = \mathbf{P}_i \cdot \mathbf{x}$
  - Don't sample densely, don't sample all the data
- Reconstruct signal assuming sparsity
  - Sparsity constraints signal space w.r.t. measurements
  - Allows for a plausible reconstruction

# What's a good measurement matrix?

- We need:
$$\mathbf{P} \cdot \mathbf{x}_1 \neq \mathbf{P} \cdot \mathbf{x}_2 \text{ for all } K\text{-sparse } \mathbf{x}_1 \neq \mathbf{x}_2$$
  - i.e. ensure that we can distinguish different inputs
- Necessary condition:  $\mathbf{P}$  must have at least  $2K$  rows
  - Assuming noiseless and well-behaved data

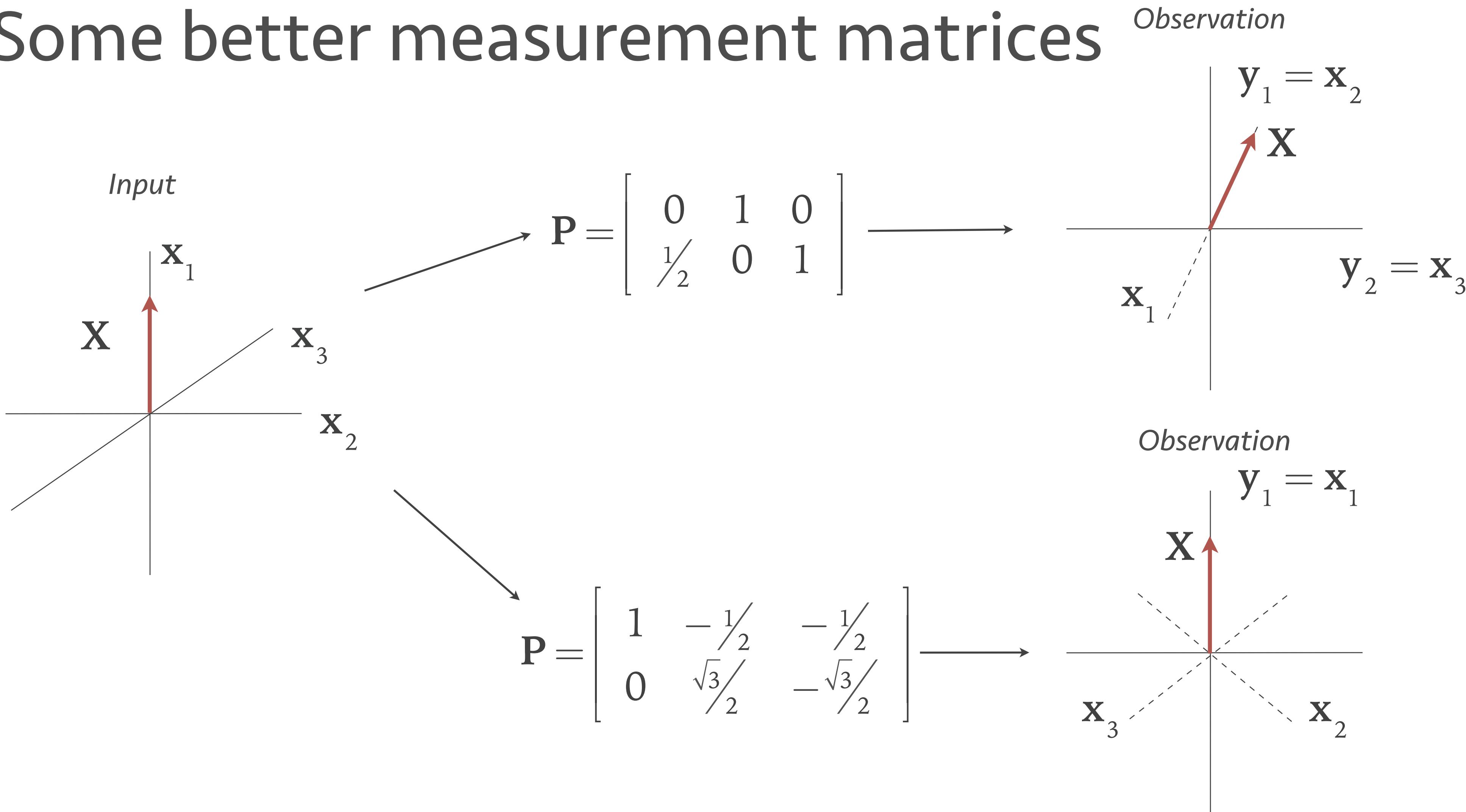
# Example 1-sparse case

- Some poor measurement matrices



# Example 1-sparse case

- Some better measurement matrices



# Restricted Isometry Property (RIP)

- A measurement matrix  $P$  has order- $K$  RIP if:

$$(1 - \delta_K) \leq \frac{\|P \cdot x\|_2^2}{\|x\|_2^2} \leq (1 + \delta_K), \forall K\text{-sparse } x$$

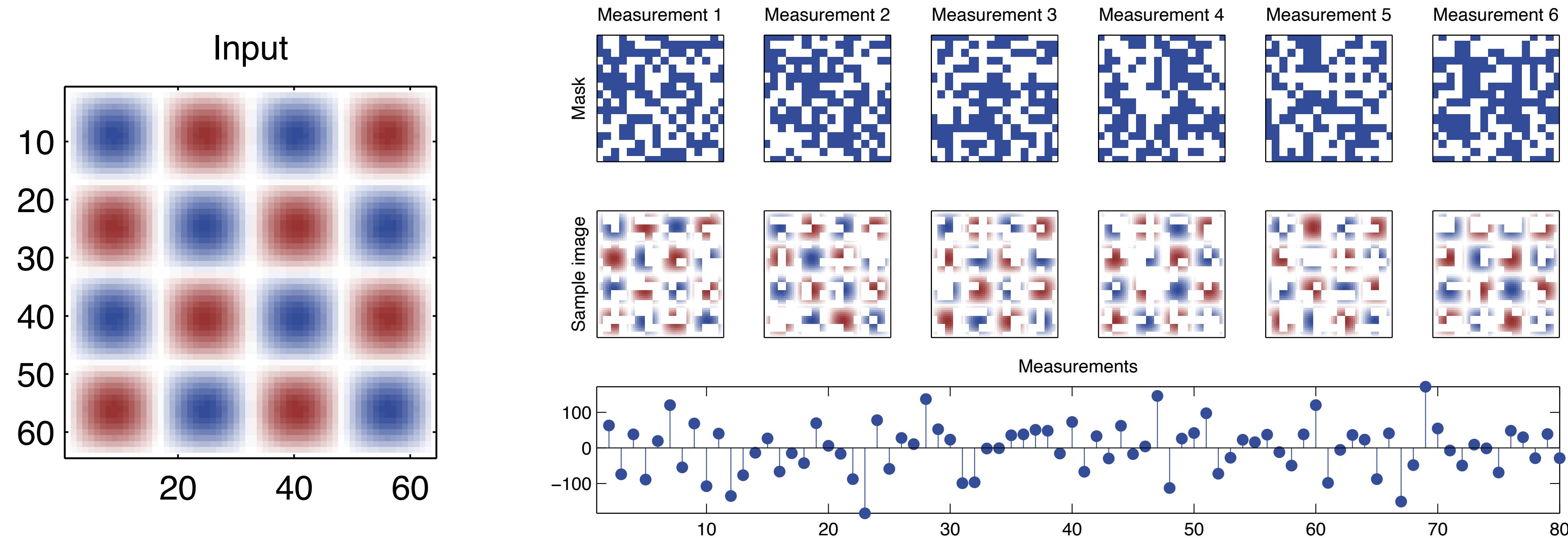
- How do we check? Difficult ...
  - But we have some good choices
    - Gaussian, Bernoulli, Fourier, etc.

