

# Chapter 6

## Alternating minimization

### Contents (class version)

---

<b>6.0 Introduction</b>	<b>6.3</b>
<b>6.1 Signal processing applications</b>	<b>6.3</b>
Compressed sensing using synthesis sparsity models	6.4
Compressed sensing with analysis regularizer	6.10
Sparse coding revisited with multi-block BCM	6.15
Sparse coding for tight frames	6.21
Patch-based regularization: analysis form	6.23
Sparsifying transform learning	6.29
Dictionary learning via two-block BCD	6.41
Joint update of atom and coefficients	6.46
<b>6.2 Machine learning applications</b>	<b>6.47</b>
Low-rank approximation for large-scale problems	6.47
Fused LASSO / generalized LASSO	6.51
Alternating minimization for 0-norm in biconvex form	6.52
<b>6.3 Convergence properties</b>	<b>6.55</b>

<b>6.4 Summary</b> . . . . .	<b>6.57</b>
------------------------------	-------------

---

---

## 6.0 Introduction

Although the proximal methods in the previous chapter are quite flexible and useful, there are many cost functions of interest that are not smooth and also not “prox friendly.” So we need additional tools.

Furthermore, all of the methods discussed so far update all elements of the optimization variable simultaneously. For many optimization problems it can be easier to update just some of the variables at a time.

This chapter develops **alternating minimization** algorithms that are suitable for such problems.

These methods have a long history in both optimization and signal/image processing. In image processing, an early approach was **iterated conditional modes (ICM)** [1].

These **coordinate descent** methods can be very useful, but also can have limitations for certain nonsmooth cost functions.

This chapter illustrates the methods by focusing on SIMPL applications, drawing from Ch. 1 and Ch. 2.

---

## 6.1 Signal processing applications

We begin with signal and image processing applications.

## Compressed sensing using synthesis sparsity models

Consider a linear measurement model assuming **synthesis** sparsity using a (usually wide) dictionary  $D$ :

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \varepsilon, \quad \mathbf{x} \approx \mathbf{D}\mathbf{z}, \quad \mathbf{z} \text{ is sparse.}$$

(If  $\mathbf{A} = \mathbf{I}$ , then this is a **denoising** problem; if  $\mathbf{A}$  is wide then it is compressed sensing.)

For these assumptions, two related but distinct ways to estimate  $\mathbf{x}$  are:

$$\hat{\mathbf{x}} = \mathbf{D}\hat{\mathbf{z}}, \quad \hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{A}\mathbf{D}\mathbf{z} - \mathbf{y}\|_2^2 + \beta \|\mathbf{z}\|_1 \quad (6.1)$$

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta R(\mathbf{x}), \quad R(\mathbf{x}) = \quad (6.2)$$

A potential advantage of (6.2) is that  $\hat{\mathbf{x}}$  need not be in the range of  $D$ .

We could write the second one (6.2) as this joint optimization problem:

$$(\hat{\mathbf{x}}, \hat{\mathbf{z}}) = \arg \min_{\mathbf{x}, \mathbf{z}} \Psi(\mathbf{x}, \mathbf{z}), \quad \Psi(\mathbf{x}, \mathbf{z}) \triangleq \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \left( \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2 + \alpha \|\mathbf{z}\|_1 \right). \quad (6.3)$$

There are at least 3 distinct ways to pursue this joint optimization problem.

- Rewrite the cost function in LASSO form (with diagonally weighted 1-norm) in terms of  $\tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{z})$

$$\Psi(\tilde{\mathbf{x}}) = \frac{1}{2} \|\mathbf{B}\tilde{\mathbf{x}} - \mathbf{b}\|_2^2 + \beta\alpha \|\mathbf{W}\tilde{\mathbf{x}}\|_1, \quad \mathbf{B} \triangleq \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \sqrt{\beta}\mathbf{I} & -\sqrt{\beta}\mathbf{D} \end{bmatrix}, \quad \mathbf{b} \triangleq \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix}, \quad \mathbf{W} \triangleq \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}, \quad (6.4)$$

and then apply any **proximal method** from Ch. 5, e.g., POGM.

- If  $\mathbf{D}$  is unitary, then we can solve analytically for the minimizer over  $\mathbf{z}$  and substitute back in to get a cost function in terms of  $\mathbf{x}$  only, where the regularizer (in this case) involves a **Huber function** (as seen in HW):

$$R(\mathbf{x}) = \mathbf{1}' \psi(\mathbf{D}'\mathbf{x}; \alpha).$$

Then we can apply any gradient-based method (such as line search OGM) to that cost function.

However, this approach is inapplicable in the general case of interest where  $\mathbf{D}$  is wide (over-complete).

- Apply **alternating minimization** or **alternating descent**, aka **block coordinate minimization (BCM)** or **block coordinate descent (BCD)** to the “two block” cost function  $\Psi(\mathbf{x}, \mathbf{z})$ .

Here, both BCM and BCD start with some initial guesses  $\mathbf{x}_0$  and  $\mathbf{z}_0$  and then alternate updates. Reasonable initial guesses are application dependent, but one option is  $\mathbf{x}_0 = c\mathbf{A}'\mathbf{y}$  for some constant  $c$  and  $\mathbf{z}_0 = \mathbf{0}$ .

**BCM**

For this “two block” cost function, the **BCM** algorithm is:

```
for  $k = 0:niter-1$ 
```

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}, \mathbf{z}_k)$$

$$\mathbf{z}_{k+1} = \arg \min_{\mathbf{z}}$$

(6.5)

If implemented as written above, then this approach decreases the cost function monotonically:

$$\Psi(\mathbf{x}_{k+1}, \mathbf{z}_{k+1}) \leq$$

For convenience, the joint cost function (6.3) is repeated here:

$$\Psi(\mathbf{x}, \mathbf{z}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \left( \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2 + \alpha \|\mathbf{z}\|_1 \right).$$

To apply **BCM** (6.5) to this cost function, the updates required are:

regularized LS :  $\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}, \mathbf{z}_k) =$

**sparse coding** :  $\mathbf{z}_{k+1} = \arg \min_{\mathbf{z}} \Psi(\mathbf{x}_{k+1}, \mathbf{z}) =$

The **sparse coding** step above involves a **proximal operation**. (?)

A: True

B: False

??

For certain applications, the  $\mathbf{x}$  update above is easy:

- image denoising:  $\mathbf{A} = \mathbf{I}$
- single-coil MRI with Cartesian sampling:  $\mathbf{A}'\mathbf{A}$  is circulant
- image inpainting:  $\mathbf{A}'\mathbf{A}$  is diagonal
- image super-resolution:  $\mathbf{A}'\mathbf{A}$  is block diagonal with small blocks

In those cases, inverting  $\mathbf{A}'\mathbf{A} + \beta\mathbf{I}$  is easy, *i.e.*,  $O(N)$  or  $O(N \log N)$ .

But for general  $\mathbf{A}$  that inverse is  $O(N^3)$  to compute exactly and infeasible if  $\mathbf{A}$  is large, so one would have to apply an iterative method like CG for the  $\mathbf{x}$  update.

However, regardless of  $\mathbf{A}$ , the  $\mathbf{z}$  update above is a LASSO problem that requires an iterative solution in general, except for certain cases like when  $\mathbf{D}$  happens to be unitary. We could use POGM to solve that inner LASSO problem but then one might ask why not just apply POGM to the joint LASSO problem (6.4)? One answer is that (6.4) involves  $\mathbf{A}$  which can be very large and expensive, whereas the  $\mathbf{z}$  update above involves only  $\mathbf{D}$  which might be much faster.

BCD

---

For any finite number of inner iterations of CG for the  $\mathbf{x}$  update, or a proximal method like POGM for the  $\mathbf{z}$  update, the BCM algorithm above does not provide an exact minimization, so the name BCM would then seem inappropriate. When using a small number (perhaps just one) inner iteration for the  $\mathbf{x}$  and/or  $\mathbf{z}$  updates, a more appropriate term is **block coordinate descent (BCD)**, which, for a two-block problem, is:

```
for  $k = 0:niter$ 
```

Find  $\mathbf{x}_{k+1}$  s.t.  $\Psi(\mathbf{x}_{k+1}, \mathbf{z}_k) \leq \Psi(\mathbf{x}_k, \mathbf{z}_k)$

Find  $\mathbf{z}_{k+1}$  s.t.  $\Psi(\mathbf{x}_{k+1}, \mathbf{z}_{k+1}) \leq \Psi(\mathbf{x}_{k+1}, \mathbf{z}_k)$ . (6.6)

Clearly this algorithm monotonically decreases the cost function by design:  $\Psi(\mathbf{x}_{k+1}, \mathbf{z}_{k+1}) \leq \Psi(\mathbf{x}_k, \mathbf{z}_k)$ .

Applying BCD (6.6) to (6.3) by using one GD update for  $\mathbf{x}$  and one PGM update for  $\mathbf{z}$  leads to the following simple algorithm:

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \frac{\gamma}{\|\mathbf{A}\|_2^2 + \beta} (\mathbf{A}'(\mathbf{A}\mathbf{x}_k - \mathbf{y}) + \beta(\mathbf{x}_k - \mathbf{D}\mathbf{z}_k)) \\ \mathbf{z}_{k+1} &= \text{soft}\left(\mathbf{z}_k - \frac{1}{L}\mathbf{D}'(\mathbf{D}\mathbf{z}_k - \mathbf{x}_{k+1}); \frac{\alpha}{L}\right), \quad L = \end{aligned} \quad (6.7)$$



For any  $0 < \gamma < 2$ , the above algorithm is appropriately named BCD. (?)

A: True

B: False

??

Depending on the relative compute effort of working with  $\mathbf{A}$  and  $\mathbf{D}$ , one could apply multiple inner GD iterations and/or PGM iterations and it will still be a BCD algorithm (for appropriate step sizes as shown above). Alternatively one could use MM updates with diagonal majorizers to avoid computing the spectral norms.

If we use multiple iterations of FGM or OGM for the  $x$  update, and/or multiple iterations of FISTA or POGM for the  $z$  update in (6.7), then the resulting algorithm is appropriately named BCD. (?)

A: True

B: False

??

The distinction between BCM and BCD is typically disregarded in the literature, but these notes will strive to use the terms appropriately for clarity.

---

## Compressed sensing with analysis regularizer

Assuming that  $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\varepsilon}$  and  $\mathbf{T}\mathbf{x}$  is sparse, for some sparsifying transform matrix  $\mathbf{T}$ , the most natural formulation of analysis regularization is the following challenging optimization problem:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \|\mathbf{T}\mathbf{x}\|_1. \quad (6.8)$$

This is simple to solve only in special cases such as when the sparsifying transform  $\mathbf{T}$  is unitary. The literature is full of *approximate* solutions to (6.8), discussed next.

## Corner rounding

---

The non-differentiability of the 1-norm is the primary challenge of (6.8), so one “approximate” approach is simply to replace the 1-norm with a smooth convex approximation  $\psi$ :

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}), \quad \Psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \mathbf{1}' \psi.(\mathbf{T}\mathbf{x}). \quad (6.9)$$

When  $\psi$  is smooth, we can apply fast methods like OGM and CG to  $\Psi$  easily. A reasonable choice for  $\psi$  is the Fair potential function, or the hyperbola  $|z| \approx \sqrt{z^2 + \epsilon^2} - \epsilon$ , and these can approximate the 1-norm very closely by taking  $\delta$  or  $\epsilon$  small enough. However, as  $\epsilon$  decreases the global Lipschitz constant of  $\Psi(\mathbf{x})$  increases, leading to slow convergence of methods like OGM that depend on the global Lipschitz constant. In contrast, CG can still work well because its step size is adaptive and depends only on the curvature along the search direction.

## Variable splitting regularizer

---

Another option is to “split” the  $\mathbf{T}$  matrix from the 1-norm term using a penalty function as follows:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \beta R(\mathbf{x}), \quad R(\mathbf{x}) = \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{T}\mathbf{x} - \mathbf{z}\|_2^2 + \alpha \|\mathbf{z}\|_1. \quad (6.10)$$

As  $\alpha$  approaches 0, the solution more strongly enforces  $\mathbf{z} \approx \mathbf{T}\mathbf{z}$  and the solution more closely approximates the original formulation (6.8). As shown in a HW problem, this formulation is mathematically equivalent to the corner rounding approach (6.9) with  $\psi$  chosen as a Huber function.

## Balanced analysis/synthesis regularization

---

(Read)

Recall that for the synthesis sparsity model, where we assume  $\hat{\mathbf{x}} \approx \mathbf{D}\hat{\mathbf{z}}$ , a typical formulation is

$$\hat{\mathbf{x}} = \mathbf{D}\hat{\mathbf{z}}, \quad \hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{ADz} - \mathbf{y}\|_2^2 + \beta \|\mathbf{z}\|_1.$$

If  $\mathbf{D}$  is a **Parseval tight frame**, for which  $\mathbf{DD}' = \mathbf{I}$ , an alternative formulation that is called the **balanced model** [2–6] is:

$$\hat{\mathbf{x}} = \mathbf{D}\hat{\mathbf{z}}, \quad \hat{\mathbf{z}} = \arg \min_{\mathbf{z}} \frac{1}{2} \|\mathbf{ADz} - \mathbf{y}\|_2^2 + \beta \|\mathbf{z}\|_1 + \alpha \frac{1}{2} \|(\mathbf{I} - \mathbf{D}'\mathbf{D})\mathbf{z}\|_2^2. \quad (6.11)$$

This cost function is essentially equivalent to (5.1).

