

A Fast Augmented Lagrangian Algorithm for Learning Low-Rank Matrices

Ryota Tomioka¹, Taiji Suzuki¹, Masashi Sugiyama²,
Hisashi Kashima¹

¹The University of Tokyo

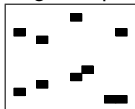
²Tokyo Institute of Technology

2010-06-22 @ ICML 2010

Learning low-rank matrices

- 1 Matrix completion [Srebro et al. 05; Abernethy et al. 09] (collaborative filtering, link prediction)

$$Y = W$$



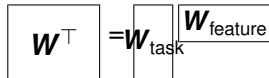
Part of Y is observed



- 2 Multi-task learning [Argyriou et al., 07]

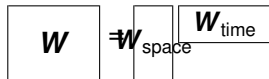
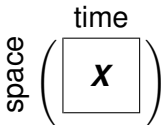
$$y = W^T x + b$$

$$W^T = \begin{bmatrix} w^{\text{task 1} T} \\ w^{\text{task 2} T} \\ \vdots \\ w^{\text{task R} T} \end{bmatrix}$$



- 3 Predicting over matrices (classification/regression) [Tomioka & Aihara, 07]

$$y = \langle W, X \rangle + b$$



Problem formulation

Primal problem

$$\underset{\mathbf{W} \in \mathbb{R}^{R \times C}}{\text{minimize}} \quad \underbrace{f_\ell(\mathcal{A}(\mathbf{W})) + \phi_\lambda(\mathbf{W})}_{=: f(\mathbf{W})}$$

- $f_\ell (\mathbb{R}^m \rightarrow \mathbb{R})$: loss function (differentiable).
- $\mathcal{A} (\mathbb{R}^{R \times C} \rightarrow \mathbb{R}^m)$: design matrix.
- ϕ_λ : regularizer (non-differentiable); for example, the trace norm:

$$\phi_\lambda(\mathbf{W}) = \lambda \|\mathbf{W}\|_* = \lambda \sum_{j=1}^r \sigma_j(\mathbf{W}) \quad (\text{linear sum of singular values}).$$

Separation of f_ℓ and $\mathcal{A} \Rightarrow$ convergence regardless of \mathcal{A} (input data).

Existing approaches and our goal

- Proximal (accelerated) gradient [Ji & Ye, 09]
 - Can keep the intermediate solution low-rank.
 - Slow for **poorly conditioned design matrix \mathcal{A}** .
 - Optimal as a first order black-box method $O(1/k^2)$.
- Interior point algorithm [Tomioka & Aihara, 07]
 - Can obtain highly precise solution in small number of iterations.
 - Only low-rank in the limit (**each step can be heavy**).

Can we keep the intermediate solution **low-rank** and still **converge rapidly** as an IP method?

- Dual Augmented Lagrangian (DAL) [Tomioka & Sugiyama, 09] for *sparse estimation*.
⇒ M-DAL: generalized to **low-rank matrix estimation**.

Existing approaches and our goal

- Proximal (accelerated) gradient [Ji & Ye, 09]
 - Can keep the intermediate solution low-rank.
 - Slow for **poorly conditioned design matrix \mathcal{A}** .
 - Optimal as a first order black-box method $O(1/k^2)$.
- Interior point algorithm [Tomioka & Aihara, 07]
 - Can obtain highly precise solution in small number of iterations.
 - Only low-rank in the limit (**each step can be heavy**).

Can we keep the intermediate solution **low-rank** and still **converge rapidly** as an IP method?

- Dual Augmented Lagrangian (DAL) [Tomioka & Sugiyama, 09] for *sparse estimation*.
⇒ M-DAL: generalized to low-rank matrix estimation.

Existing approaches and our goal

- Proximal (accelerated) gradient [Ji & Ye, 09]
 - Can keep the intermediate solution low-rank.
 - Slow for **poorly conditioned design matrix \mathcal{A}** .
 - Optimal as a first order black-box method $O(1/k^2)$.
- Interior point algorithm [Tomioka & Aihara, 07]
 - Can obtain highly precise solution in small number of iterations.
 - Only low-rank in the limit (**each step can be heavy**).