# Embedding viewpoint

- This is similar to the subspace/manifold methods
  - Find low-rank projection that preserves sample distances
- Special case: Random projection
  - For  $n$  points in  $p$ -dims there exists a  $q$ -dim projection that preserves distances by a factor of  $1 + \varepsilon$ ,  $q \geq O(\varepsilon^{-2} \log n)$

# A simple compressed sensing case

- Simple  $64 \times 64$  input
  - Take a set of masked measurements, store only average value
  - Measure using:  $y_i = \mathbf{p}_i^\top \cdot \text{vec}(\mathbf{x})$ ,  $\mathbf{p}_i \sim \text{Bern}(0.5)$



# Resolving via the DCT

- Model is:

$$\mathbf{y} = \mathbf{P}^\top \cdot \mathbf{x} = \mathbf{P}^\top \cdot (\mathbf{C}^\top \cdot \mathbf{z})$$

*Measure in the DCT domain instead*

- Assume sparsity in  $\mathbf{z}$

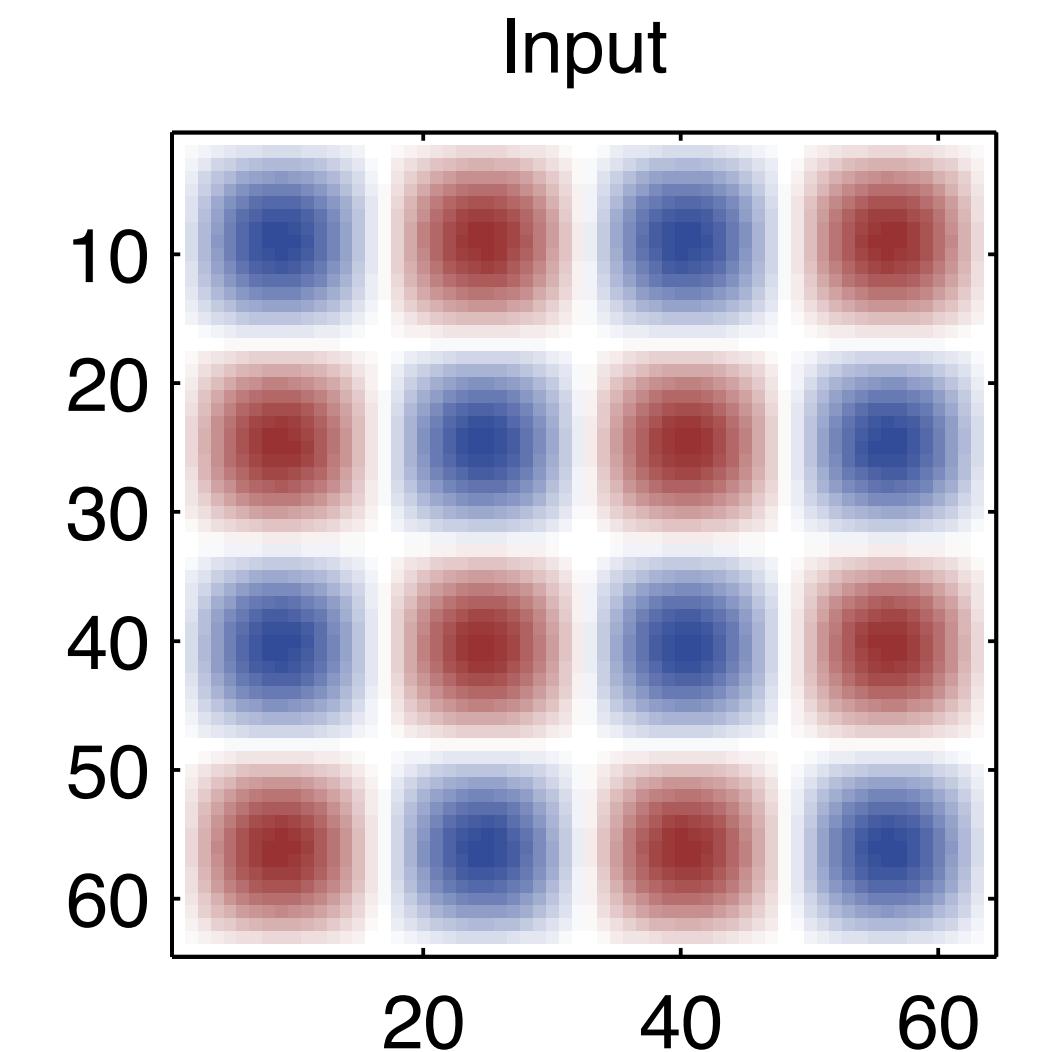
- Resolve:

$$\min \|\mathbf{z}\|_1 + \|\mathbf{y} - \mathbf{P}^\top \cdot \mathbf{C}^\top \cdot \mathbf{z}\|$$