- When  $\alpha = 0$ , this reverts to the synthesis model.

- As  $\alpha \rightarrow \infty$ , the additional term enforces the constraint  $\mathbf{z} = \mathbf{D}'\mathbf{D}\mathbf{z} \implies \|\mathbf{z}\|_1 = \|\mathbf{D}'\mathbf{D}\mathbf{z}\|_1 = \|\mathbf{D}'\mathbf{x}\|_1$ , which is a particular type of analysis regularizer  $\|\mathbf{T}\mathbf{x}\|_1$ , where  $\mathbf{T} = \mathbf{D}'$ .
- For any finite  $\alpha$ , this formulation is yet another approximation to the analysis regularizer formulation, and even then only in the case of a Parseval tight frame.

Results in [6] for compressed sensing MRI using shift-invariant wavelets for  $\mathbf{D}$  showed no empirical performance advantages of the “balanced” approach over an analysis regularizer.

If  $\mathbf{D}$  above is unitary, then the “balanced” formulation (6.11) is equivalent to a usual synthesis (?) and/or analysis (?) formulations?

A: T,T

B: T,F

C: F,T

D: F,F

??

## Optimization strategies

We have at least three viable options for optimizing the variable split form (6.10).

- Rewrite as a **joint cost function** in terms of  $\tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{z})$  and apply any proximal method like **POGM** to that joint cost function:

$$\Psi(\tilde{\mathbf{x}}) = \frac{1}{2} \|\mathbf{B}\tilde{\mathbf{x}} - \mathbf{b}\|_2^2 + \beta\alpha \|\mathbf{W}\tilde{\mathbf{x}}\|_1, \quad \mathbf{b} \triangleq \begin{bmatrix} \mathbf{y} \\ 0 \end{bmatrix}, \quad \mathbf{W} \triangleq \begin{bmatrix} 0 & \mathbf{I} \end{bmatrix}, \quad \mathbf{B} \triangleq$$

- Minimize over  $\mathbf{z}$  and plug back in leading to a Huber function regularizer in  $\mathbf{x}$ . Then apply any gradient-based method like line-search **OGM** to optimize over  $\mathbf{x}$ .

Letting  $\psi$  denote that Huber function, so  $R(\mathbf{x}) = \mathbf{1}' \psi(\mathbf{T}\mathbf{x}; \alpha)$ , the gradient here is

$$\nabla \Psi(\mathbf{x}) = \mathbf{A}'(\mathbf{A}\mathbf{x} - \mathbf{y}) + \beta \mathbf{T}' \dot{\psi}(\mathbf{T}\mathbf{x}; \alpha),$$

for which there is no non-iterative solution if set to zero, unless  $\alpha = 0$ , so iterative method are required.

- Write a two-block cost function and apply **BCD** or **BCM** to it; here is that cost function:

$$\Psi(\mathbf{x}, \mathbf{z}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \left( \frac{1}{2} \|\mathbf{T}\mathbf{x} - \mathbf{z}\|_2^2 + \alpha \|\mathbf{z}\|_1 \right).$$

The  $\mathbf{z}$  update is “easy,” *i.e.*, non-iterative even for large problems. (?)

A: True

B: False

??

In general the  $\mathbf{z}$  update is

$$\mathbf{z}_{k+1} =$$

For the denoising case where  $\mathbf{A} = \mathbf{I}$ , the  $\mathbf{x}$  update is typically “easy,” *i.e.*, non-iterative even for large problems. (?)

A: True

B: False

??

In general the  $\mathbf{x}$  update is

$$\mathbf{x}_{k+1} =$$

One case where this update is easy is when both  $\mathbf{A}'\mathbf{A}$  and  $\mathbf{T}'\mathbf{T}$  are *circulant*, because then we can use FFT operations for the inverse.  $\mathbf{A}'\mathbf{A}$  is circulant, *e.g.*, for denoising, for single-coil Cartesian MRI, for deblurring with periodic boundary conditions.  $\mathbf{T}'\mathbf{T}$  circulant when  $\mathbf{T}$  is unitary and when  $\mathbf{T}$  corresponds to finite differences with periodic boundary conditions. Often the Hessian is Toeplitz, *i.e.*, approximately circulant, and we can use a circulant preconditioner for CG.

## Sparse coding revisited with multi-block BCM

Recall that the **sparse coding** step of **BCM** for synthesis-based regularization is a LASSO problem:

$$\hat{\mathbf{z}} = \arg \min_{\mathbf{z} \in \mathbb{F}^K} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2 + \alpha \|\mathbf{z}\|_1,$$

where  $\mathbf{D} \in \mathbb{F}^{N \times K}$ , for which there is no closed-form solution, so iterative methods like POGM are required.

Recall this inner minimization arose from a “two-block” joint cost function  $\Psi(\mathbf{x}, \mathbf{z})$  with corresponding two-way alternating minimization. Another approach is to think of the cost function as having  $K + 1$  blocks:  $\Psi(\mathbf{x}, z_1, \dots, z_K)$ , and to implement BCM or BCD by updating one of these (now much smaller) blocks at a time:

for  $t = 0:\text{niter}$  (outer loop over iteration)

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x} \in \mathbb{F}^N} \Psi\left(\mathbf{x}, z_1^{(t)}, \dots, z_K^{(t)}\right)$$

for  $k = 1:K$  (inner loop over coefficients)

$$z_k^{(t+1)} = \arg \min_{z_k \in \mathbb{F}} \Psi\left(\mathbf{x}^{(t+1)}, z_1^{(t+1)}, z_2^{(t+1)}, \dots, z_{k-1}^{(t+1)}, z_k, z_{k+1}^{(t)}, \dots, z_K^{(t)}\right).$$

As written, this algorithm is guaranteed to monotonically decrease  $\Psi$ . (?)

A: True

B: False

??

Note that  $\mathbf{D}\mathbf{z} = \sum_{j=1}^K \mathbf{d}_j z_j$  where  $\mathbf{d}_j = D[:, j]$ .

Now focus on updating one  $z_k$  coefficient by defining

$$f_k(z_k) = g(z_1, \dots, z_{k-1}, z_k, z_{k+1}, \dots, z_K) = \frac{1}{2} \left\| \mathbf{x} - \sum_{j \neq k} \mathbf{d}_j z_j - \mathbf{d}_k z_k \right\|_2^2 + \alpha \left( \sum_{j \neq k} |z_j| + |z_k| \right).$$

Define  $\mathbf{r} \triangleq \mathbf{x} - \sum_{k=1}^K \mathbf{d}_k z_k^{(t)}$  and  $\mathbf{r}_k \triangleq \mathbf{r} + \mathbf{d}_k z_k^{(t)}$ , then ignoring constants independent of  $z_k$ :

$$\begin{aligned} f_k(z_k) &= \frac{1}{2} \|\mathbf{r}_k - \mathbf{d}_k z_k\|_2^2 + \alpha |z_k| + c_1 \\ &= \frac{1}{2} \|\mathbf{r}_k\|_2^2 - \text{real}\{\mathbf{r}'_k \mathbf{d}_k z_k\} + \frac{1}{2} |z_k|^2 \|\mathbf{d}_k\|_2^2 + \alpha |z_k| \\ &= \end{aligned}$$

Thus the update for  $z_k$  is

$$z_k^{(t+1)} = \text{soft}\left(\tilde{\mathbf{d}}'_k \mathbf{r}_k, \alpha / \|\mathbf{d}_k\|_2^2\right).$$

For an efficient implementation we keep  $\mathbf{r}$  updated as a state vector, something like this:

$$\mathbf{r} = \mathbf{r}_k - \mathbf{d}_k z_k^{(t+1)}.$$

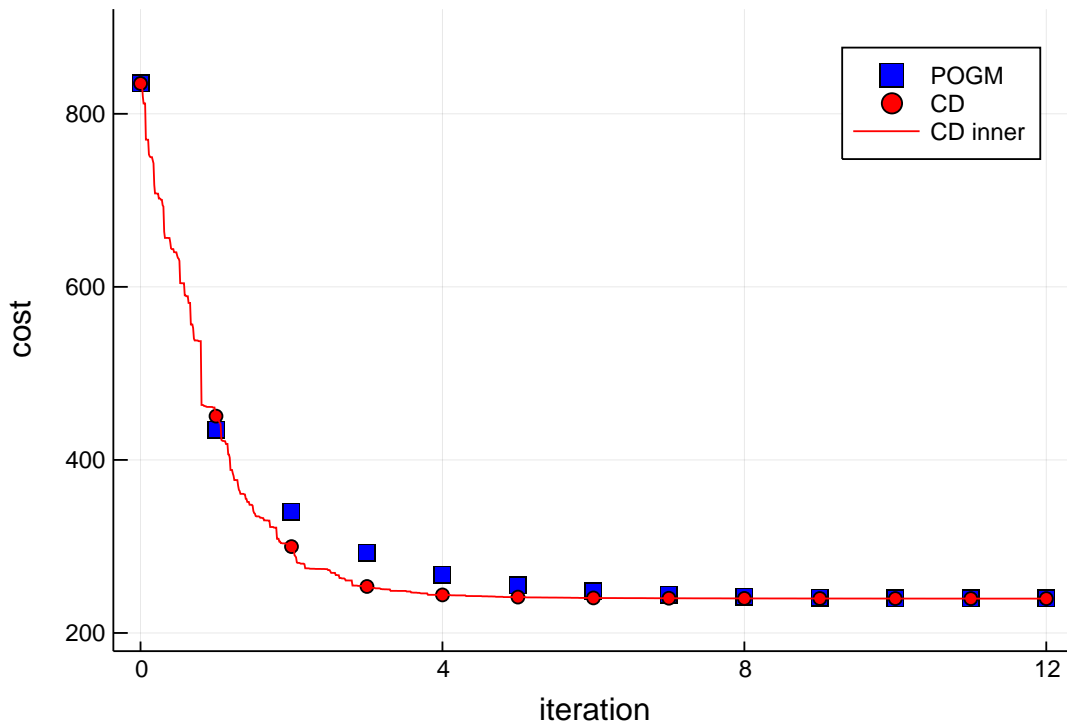


Here is JULIA code

```
# argmin_x 0.5*|A x - y|^2 + reg |x|_1
function sparse_code_cd(y, A, x0::AbstractVector{<:Number}, reg::Real)
    N = length(x0)
    vr = y - A * x0 # residual vector
    norm2 = [norm(A[:,n])^2 for n=1:N] # normalize
    D = hcat([A[:,n]/norm2[n] for n=1:N]...) # normalize
    x = copy(x0)

    for iter=1:niter # outer loop over iteration
        for n=1:N # inner loop over elements
            an = A[:,n]
            vr += an * x[n] # r_k
            x[n] = soft(D[:,n]'*vr, reg/norm2[n])
            vr -= an * x[n] # full residual again
        end
    end
    return x
end
```

Example. Same data as HW that compared POGM and CLS, where  $N = 50$ ,  $K = 99$  and 32/99 coefficients are zero. CD converges faster than POGM in terms of reducing cost per iteration. But wall time?