Can we keep the intermediate solution **low-rank** and still **converge rapidly** as an IP method?

- Dual Augmented Lagrangian (DAL) [Tomioka & Sugiyama, 09] for *sparse estimation*.
⇒ M-DAL: generalized to **low-rank matrix estimation**.

Outline

- 1 Introduction
 - Why learn *low-rank* matrices?
 - Existing algorithms
- 2 Proposed algorithm
 - **Augmented Lagrangian** \leftrightarrow Proximal minimization
 - Super-linear convergence
 - Generalizations
- 3 Experiments
 - Synthetic $10,000 \times 10,000$ matrix completion.
 - BCI classification problem (learning multiple matrices).
- 4 Summary

M-DAL algorithm

- 1 Initialize \mathbf{W}^0 . Choose $\eta_1 \leq \eta_2 \leq \dots$
- 2 Iterate until relative duality gap $< \epsilon$

$$\alpha^t := \operatorname{argmin}_{\alpha \in \mathbb{R}^m} \varphi_t(\alpha) \quad (\text{inner objective})$$

$$\mathbf{W}^{t+1} = \operatorname{ST}_{\lambda \eta_t} \left(\mathbf{W}^t + \eta_t \mathcal{A}^\top(\alpha^t) \right)$$

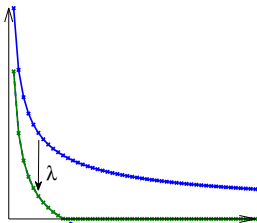
Choosing $\alpha = -\nabla f_\ell^t$ yields the proximal gradient (forward-backward) method.

$$\operatorname{ST}_\lambda(\mathbf{W}) := \operatorname{argmin}_{\mathbf{X} \in \mathbb{R}^{R \times C}} \left(\phi_\lambda(\mathbf{X}) + \frac{1}{2} \|\mathbf{X} - \mathbf{W}\|_{\text{fro}}^2 \right).$$

For the trace-norm $\phi_\lambda(\mathbf{W}) = \lambda \|\mathbf{W}\|_*$,

$$\operatorname{ST}_\lambda(\mathbf{W}) = \mathbf{U} \max(\mathbf{S} - \lambda \mathbf{I}, 0) \mathbf{V}^\top,$$

where $\mathbf{W} = \mathbf{U} \mathbf{S} \mathbf{V}^\top$ (singular-value decomposition).



Proximal minimization [Rockafellar 76a, 76b]

1 Initialize \mathbf{W}^0 .

2 Iterate:

$$\begin{aligned} \mathbf{W}^{t+1} &= \operatorname{argmin}_{\mathbf{W}} \left(f(\mathbf{W}) + \overbrace{\frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}^{\text{proximal term}} \right) \\ &= \operatorname{argmin}_{\mathbf{W}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \underbrace{\phi_{\lambda}(\mathbf{W}) + \frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}_{=:\frac{1}{\eta_t} \Phi_{\lambda\eta_t}(\mathbf{W}; \mathbf{W}^t)} \right) \end{aligned}$$

Fenchel dual (see Rockafellar 1970):

$$\begin{aligned} \min_{\mathbf{W}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \frac{1}{\eta_t} \Phi_{\lambda\eta_t}(\mathbf{W}; \mathbf{W}^t) \right) &= \max_{\alpha} \left(-f_{\ell}^*(-\alpha) - \frac{1}{\eta_t} \Phi_{\lambda\eta_t}^*(\eta_t \mathcal{A}^T(\alpha); \mathbf{W}^t) \right) \\ &=: -\min_{\alpha} \varphi_t(\alpha) \end{aligned}$$

Proximal minimization [Rockafellar 76a, 76b]

