

Analyzing Non-convex Optimization for Sparse Coding

Tengyu Ma

Princeton

Sanjeev Arora
Princeton

Rong Ge
Microsoft Research

Ankur Moitra
MIT

Sparse Coding (a.k.a Dictionary Learning)

- Data has sparse representation in an appropriate basis:

Sparse Coding (a.k.a Dictionary Learning)

- Data has sparse representation in an appropriate basis:

$$\begin{bmatrix} y \end{bmatrix} \approx \begin{bmatrix} \vdots \\ A_1 & \dots & \dots \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_m \end{bmatrix}$$

Sparse Coding (a.k.a Dictionary Learning)

- Data has sparse representation in an appropriate basis:

$$\begin{bmatrix} y \end{bmatrix} \approx \begin{bmatrix} \vdots & & & \vdots \\ A_1 & \dots & \dots & A_m \\ \vdots & & & \vdots \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_m \end{bmatrix}$$

unknown dictionary
($n \times m, m > n$)

Sparse Coding (a.k.a Dictionary Learning)

- Data has sparse representation in an appropriate basis:

$$\begin{bmatrix} y \end{bmatrix} \approx \begin{bmatrix} \vdots & & & & \vdots \\ A_1 & \dots & \dots & & A_m \\ \vdots & & & & \vdots \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_i \\ \vdots \\ x_m \end{bmatrix}$$

unknown dictionary representation
($n \times m, m > n$) **k** non-zeros

Sparse Coding (a.k.a Dictionary Learning)

- Data has sparse representation in an appropriate basis:

$$\begin{bmatrix} y \end{bmatrix} \approx \begin{bmatrix} \vdots \\ A_1 & \dots & \dots \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} A_m \\ \vdots \\ x_1 \\ \vdots \\ x_i \\ \vdots \\ x_m \end{bmatrix}$$

unknown dictionary representation
($n \times m, m > n$) **k** non-zeros

- In compressed sensing, dict. is known and designed
- Here: **learn** unknown dict. (designed by nature)

Sparse Coding (a.k.a Dictionary Learning)

- Learning the basis from multiple samples

$$\begin{bmatrix} y^1 & \dots & \dots & y^p \end{bmatrix} \approx \begin{bmatrix} \vdots & \dots & \dots & \vdots \\ A_1 & \dots & \dots & A_m \\ \vdots & \dots & \dots & \vdots \end{bmatrix} \cdot \begin{bmatrix} x^1 & \dots & \dots & x^p \end{bmatrix}$$

unknown dictionary representations
($n \times m, m > n$) k non-zeros each col.

Sparse Coding (a.k.a Dictionary Learning)

- Learning the basis from multiple samples

$$\begin{bmatrix} y^1 & \dots & \dots & y^p \end{bmatrix} \approx \begin{bmatrix} \vdots & \dots & \dots & \vdots \\ A_1 & \dots & \dots & A_m \\ \vdots & \dots & \dots & \vdots \end{bmatrix} \cdot \begin{bmatrix} x^1 & \dots & \dots & x^p \end{bmatrix}$$

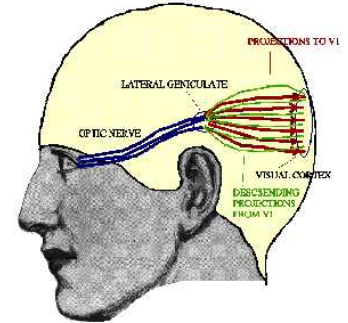
unknown dictionary representations
($n \times m, m > n$) k non-zeros each col.

notation: $Y \approx A \cdot X$

Origins of Sparse Coding

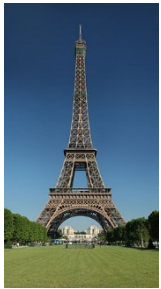
Olshausen and Field'97:

Primal visual cortex V1

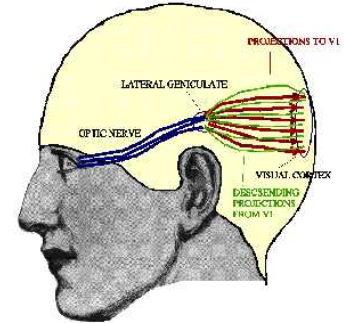


Origins of Sparse Coding

Olshausen and Field'97:



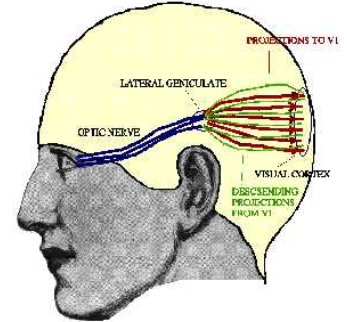
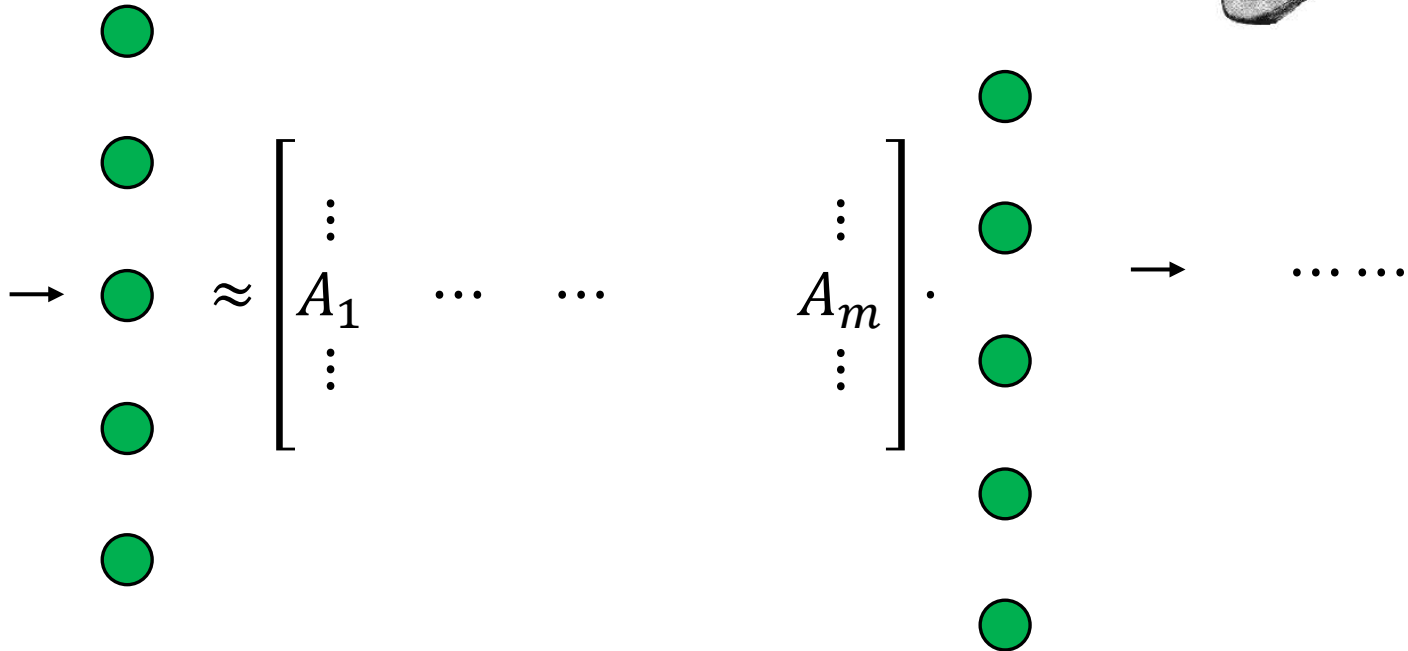
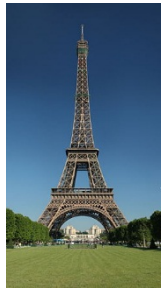
Primal visual cortex V1