## CD approach to $\mathbf{x}$ update

---

Recall that the  $\mathbf{x}$  update in general required inverting a large matrix involving  $\mathbf{A}'\mathbf{A}$  for the synthesis form:

$$\Psi(x_1, \dots, x_N, \mathbf{z}) = \Psi(\mathbf{x}, \mathbf{z}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \left( \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2 + \alpha \|\mathbf{z}\|_1 \right).$$

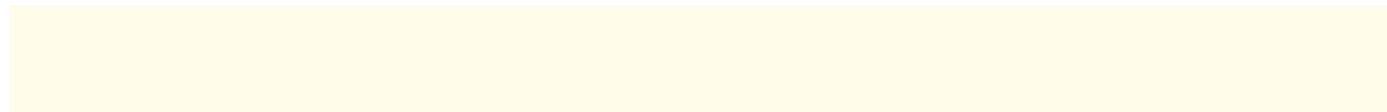
A multi-block approach, aka **coordinate descent**, can also avoid this matrix inverse. The update for  $x_n$  is:

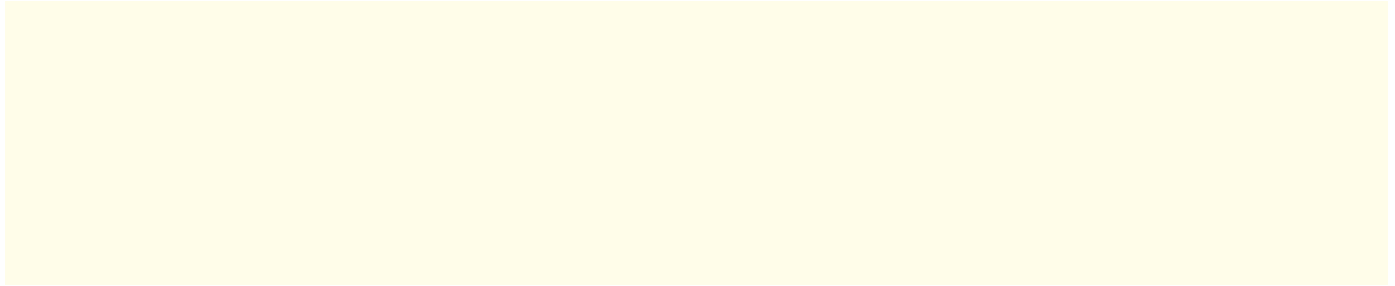
$$x_n^{(t+1)} = \arg \min_{x_n \in \mathbb{F}} \Psi \left( x_1^{(t+1)}, \dots, x_{n-1}^{(t+1)}, x_n, x_{n+1}^{(t)}, \dots, x_N^{(t)} \right) = \arg \min_{x_n \in \mathbb{F}} f_n(x_n)$$

$$f_n(x_n) = \frac{1}{2} \left\| \mathbf{A}\mathbf{x}^{(t,n)} + \mathbf{a}_n (x_n - x_n^{(t)}) - \mathbf{y} \right\|_2^2 + \beta \frac{1}{2} \left\| \mathbf{x}^{(t,n)} + \mathbf{a}_n (x_n - x_n^{(t)}) - \mathbf{D}\mathbf{z} \right\|_2^2,$$

where  $\mathbf{a}_n = \mathbf{A}[:, n]$  and  $\mathbf{x}^{(t,n)} = \left( x_1^{(t+1)}, \dots, x_{n-1}^{(t+1)}, x_n^{(t)}, x_{n+1}^{(t)}, \dots, x_N^{(t)} \right)$  contains all the most recent values.

In-class **group work** on  $\mathbf{x}$  update:





## Sparse coding for tight frames

The **sparse coding** problem is:

$$\hat{\mathbf{z}} = \arg \min_{\mathbf{z} \in \mathbb{F}^K} \frac{1}{2} \|\mathbf{x} - \mathbf{D}\mathbf{z}\|_2^2 + R(\mathbf{z}),$$

for some regularizer such as  $R(\mathbf{z}) = \alpha \|\mathbf{z}\|_1$  or  $R(\mathbf{z}) = \alpha \|\mathbf{z}\|_0$ .

So far we have discussed 2 ways to approach this optimization problem:

- proximal methods like **POGM** that update *all* coefficients  $\mathbf{z}$  simultaneously,
- multi-block **BCM** where we update *one* coefficient  $z_k$  at a time, sequentially.

These two options represent two extremes of parallel versus sequential; there are also “in-between” options.

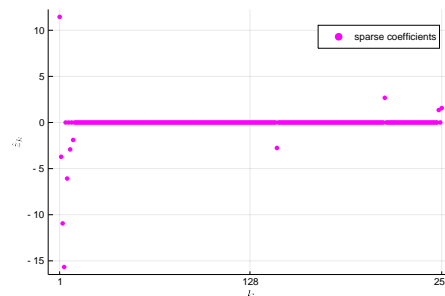
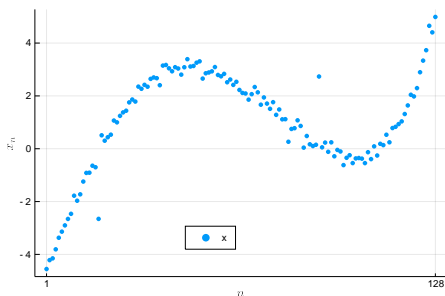
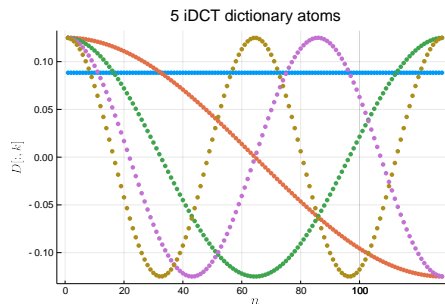
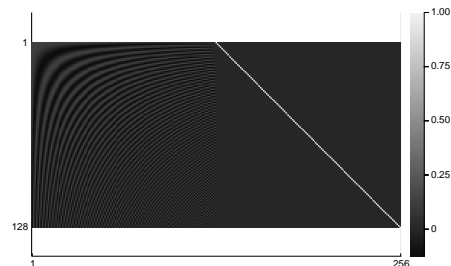
Consider the case where the dictionary  $\mathbf{D}$  is a **tight frame** consisting of two  $N \times N$  unitary matrices:  $\mathbf{D} = [\mathbf{D}_1 \ \mathbf{D}_2]$ . In the usual case where  $R$  is **additively separable**, we can rewrite the sparse coding problem as:

$$(\hat{\mathbf{z}}_1, \hat{\mathbf{z}}_2) = \arg \min_{\mathbf{z}_1, \mathbf{z}_2 \in \mathbb{F}^N}$$

There is still no joint closed-form solution here. But because  $\mathbf{D}_1$  and  $\mathbf{D}_2$  are unitary, it is very easy to perform two-block BCM where we alternate between updating  $\mathbf{z}_1$  and  $\mathbf{z}_2$ . The  $\mathbf{z}_1$  update is  $\text{prox}_R(\mathbf{D}_1'(\mathbf{x} - \mathbf{D}_2\mathbf{z}_2))$ .

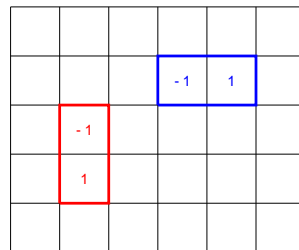
Example.

An **in-class task** will focus on a waves+spikes application for signals that are smooth + some impulses, where  $D_1$  is the (inverse) **DCT** matrix and  $D_2 = I$ , both of which are unitary matrices.



## Patch-based regularization: analysis form

Using **TV** regularizer  $R(\mathbf{x}) = \|\mathbf{T}\mathbf{x}\|_1$   
 where  $\mathbf{T}$  is 1st-order finite-differences  
 $\equiv$  **patches** of size  $2 \times 1$ .

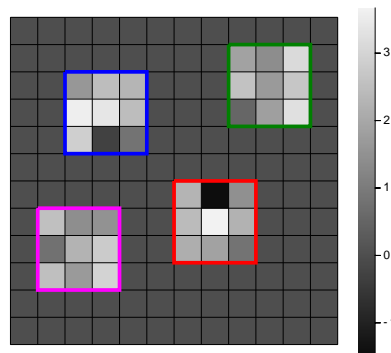


Larger patches provide more context  
 for distinguishing signal from noise.

*cf.* CNN approaches

Patch-based regularizers:

- **synthesis** models
- **analysis** methods



Especially for data-driven models, often it is more appropriate to analyze / regularize each **patch** of an image rather than trying to model the entire image.

For the model “ $\mathbf{T}\mathbf{R}_p\mathbf{x}$  tends to be sparse,” a typical patch-based analysis (or sparsifying transform) regularizer is:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{F}^N} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta R(\mathbf{x}), \quad R(\mathbf{x}) = \min_{\mathbf{Z} \in \mathbb{F}^{K \times P}} \sum_{p=1}^P \frac{1}{2} \|\mathbf{T}\mathbf{R}_p\mathbf{x} - \mathbf{z}_p\|_2^2 + \alpha \phi(\mathbf{z}_p), \quad (6.12)$$

where  $\mathbf{T}$  is a  $K \times d$  sparsifying transform matrix for **vectorized** patches and  $\mathbf{Z} = [\mathbf{z}_1 \dots \mathbf{z}_K]$ . Often  $K \approx d$ . Here  $\phi(\cdot)$  is a sparsity regularizer such as  $\|\cdot\|_0$  or  $\|\cdot\|_1$  or  $\|\mathbf{W}\cdot\|_0$  for some diagonal weighting matrix  $\mathbf{W}$  like we used with the wavelet transform.

**Example.** The most minimalist version would be  $d = 2$  and  $K = 1$  and  $\mathbf{T} = \begin{bmatrix} -1 & 1 \end{bmatrix}$ , which essentially ends up being very similar (but not identical) to a 1D **TV** regularizer. When we use (6.12), we are hoping to outperform methods like TV. (In 2D we would need both  $2 \times 1$  and  $1 \times 2$  patches.)

We write  $\mathbf{R}_p$  as a matrix above, but  $\mathbf{R}_p\mathbf{x}$  is yet another linear operation that we implement efficiently in code, not as matrix-vector multiplication.

If  $\mathbf{X}$  is a 2D image of size  $N_x \times N_y$  and  $\mathbf{x} = \text{vec}(\mathbf{X})$  and we want to use a patch size  $p_x \times p_y$ , for which  $N = N_x N_y$  and  $d = d_x d_y$ , the code for computing  $\mathbf{R}_1\mathbf{x}$  is `reshape(x, Nx, Ny) [1:px, 1:py]` and for  $\mathbf{R}_2\mathbf{x}$  is `reshape(x, Nx, Ny) [ (1:px)+1, 1:py]` etc. (See Ch. 2.)



In practice one could use something like MATLAB's `im2col` function to extract all the patches. (See next page for more memory efficient way.) There are many versions online for JULIA:

<https://discourse.julialang.org/t/what-is-julias-im2col/14066>

<https://github.com/pluskid/Mocha.jl/blob/master/benchmarks/native-im2col/im2col-bm.jl>

<https://github.com/outyang/MatlabFun.jl/blob/master/im2col.jl>

## BCD/BCM for patch-based analysis regularization

---

To perform the joint optimization problem (6.12), BCD/BCM are natural algorithm choices.

We alternate between updating the image  $x$  and updating the sparse coefficients  $z$ .

- The Hessian of (6.12) w.r.t.  $x$  is  $A'A + \beta D$  where  $D \triangleq \sum_{p=1}^P R_p' T' T R_p$ .

If that Hessian does not happen to have an easy inverse, what algorithm is the most natural choice for updating  $x$ ?

A: GD

B: (P)SD

C: (P)CG

D: OGM

E: POGM

??

If  $T$  is unitary, then  $D = \sum_{p=1}^P R_p' R_p$  is a  $N \times N$  diagonal, where the  $n$ th diagonal element is the number of patches that contain the  $n$ th pixel. If we choose patches with stride=1 and periodic boundary conditions, then that number is always  $d$ , the patch size, so  $D = dI$ . Otherwise it is at most  $d$ , and  $D \preceq dI$ .

So if  $A'A$  is also diagonal (e.g., denoising, inpainting, single-coil Cartesian MRI), then the  $x$  update is an exact minimization.

- The  $z_p$  updates are an **embarrassingly parallel** proximal operation:

$$z_p^{(t+1)} =$$

In JULIA this operation is simply: `Z = mapslices(prox, T*Xpatch, dims=1)`

where `Xpatch` =  $[R_1x \ \dots R_Px] \in \mathbb{F}^{d \times P}$  and `prox` is a `Function` that computes the proximity operator of  $\alpha\phi(\cdot)$ .

If  $\phi$  is additively separable, then the code simplifies further to `Z = prox.(T*Xpatch)`

## Practical implementation (Read)

As written, this BCD/BCM approach would be memory intensive because it stores  $Z = [z_1 \ \dots z_P]$ . Careful implementation can reduce the memory greatly, at least when  $T$  is unitary.

Consider the part of the regularizer in (6.12) involving  $x$  at iteration  $t$ :

$$f(x) \triangleq \sum_{p=1}^P \frac{1}{2} \|TR_p x - z_p^{(t)}\|_2^2 \implies \nabla f(x) = \sum_{p=1}^P R_p' T' (TR_p x - z_p^{(t)}) = D(x - \tilde{x}^{(t)})$$

$$\tilde{x}^{(t)} \triangleq D^{-1} \sum_{p=1}^P R_p' T' z_p^{(t)} = D^{-1} \sum_{p=1}^P R_p' T' \text{prox}_{\alpha\phi}(TR_p x^{(t-1)}) .$$

This summation form means that we can extract one (or several) of the patches from  $x^{(t-1)}$  at time, apply the

transform, threshold, and inverse transform, and put the result into an **accumulator**  $\tilde{\mathbf{x}}^{(t)}$  that is the same size as  $\mathbf{x}$ , and finally apply  $\mathbf{D}^{-1}$  (typically just  $1/d$ ). We never need to store all of  $\mathbf{Z}$ .