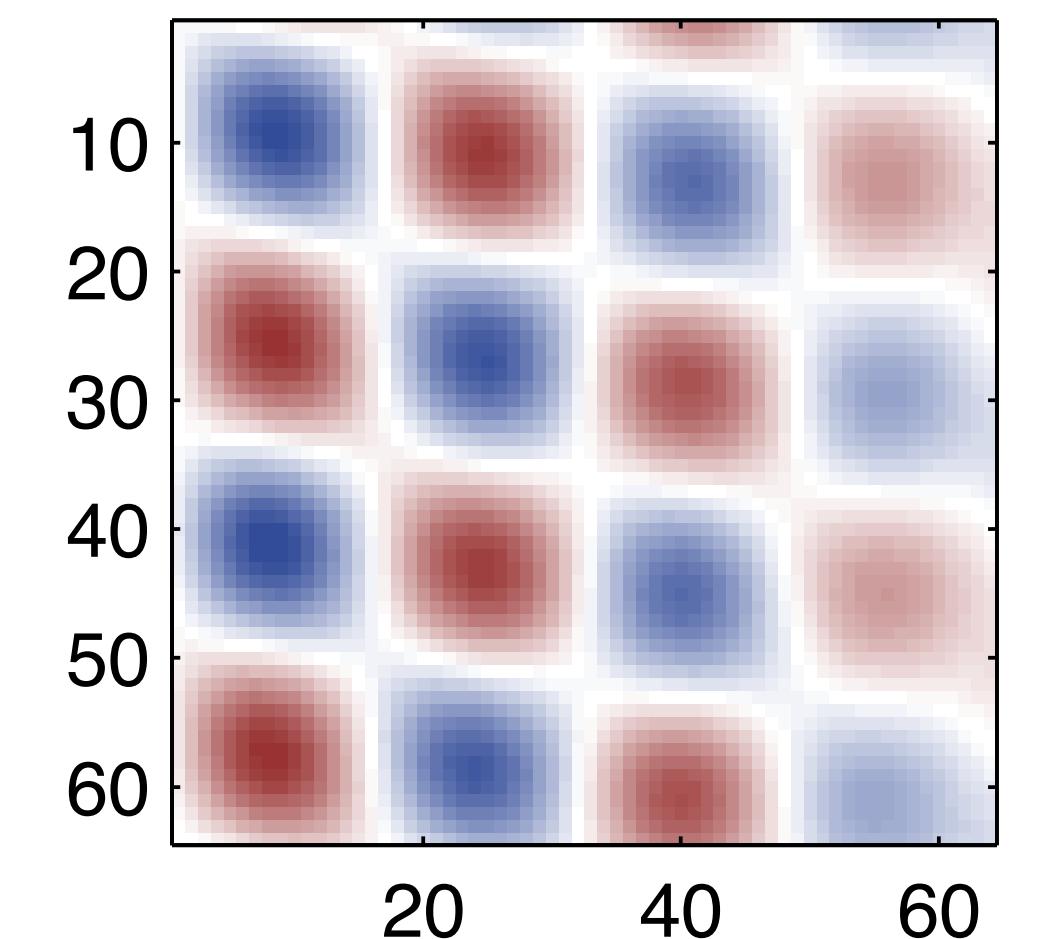
- Reconstruct:

$$\hat{\mathbf{x}} = \mathbf{C}^\top \cdot \mathbf{z}$$

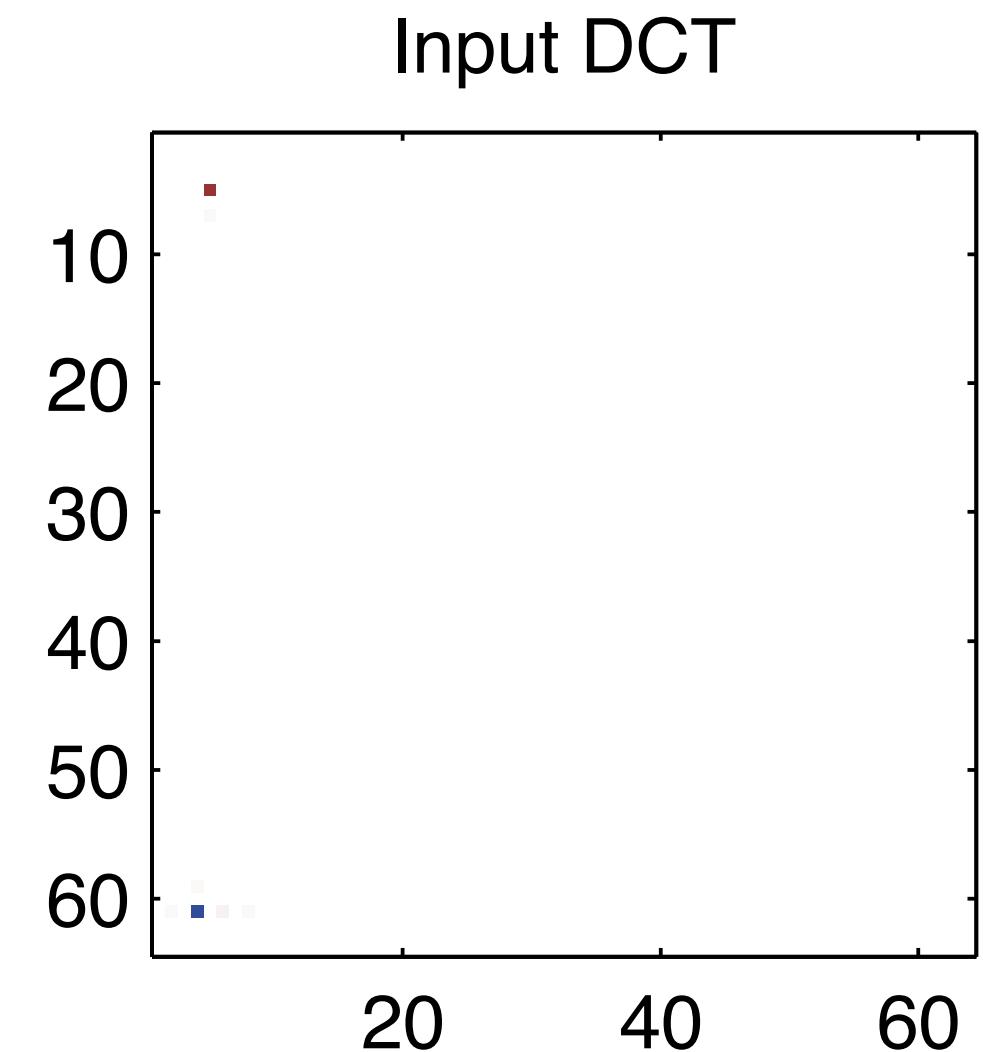
- Using ~2% of samples!