1 Initialize \mathbf{W}^0 .

2 Iterate:

$$\begin{aligned} \mathbf{W}^{t+1} &= \operatorname{argmin}_{\mathbf{W}} \left(f(\mathbf{W}) + \overbrace{\frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}^{\text{proximal term}} \right) \\ &= \operatorname{argmin}_{\mathbf{W}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \underbrace{\phi_{\lambda}(\mathbf{W}) + \frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}_{=:\frac{1}{\eta_t} \Phi_{\lambda\eta_t}(\mathbf{W}; \mathbf{W}^t)} \right) \end{aligned}$$

Fenchel dual (see Rockafellar 1970):

$$\begin{aligned} \min_{\mathbf{W}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \frac{1}{\eta_t} \Phi_{\lambda\eta_t}(\mathbf{W}; \mathbf{W}^t) \right) &= \max_{\alpha} \left(-f_{\ell}^*(-\alpha) - \frac{1}{\eta_t} \Phi_{\lambda\eta_t}^*(\eta_t \mathcal{A}^T(\alpha); \mathbf{W}^t) \right) \\ &=: -\min_{\alpha} \varphi_t(\alpha) \end{aligned}$$

Proximal minimization [Rockafellar 76a, 76b]

1 Initialize \mathbf{W}^0 .

2 Iterate:

$$\begin{aligned}
 \mathbf{W}^{t+1} &= \operatorname{argmin}_{\mathbf{W}} \left(f(\mathbf{W}) + \overbrace{\frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}^{\text{proximal term}} \right) \\
 &= \operatorname{argmin}_{\mathbf{W}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \underbrace{\phi_{\lambda}(\mathbf{W}) + \frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}_{=:\frac{1}{\eta_t} \Phi_{\lambda\eta_t}(\mathbf{W}; \mathbf{W}^t)} \right)
 \end{aligned}$$

Fenchel dual (see Rockafellar 1970):

$$\begin{aligned}
 \min_{\mathbf{W}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \frac{1}{\eta_t} \Phi_{\lambda\eta_t}(\mathbf{W}; \mathbf{W}^t) \right) &= \max_{\alpha} \left(-f_{\ell}^*(-\alpha) - \frac{1}{\eta_t} \Phi_{\lambda\eta_t}^*(\eta_t \mathcal{A}^{\top}(\alpha); \mathbf{W}^t) \right) \\
 &=: -\min_{\alpha} \varphi_t(\alpha)
 \end{aligned}$$

Definition

- \mathbf{W}^* : the unique minimizer of the objective $f(\mathbf{W})$.
- \mathbf{W}^t : sequence generated by the M-DAL algorithm with

$$\|\nabla\varphi_t(\boldsymbol{\alpha}^t)\| \leq \sqrt{\frac{\gamma}{\eta_t}} \|\mathbf{W}^{t+1} - \mathbf{W}^t\|_{\text{fro}} \quad \left(\begin{array}{l} 1/\gamma: \text{ Lipschitz con-} \\ \text{stant of } \nabla f_\ell. \end{array} \right)$$

Assumption

There is a constant $\sigma > 0$ such that

$$f(\mathbf{W}^{t+1}) - f(\mathbf{W}^*) \geq \sigma \|\mathbf{W}^{t+1} - \mathbf{W}^*\|_{\text{fro}}^2 \quad (t = 0, 1, 2, \dots).$$

Theorem: Super-linear convergence

$$\|\mathbf{W}^{t+1} - \mathbf{W}^*\|_{\text{fro}} \leq \frac{1}{\sqrt{1 + 2\sigma\eta_t}} \|\mathbf{W}^t - \mathbf{W}^*\|_{\text{fro}}.$$

I.e., \mathbf{W}^t converges super-linearly to \mathbf{W}^* if η_t is increasing.

Why is M-DAL efficient?