Clearly  $\nabla^2 f = \mathbf{D}$  so we can also write

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \tilde{\mathbf{x}}^{(t)}\|_{\mathbf{D}}^2 + c,$$

so the  $\mathbf{x}$  update is simply

$$\mathbf{x}^{(t+1)} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \frac{1}{2} \|\mathbf{x} - \tilde{\mathbf{x}}^{(t)}\|_{\mathbf{D}}^2 = (\mathbf{A}'\mathbf{A} + \beta\mathbf{D})^{-1}(\mathbf{A}'\mathbf{y} + \beta\mathbf{D}\tilde{\mathbf{x}}^{(t)}).$$

When  $\mathbf{D} = d\mathbf{I}$  this simplifies to

$$\mathbf{x}^{(t+1)} = (\mathbf{A}'\mathbf{A} + \beta d\mathbf{I})^{-1}(\mathbf{A}'\mathbf{y} + \beta d\tilde{\mathbf{x}}^{(t)}). \quad (6.13)$$

(Read)

Example. For single-coil Cartesian MRI, where  $\mathbf{A} = \mathbf{P}\mathbf{F}$  where  $\mathbf{F}^{-1} = \frac{1}{N}\mathbf{F}'$  and where  $\mathbf{P}$  is the  $K \times N$  sample selection matrix, for which  $\mathbf{P}'\mathbf{P}$  is diagonal, the update (6.13) further simplifies to

$$\begin{aligned}\mathbf{x}^{(t+1)} &= (\mathbf{F}'\mathbf{P}'\mathbf{P}\mathbf{F} + \beta d\mathbf{I})^{-1}(\mathbf{F}'\mathbf{P}'\mathbf{y} + \beta d\tilde{\mathbf{x}}^{(t)}) = (\mathbf{F}'\mathbf{P}'\mathbf{P}\mathbf{F} + \beta \frac{d}{N}\mathbf{F}'\mathbf{F})^{-1}(\mathbf{F}'\mathbf{P}'\mathbf{y} + \beta d\tilde{\mathbf{x}}^{(t)}) \\ &= \mathbf{F}^{-1}(\mathbf{P}'\mathbf{P} + \beta \frac{d}{N}\mathbf{I})^{-1}(\mathbf{F}')^{-1}(\mathbf{F}'\mathbf{P}'\mathbf{y} + \beta d\tilde{\mathbf{x}}^{(t)}) = \mathbf{F}^{-1}(\mathbf{P}'\mathbf{P} + \beta \frac{d}{N}\mathbf{I})^{-1}(\mathbf{P}'\mathbf{y} + \beta \frac{d}{N}\mathbf{F}\tilde{\mathbf{x}}^{(t)}) \\ &= \mathbf{F}^{-1}\mathbf{v}, \quad v_k = \begin{cases} \frac{1}{1+\beta d/N}([\mathbf{P}'/\mathbf{y}]_k + (\beta d/N)[\mathbf{F}\tilde{\mathbf{x}}^{(t)}]_k), & k \in \Omega \\ [\mathbf{F}\tilde{\mathbf{x}}^{(t)}]_k, & k \notin \Omega. \end{cases}\end{aligned}$$

From this expression, small  $\beta$  seems desirable, *i.e.*,  $\beta d/N \ll 1$ .

(Read)

If  $\mathbf{T}$  is not unitary, then in general the matrix  $\mathbf{D}$  above is not diagonal but we still have that

$$\mathbf{D}\mathbf{x} = \sum_{p=1}^P \mathbf{R}'_p \mathbf{T}' \mathbf{T} \mathbf{R}_p \mathbf{x}$$

and this summation could be done incrementally (one or several patches at a time) instead of extracting all patches at once, to save memory.

Nonlinear models for patches based on artificial **neural networks** are a recent trend [7].

## Sparsifying transform learning

So far we have considered a dictionary  $\mathbf{D}$  or a sparsifying transform  $\mathbf{T}$  to be “given,” but often we want to **learn**  $\mathbf{T}$  from training data. Given a set of training examples (typically image patches)  $\mathbf{X} \triangleq [\mathbf{x}_1 \ \dots \ \mathbf{x}_L] \in \mathbb{F}^{d \times L}$ , we want to find a transform  $\mathbf{T} \in \mathbb{F}^{K \times d}$  such that the transform coefficients  $\{\mathbf{z}_l = \mathbf{T}\mathbf{x}_l\}$  are typically **sparse**. This process is called **sparsifying transform learning**. Often  $K = d$  but we also consider other cases here. Let  $\mathbf{Z} = [\mathbf{z}_1 \ \dots \ \mathbf{z}_L] \in \mathbb{F}^{K \times L}$  denote the transform coefficient matrix. A typical transform learning optimization problem is [8, 9]:

$$\hat{\mathbf{T}} = \arg \min_{\mathbf{T} \in \mathcal{T}} \min_{\mathbf{Z} \in \mathbb{F}^{K \times L}} \Psi(\mathbf{T}, \mathbf{Z}), \quad \Psi(\mathbf{T}, \mathbf{Z}) \triangleq \quad (6.14)$$

where the “arg min” and “min” above are deliberately different and  $\phi$  is some sparsity regularizer like  $\|\cdot\|_0$ . An alternative formulation that looks simpler (no  $\alpha$  choice), but perhaps harder to optimize, is:

$$\hat{\mathbf{T}} = \arg \min_{\mathbf{T} \in \mathcal{T}} \Psi(\mathbf{T}), \quad \Psi(\mathbf{T}) = \sum_{l=1}^L \phi(\mathbf{T}\mathbf{x}_l).$$

For transform learning, we need to avoid a scale ambiguity *and* we would like to ensure that the rows of  $\mathbf{T}$  are not redundant. So one natural approach is to use the following (row) orthonormality constraint:

$$\mathcal{T} = \{\mathbf{T} \in \mathbb{F}^{K \times d} : \mathbf{T}\mathbf{T}' = \mathbf{I}_K\}. \quad (6.15)$$

With this choice of  $\mathcal{T}$  the problem (6.14) is always nonconvex, even if  $\phi$  is the convex 1-norm.

For the constraint (6.15) to hold, we must have (choose most general correct condition):

A:  $K \leq d$

B:  $K \geq d$

C:  $K = d$

D:  $K \neq d$

E:  $d, K \in \mathbb{N}$

??

Some authors consider tall  $\mathbf{T}$ , called an “over-complete” transform [10].

---

Example. The following matrix  $\mathbf{T} \in \mathcal{T}$  satisfies the constraint for  $K = 2$  and  $d = 5$ :

$$\mathbf{T}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 \end{bmatrix}.$$

Think of each row as a filter that we hope can sparsify patches extracted from images.

This example  $\mathbf{T}_2$  is in  $\mathcal{T}$ , but still has some redundancy in it from a filtering perspective. If we use this  $\mathbf{T}_2$  as part of a regularizer where we extract signal patches of size  $5 \times 1$  with a stride of 1 pixel (see Ch. 2), then we would get the same results using this simpler sparsifying transform

$$\mathbf{T}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \end{bmatrix},$$

with  $\beta$  adjusted by a factor of 2.

How to design  $\mathcal{T}$  to encourage less redundancy is an active research topic. See [11] for a Fourier approach.

## Two-block BCM for transform learning

There is no closed-form joint solution for  $\mathbf{T}$  and  $\mathbf{Z}$  in the transform learning problem (6.14), but its form suggests applying two-block **BCM**. Repeating (6.14) here for convenience:

$$\Psi(\mathbf{T}, \mathbf{Z}) \triangleq \sum_{l=1}^L \frac{1}{2} \|\mathbf{T} \mathbf{x}_l - \mathbf{z}_l\|_2^2 + \alpha \phi(\mathbf{z}_l), \text{ s.t. } \mathbf{T} \in \mathcal{T}.$$

- The  $\mathbf{Z}$  update is an **embarrassingly parallel** proximal operation:

$$\mathbf{z}_l^{(t+1)} =$$

In JULIA this operation is simply: `Z = mapslices(prox, T*X, dims=1)`

where `prox` is a `Function` that computes the proximity operator of  $\alpha\phi(\cdot)$ .

If  $\phi$  is additively separable, then the code simplifies further to `Z = prox.(T*X)`

- Now consider the  $\mathbf{T}$  update:

$$\hat{\mathbf{T}} = \arg \min_{\mathbf{T} \in \mathcal{T}} \sum_{l=1}^L \|\mathbf{T} \mathbf{x}_l - \mathbf{z}_l\|_2^2 = \quad (6.16)$$

Because of the  $\mathbf{X}$  multiplying  $\mathbf{T}$ , this is not a proximal operator.

It is *almost* a problem you have solved already!

HN 199

What is the name of this problem? ??

Why   ?? Why almost? ??

### Square transform learning: $T$ update

---

The solution for the  $T$  update (6.16) in the square case ( $K = d$ ) uses the following **SVD**:

$$\hat{T} = UV', \quad \underbrace{\mathbf{Z}}_{d \times L} \underbrace{\mathbf{X}'}_{L \times d} = \underbrace{\mathbf{U}}_{d \times d} \underbrace{\mathbf{\Sigma}}_{d \times d} \underbrace{\mathbf{V}'}_{d \times d}. \quad (6.17)$$

As a sanity check:  $\hat{T}\hat{T}' = UV'VU' = I_d$ .

Does the use of an SVD here preclude large-scale problems?

Typically  $K = d \ll L$ , e.g.,  $K = d = 8^2$  for  $8 \times 8$  patches whereas  $L$  can greatly exceed  $10^6$ .

The most expensive part is the matrix multiplication  $\mathbf{Z}\mathbf{X}'$  that is  $O(d^2L)$ , whereas the  $d \times d$  SVD is  $O(d^3)$ .

For  $8 \times 8$  patches doing a  $8^2 \times 8^2$  SVD is trivial.

---

Here is a summary of the BCM algorithm for (square) transform learning with  $K = d$ :

- Apply the **proximal operator** (e.g., soft or hard thresholding) to each column of  $\mathbf{T}\mathbf{X}$  to get new  $\mathbf{Z}$ .
- Apply orthogonal **Procrustes method** using SVD of the product  $\mathbf{Z}\mathbf{X}'$  to get new  $\mathbf{T}$ .
- Repeat until convergence.

See [8] for some convergence theory for this alternating method, even for  $\phi(\mathbf{z}) = \|\mathbf{z}\|_0$ .

This method for updating  $T, Z$  should be named **BCD** instead of **BCM**. (?)

A: True

B: False

??



## Non-square transform learning

---

In the non-square case, specifically where  $K < d$ , the generalized orthogonal Procrustes solution from EECS 551 and HW1#9 (stiefl) is inapplicable. Those solutions were for the case where  $\mathbf{T}$  is tall ( $K \geq d$  here), with the constraint  $\mathbf{T}'\mathbf{T} = \mathbf{I}_d$ .

The derivation there does not generalize readily to handle  $\mathbf{T}\mathbf{T}' = \mathbf{I}_K$ . To synopsise the issue, consider:

$$\begin{aligned}\|\mathbf{T}\mathbf{X} - \mathbf{Z}\|_{\text{F}}^2 &= \text{trace}\{(\mathbf{T}\mathbf{X} - \mathbf{Z})'(\mathbf{T}\mathbf{X} - \mathbf{Z})\} \\ &= \text{trace}\{\mathbf{X}'\mathbf{T}'\mathbf{T}\mathbf{X}\} - 2\text{real}\{\text{trace}\{\mathbf{T}\mathbf{X}\mathbf{Z}'\}\} + \text{trace}\{\mathbf{Z}'\mathbf{Z}\}.\end{aligned}$$

The previous solution used the constraint  $\mathbf{T}'\mathbf{T} = \mathbf{I}_d$  to simplify the first term. Here we have the constraint  $\mathbf{T}\mathbf{T}' = \mathbf{I}_K$  that does not help simplify the first term in general.

In the square case,  $\mathbf{T}'\mathbf{T} = \mathbf{I} \iff \mathbf{T}\mathbf{T}' = \mathbf{I}$ , so the previous solution applies, but not more generally.

One possible approach is to use the trace circular commutative property to write the first term as

$$\begin{aligned}\text{trace}\{\mathbf{X}'\mathbf{T}'\mathbf{T}\mathbf{X}\} &= \text{trace}\{\mathbf{T}\mathbf{X}\mathbf{X}'\mathbf{T}'\} \leq \text{trace}\{\mathbf{T}\mathbf{\Pi}\mathbf{T}'\} + \text{trace}\{(\mathbf{T} - \mathbf{T}_k)(\rho\mathbf{I}_d - \mathbf{\Pi})(\mathbf{T} - \mathbf{T}_k)'\} \\ &= -2\text{real}\{\text{trace}\{\mathbf{T}(\rho\mathbf{I}_d - \mathbf{\Pi})\mathbf{T}_k'\}\} + \text{trace}\{\mathbf{T}_k(\rho\mathbf{I}_d - \mathbf{\Pi})\mathbf{T}_k'\} + \underbrace{\rho\text{trace}\{\mathbf{T}\mathbf{T}'\}}_K,\end{aligned}$$

where  $\mathbf{\Pi} \triangleq \mathbf{X}\mathbf{X}'$  is the  $d \times d$  data covariance matrix and  $\rho$  is its spectral radius and  $\mathbf{T}_k$  is the current transform estimate. This inequality can be the basis for a (possibly novel) **MM** approach.

