

Chapter 9

Misc topics

Contents (class version)

9.0 Review	9.2
Binary classifier design: Review	9.4
9.1 If time had permitted...	9.13
9.2 Towards CNN methods	9.14
Review of denoising by soft thresholding	9.14
Review of compressed sensing with transform sparsity	9.15
Review of patch transform sparsity	9.16

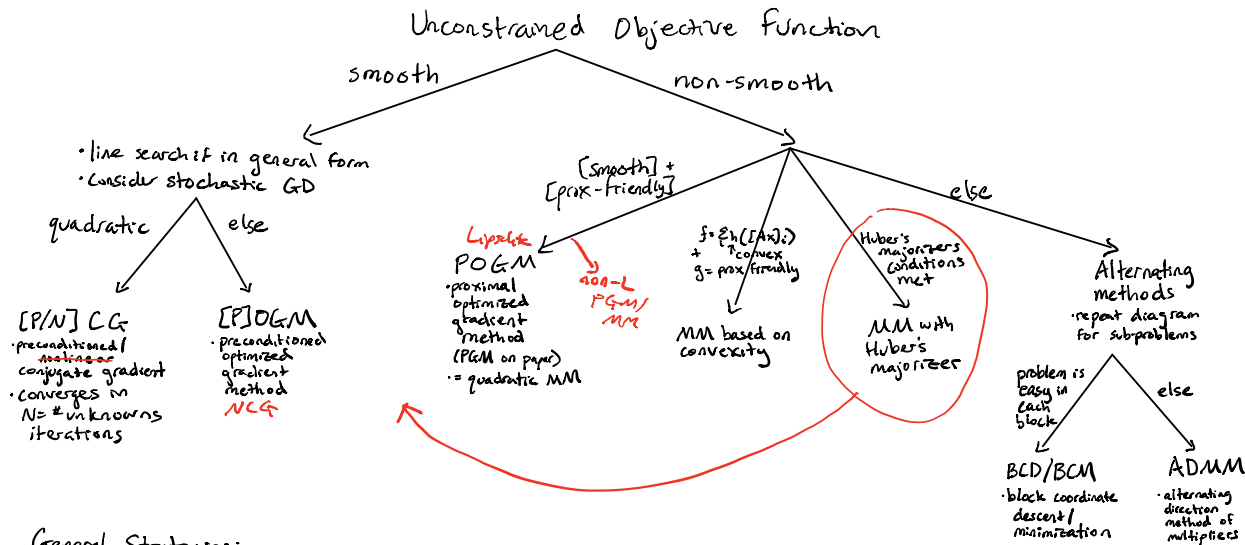
9.0 Review

Final course notes on gdrive.

Two principal kinds of questions.

- **Procedural.** Given $\Psi(\mathbf{x})$, given algorithm (by name), then implement/analyze it
Interesting Exam2 sample questions (by pdf file page)
 - p1: composite with $\|\mathbf{x}\|_1$ and box constraint
 - p20: LS + hinge regularizer
 - p24: sigmoid as smooth 0-1 loss
 - p43: prox for leaky ReLU
 - p44: hinge + 0-norm
- **Conceptual.** Given application, then determine $\Psi(\mathbf{x})$, choose algorithm, and implement/analyze it
 - p2: $\|\mathbf{x}\|_0$ constraint
 - p6: L+S
 - p45: compressed sensing with transform/analysis sparsity

Caroline's review diagram

General Strategies:

- Rewrite as a joint cost function $\tilde{x} = \begin{bmatrix} x \\ z \end{bmatrix}$
- Write constraints using the characteristic function
- Rewrite to make convex, e.g., $\|x\|_1 = u + v$
or linear: $\|x\|_1 \leq t \Rightarrow -t \leq x \leq t$

• Things never appearing here:

- PSD / PGD \rightarrow use $\{P\}$ OGM?
- FGM \rightarrow use OGM
- FPGM \rightarrow use POGM

Binary classifier design: Review

General cost function for binary classifier design for $\text{sign}(\mathbf{v}'\mathbf{x})$ with $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, $\beta \geq 0$:

$$\Psi(\mathbf{x}) = f(\mathbf{x}) + \beta \|\mathbf{x}\|_p^p, \quad f(\mathbf{x}) = \mathbf{1}_M' h(\mathbf{A}\mathbf{x}), \quad p \in \{1, 2\}$$

Quadratic

$$h(t) = \frac{1}{2}(t - 1)^2 \implies f(\mathbf{x}) =$$

Method for $p = 2$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

Method for $p = 1$?