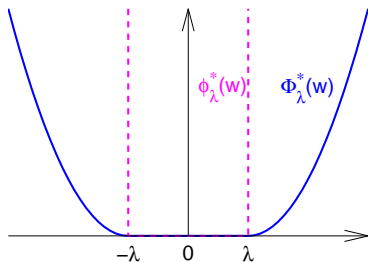
(1) Proximation wrt ϕ_λ is analytic (though non-smooth):

$$\mathbf{W}^{t+1} = \text{ST}_{\eta_t \lambda} \left(\mathbf{W}^t + \eta_t \mathcal{A}^\top(\boldsymbol{\alpha}^t) \right)$$

(2) Inner minimization is smooth:

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{argmin}} \left(\underbrace{f_\ell^*(-\boldsymbol{\alpha})}_{\text{Differentiable}} + \frac{1}{2\eta_t} \underbrace{\|\text{ST}_{\eta_t \lambda}(\mathbf{W}^t + \eta_t \mathcal{A}^\top(\boldsymbol{\alpha}))\|_{\text{fro}}^2}_{= \Phi_\lambda^*(\cdot)} \right)$$

(linear to the estimated rank)



Generalizations

1 Learning multiple matrices

$$\phi_\lambda(\mathbf{W}) = \lambda \sum_{k=1}^K \|\mathbf{W}^{(k)}\|_* = \lambda \left\| \begin{pmatrix} \mathbf{W}^{(1)} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{W}^{(K)} \end{pmatrix} \right\|_*$$

- No need to form the big matrix.
- Can be used to select informative data sources and learn feature extractors simultaneously.

2 General spectral regularization

$$\phi_\lambda(\mathbf{W}) = \sum_{j=1}^r g_\lambda(\sigma_j(\mathbf{W}))$$

for any convex function g_λ for which the proximal operator:

$$\text{ST}_\lambda^g(\sigma_j) = \underset{x \in \mathbb{R}}{\text{argmin}} \left(g_\lambda(x) + \frac{1}{2}(x - \sigma_j)^2 \right)$$

can be computed in closed form.

Synthetic experiment 1: low-rank matrix completion

Large scale & structured.

- True matrix \mathbf{W}^* : $10,000 \times 10,000$ (100M elements), low rank.
- Observation: randomly chosen m elements (sparse).

⇒ **Quasi-Newton method** for the minimization of $\varphi_t(\alpha)$

⇒ No need to form the full matrix!

$m=1,200,000$

$m=2,400,000$

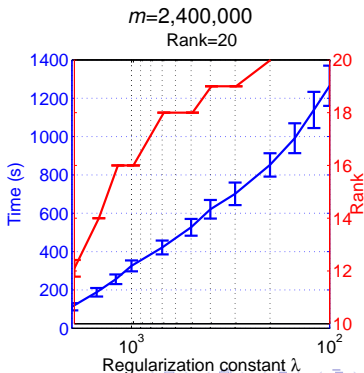
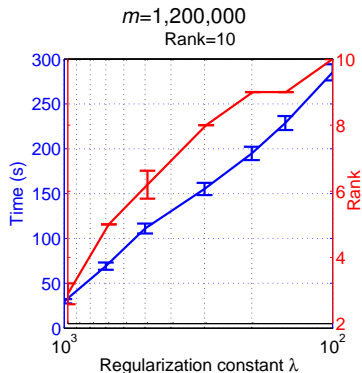
Synthetic experiment 1: low-rank matrix completion

Large scale & structured.

- True matrix \mathbf{W}^* : $10,000 \times 10,000$ (100M elements), low rank.
- Observation: randomly chosen m elements (sparse).

⇒ **Quasi-Newton method** for the minimization of $\varphi_t(\alpha)$

⇒ No need to form the full matrix!