Here is an alternative approach that seems simpler.

Suppose we want to learn  $K < d$  filters. We can still estimate a  $d \times d$  matrix  $\mathbf{T}$ , but ignore the last  $d - K$  “dummy” rows of the matrix. To truly ignore those rows, we need to ignore the corresponding rows of  $\mathbf{Z}$  that correspond to the products of the dummy rows of  $\mathbf{T}[(K + 1) : d, :]$  with  $\mathbf{X}$ , *i.e.*, we want the first  $K$  rows of  $\mathbf{Z}$  to be sparse, but we do not care about the remaining rows. Thus, we define the following (possibly novel) cost function

$$\Psi(\mathbf{T}, \mathbf{Z}) = \frac{1}{2} \|\mathbf{T}\mathbf{X} - \mathbf{Z}\|_{\text{F}}^2 + \alpha \quad (6.18)$$

(The idea here is somewhat similar to the **HW** problem where we apply sparsity regularization to the wavelet detail coefficients only.)

With this approach, the  $\mathbf{Z}$  update applies hard thresholding to the first  $K < d$  rows of  $\mathbf{T}\mathbf{X}$  and leaves untouched the remaining rows: `Z[1:K, :] = hard.(T[1:K, :] * X)`

Then the  $\mathbf{T}$  update is simply the standard orthogonal Procrustes solution in (6.17).

The potential advantage of the MM approach is that the matrix  $\mathbf{Z}$  requires storing only a  $K \times L$  matrix whereas the dummy rows approach appears to require storing a  $d \times L$  matrix. If  $d = 8^3$  for a 3D imaging problem and we are content learning, say, 32, filters, that is a 16-fold difference in storage that could be significant. The advantage of the dummy rows approach is that we can easily reuse the `stief1` code instead of writing a new MM algorithm.

Convergence properties of both approaches would require investigation, but most likely the proofs in [8] could be adapted.

## BCM for non-square transform learning

---

Here is a summary of a **BCM** algorithm for transform learning with  $K < d$ .

- Apply the **proximal operator** (e.g., soft or hard thresholding) to the first  $K$  rows of  $\mathbf{T}\mathbf{X}$  to update  $\mathbf{Z}$ .
  - Apply the **orthogonal Procrustes method** using SVD of the product  $\mathbf{Z}\mathbf{X}'$  to update  $\mathbf{T}$ .
  - Repeat until convergence, then perhaps keep just  $\hat{\mathbf{T}}[1 : K, :]$
- 

Why perhaps? Because to use  $\hat{\mathbf{T}}$  as a patch-based regularizer, we might want to use the same trick:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta R(\mathbf{x}), \quad R(\mathbf{x}) = \min_{\mathbf{z} \in \mathbb{R}^{d \times L}} \sum_{l=1}^L \frac{1}{2} \left\| \hat{\mathbf{T}} \mathbf{R}_l \mathbf{x} - \mathbf{z}_l \right\|_2^2 + \alpha \phi(\mathbf{W} \mathbf{z}_l), \quad (6.19)$$

where  $\mathbf{W} = \text{Diag}\{\mathbf{w}\}$ ,  $\mathbf{w} = \begin{bmatrix} \mathbf{1}_K \\ \mathbf{0}_{d-K} \end{bmatrix}$ .

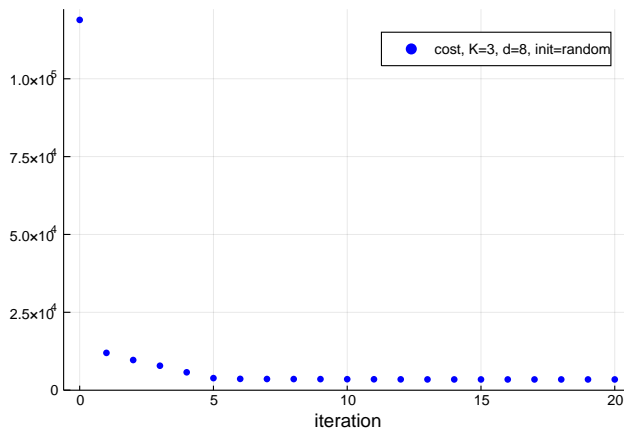
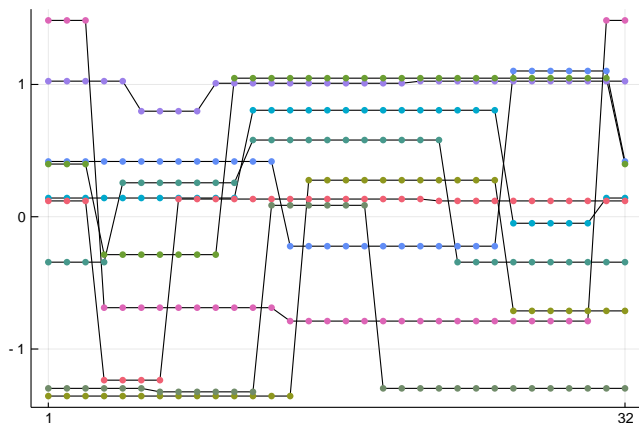
This way  $\hat{\mathbf{T}}' \hat{\mathbf{T}} = \mathbf{I}$ , simplifying the update, while we essentially ignore the last  $d - K$  rows of  $\mathbf{Z}$  by not enforcing sparsity there. To be explicit, for the above  $\mathbf{W}$  we have:

$$R(\mathbf{x}) = \min_{\mathbf{z} \in \mathbb{R}^{d \times L}} \sum_{l=1}^L \frac{1}{2} \left\| \hat{\mathbf{T}} \mathbf{R}_l \mathbf{x} - \mathbf{z}_l \right\|_2^2 + \alpha \phi(\mathbf{W} \mathbf{z}_l) = \min_{\mathbf{z} \in \mathbb{R}^{K \times L}} \sum_{l=1}^L \frac{1}{2} \left\| \hat{\mathbf{T}}[1 : K, :] \mathbf{R}_l \mathbf{x} - \mathbf{z}_l \right\|_2^2 + \alpha \phi(\mathbf{z}_l).$$

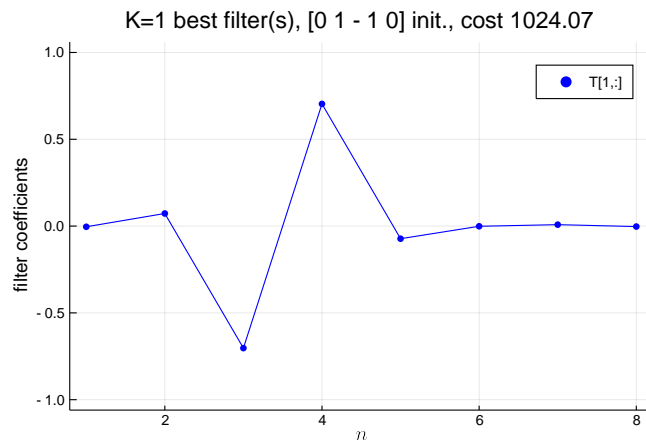
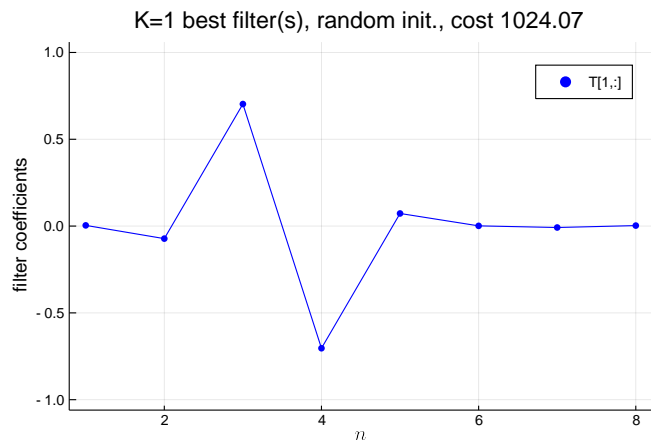
This dummy-row trick may be novel. (If you see it in the literature please let me know.)

Example. Here we consider an ensemble of  $2^{11}$  1D piecewise-constant signals of length  $N = 32$ , several of which are illustrated in the left figure below.

From these signals I extracted all  $5 \cdot 10^4$  patches of size  $d = 8 \times 1$ , and then discarded all patches that are completely uniform because they seem to contain little useful information. There were about  $L = 3 \cdot 10^4$  “interesting” patches for training. The right figure shows the cost function (6.18) decreasing with each BCM update, for the case  $K = 1$ .



Here are learned filters  $\hat{\mathbf{T}}$  for  $K = 1$ . For the left figure  $\mathbf{T}_0$  was randomly initialized, and for the right figure the first row of  $\mathbf{T}_0$  was  $[-1 \ 1 \ 0 \ \dots \ 0]$ . The consistency of the shape (up to a sign flip) is interesting.

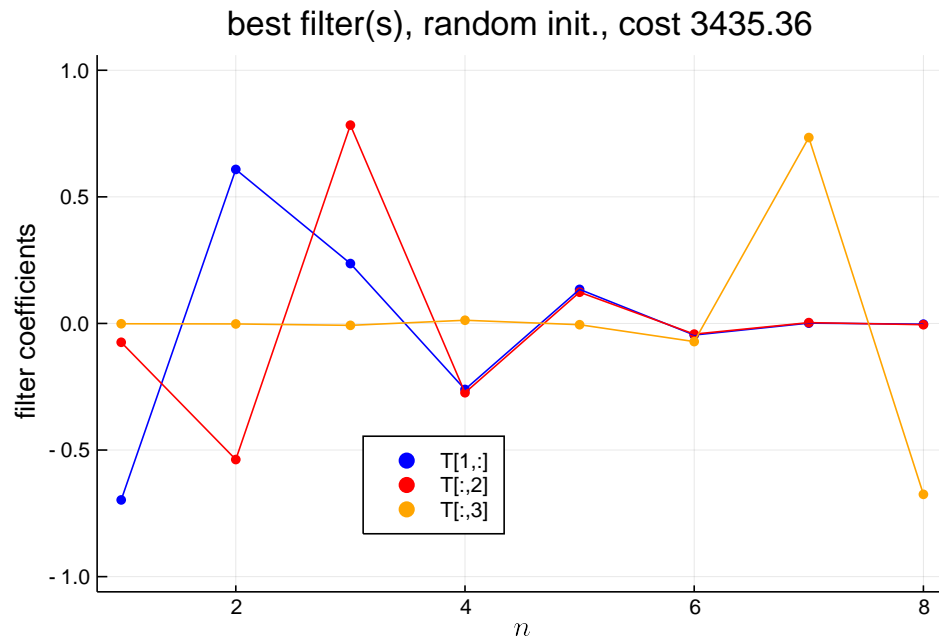


The filter values  $\hat{\mathbf{T}}$  are  $[-0.0 \ -0.05 \ 0.7 \ -0.71 \ 0.07 \ 0.0 \ -0.01 \ 0.0]$

I used  $\alpha = 0.4$ , for which about 5% of the  $\mathbf{Z}[1, :]$  values were nonzero.

Because all the training signals are piecewise-constant (by design), it is unsurprising that the learned filter is close to a finite difference filter. But it is not exactly  $[1 \ -1]$ , so it would be interesting to compare denoising using a regularizer based on this learned filter to that based on TV. Increasing  $L$  did not change the filter estimates.

Here is the case where we learn  $K = 3$  filters:



## Memory efficient implementation of transform learning (Read)

The BCM algorithm for **transform learning** on p. 6.32 is elegant in its simplicity, but as written it is memory inefficient because it must store both the  $d \times L$  matrix  $\mathbf{X}$  and the  $K \times L$  matrix  $\mathbf{Z}$ , where  $L$  can be enormous. Furthermore, typically we form the columns of  $\mathbf{X}$  from overlapping patches from training images, and the patch extraction process (with stride=1) increases the memory of an image with  $N$  pixels to a set of patches with  $dN$  elements, so by a factor of  $d$ . Careful implementation can avoid these memory issues. First we rewrite the two BCM steps:

$$\begin{aligned} z_l^{(t+1)} &= \text{prox}_{\alpha\phi}(\mathbf{T}^{(t)} \mathbf{x}_l), \quad l = 1, \dots, L \\ \mathbf{T}^{(t+1)} &= \arg \min_{\mathbf{T} \in \mathcal{T}} \left\| \mathbf{T} \mathbf{X} - \mathbf{Z}^{(t+1)} \right\|_{\text{F}}^2 = \mathbf{U} \mathbf{V}', \quad \mathbf{Z}^{(t+1)} \mathbf{X}' = \mathbf{U} \mathbf{\Sigma} \mathbf{V}'. \end{aligned}$$

The key is to implement the matrix product  $\mathbf{Z}^{(t+1)} \mathbf{X}'$  with an **accumulator**:

$$\mathbf{Z}^{(t+1)} \mathbf{X}' = \sum_{l=1}^L z_l^{(t+1)} \mathbf{x}'_l = \sum_{l=1}^L \text{prox}_{\alpha\phi}(\mathbf{T}^{(t)} \mathbf{x}_l) \mathbf{x}'_l. \quad (6.20)$$

In this form, we never need to store all of the  $\{z_l\}$  coefficients. Furthermore, if we are extracting each patch  $\mathbf{x}_l$  from a set of training images, then we never need to store the individual patches; we just loop over each training image, then loop over each patch in those images; we extract one patch, multiply it by  $\mathbf{T}^{(t)}$ , apply the proximal operator, make the outer product with the patch, and accumulate using `+=`.