Origins of Sparse Coding

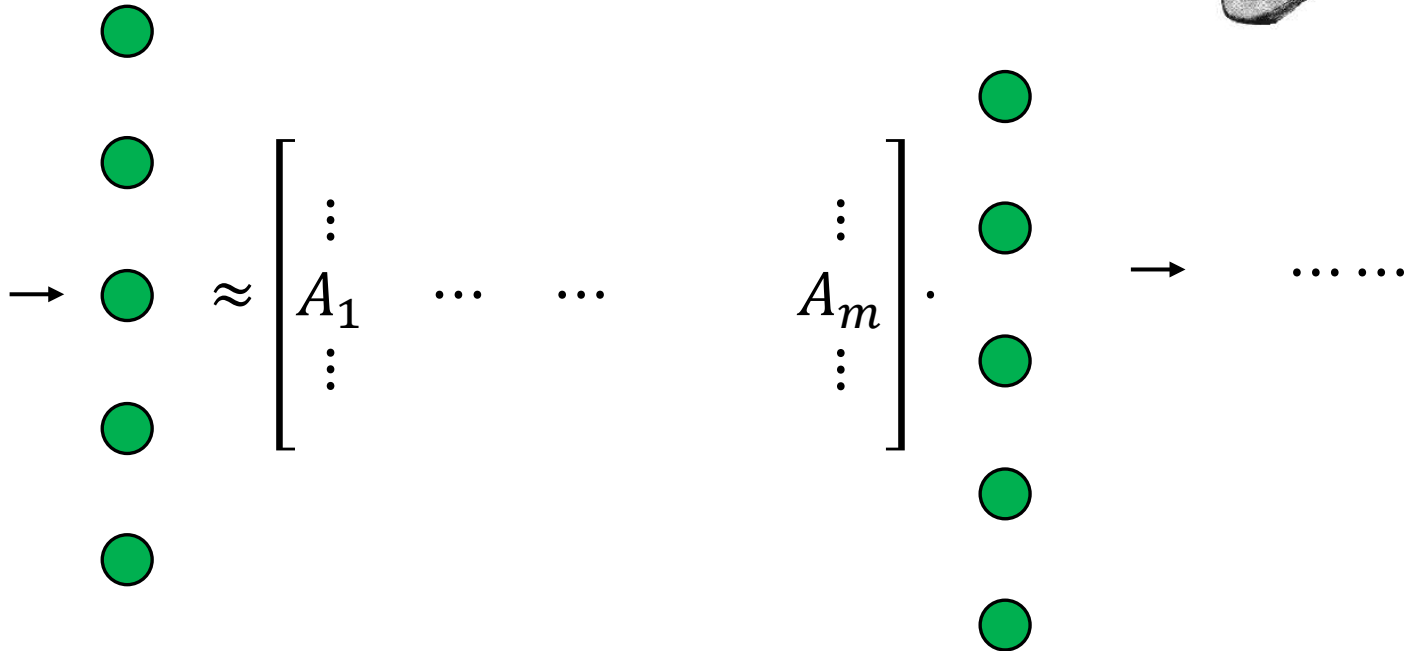
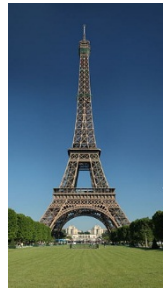
Olshausen and Field'97:

Primal visual cortex V1



Origins of Sparse Coding

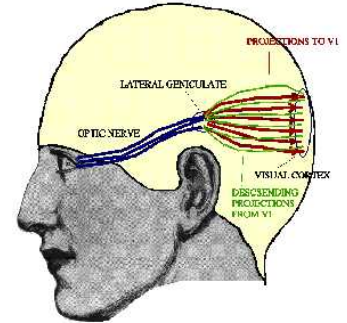
Olshausen and Field'97:



neurons

neurons

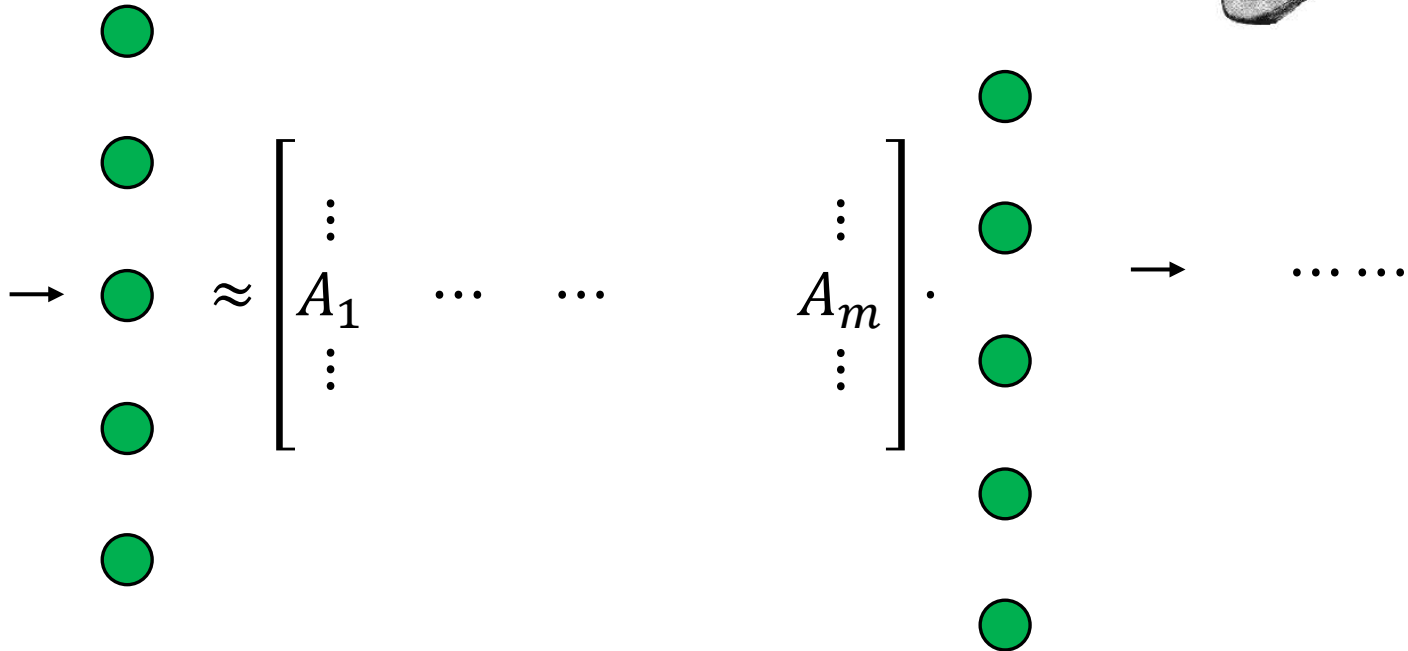
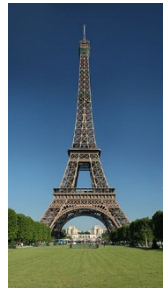
Primal visual cortex V1



Origins of Sparse Coding

Olshausen and Field'97:

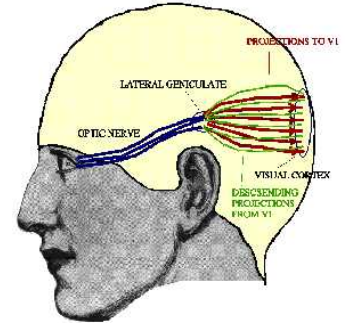
Primal visual cortex V1



neurons

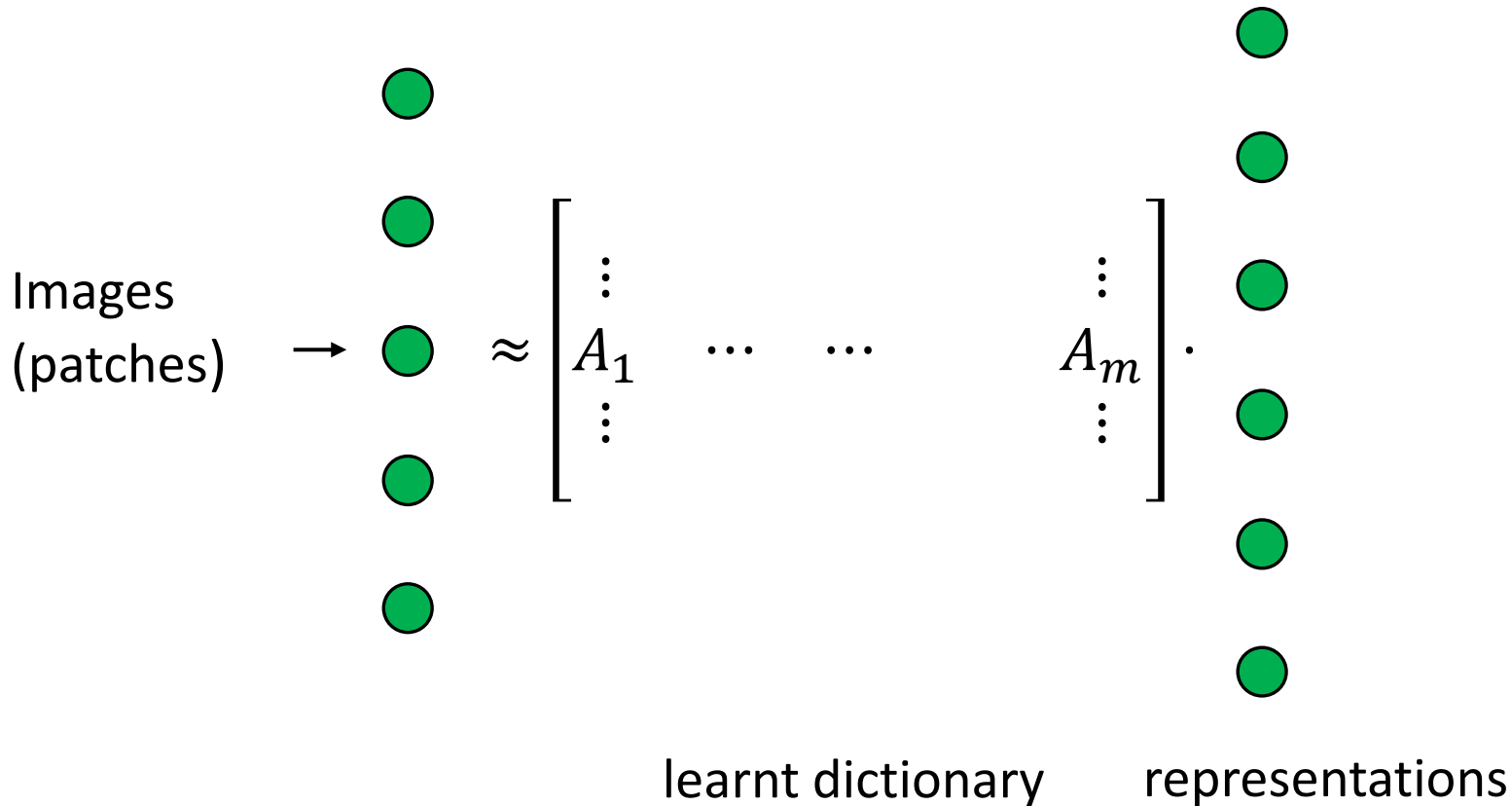
A : weights on synapses

neurons



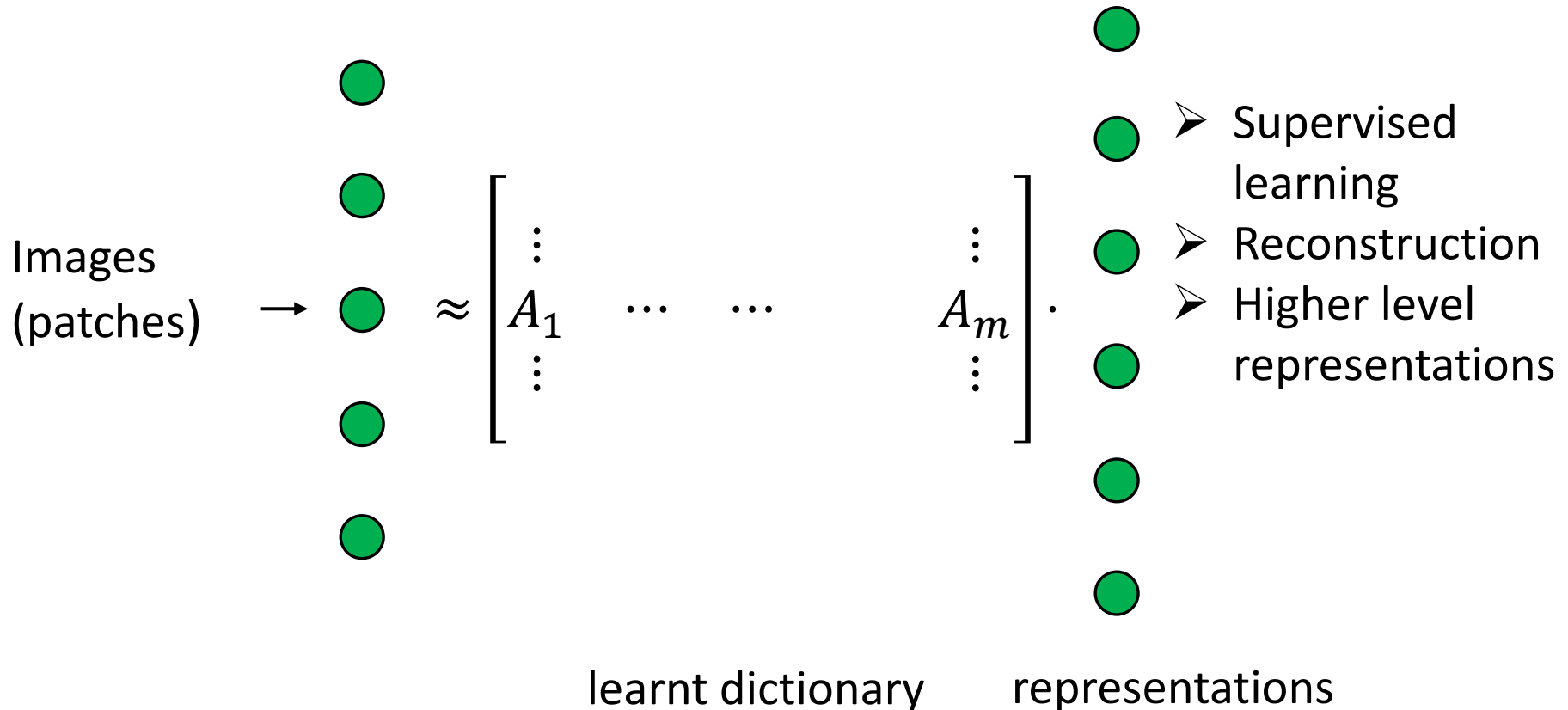
Applications to Machine Learning

- Denoising, Edge Detection, Super-resolution, Deep Learning



Applications to Machine Learning

- Denoising, Edge Detection, Super-resolution, Deep Learning



Outline

Efficient and provable algorithms for sparse coding?

- Introduction
 - Non-convex heuristics used in practice
 - Generative model and main results
- Analyzing alternating minimization

Non-convex Heuristics

- Recall

$$\begin{bmatrix} y^1 & \dots & \dots & y^p \end{bmatrix} \approx \begin{bmatrix} \vdots & & & \\ A_1 & \dots & \dots & \\ \vdots & & & \\ A_m & & & \\ \vdots & & & \end{bmatrix} \cdot \begin{bmatrix} x^1 & \dots & \dots & x^p \end{bmatrix}$$

data

dictionary

representation

- Energy function

$$\mathcal{E}(A, X) = \sum_{j=1}^p \|y^j - Ax^j\|^2 + \sum_{j=1}^p \ell(x^j)$$

Non-convex Heuristics

- Recall

$$\begin{bmatrix} y^1 & \dots & \dots & y^p \end{bmatrix} \approx \begin{bmatrix} \vdots & & & \\ A_1 & \dots & \dots & \\ \vdots & & & \end{bmatrix} \cdot \begin{bmatrix} \vdots \\ A_m \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} x^1 & \dots & \dots & x^p \end{bmatrix}$$

data

dictionary

representation

- Energy function

$$\mathcal{E}(A, X) = \sum_{j=1}^p \|y^j - Ax^j\|^2 + \sum_{j=1}^p \ell(x^j)$$

Sparse penalty, e.g.

$$\ell \cdot = \begin{bmatrix} |\cdot & 0, & |\cdot & 1 \\ 1 & 1 & 1 & 0 & 0 & |\cdot & | \\ & & & & & & 1 \end{bmatrix}$$

Minimizing Energy Function

$$\min_{A, X} \mathcal{E}(A, X) = \sum_{j=1}^p \|y^j - Ax^j\|^2 + \sum_{j=1}^p \ell(x^j)$$

Minimizing Energy Function

$$\min_{A, X} \mathcal{E}(A, X) = \sum_{j=1}^p \|y^j - Ax^j\|^2 + \sum_{j=1}^p \ell(x^j)$$

- **Alternating** Update Algorithms (Olshausen and Field'97)
- Repeat until convergence
 - **Decoding Step:** Fix A , update X

$$x^j \leftarrow \operatorname{argmin}_{x^j} \|y^j - Ax^j\|^2 + \ell(x^j)$$

- **Learning Step:** Fix X , update A

$$A \leftarrow A - \eta \nabla_A \mathcal{E}(A, X)$$

Previous Works

| | sparsity | Over-cmplt. | #samples |
|------------------------------------|---------------------------|-------------|---------------------------------------|
| LP-based Algo [SWW'12] | $k \lesssim \sqrt{n}$ | No | n^2 |
| Combinatorial [AGM'13]][AAJNT'13] | $k \lesssim n^{1/2-\eta}$ | Yes | m^2 |
| Lassere/SoS relaxation [BKS'14] | $k \lesssim n^{1-\eta}$ | Yes | $m^{\text{poly}(\frac{1}{\epsilon})}$ |

Previous Works

| | sparsity | Over-cmplt. | #samples |
|-----------------------------------|---------------------------|-------------|---------------------------------------|
| LP-based Algo [SWW'12] | $k \lesssim \sqrt{n}$ | No | n^2 |
| Combinatorial [AGM'13] [AAJNT'13] | $k \lesssim n^{1/2-\eta}$ | Yes | m^2 |
| Lassere/SoS relaxation [BKS'14] | $k \lesssim n^{1-\eta}$ | Yes | $m^{\text{poly}(\frac{1}{\epsilon})}$ |

These provable approaches are highly noncompetitive with non-convex optimization via alternating updates!