Rank=10

Rank=10, #observations $m=1,200,000$

λ	time (s)	#outer	#inner	rank	S-RMSE
1000	33.1 (± 2.0)	5 (± 0)	8 (± 0)	2.8 (± 0.4)	0.0158 (± 0.0024)
700	77.1 (± 5.6)	11 (± 0)	18 (± 0)	5 (± 0)	0.0133 (± 0.0008)
500	124 (± 7.2)	17 (± 0)	28 (± 0)	6.4 (± 0.5)	0.0113 (± 0.0015)
300	174 (± 8.0)	23 (± 0)	38.4 (± 0.84)	8 (± 0)	0.00852 (± 0.00039)
200	220 (± 9.9)	29 (± 0)	48.4 (± 0.84)	9 (± 0)	0.00767 (± 0.00031)
150	257 (± 9.9)	35 (± 0)	58.4 (± 0.84)	9 (± 0)	0.00498 (± 0.00026)
100	319 (± 11)	41 (± 0)	70 (± 0.82)	10 (± 0)	0.00743 (± 0.00013)

- #inner iterations is roughly 2 times #outer iterations.
 ← Because we don't need to solve the inner problem very precisely!

Rank=20

Rank=20, #observations $m=2,400,000$

λ	time (s)	#outer	#inner	rank	S-RMSE
2000	112 (± 19)	6 (± 0)	15.1 (± 1.0)	12.1 (± 0.3)	0.011 (± 0.002)
1500	188 (± 22)	11 (± 0)	24.1 (± 1.0)	14 (± 0)	0.0094 (± 0.001)
1200	256 (± 25)	15 (± 0)	31.1 (± 1.0)	16 (± 0)	0.0090 (± 0.0008)
1000	326 (± 29)	19 (± 0)	38.1 (± 1.0)	16 (± 0)	0.0073 (± 0.0007)
700	421 (± 36)	24 (± 0)	48.1 (± 1.0)	18 (± 0)	0.0065 (± 0.0004)
500	527 (± 44)	29 (± 0)	57.1 (± 1.0)	18 (± 0)	0.0042 (± 0.0003)
400	621 (± 48)	34 (± 0)	66.1 (± 1.0)	19 (± 0)	0.0044 (± 0.0002)
300	702 (± 59)	38.5 (± 0.5)	74.1 (± 1.5)	19 (± 0)	0.0030 (± 0.0003)
200	852 (± 61)	43.6 (± 0.5)	83.9 (± 2.3)	20 (± 0)	0.0039 (± 0.0001)
150	992 (± 78)	48.4 (± 0.7)	92.5 (± 1.5)	20 (± 0)	0.0024 (± 0.0002)
120	1139 (± 94)	53.4 (± 0.7)	102 (± 1.5)	20 (± 0)	0.0016 ($\pm 6 \times 10^{-5}$)
100	1265 (± 105)	57.7 (± 0.8)	109 (± 2.4)	20 (± 0)	0.0013 ($\pm 8 \times 10^{-5}$)

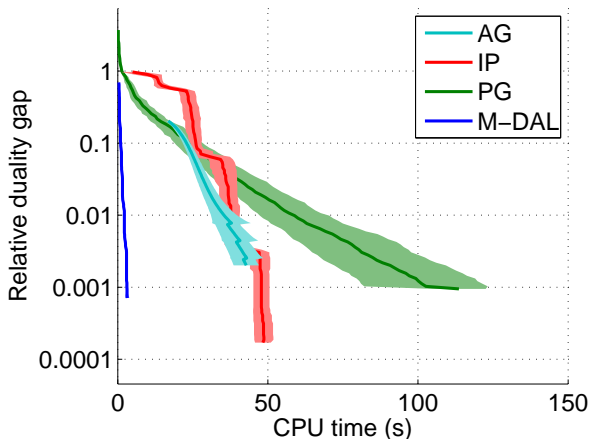
Synthetic experiment 2: classifying matrices