One small drawback of using (6.20) is that without all of  $\mathbf{X}$  and  $\mathbf{Z}$  available, one cannot compute the cost function  $\Psi(\mathbf{T}, \mathbf{Z})$  for any  $\mathbf{T}$  and  $\mathbf{Z}$ . However, if  $\mathbf{T}$  is unitary, then the key term in cost is

$$\begin{aligned} f(\mathbf{T}, \mathbf{Z}) &= -2 \operatorname{real}\{\operatorname{trace}\{\mathbf{Z}(\mathbf{T}\mathbf{X})'\}\} + \operatorname{trace}\{\mathbf{Z}\mathbf{Z}'\} \\ &= -2 \operatorname{real}\left\{\operatorname{trace}\left\{\sum_{l=1}^L \mathbf{z}_l(\mathbf{T}\mathbf{x}_l)'\right\}\right\} + \operatorname{trace}\left\{\sum_{l=1}^L \mathbf{z}_l\mathbf{z}_l'\right\} \\ &= -2 \operatorname{real}\left\{\sum_{l=1}^L \langle \mathbf{z}_l, \mathbf{T}\mathbf{x}_l \rangle\right\} + \sum_{l=1}^L \|\mathbf{z}_l\|_2^2 = \sum_{l=1}^L \|\mathbf{z}_l - \mathbf{T}\mathbf{x}_l\|_2^2 + c. \end{aligned}$$

If needed, we can use this expression to evaluate  $f(\mathbf{T}^{(t)}, \mathbf{Z}^{(t+1)})$  while computing (6.20).

For an extension to multi-layer transform learning see [12].



## Dictionary learning via two-block BCD

Problem statement: Given training data  $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_L] \in \mathbb{F}^{N \times L}$ , typically patches extracted from images, find a **dictionary**  $\mathbf{D} \in \mathbb{F}^{d \times K}$  such that  $\mathbf{x}_l \approx \mathbf{D} \mathbf{z}_l$  where  $\mathbf{z}_l \in \mathbb{F}^K$  is (typically) a **sparse** coefficient vector. Optimization formulation:

$$\hat{\mathbf{D}} = \arg \min_{\mathbf{D} \in \mathcal{D}} \min_{\mathbf{Z} \in \mathbb{F}^{K \times L}} \Psi(\mathbf{D}, \mathbf{Z}), \quad \Psi(\mathbf{D}, \mathbf{Z}) \triangleq$$

where  $\|\mathbf{Z}\|_0 = \sum_{l=1}^L \|\mathbf{z}_l\|_0$  is the *aggregate* non-sparsity, *i.e.*, the total number of nonzero coefficients. This cost function  $\Psi$  is nonconvex due to the product  $\mathbf{D}\mathbf{Z}$ .

(It would be **biconvex** if we used  $\|\text{vec}(\mathbf{Z})\|_1$ .)

To avoid the scale ambiguity, we focus on this typical choice for the set admissible dictionaries:

$$\mathcal{D} = \{ \mathbf{D} \in \mathbb{F}^{d \times K} : \|\mathbf{d}_k\|_2 = 1, \ k = 1, \dots, K \}.$$

To minimize  $\Psi$ , one could attempt two-block **BCD**, alternating between updating  $\mathbf{D}$  and  $\mathbf{Z}$ .

- The  $\mathbf{D}$  update is a (nonconvex) constrained problem and one could apply **gradient projection** for it (see **HW**). If we replace  $\|\mathbf{d}_k\|_2 = 1$  with  $\|\mathbf{d}_k\|_2 \leq 1$  then the  $\mathbf{D}$  update would be a convex constrained problem.
- The  $\mathbf{Z}$  update is a set of  $L$  separate sparse coding problems, If we used the 1-norm, then it would be convex and one could apply **POGM** in parallel to each column of  $\mathbf{Z}$ . This is a simple use of parallel

processing. For the 0-norm, there is no convergence guarantee of momentum methods like POGM (to my knowledge), so it seems safer to use **PGM**, corresponding to iterative hard thresholding:

$$\mathbf{z}_l^{(t+1)} = \text{hard} \left( \mathbf{z}_l^{(t)} - \frac{1}{\|D\|_2^2} D' (D \mathbf{z}_l - \mathbf{x}_l), \frac{\beta}{\|D\|_2^2} \right), \quad l = 1, \dots, L.$$

Both of these updates require 1 or more inner iterations, so overall it is a **BCD** approach.

### Dictionary learning via multi-block BCM \_\_\_\_\_ (**SOUP-DIL**)

Instead of updating the entire dictionary  $D$  and the entire coefficient vector  $Z$  simultaneously, we can instead think of the multi-block cost function  $\Psi(\mathbf{d}_1, \dots, \mathbf{d}_K, \mathbf{c}_1, \dots, \mathbf{c}_K)$  where  $C = Z' = [\mathbf{c}_1 \dots \mathbf{c}_K] \in \mathbb{F}^{L \times K}$ , and write the product in the following sum-of-outer-products (**SOUP**) form:

$$[D \mathbf{z}_1 \quad \dots \quad D \mathbf{z}_L] = DZ = DC' =$$

With this formulation, a useful multi-block **BCM** approach is to update one atom  $\mathbf{d}_k$  at a time, and then update the corresponding coefficient vector  $\mathbf{c}_k$ , and loop sequentially through  $k = 1, \dots, K$ .

Updating the variables in matched pairs, *e.g.*, in the order:  $\mathbf{d}_1, \mathbf{c}_1, \mathbf{d}_2, \mathbf{c}_2, \dots, \mathbf{d}_K, \mathbf{c}_K$ , seems to accelerate convergence. The next pages describe those updates.

**Dictionary atom update** 

---

To update  $\mathbf{d}_k$ , define the residual matrix

$$\mathbf{R} = \mathbf{X} - \sum_{j \neq k} \mathbf{d}_j \mathbf{c}'_j \quad (6.21)$$

then the update is

$$\mathbf{d}_k^{(t+1)} = \arg \min_{\mathbf{d}_k \in \mathbb{R}^d : \|\mathbf{d}_k\|_2=1} \frac{1}{2} \|\mathbf{R} - \mathbf{d}_k \mathbf{c}'_k\|_F^2$$

where



**Sparse coefficient update**

The update for  $\mathbf{c}_k$  is

group work:

$$\mathbf{c}_k^{(t+1)} = \arg \min_{\mathbf{c}_k \in \mathbb{R}^L} \frac{1}{2} \|\mathbf{R} - \mathbf{d}_k \mathbf{c}_k'\|_{\text{F}}^2 + \beta \|\mathbf{c}_k\|_0.$$

Because  $\|\mathbf{d}_k\|_2 = 1$ :

$$\|\mathbf{R} - \mathbf{d}_k \mathbf{c}_k'\|_{\text{F}}^2 = \text{trace}\{(\mathbf{R} - \mathbf{d}_k \mathbf{c}_k')(\mathbf{R} - \mathbf{d}_k \mathbf{c}_k')'\} = \|\mathbf{c}_k\|_2^2 - 2 \text{real}\{\mathbf{d}_k' \mathbf{R} \mathbf{c}_k\} + c_1 = \|\mathbf{c}_k - \mathbf{R}' \mathbf{d}_k\|_2^2 + c_2$$

$$\mathbf{c}_k^{(t+1)} = \text{hard.}(\mathbf{R}' \mathbf{d}_k, \beta).$$

The main drawback of this SOUP approach is that it is less parallelizable than the two-block BCD approach.

This SOUP alternating approach is:

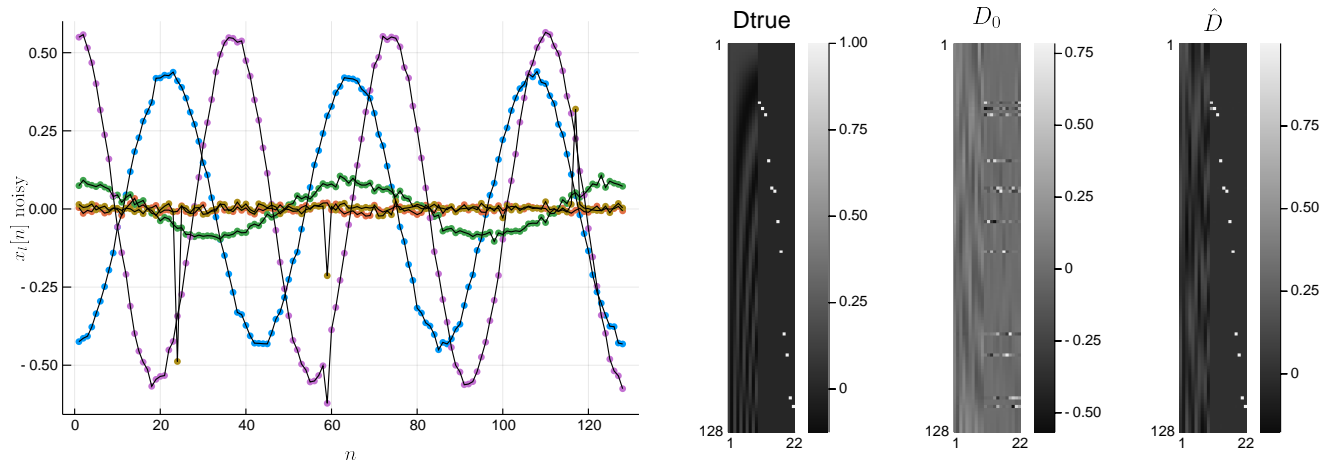
A: BCD

B: BCM

C: Neither

??

Example. An in-class **group activity** will apply this SOUP **dictionary learning** method to signals consisting of a smooth part and some added spikes. The left plot shows some example signals. The right plot shows the true dictionary, the initial dictionary estimate  $D_0$  (from PCA) and the estimated dictionary  $\hat{D}$ .



## Joint update of atom and coefficients

(Read)

Previously we focused on alternating between updating one atom  $\mathbf{d}_k$  and the corresponding coefficient vector  $\mathbf{c}_k$ , ala SOUP. Can we update them concurrently? Starting with the residual (6.21) and dropping subscript  $k$  for simplicity, we want to solve (or descend)

$$\arg \min_{\mathbf{c} \in \mathbb{F}^L, \mathbf{d} \in \mathbb{F}^d: \|\mathbf{d}\|_2 \leq 1} f(\mathbf{d}, \mathbf{c}), \quad f(\mathbf{d}, \mathbf{c}) \triangleq \frac{1}{2} \|\mathbf{d}\mathbf{c}' - \mathbf{R}\|_{\text{F}}^2 \stackrel{\text{c}}{=} \frac{1}{2} \|\mathbf{d}\|_2^2 \|\mathbf{c}\|_2^2 - \mathbf{d}' \mathbf{R} \mathbf{c}.$$

The Hessian of this cost function is

$$\nabla^2 f = \begin{bmatrix} \|\mathbf{c}\|^2 \mathbf{I}_d & 2\mathbf{d}\mathbf{c}' - \mathbf{R} \\ 2\mathbf{c}\mathbf{d}' - \mathbf{R}' & \|\mathbf{d}\|^2 \mathbf{I}_L \end{bmatrix}. \quad (6.22)$$

Consider the scalar case ( $d = L = 1$ ) and suppose  $R = 0$ . Then  $\nabla^2 f = \begin{bmatrix} c^2 & 2dc \\ 2cd & d^2 \end{bmatrix}$  which has eigenvalues  $\lambda = (d^2 + c^2 \pm \sqrt{(d^2 - c^2)^2 + 16d^2c^2})/2$ , so  $\rho(\nabla^2 f)$  is unbounded as  $c$  increases. Thus the claim in [13] that  $\nabla f$  is jointly global Lipschitz is incorrect in general.

However, in practical problems it seems reasonable to assume (or impose) that the sparse coefficients are bounded by some finite maximum value:  $\|\mathbf{c}\|_\infty \leq \bar{z}$ , cf. [13, 14]. In the scalar case this leads to a simple upper bound on the spectral radius of the Hessian.

Challenge: using  $\|\mathbf{d}\|_2 \leq 1$  and  $\|\mathbf{c}\|_\infty \leq \bar{z}$ , find an upper bound on the spectral radius of the Hessian in (6.22). It is fine to assume  $1 \leq \bar{z}$ . Perhaps using  $\|\cdot\|_1$  could simplify.