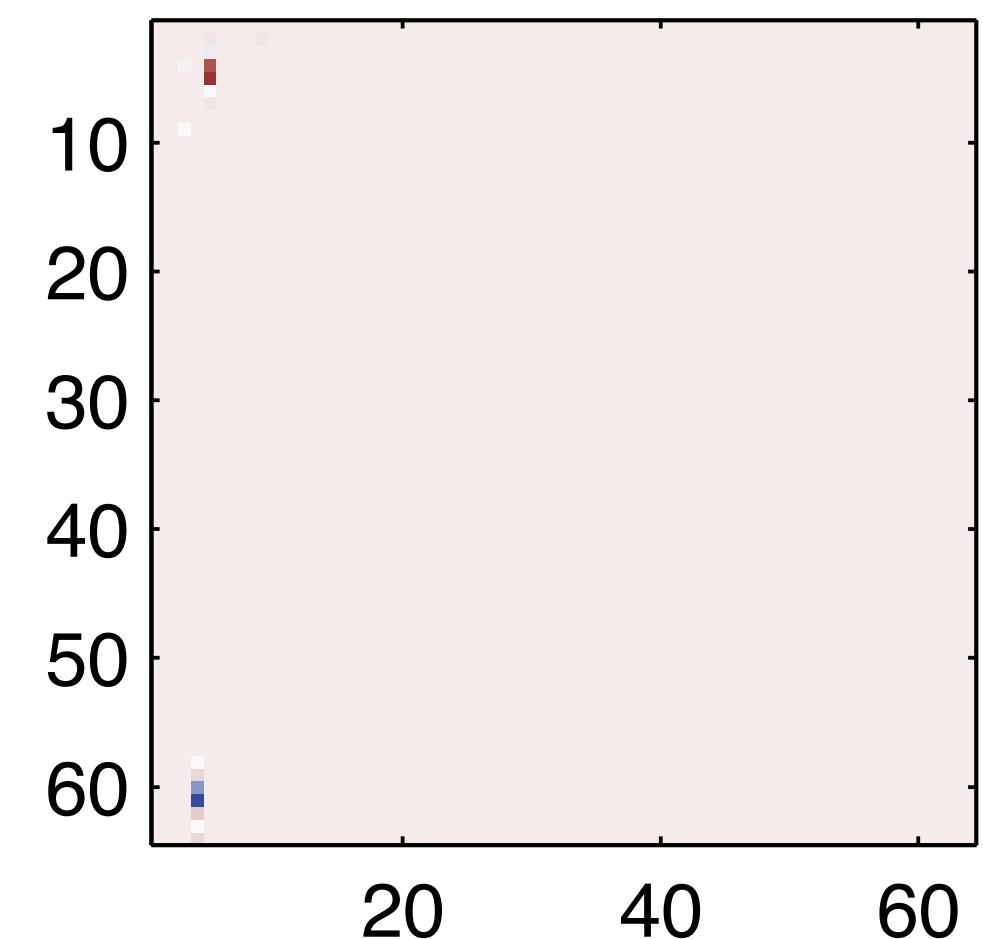
Input



Using 1.95% measurements



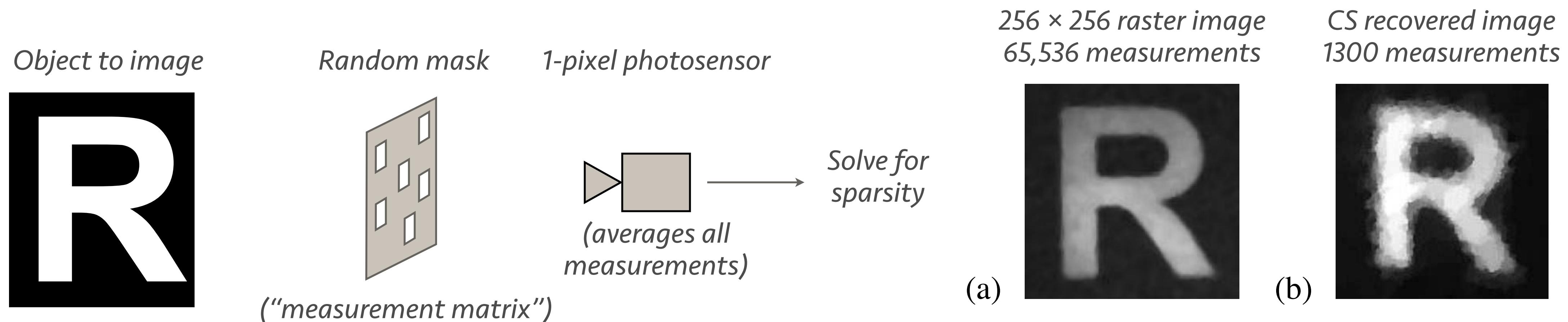
Input DCT



Extracted DCT

# Putting compressive sensing to work

- The single-pixel camera
  - Measure image intensity using multiple random masks
  - Reconstruct assuming sparsity
    - In some domain ...



# Other kinds of sparsity

- Up to now we talked about a simple form of sparsity
  - Sparsity over all coefficients separately
- Some problems need more elaborate definitions
  - e.g. block structure, joint structure, temporal structure, etc.