- True matrix \mathbf{W}^* : 64×64 , rank=16.
- Classification problem (# samples $m = 1000$):

$$f_\ell(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-y_i z_i))$$

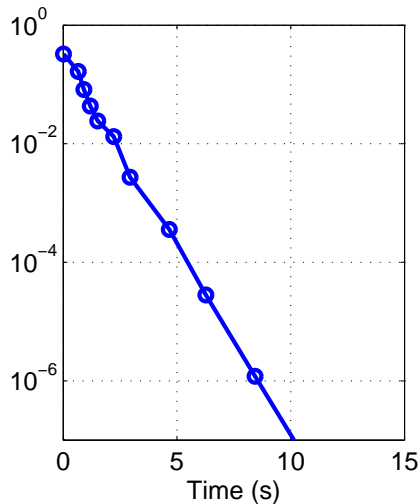
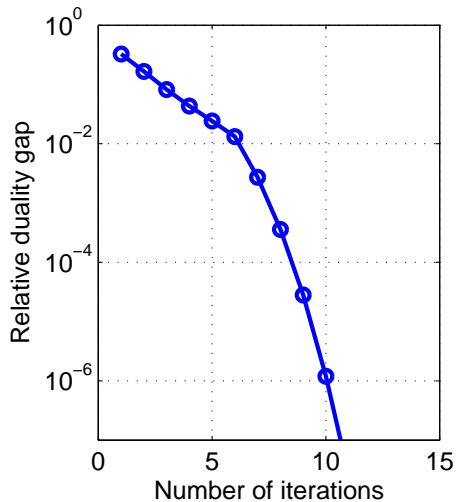
- Design matrix \mathcal{A} : 1000×64^2 (dense) each example drawn from Wishart distribution.
- $\lambda = 800$ (chosen to roughly reproduce rank=16).
- Medium scale & dense \Rightarrow Newton method for minimizing $\varphi_t(\alpha)$ (works better when the condition is poor).
- Methods:
 - M-DAL (proposed)
 - IP (interior point method [T& A 2007])
 - PG (projected gradient method [T& S 2008])
 - AG (accelerated gradient method [Ji & Ye 2009])

Comparison



- **M-DAL** (proposed)
- **IP** (interior point method [T& A 2007])
- **PG** (projected gradient method [T& S 2008])
- **AG** (accelerated gradient method [Ji & Ye 2009])

Super-linear convergence



Benchmark experiment: BCI dataset

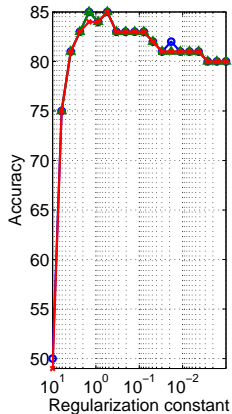
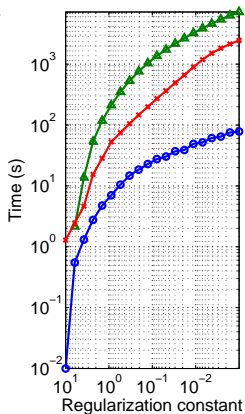
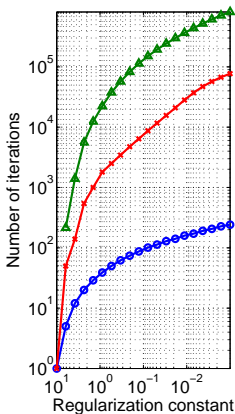
- BCI competition 2003 dataset IV.
- Task: predict the upcoming finger movement is **right or left**.
- Original data: multichannel EEG time-series 28 channels \times 50 time-points (500ms long)
- Each training example is preprocessed into **three matrices** and whitened [T&M 2010]:
 - First order (<20Hz) component: 28×50 matrix.
 - Second order (alpha 7-15Hz) component: 28×28 (covariance)
 - Second order (beta 15-30Hz) component: 28×28 (covariance)
- 316 training examples. 100 test samples.

BCI dataset (regularization path)

-o- M-DAL (proposed)

-△- PG (T&S 08)

-x- AG (Ji&Ye 09)



- Stopping criterion: $RDG \leq 10^{-3}$.
- Note: these are costs for obtaining the entire regularization path.

Summary

- M-DAL: Dual Augmented Lagrangian algorithm for **learning low-rank matrices**.
- Theoretically shown that M-DAL converges **superlinearly**
⇒ requires small number of steps.
- Separable regularizer + differentiable loss function
⇒ each step can be efficiently solved.
- Empirically we can solve matrix completion problem of size $10,000 \times 10,000$ in roughly 5min.
- Training a low-rank classifier that automatically combines multiple data sources (matrices) is sped up by a factor of 10–100.