A: GD

B: CG

C: OGM

D: POGM

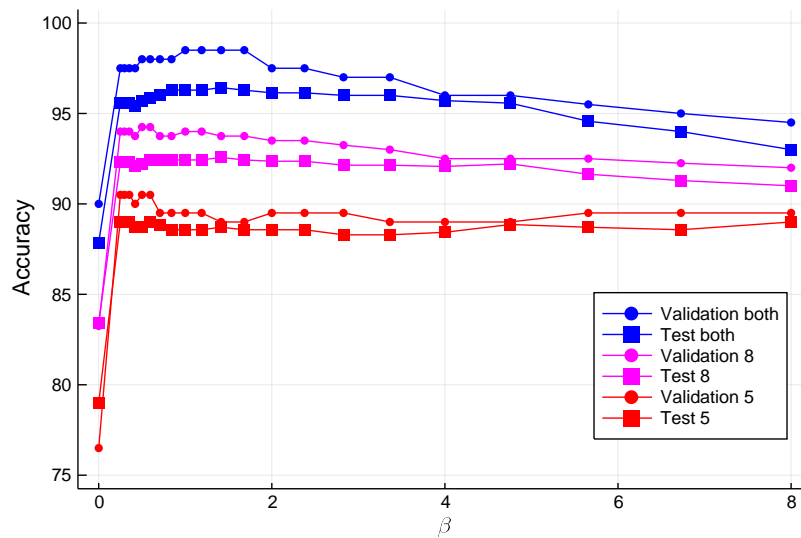
E: ADMM

??

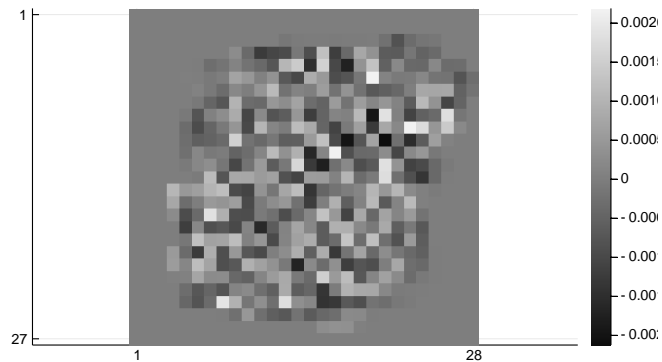
Other options?

Example

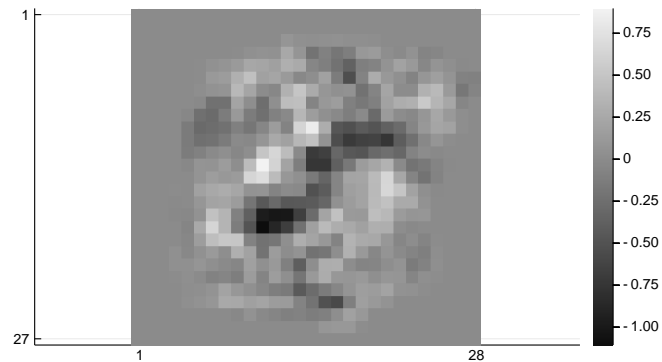
Same 5/8 handwritten digit data as in HW



LS classifier weights



Best regularized LS weights



Logistic / Huber hinge

Method for $p = 2$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

Method for $p = 1$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

Hinge _____Method for $p = 2$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

Method for $p = 1$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

Sigmoid (smooth version of 0-1) loss _____Method for $p = 2$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

Method for $p = 1$?

A: GD

B: CG

C: OGM

D: POGM

E: ADMM

??

Other options?

0-1 loss



(blank)

(blank)

9.1 If time had permitted...

recommender systems

neural networks

parallel computing

semidefinite programming

non-convex regularizers that lead to convex cost functions [9]

relationships of image models and priors for restoration problems [11]

Regularization parameter selection

SURE etc.

CNN training using SURE without ground-truth images [1] [2]

Optimization on manifolds

Minimization subject to constraints like a matrix being unitary (Stiefel manifold) [5–8]

9.2 Towards CNN methods

This section reviews some iterative algorithms based on sparsity models and summarizes how those algorithms provide a foundation for “variational neural networks” when “unrolled.”