Previous Works

| | sparsity | Over-cmplt. | #samples |
|------------------------------------|---------------------------|-------------|---------------------------------------|
| LP-based Algo [SWW'12] | $k \lesssim \sqrt{n}$ | No | n^2 |
| Combinatorial [AGM'13]][AAJNT'13] | $k \lesssim n^{1/2-\eta}$ | Yes | m^2 |
| Lassere/SoS relaxation [BKS'14] | $k \lesssim n^{1-\eta}$ | Yes | $m^{\text{poly}(\frac{1}{\epsilon})}$ |

These provable approaches are highly **noncompetitive** with non-convex optimization via alternating updates!

Can we analyze non-convex optimization??

Previous Works

| | sparsity | Over-cmplt. | #samples |
|-----------------------------------|---------------------------|-------------|---------------------------------------|
| LP-based Algo [SWW'12] | $k \lesssim \sqrt{n}$ | No | n^2 |
| Combinatorial [AGM'13] [AAJNT'13] | $k \lesssim n^{1/2-\eta}$ | Yes | m^2 |
| Lassere/SoS relaxation [BKS'14] | $k \lesssim n^{1-\eta}$ | Yes | $m^{\text{poly}}(\frac{1}{\epsilon})$ |

These provable approaches are highly **noncompetitive** with non-convex optimization via alternating updates!

Can we analyze non-convex optimization??

Local convergence is known (if start with dictionary that is $1/k$ – close to optimal) [AGM'13][AAJNT'13]

- This paper: convergence starting from $1/\log n$ -close.
+ **better** runtime/sample complexity

Model Assumptions

$$\begin{array}{c} n\text{-dim} \left\{ \begin{array}{l} \left[\begin{array}{cccc} y^1 & \dots & \dots & y^p \end{array} \right] \\ \text{data} \end{array} \right. = \begin{array}{l} \left[\begin{array}{ccc} \vdots & \dots & \dots \\ A_1^* & \dots & \dots \\ \vdots & \dots & \dots \end{array} \right] \\ \text{dictionary} \end{array} \cdot \begin{array}{l} \left[\begin{array}{cccc} \vdots & \dots & \dots & \vdots \\ A_m^* & \dots & \dots & \dots \\ \vdots & \dots & \dots & \vdots \end{array} \right] \cdot \left[\begin{array}{cccc} x^{*1} & \dots & \dots & x^{*p} \end{array} \right] \\ \text{representation} \end{array} \end{array}$$

Model Assumptions

$$\begin{array}{c} n\text{-dim} \left\{ \begin{array}{l} \left[\begin{array}{cccc} y^1 & \dots & \dots & y^p \end{array} \right] \\ \text{data} \end{array} \right. = \begin{array}{c} \left[\begin{array}{ccc} \vdots & & \\ A_1^* & \dots & \dots \\ \vdots & & \end{array} \right] \\ \text{dictionary} \end{array} \cdot \begin{array}{c} \left[\begin{array}{cccc} \vdots & & & \\ A_m^* & & & \\ \vdots & & & \end{array} \right] \cdot \left[\begin{array}{cccc} x^{*1} & \dots & \dots & x^{*p} \end{array} \right] \\ \text{representation} \end{array} \end{array}$$

- $m = O(n)$;
- $\|A_i^*\| = 1$;
- $\langle A_i^*, A_j^* \rangle \lesssim 1/\sqrt{n}$;
- $\|A^*\| \leq \sqrt{m/n}$;

Model Assumptions

$$\begin{array}{c} n\text{-dim} \end{array} \left\{ \begin{array}{c} \left[\begin{array}{cccc} y^1 & \dots & \dots & y^p \end{array} \right] = \left[\begin{array}{ccc} \vdots & & \\ A_1^* & \dots & \dots \\ \vdots & & \\ & & A_m^* \\ & & \vdots \end{array} \right] \cdot \left[\begin{array}{cccc} x^{*1} & \dots & \dots & x^{*p} \end{array} \right]
 \end{array} \right.$$

data

dictionary

representation

- $m = O(n)$;
- $\|A_i^*\| = 1$;
- $\langle A_i^*, A_j^* \rangle \lesssim 1/\sqrt{n}$;
- $\|A^*\| \leq \sqrt{m/n}$;
- $k \lesssim \sqrt{n}$
- i.i.d k -sparse columns
- entries of x^{*j} has bounded correlation

Our results

Theorem 1 (Convergence):

- Given A^0 column-wise $1/\log n$ – close to A^*
- Olshausen and Field algo (with a simpler decoding) returns A^s after s iteration s.t.

$$\|A^s - A^*\|_F^2 \leq 0.99^s \cdot \|A^0 - A^*\|_F^2 + O(mk^2/n^2)$$

Our results

Theorem 1 (Convergence):

$$\triangleright \|A_i^0 - A_i^*\| \leq 1/\log n$$

- Given A^0 column-wise $1/\log n$ – close to A^*
- Olshausen and Field algo (with a simpler decoding) returns A^s after s iteration s.t.

$$\|A^s - A^*\|_F^2 \leq 0.99^s \cdot \|A^0 - A^*\|_F^2 + O(mk^2/n^2)$$

Our results

Theorem 1 (Convergence):

$$\triangleright \|A_i^0 - A_i^*\| \leq 1/\log n$$

- Given A^0 column-wise $1/\log n$ – close to A^*
- Olshausen and Field algo (with a simpler decoding) returns A^s after s iteration s.t.