Future work:

- Other AL algorithms.
- Large scale and unstructured (or badly conditioned) problems.

Background

- Convex conjugate of f

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x}))$$

- Convex conjugate of a sum:

$$(f + g)^*(\mathbf{y}) = (f^* \oplus g^*)(\alpha) = \inf_{\alpha} (f^*(\alpha) + g^*(\mathbf{y} - \alpha))$$

- Fenchel dual:

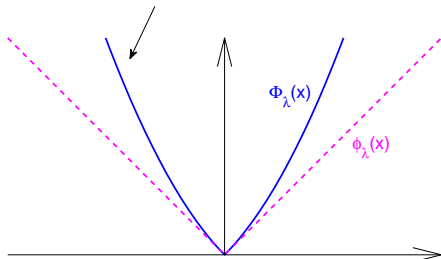
$$\inf_{\mathbf{x} \in \mathbb{R}^n} (f(\mathbf{x}) + g(\mathbf{x})) = \sup_{\alpha \in \mathbb{R}^n} (-f^*(\alpha) - g^*(-\alpha))$$

Inf-convolution

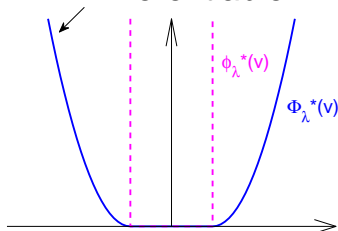
Inf-convolution (envelope function):

$$\begin{aligned}\Phi_{\lambda}^*(\mathbf{V}; \mathbf{W}^t) &= (\phi_{\lambda} + \frac{1}{2} \|\cdot - \mathbf{W}^t\|)^*(\mathbf{V}) \\ &= \inf_{\mathbf{Y} \in \mathbb{R}^{R \times C}} \left(\phi_{\lambda}^*(\mathbf{Y}) + \frac{1}{2} \|\mathbf{W}^t + \mathbf{V} - \mathbf{Y}\|^2 \right) \\ &= \Phi_{\lambda}^*(\mathbf{W}^t + \mathbf{V}; \mathbf{0}) = \Phi_{\lambda}^*(\mathbf{W}^t + \mathbf{V})\end{aligned}$$

Nondifferentiable



Differentiable



M-DAL algorithm for the trace-norm

$$\mathbf{W}^{t+1} = \text{ST}_{\lambda\eta_t} \left(\mathbf{W}^t + \eta_t \mathcal{A}^\top(\boldsymbol{\alpha}^t) \right)$$

where

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha}}{\text{argmin}} \underbrace{\left(f_\ell^*(-\boldsymbol{\alpha}) + \frac{1}{2\eta_t} \left\| \text{ST}_{\lambda\eta_t} \left(\mathbf{W}^t + \eta_t \mathcal{A}^\top(\boldsymbol{\alpha}) \right) \right\|^2 \right)}_{\varphi_t(\boldsymbol{\alpha})}$$

- $\varphi_t(\boldsymbol{\alpha})$ is differentiable:

$$\nabla \varphi_t(\boldsymbol{\alpha}) = -\nabla f_\ell^*(-\boldsymbol{\alpha}) + \mathcal{A}(\text{ST}_{\lambda\eta_t}(\mathbf{W}^t + \eta_t \mathcal{A}^\top(\boldsymbol{\alpha})))$$

$\nabla^2 \varphi_t(\boldsymbol{\alpha})$: a bit more involved but can be computed (Wright 92)

Minimizing the inner objective $\varphi_t(\alpha)$

- $\varphi_t(\alpha)$, $\nabla\varphi_t(\alpha)$ can be computed using only the **singular values** $\sigma_j(\mathbf{W}_\alpha^t) \geq \lambda\eta_t$ (and the corresponding SVs).
 \Rightarrow no need to compute full SVD of $\mathbf{W}_\alpha^t := \mathbf{W}^t + \eta_t\mathcal{A}^\top(\alpha)$.
- But, computation of $\nabla^2\varphi_t(\alpha)$ requires full SVD.