# Block sparsity

- Have sparsity appear in non-overlapping blocks
  - Useful for some imaging operations

$$\mathbf{y} = \Phi \mathbf{x}$$

Diagram illustrating the relationship between measurements  $\mathbf{y}$ , measurement matrix  $\Phi$ , and sparse representation  $\mathbf{x}$ .

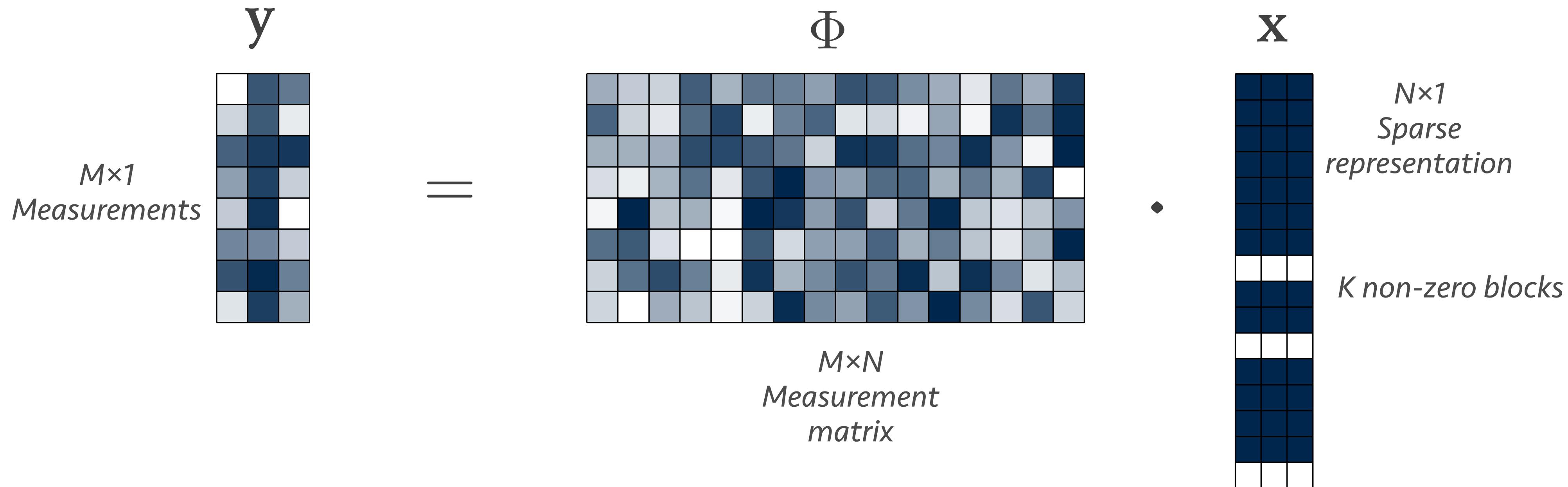
$\mathbf{y}$ :  $M \times 1$  Measurements

$\Phi$ :  $M \times N$  Measurement matrix

$\mathbf{x}$ :  $N \times 1$  Sparse representation with  $K$  non-zero blocks.

# Joint sparsity

- Obtain multiple solutions with the same sparsity
  - Useful for multimodal/multichannel data



# Learning and sparsity

- How about other models?
  - Add sparsity as an extra “regularizer”
- Sparse decompositions
  - Sparse NMF, sparse PCA, etc ...
- ICA is already (sort of) sparse!

# Recap

- Sparsity and  $\ell_p$  norms
  - Different definition of sparsity
- Minimum- $\ell_1$  coefficient algorithms
  - Linear programming
  - Greedy methods
- Compressive sensing and random projections
  - 1-pixel camera

# Reading material

- Compressive Sensing page
  - <http://dsp.rice.edu/cs>
- Experiments with Random Projections
  - <http://dimacs.rutgers.edu/Research/MMS/PAPERS/rp.pdf>

# Next lecture

- Deep learning!
  - aka neural nets v2.0
- Also Problem Set 4 will come out tomorrow
  - *Optional*, due at last day of classes
  - Use it to perk up your grade if you need to
    - Would also help if your final project is floundering