$$\|A^s - A^*\|_F^2 \leq 0.99^s \cdot \|A^0 - A^*\|_F^2 + O(mk^2/n^2)$$

Theorem 2 (Initialization):

spectral method based algo returns A^0 that is $1/\log n$ – close to A^*

| | sparsity | Over-cmplt. | #samples |
|------------------------------------|---------------------------|-------------|---------------------------------------|
| LP-based Algo [SWW'12] | $k \lesssim \sqrt{n}$ | No | n^2 |
| Combinatorial [AGM'13]][AAJNT'13] | $k \lesssim n^{1/2-\eta}$ | Yes | $m^2 \cdot 1/\epsilon$ |
| Lassere/SoS relaxation [BKS'14] | $k \lesssim n^{1-\eta}$ | Yes | $m^{\text{poly}(\frac{1}{\epsilon})}$ |
| This paper | $k \lesssim \sqrt{n}$ | Yes | $\tilde{O}(mk \log 1/\epsilon)$ |

| | sparsity | Over-cmplt. | #samples |
|------------------------------------|---------------------------|-------------|---------------------------------------|
| LP-based Algo [SWW'12] | $k \lesssim \sqrt{n}$ | No | n^2 |
| Combinatorial [AGM'13]][AAJNT'13] | $k \lesssim n^{1/2-\eta}$ | Yes | $m^2 \cdot 1/\epsilon$ |
| Lassere/SoS relaxation [BKS'14] | $k \lesssim n^{1-\eta}$ | Yes | $m^{\text{poly}(\frac{1}{\epsilon})}$ |
| This paper | $k \lesssim \sqrt{n}$ | Yes | $\tilde{O}(mk \log 1/\epsilon)$ |

Non-convex approach is powerful!

- Aside: new initialization also heavily inspired by alt. update algorithm.

Outline

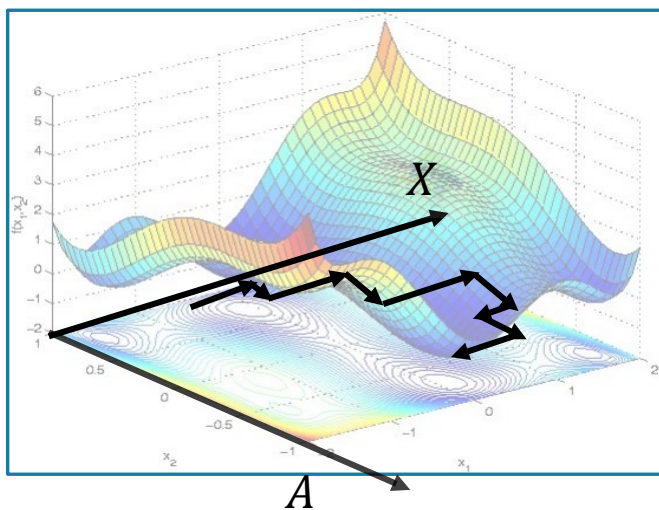
Efficient and provable algorithms for sparse coding?

- Introduction
 - Non-convex heuristics used in practice
 - Generative model and main results
- Analyzing alternating minimization

Alternating Updates \rightarrow Updates on A

- $\min_{A, X} \mathcal{E}(A, X)$
- **Decoding Step:** Fixing A , update
- $x^j \leftarrow \operatorname{argmin}_{x^j} \|y^j - Ax^j\|^2 + \ell(x^j)$
- **Learning Step:** Fixing X , update A

$$A \leftarrow A - \eta \nabla_A \mathcal{E}(A, X)$$

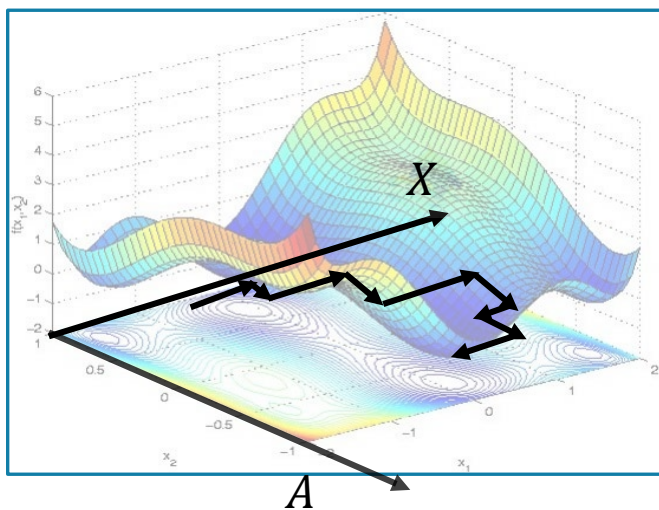


Alternating Updates \rightarrow Updates on A

- $\min_{A, X} \mathcal{E}(A, X)$
- **Decoding Step:** Fixing A , update
- $x^j \leftarrow \operatorname{argmin}_{x^j} \|y^j - Ax^j\|^2 + \ell(x^j)$
- **Learning Step:** Fixing X , update A

$$A \leftarrow A - \eta \nabla_A \mathcal{E}(A, X)$$

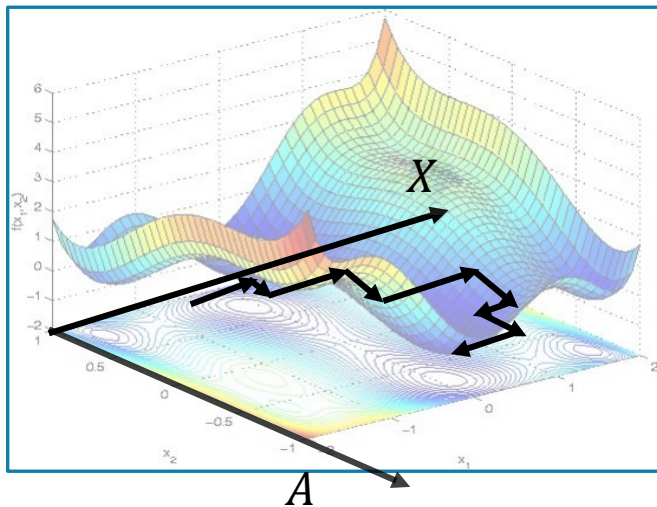
$$A^{s+1} = A^s - \eta g^s$$



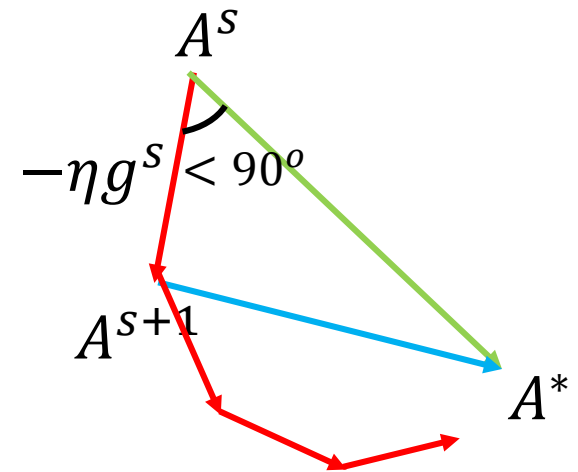
Alternating Updates \rightarrow Updates on A