Consequence:

- When $\mathcal{A}^\top(\alpha)$ is **structured**, computing the above SVD is cheap.
 E.g., matrix completion: $\mathcal{A}^\top(\alpha)$ is sparse with only m non-zeros.

$$\boxed{\mathbf{W}_\alpha^t} = \boxed{\mathbf{W}_L^t} \boxed{\mathbf{W}_R^t} + \eta_t \begin{array}{|c|c|c|c|} \hline \cdot & & \cdot & \\ \hline & \cdot & & \cdot \\ \hline \cdot & & \cdot & \\ \hline & \cdot & & \cdot \\ \hline \cdot & & \cdot & \\ \hline & \cdot & & \cdot \\ \hline \cdot & & \cdot & \\ \hline & \cdot & & \cdot \\ \hline \end{array}$$

\Rightarrow Quasi-Newton method (maintains factorized \mathbf{W}^t and scalable)

- When $\mathcal{A}^\top(\alpha)$ is **not structured**,
 \Rightarrow Full Newton method (converges faster and more stable)

Primal Proximal Minimization A (M-DAL)

$$\mathbf{W}^{t+1} = \underset{\mathbf{W}}{\operatorname{argmin}} \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \left(\phi_{\lambda}(\mathbf{W}) + \frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|_{\text{fro}}^2 \right) \right)$$

- Primal Proximal Minimization B

$$\mathbf{W}^{t+1} = \underset{\mathbf{W}}{\operatorname{argmin}} \left(\phi_{\lambda}(\mathbf{W}) + \left(f_{\ell}(\mathcal{A}(\mathbf{W})) + \frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|_{\text{fro}}^2 \right) \right)$$

- Dual Proximal Minimization A (~ split Bregman iteration)

$$\alpha^{t+1} = \underset{\alpha}{\operatorname{argmin}} \left(f_{\ell}^*(-\alpha) + \left(\phi_{\lambda}^*(\mathcal{A}^{\top}(\alpha)) + \frac{1}{2\eta_t} \|\alpha - \alpha^t\|^2 \right) \right)$$

- Dual Proximal Minimization B

$$\alpha^{t+1} = \underset{\alpha}{\operatorname{argmin}} \left(\phi_{\lambda}^*(\mathcal{A}^{\top}(\alpha)) + \left(f_{\ell}^*(-\alpha) + \frac{1}{2\eta_t} \|\alpha - \alpha^t\|^2 \right) \right)$$

References

- Abernethy et al. A new approach to collaborative filtering: Operator estimation with spectral regularization. JMLR 10, 2009.
- Argyriou et al. Multi-task feature learning In NIPS 19, 2007.
- Ji & Ye. An accelerated gradient method for trace norm minimization. In ICML 2009.
- Rockafellar. Monotone operators and the proximal point algorithm. SIAM J. on Control and Optimization 14, 1976a.
- Rockafellar. Augmented Lagrangians and applications of the proximal point algorithm in convex programming. Math. of Oper. Res. 1, 1976b.
- Srebro et al. Maximum-margin matrix factorization. In NIPS 17, 2005.
- Tomioka & Aihara. Classifying matrices with a spectral regularization. In ICML 2007.
- Tomioka et al. Super-Linear Convergence of Dual Augmented-Lagrangian Algorithm for Sparsity Regularized Estimation. Arxiv:0911.4046, 2010.
- Tomioka & Sugiyama. Sparse learning with duality gap guarantee. In NIPS workshop OPT 2008 Optimization for Machine Learning, 2008.
- Tomioka & Müller. A regularized discriminative framework for EEG analysis with application to brain-computer interface. Neuroimage 49, 2010.
- Wright. Differential equations for the analytic singular value decomposition of a matrix. Numer. Math. 63, 1992.