Review of denoising by soft thresholding

Consider the measurement model $\mathbf{y} = \mathbf{x} + \boldsymbol{\varepsilon}$ and the signal model that assumes $\mathbf{T}\mathbf{x}$ is sparse for some unitary transform \mathbf{T} . The natural optimization problem for estimating \mathbf{x} is

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

The non-iterative solution is the following **denoising** operation:

$$\hat{\mathbf{x}} =$$

The essential ingredients of a **convolutional neural network (CNN)** are present in this simple form:

Review of compressed sensing with transform sparsity

For the compressed sensing measurement model $\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\varepsilon}$, again assuming $\mathbf{T}\mathbf{x}$ is sparse for some unitary transform \mathbf{T} , the natural optimization problem for estimating \mathbf{x} is

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

In this case there is no closed-form solution and iterative algorithms are needed. The simplest algorithm is the **proximal gradient method (PGM)**, aka **ISTA**, which is an MM update based on the majorizer

$$\begin{aligned} \tilde{\mathbf{x}}_k &= \mathbf{x}_k - \frac{1}{L} \mathbf{A}'(\mathbf{A}\mathbf{x}_k - \mathbf{y}) \\ \phi_k(\mathbf{x}) &= \frac{L}{2} \|\mathbf{x} - \tilde{\mathbf{x}}_k\|_2^2 + \beta \|\mathbf{T}\mathbf{x}\|_1, \end{aligned}$$

where $L = \|\mathbf{A}\|_2^2$, for which the minimization step is a **denoising** operation:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \phi_k(\mathbf{x}) =$$

The “unrolled loop” block diagram for this algorithm is the basis for **learned ISTA (LISTA)** [3]:

$$\mathbf{y} \rightarrow \mathbf{x}_0 \rightarrow \boxed{\text{data}} \rightarrow \tilde{\mathbf{x}}_0 \rightarrow \boxed{\text{denoise}} \rightarrow \mathbf{x}_1 \rightarrow \boxed{\text{data}} \rightarrow \tilde{\mathbf{x}}_1 \rightarrow \boxed{\text{denoise}} \rightarrow \mathbf{x}_2 \rightarrow \boxed{\text{data}} \rightarrow \tilde{\mathbf{x}}_2 \cdots$$

Denoise options:

Can learn:

Early SURE-LET work: [4].

Review of patch transform sparsity

The natural cost function for a patch transform sparse model is

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta R(\mathbf{x}), \quad R(\mathbf{x}) = \min_{\mathbf{Z}} \sum_{p=1}^P \frac{1}{2} \|\mathbf{T}\mathbf{P}_p\mathbf{x} - \mathbf{z}_p\|_2^2 + \alpha \|\mathbf{z}_p\|_1.$$

For a BCD approach, the \mathbf{Z} update is simple:

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

To better understand the \mathbf{x} update, it is helpful to expand the regularizer:

$$\sum_{p=1}^P \frac{1}{2} \|\mathbf{T}\mathbf{P}_p\mathbf{x} - \mathbf{z}_p^{(t+1)}\|_2^2 =$$

$$= \frac{1}{2} \mathbf{x}' \mathbf{H} \mathbf{x} - \mathbf{x}' \tilde{\mathbf{x}}_t + c_1,$$

If \mathbf{T} has orthonormal columns (e.g., is unitary), and if the patches have d pixels and are chosen with stride=1

and periodic boundary conditions, then

$$\mathbf{H} =$$

Completing the square for the regularizer term yields

$$\dots = \frac{d}{2} \mathbf{x}' \mathbf{x} - \mathbf{x}' \tilde{\mathbf{x}}_t + c_1 = \quad \quad \quad \tilde{\mathbf{x}}_t \triangleq$$

Thus the \mathbf{x} update for the **BCD** algorithm is

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \frac{d}{2} \|\mathbf{x} - \tilde{\mathbf{x}}_t\|_2^2.$$

Here, $\tilde{\mathbf{x}}_t$ acts like a prior for the update. For some cases there is a closed-form solution (like single-coil Cartesian MRI). Otherwise, one or more iterations are needed.