- $\min_{A, X} \mathcal{E}(A, X)$
- **Decoding Step:** Fixing A , update
- $x^j \leftarrow \operatorname{argmin}_{x^j} \|y^j - Ax^j\|^2 + \ell(x^j)$
- **Learning Step:** Fixing X , update A

$$A \leftarrow A - \eta \nabla_A \mathcal{E}(A, X)$$



$$A^{s+1} = A^s - \eta g^s$$



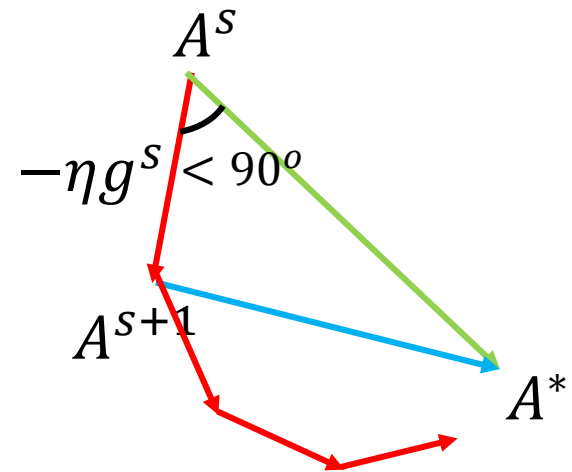
Proof Framework

- Suppose we have an update rule

$$A^{s+1} = A^s - \eta g^s$$

and a desired point A^*

(g doesn't have to be grad; and there may be no “objective fn”)



Proof Framework

- Suppose we have an update rule

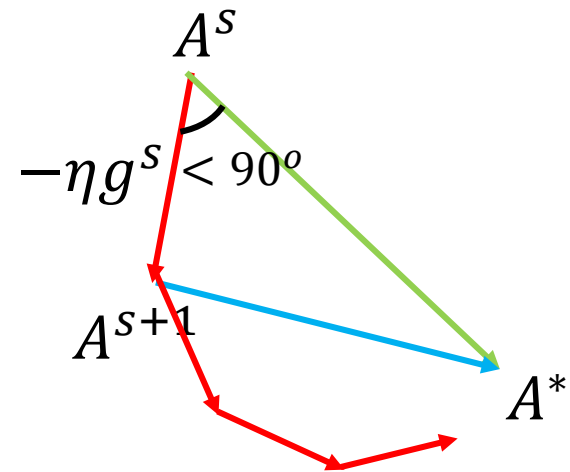
$$A^{s+1} = A^s - \eta g^s$$

and a desired point A^*

(g doesn't have to be grad; and there may be no “objective fn”)

Condition:

$$\langle g^s, A^s - A^* \rangle \geq \alpha \|A^s - A^*\|_F^2 + \beta \|g^s\|_F^2 - \epsilon_s$$



Proof Framework

- Suppose we have an update rule

$$A^{s+1} = A^s - \eta g^s$$

and a desired point A^*

(g doesn't have to be grad; and there may be no “objective fn”)

Condition:

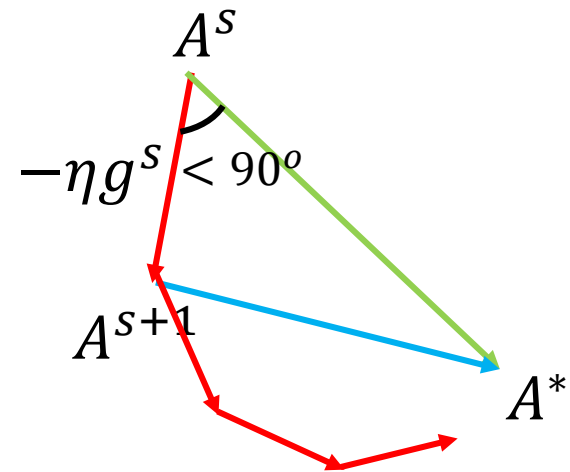
$$\langle g^s, A^s - A^* \rangle \geq \alpha \|A^s - A^*\|_F^2 + \beta \|g^s\|_F^2 - \epsilon_s$$

⇓

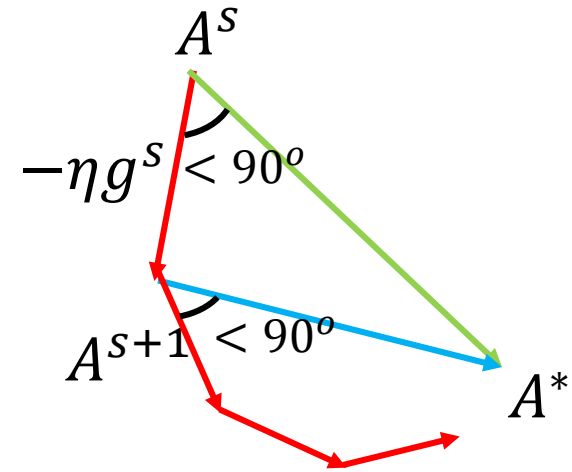
geometrical convergence:

$$\|A^s - A^*\|_F^2 \leq (1 - 2\alpha\eta)^s \cdot \|A^0 - A^*\|_F^2 + \epsilon/\alpha$$

(where $\eta \leq 2\beta, \epsilon = \max \epsilon_s$)