---

## 6.2 Machine learning applications

---

### Low-rank approximation for large-scale problems

Given  $\mathbf{Y} = \mathbf{X} + \varepsilon \in \mathbb{F}^{M \times N}$  where  $\varepsilon$  denotes a  $M \times N$  noise matrix and we believe  $\text{rank}(\mathbf{X}) \leq K$ .

EECS 551 used truncated/thresholded SVD methods requiring  $O(MN^2)$  operations that do not scale to large problems where both  $M$  and  $N$  are large, even if the rank  $K$  is very small. We can overcome this problem using **alternating minimization** (two-block BCD/BCM) methods.

### Matrix factorization approach

---

To overcome this limitation of SVD-based approaches, one can take a **matrix factorization** approach by choosing a desired (maximum) rank  $K$  and expressing  $\hat{\mathbf{X}}$  directly as

$$\hat{\mathbf{X}} = \underbrace{\mathbf{U}}_{M \times K} \underbrace{\mathbf{V}}_{K \times N},$$

where now  $\mathbf{U}$  and  $\mathbf{V}$  should be simply interpreted as factors, not as matrices with singular vectors.

With this formulation, a typical optimization problem for finding  $U$  and  $V$ , and hence  $\hat{X}$ , looks like:

$$\begin{aligned}\hat{X} &= \hat{U}\hat{V} \\ (\hat{U}, \hat{V}) &= \arg \min_{U \in \mathbb{R}^{M \times K}, V \in \mathbb{R}^{K \times N}} \Psi(U, V) \\ \Psi(U, V) &\triangleq \frac{1}{2} \|Y - UV\|_F^2 + \beta_1 R_1(U) + \beta_2 R_2(V),\end{aligned}$$

where one must select appropriate regularizers for  $U$  and  $V$ . The data term here is **biconvex**.

## Ambiguities

---

Scale ambiguity:

Factorization ambiguity for invertible  $P$ :

- These ambiguities might not matter if we just want the final  $\hat{X}$ .
- They might matter for optimization where if the iterates do not converge.

Constraining  $U$  to have orthonormal columns resolves the above scale (?) and factorization (?) ambiguity.

A: F,F

B: F,T

C: T,F

D: T,T

??



## Two-block BCM with unitary constraint

---

One way to resolve the ambiguity is to constrain  $\mathbf{U}$  to have orthonormal columns:

$$R_1(\mathbf{U}) = \chi_{\mathcal{V}_K(\mathbb{F}^M)}(\mathbf{U}),$$

where the **Stiefel manifold** is  $\mathcal{V}_K(\mathbb{F}^M) = \{\mathbf{Q} \in \mathbb{F}^{M \times K} : \mathbf{Q}'\mathbf{Q} = \mathbf{I}_K\}$ . The optimization problem becomes:

$$\hat{\mathbf{X}} = \hat{\mathbf{U}}\hat{\mathbf{V}}, \quad (\hat{\mathbf{U}}, \hat{\mathbf{V}}) = \arg \min_{\mathbf{U} \in \mathbb{F}^{M \times K}, \mathbf{V} \in \mathbb{F}^{K \times N}} \Psi(\mathbf{U}, \mathbf{V}), \quad \Psi(\mathbf{U}, \mathbf{V}) \triangleq \frac{1}{2} \|\mathbf{Y} - \mathbf{U}\mathbf{V}\|_{\text{F}}^2 + \chi_{\mathcal{V}_K(\mathbb{F}^M)}(\mathbf{U}).$$

A two-block **BCM** for this problem is:

$$\mathbf{U}_{t+1} = \arg \min_{\mathbf{U} : \mathbf{U}'\mathbf{U} = \mathbf{I}_K} \|\mathbf{Y} - \mathbf{U}\mathbf{V}_t\|_{\text{F}}^2 = \tilde{\mathbf{U}}\tilde{\mathbf{V}}', \quad \underbrace{\mathbf{Y}}_{M \times N} \underbrace{\mathbf{V}_t'}_{N \times K} = \underbrace{\tilde{\mathbf{U}}_K}_{M \times K} \boldsymbol{\Sigma}_K \underbrace{\tilde{\mathbf{V}}'}_{K \times K}. \quad O(MK^2)$$

$$\mathbf{V}_{t+1} = \arg \min_{\mathbf{V} \in \mathbb{F}^{K \times N}} \|\mathbf{Y} - \mathbf{U}_{t+1}\mathbf{V}\|_{\text{F}}^2 \implies \mathbf{V}_{t+1}[:, n] = \mathbf{U}_{t+1}' \mathbf{Y}[:, n] \implies \mathbf{V}_{t+1} = \mathbf{U}_{t+1}' \mathbf{Y}, \quad O(MNK)$$

by solving separate LS problems for each column of  $\mathbf{V}$ , because

$$\|\mathbf{Y} - \mathbf{U}_{t+1}\mathbf{V}\|_{\text{F}}^2 = \sum_{n=1}^N \|\mathbf{Y}[:, n] - \mathbf{U}_{t+1}\mathbf{V}[:, n]\|_2^2.$$

- Convergence?

Cost function decreases every iteration, bounded below,  $\implies$  cost converges.

Convergence of iterates is an active research area.

Yes, Under some RIP conditions [15]

- Other regularizers or constraints?

Sparsity (even more structure than low-rank!)  $R_2(\mathbf{V}) = \|\text{vec}(\mathbf{V})\|_1$

group work:  $\mathbf{V}$  update using **PGM** (because we have so many tools now)

---

## Fused LASSO / generalized LASSO

(Read)

The **fused LASSO** is a generalization of the LASSO problem used in machine learning for regression problems where some features are known to be related [17]. The cost function has the form

$$\Psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{Ax} - \mathbf{y}\|_2^2 + \beta \|\mathbf{x}\|_1 + \gamma \|\mathbf{T}\mathbf{x}\|_1,$$

for some special matrix  $\mathbf{T}$  that involves finite differences between correlated features.

This is a challenging optimization problem.

A “dual path algorithm” [18] [19] [20] is available in R:

<https://rdrr.io/cran/genlasso/man/fusedlasso.html>

However, the documentation states: “Hence it is not advisable to run fusedlasso2d on image denoising problems of large scale, as the dual solution path is computationally infeasible. It should be noted that a faster algorithm for the 2d fused LASSO solution path (when the predictor matrix is the identity), which begins at the dense end of the path, is available in the `flsa` package.”

A related cost function called **generalized LASSO** simply uses  $\beta = 0$  above.

<https://rdrr.io/cran/genlasso/man/genlasso.html>

These challenges motivate the AL/ADMM methods described in Ch. 7.

## Alternating minimization for 0-norm in biconvex form

(Read)

Consider the very challenging sparsity-constrained optimization problem

$$\arg \min_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x}) \text{ s.t. } \|\mathbf{T}\mathbf{x}\|_0 \leq K. \quad (6.23)$$

Nothing we have covered so far solves this type of problem for a general matrix  $\mathbf{T}$ .

If  $f$  has a Lipschitz gradient, then when  $\mathbf{T} = \mathbf{I}$  we could make a proximal gradient method for (6.23), which would involve a form of **iterative hard thresholding**.

For  $\mathbf{z} \in \mathbb{R}^N$ , we can write the 0-norm as the following (**convex**!) optimization problem solution [21]:

$$\|\mathbf{z}\|_0 = \min_{-\mathbf{1}_N \leq \mathbf{u} \leq \mathbf{1}_N} \|\mathbf{u}\|_1 \text{ s.t. } \|\mathbf{z}\|_1 = \langle \mathbf{u}, \mathbf{z} \rangle.$$

The unique solution is  $\mathbf{u}_* = \mathbf{u}_*(\mathbf{z}) = \text{sign}(\mathbf{z})$ , for which  $\|\mathbf{u}_*\|_1 = \|\mathbf{z}\|_0$ .

More generally (thanks to Katherine Banas for helping me see this), for  $\mathbf{z} \in \mathbb{R}^N$ , we can write the 0-norm as the following (**convex**!) optimization problem solution:

$$\|\mathbf{z}\|_0 = \min_{\mathbf{u} \in \mathbb{R}^N : \|\mathbf{u}\|_\infty \leq 1} \|\mathbf{u}\|_1 \text{ s.t. } \|\mathbf{z}\|_1 = \langle \mathbf{u}, \mathbf{z} \rangle. \quad (6.24)$$

The unique solution is  $\mathbf{u}_* = \mathbf{u}_*(\mathbf{z}) = \text{sign}(\mathbf{z})$ , for which  $\|\mathbf{u}_*\|_1 = \|\mathbf{z}\|_0$ .

The form (6.24) is general enough to cover both  $\mathbb{R}$  and  $\mathbb{C}$  cases.

Thus we can rewrite the original problem (6.23) as:

$$\arg \min_{\mathbf{x} \in \mathbb{F}^N} \min_{\mathbf{u} \in \mathbb{F}^N} f(\mathbf{x}) + \chi_{\mathcal{C}_K}(\mathbf{u}) \text{ s.t. } \|\mathbf{T}\mathbf{x}\|_1 = \langle \mathbf{u}, \mathbf{T}\mathbf{x} \rangle, \quad \mathcal{C}_K = \{\mathbf{u} \in \mathbb{F}^N : \|\mathbf{u}\|_\infty \leq 1, \|\mathbf{u}\|_1 \leq K\}. \quad (6.25)$$

One can verify that the constraint set  $\|\mathbf{z}\|_1 = \langle \mathbf{u}, \mathbf{z} \rangle$  is convex, so (6.25) is **biconvex** [22].

When  $\|\mathbf{u}\|_\infty \leq 1$ , we have  $\langle \mathbf{u}, \mathbf{z} \rangle \leq \sum_{n=1}^N |z_n| = \|\mathbf{z}\|_1$ .

Thus  $\{\mathbf{z} : \|\mathbf{z}\|_1 = \langle \mathbf{u}, \mathbf{z} \rangle\} = \{\mathbf{z} : \|\mathbf{z}\|_1 \leq \langle \mathbf{u}, \mathbf{z} \rangle\}$ .

Now if  $\mathbf{z}$  and  $\mathbf{w}$  are both in this set then  $\alpha\mathbf{z} + \beta\mathbf{w}$  is also in this set for  $0 \leq \alpha \leq 1$  and  $\beta = 1 - \alpha$ , because  $\|\alpha\mathbf{z} + \beta\mathbf{w}\|_1 \leq \alpha\|\mathbf{z}\|_1 + \beta\|\mathbf{w}\|_1 \leq \alpha\langle \mathbf{u}, \mathbf{z} \rangle + \beta\langle \mathbf{u}, \mathbf{w} \rangle = \langle \alpha\mathbf{u} + \beta\mathbf{u}, \mathbf{z} \rangle$ .

Still, that constraint seems challenging, so one can replace the constraint in (6.25) with a penalty on the gap:

$$\arg \min_{\mathbf{x} \in \mathbb{F}^N} \min_{\mathbf{u} \in \mathbb{F}^N} f(\mathbf{x}) + \chi_{\mathcal{C}_K}(\mathbf{u}) + \mu (\|\mathbf{T}\mathbf{x}\|_1 - \langle \mathbf{u}, \mathbf{T}\mathbf{x} \rangle). \quad (6.26)$$

This is a **biconvex** problem. One can increase  $\mu$  as the iterations proceed to (asymptotically) enforce the constraint. The obvious approach here is **alternating minimization** and there are convergence results in [21].

**Interpretation**

(Read)

For any finite  $\mu > 0$ , we can rewrite the penalized formulation (6.26) as:

$$\arg \min_{\mathbf{x}} f(\mathbf{x}) + \mu R(\mathbf{T}\mathbf{x}), \quad R(\mathbf{z}) \triangleq \min_{\mathbf{u} \in \mathcal{C}_K} (\|\mathbf{z}\|_1 - \langle \mathbf{u}, \mathbf{z} \rangle).$$

We can simplify the expression for  $R$  as follows:

$$R(\mathbf{z}) = \|\mathbf{z}\|_1 - g(\mathbf{z}), \quad g(\mathbf{z}) \triangleq \max_{\mathbf{u} \in \mathcal{C}_K} \langle \mathbf{u}, \mathbf{z} \rangle = \text{sum}(\text{sort}(\text{abs}(\mathbf{z})) [1:K]).$$

(This is related to the “order weighted L1” (**OWL**) regularizer [23].) Thus

$$R(\mathbf{z}) = \mu \text{sum}(\text{sort}(\text{abs}(\mathbf{z})) [K+1:\text{end}]). \quad (6.27)$$

So this regularizer penalizes the  $N - K$  smallest (in magnitude) values of  $\mathbf{z} = \mathbf{T}\mathbf{x}$ , so as  $\mu$  increases those values will be thresholded to zero, which is what (6.23) requires! To help see this, note that

$$\|\mathbf{z}\|_0 \leq K \iff \text{norm}(\text{sort}(\text{abs}(\mathbf{z})) [K+1:\text{end}], 1) == 0$$

This analysis is helpful for insight, but the sorting operation in (6.27) is very nonconvex and hard for optimization. In contrast, the biconvex set up (6.26) seems much simpler for optimization.