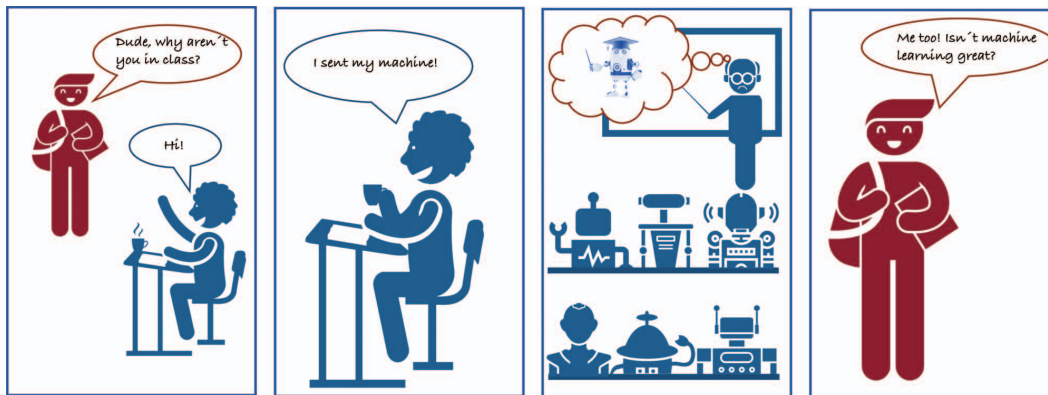
Unrolling the BCD loop in block diagram form:

$$\mathbf{y} \rightarrow \mathbf{x}_0 \rightarrow \boxed{\text{denoise}} \rightarrow \tilde{\mathbf{x}}_0 \rightarrow \boxed{\text{data}} \rightarrow \mathbf{x}_1 \rightarrow \boxed{\text{denoise}} \rightarrow \tilde{\mathbf{x}}_1 \rightarrow \boxed{\text{data}} \rightarrow \mathbf{x}_2 \rightarrow \boxed{\text{denoise}} \rightarrow \tilde{\mathbf{x}}_2 \cdots$$

Denoise options:

Machine-Learning Class

by Robert W. Heath, Jr. and Nuria González-Prelcic



Note: This cartoon was created entirely by real humans, no machine learning was involved with conception or realization to practice.

(c) Robert W. Heath Jr. and Nuria González-Prelcic, 2019

Bibliography

- [1] C. A. Metzler, A. Mousavi, R. Heckel, and R. G. Baraniuk. *Unsupervised learning with Stein's unbiased risk estimator*. 2018 (cit. on p. 9.13).
- [2] S. Soltanayev and S. Y. Chun. "Training deep learning based denoisers without ground truth data". In: *Neural Info. Proc. Sys.* 2018, 3257–67 (cit. on p. 9.13).
- [3] K. Gregor and Y. LeCun. "Learning fast approximations of sparse coding". In: *Proc. Intl. Conf. Mach. Learn.* 2010 (cit. on p. 9.15).
- [4] T. Blu and F. Luisier. "The SURE-LET approach to image denoising". In: *IEEE Trans. Im. Proc.* 16.11 (Nov. 2007), 2778–86 (cit. on p. 9.15).
- [5] A. Edelman, Tomás A Arias, and S. T. Smith. "The geometry of algorithms with orthogonality constraints". In: *SIAM J. Matrix. Anal. Appl.* 20.2 (1998), 303–53 (cit. on p. 9.13).
- [6] T. E. Abrudan, J. Eriksson, and V. Koivunen. "Steepest descent algorithms for optimization under unitary matrix constraint". In: *IEEE Trans. Sig. Proc.* 56.3 (Mar. 2008), 1134–47 (cit. on p. 9.13).
- [7] B. Gao, X. Liu, X. Chen, and Y. Yuan. "A new first-order algorithmic framework for optimization problems with orthogonality constraints". In: *SIAM J. Optim.* 28.1 (Jan. 2018), 302–332 (cit. on p. 9.13).
- [8] S. Chen, S. Ma, A. M-C. So, and T. Zhang. "Proximal gradient method for nonsmooth optimization over the Stiefel manifold". In: *SIAM J. Optim.* 30.1 (Jan. 2020), 210–39 (cit. on p. 9.13).
- [9] A. Lanza, S. Morigi, I. W. Selesnick, and F. Sgallari. "Sparsity-inducing nonconvex nonseparable regularization for convex image processing". In: *SIAM J. Imaging Sci.* 12.2 (Jan. 2019), 1099–34 (cit. on p. 9.13).
- [10] R. T. Rockafellar. "Lagrange multipliers and optimality". In: *SIAM Review* 35.2 (June 1993), 183–238.
- [11] B. Wen, Y. Li, Y. Li, and Y. Bresler. *A set-theoretic study of the relationships of image models and priors for restoration problems*. 2020 (cit. on p. 9.13).
- [12] W. Zuo and Z. Lin. "A generalized accelerated proximal gradient approach for total-variation-based image restoration". In: *IEEE Trans. Im. Proc.* 20.10 (Oct. 2011), 2748–59.