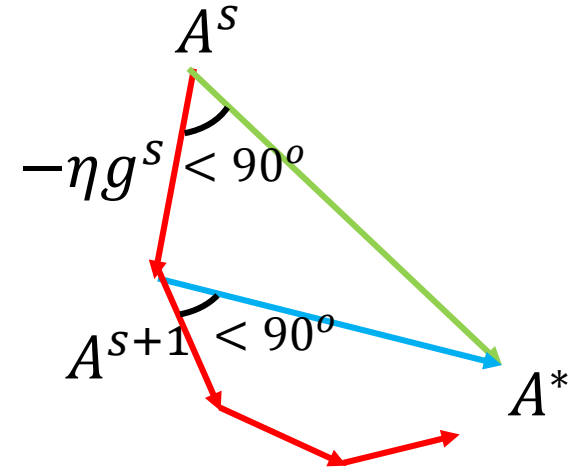


Analysis of Sparse Coding



Analysis of Sparse Coding

- $\min_{A, X} \mathcal{E}(A, X)$



Analysis of Sparse Coding

◦ $\arg\min_{x_j} \|x_j\|_1 \leftarrow$

$$\arg\min_{x_j} \|y_j - Ax_j\|_2^2 +$$

$$\ell(x_j) \quad y_j - Ax_j$$

$$y_j -$$

jjyyjjxxnimgr $\arg\min_{x_j} \|a\|_1 \leftarrow$

$$\arg\min_{x_j} \|y_j - Ax_j\|_2^2 +$$

$$\ell(x_j) \quad y_j - Ax_j$$

$$y_j -$$

jjyyjjxxnimgr $\arg\min_{x_j} \|a\|_1 \leftarrow$

$$\arg\min_{x_j} \|y_j - Ax_j\|_2^2 +$$

$$\ell(x_j) \quad y_j - Ax_j$$

$$y_j -$$

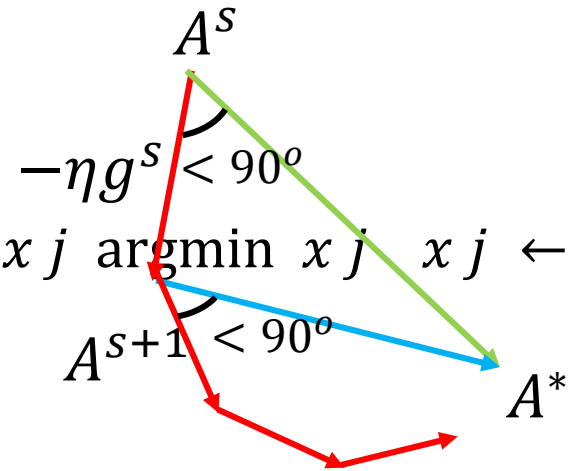
jjyyjjxxnimgr $\arg\min_{x_j} \|a\|_1 \leftarrow$

$$\arg\min_{x_j} \|y_j - Ax_j\|_2^2 +$$

$$\ell(x_j) \quad y_j - Ax_j$$

$$y_j -$$

jjyyjjxxnimgr $\arg\min_{x_j} \|a\|_1 \leftarrow$



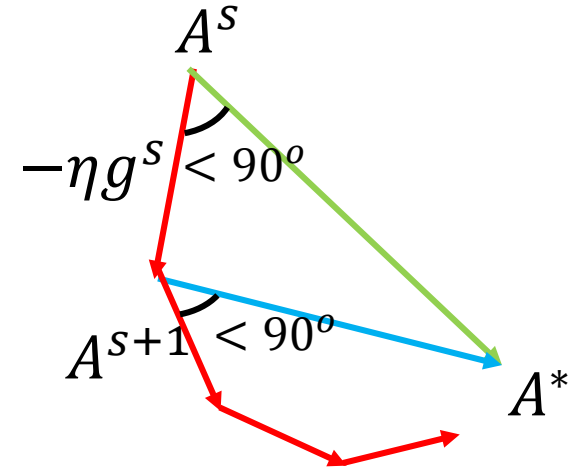
◦ **Decoding Step:** Fixing A , update

Analysis of Sparse Coding

- $\min_{A, X} \mathcal{E}(A, X)$
- **Decoding Step:** Fixing A , update
- $x^j \leftarrow \text{threshold}(A^T y^j)$
- **Learning Step:** Fixing X , update A

$$x^j \leftarrow \text{threshold}(A^T y^j)$$

$$A \leftarrow A - \eta g$$



Proof sketch for convergence Thm:

- Get a “closed” form for $g \approx B_A A - C_A A^* + \epsilon$ with $B_A, C_A \approx I$
 - using the randomness of data y^j
 - using incoherence and isotropy of A^*
 - + other tricks
- Check that $\langle g, A - A^* \rangle \geq \alpha \|A - A^*\|^2 + \beta \|g\|^2 - \epsilon_s$

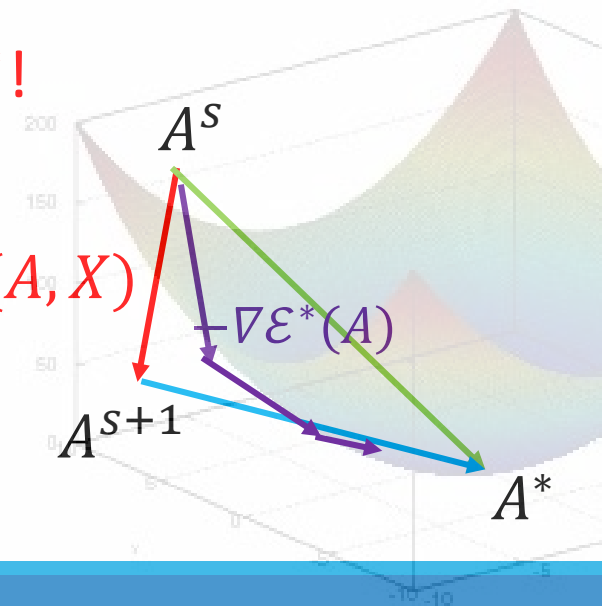
Why should we expect that $g = \nabla_A \mathcal{E}(A, X)$ is correlated with $A - A^*$?

- Consider **unknown** convex function

$$\mathcal{E}^*(A) := \mathcal{E}(A, X^*) = \sum_{j=1}^p \|y^j - Ax^{*j}\|^2 + C$$

- $X \approx X^* \Rightarrow \nabla_A \mathcal{E}(A, X) \approx \nabla_A \mathcal{E}(A, X^*) = \nabla \mathcal{E}^*(A)$
- $\mathcal{E}^*(\cdot)$ is convex $\Rightarrow \nabla \mathcal{E}^*(A)$ correlated with $A - A^*$
- Alt. update is simulating grad descent on \mathcal{E}^* !**
- (Actual proof doesn't use this explicitly)

$$-\eta \nabla_A \mathcal{E}(A, X)$$



Discussion and Open Questions

- Generalizable to alternating updates algo. for other hidden variable models?
 - Key technical difficulty: what if the decoding doesn't have a simple closed form?
- Limited to “local” convergence analysis
 - Analyzing global convergence from random initialization?
- Practical algo. (initialization) beyond $k = \sqrt{n}$?