### 6.3 Convergence properties

BCD converges in 1 iteration in the rate cases where the cost function is block separable (decoupled):

$$\Psi(\mathbf{x}_1, \dots, \mathbf{x}_B) = \sum_{b=1}^B \Psi_b(\mathbf{x}_b).$$

Convergence results (for limit points) under weak assumptions are given in [24, p. 268]. See also [25–34].

Convergence of the coordinate descent method for strictly convex, twice-differentiable cost functions is analyzed in detail in [35], including consideration of **box constraints**. Powell demonstrates that uniqueness of the “arg min” step is important to ensure convergence [36].

Convergence rate analysis, including constrained cases, is given in [31, 37].

For “pure” **coordinate descent (CD)**, where we update one parameter at a time in sequence, if the cost function is twice differentiable then there is an asymptotic linear convergence rate. Specifically, when  $\mathbf{x}_n$  is near the minimizer  $\hat{\mathbf{x}}$ :

$$\|\mathbf{x}_{n+1} - \hat{\mathbf{x}}\|_{\mathbf{H}^{1/2}} \leq \rho(\mathbf{M}) \|\mathbf{x}_n - \hat{\mathbf{x}}\|_{\mathbf{H}^{1/2}}, \quad \mathbf{M} = \mathbf{I} - [\mathbf{D} + \mathbf{L}]^{-1} \mathbf{H},$$

where  $\mathbf{H} = \nabla^2 \Psi(\hat{\mathbf{x}}) = \mathbf{L} + \mathbf{D} + \mathbf{L}'$  where  $\mathbf{D}$  is diagonal and  $\mathbf{L}$  is lower triangular. The analysis is similar to that of the **Gauss-Seidel method** for solving a linear system of equations.

Analysis of the convergence rate of “pure” CD is a special case of the analysis of **SAGE** in [38].

For a BPGM version with **momentum** see [34, 39].

For analysis of inexact proximal BCD see [40].

Generalizations for non-smooth and non-convex functions are in [29, 40, 41]. This is an evolving area because of growing interest in BCD methods.

In particular, there is considerable recent focus on **biconvex** problems like those involving **matrix factorization**  $\mathbf{X} = \mathbf{UV}$ . One can show that that some such problems have no spurious local minima [42–44].

Global convergence guarantees for dictionary learning appear in [45], despite many saddle points, for random initialization.

For randomized block selection, see [46].



---

### 6.4 Summary

Alternating minimization methods have a multitude of applications. They also provide a nice context for course review because they use many previous methods (gradient, MM, proximal) as intermediate steps.

**Limitations of BCD** methods are difficulty with **parallelism** (if one uses too many blocks with too small sizes) and getting stuck at non-stationary points for non-smooth cost functions.

### Bibliography

---

- [1] J. Besag. “On the statistical analysis of dirty pictures”. In: *J. Royal Stat. Soc. Ser. B* 48.3 (1986), 259–302 (cit. on p. 6.3).
- [2] J-F. Cai, R. H. Chan, and Z. Shen. “A framelet-based image inpainting algorithm”. In: *Applied and Computational Harmonic Analysis* 24.2 (Mar. 2008), 131–49 (cit. on p. 6.11).
- [3] J-F. Cai, R. Chan, L. Shen, and Z. Shen. “Restoration of chopped and noded images by framelets”. In: *SIAM J. Sci. Comp.* 30.3 (2008), 1205–27 (cit. on p. 6.11).
- [4] J-F. Cai, R. H. Chan, L. Shen, and Z. Shen. “Convergence analysis of tight framelet approach for missing data recovery”. In: *Applied and Computational Harmonic Analysis* 31.1 (Oct. 2009), 87–113 (cit. on p. 6.11).
- [5] J-F. Cai and Z. Shen. “Framelet based deconvolution”. In: *J. Comp. Math.* 28.3 (May 2010), 289–308 (cit. on p. 6.11).
- [6] Y. Liu, J-F. Cai, Z. Zhan, D. Guo, J. Ye, Z. Chen, and X. Qu. “Balanced sparse model for tight frames in compressed sensing magnetic resonance imaging”. In: *PLoS One* 10.4 (2015), 1–19 (cit. on pp. 6.11, 6.12).
- [7] D. Gilton, G. Ongie, and R. Willett. “Learned patch-based regularization for inverse problems in imaging”. In: *Proc. Intl. Wkshp. Comp. Adv. Multi-Sensor Adapt. Proc.* 2019, 211–5 (cit. on p. 6.28).
- [8] S. Ravishankar and Y. Bresler. “ $l_0$  sparsifying transform learning with efficient optimal updates and convergence guarantees”. In: *IEEE Trans. Sig. Proc.* 63.9 (May 2015), 2389–404 (cit. on pp. 6.29, 6.32, 6.34).

- [9] B. Wen, S. Ravishankar, L. Pfister, and Y. Bresler. “Transform learning for magnetic resonance image reconstruction: from model-based learning to building neural networks”. In: *IEEE Sig. Proc. Mag.* 37.1 (Jan. 2020), 41–53 (cit. on p. 6.29).
- [10] Z. Li, S. Xie, W. Chen, and Z. Yang. “Overcomplete transform learning with the log regularizer”. In: *IEEE Access* 6 (2018), 65239–49 (cit. on p. 6.30).
- [11] L. Pfister and Y. Bresler. “Learning filter bank sparsifying transforms”. In: *IEEE Trans. Sig. Proc.* 67.2 (Jan. 2019), 504–19 (cit. on p. 6.30).
- [12] S. Ravishankar and B. Wohlberg. “Learning multi-layer transform models”. In: *Allerton Conf. on Comm., Control, and Computing*. 2018 (cit. on p. 6.40).
- [13] G-J. Peng. “Joint and direct optimization for dictionary learning in convolutional sparse representation”. In: *IEEE Trans. Neural Net. Learn. Sys.* (2019) (cit. on p. 6.46).
- [14] G-J. Peng. “Adaptive ADMM for dictionary learning in convolutional sparse representation”. In: *IEEE Trans. Im. Proc.* 28.7 (July 2019), 3408–422 (cit. on p. 6.46).
- [15] P. Jain, P. Netrapalli, and S. Sanghavi. “Low-rank matrix completion using alternating minimization”. In: *ACM Symp. Theory Comp.* 2013, 665–74 (cit. on p. 6.50).
- [17] R. Tibshirani, M. Saunders, S. Rosset, J. Zhu, and K. Knight. “Sparsity and smoothness via the fused LASSO”. In: *J. Royal Stat. Soc. Ser. B* 67.1 (Feb. 2005), 91–108 (cit. on p. 6.51).
- [18] H. Hoefling. *A path algorithm for the Fused Lasso signal approximator*. 2009 (cit. on p. 6.51).
- [19] R. J. Tibshirani and J. Taylor. “The solution path of the generalized LASSO”. In: *Ann. Stat.* 39.3 (June 2011), 1335–71 (cit. on p. 6.51).
- [20] T. B. Arnold and R. J. Tibshirani. “Efficient implementations of the generalized LASSO dual path algorithm”. In: *J. Computational and Graphical Stat.* 25.1 (2016), 1–27 (cit. on p. 6.51).
- [21] G. Yuan and B. Ghanem. *Sparsity constrained minimization via mathematical programming with equilibrium constraints*. 2018 (cit. on pp. 6.52, 6.53).
- [22] A. Bechensteen, L. Blanc-Feraud, and G. Aubert. “New methods for l2-l0 minimization and their applications to 2D single-molecule localization microscopy”. In: *Proc. IEEE Intl. Symp. Biomed. Imag.* 2019, 1377–81 (cit. on p. 6.53).
- [23] M. A. T. Figueiredo and R. D. Nowak. “Ordered weighted l1 regularized regression with strongly correlated covariates: Theoretical aspects”. In: *aistats*. 2016, 930–8 (cit. on p. 6.54).
- [24] D. P. Bertsekas. *Nonlinear programming*. 2nd ed. Belmont: Athena Scientific, 1999 (cit. on p. 6.55).
- [25] P. Tseng. “Convergence of a block coordinate descent methods for nondifferentiable minimization”. In: *J. Optim. Theory Appl.* 109 (2001), 475–94 (cit. on p. 6.55).

- [26] Y. Nesterov. “Efficiency of coordinate descent methods on huge-scale optimization problems”. In: *SIAM J. Optim.* 22.2 (2012), 341–62 (cit. on p. 6.55).
- [27] M. W. Jacobson and J. A. Fessler. “An expanded theoretical treatment of iteration-dependent majorize-minimize algorithms”. In: *IEEE Trans. Im. Proc.* 16.10 (Oct. 2007), 2411–22 (cit. on p. 6.55).
- [28] A. Beck and L. Tetruashvili. “On the convergence of block coordinate descent type methods”. In: *SIAM J. Optim.* 23.4 (2013), 2037–60 (cit. on p. 6.55).
- [29] Y. Xu and W. Yin. “A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion”. In: *SIAM J. Imaging Sci.* 6.3 (2013), 1758–89 (cit. on pp. 6.55, 6.56).
- [30] J. Bolte, S. Sabach, and M. Teboulle. “Proximal alternating linearized minimization for nonconvex and nonsmooth problems”. In: *Mathematical Programming* 146.1 (Aug. 2014), 459–94 (cit. on p. 6.55).
- [31] S. Yun. “On the iteration complexity of cyclic coordinate gradient descent methods”. In: *SIAM J. Optim.* 24.3 (2014), 1567–80 (cit. on p. 6.55).
- [32] K. Khare and B. Rajaratnam. *Convergence of cyclic coordinatewise l1 minimization*. 2015 (cit. on p. 6.55).
- [33] Z. Shi and R. Liu. *A better convergence analysis of the block coordinate descent method for large scale machine learning*. 2016 (cit. on p. 6.55).
- [34] Y. Xu and W. Yin. “A globally convergent algorithm for nonconvex optimization based on block coordinate update”. In: *J. of Scientific Computing* 72.2 (Aug. 2017), 1–35 (cit. on pp. 6.55, 6.56).
- [35] Z. Q. Luo and P. Tseng. “On the convergence of the coordinate descent method for convex differentiable minimization”. In: *J. Optim. Theory Appl.* 72.1 (Jan. 1992), 7–35 (cit. on p. 6.55).
- [36] M. J. D. Powell. “On search directions for minimization algorithms”. In: *Mathematical Programming* 4.1 (1973), 193–201 (cit. on p. 6.55).
- [37] Z-Q. Luo and P. Tseng. “On the convergence rate of dual ascent methods for linearly constrained convex minimization”. In: *Math. Oper. Res.* 18.4 (Nov. 1993), 846–67 (cit. on p. 6.55).
- [38] J. A. Fessler and A. O. Hero. “Space-alternating generalized expectation-maximization algorithm”. In: *IEEE Trans. Sig. Proc.* 42.10 (Oct. 1994), 2664–77 (cit. on p. 6.56).
- [39] I. Y. Chun and J. A. Fessler. “Convolutional analysis operator learning: acceleration and convergence”. In: *IEEE Trans. Im. Proc.* 29.1 (Jan. 2020), 2108–22 (cit. on p. 6.56).
- [40] E. Chouzenoux, J-C. Pesquet, and A. Repetti. “A block coordinate variable metric forward-backward algorithm”. In: *J. of Global Optimization* 66.3 (Nov. 2016), 457–85 (cit. on p. 6.56).

- [41] M. Razaviyayn, M. Hong, and Z. Luo. “A unified convergence analysis of block successive minimization methods for nonsmooth optimization”. In: *SIAM J. Optim.* 23.2 (2013), 1126–53 (cit. on p. 6.56).
- [42] M. Hardt. “Understanding alternating minimization for matrix completion”. In: *Foundations of Computer Science (FOCS)*. 2014, 651–60 (cit. on p. 6.56).
- [43] R. Ge, J. D. Lee, and T. Ma. “Matrix completion has no spurious local minimum”. In: *Neural Info. Proc. Sys.* 2016, 2973–81 (cit. on p. 6.56).
- [44] R. Ge, C. Jin, and Y. Zheng. “No spurious local minima in nonconvex low rank problems: A unified geometric analysis”. In: *Proc. Intl. Conf. Mach. Learn.* Vol. 70. 2017, 1233–42 (cit. on p. 6.56).
- [45] D. Gilboa, S. Buchanan, and J. Wright. *Efficient dictionary learning with gradient descent*. 2018 (cit. on p. 6.56).
- [46] J. Diakonikolas and L. Orecchia. “Alternating randomized block coordinate descent”. In: *Proc. Intl. Conf. Mach. Learn.* Vol. 80. 2018, 1224–32 (cit. on p. 6.56).