

Projection-based Dual Averaging For Stochastic Sparse Optimization

Asahi Ushio, Masahiro Yukawa

Dept. Electronics and Electrical Engineering, Keio University, Japan.
E-mail: ushio@ykw.elec.keio.ac.jp, yukawa@elec.keio.ac.jp

The 42nd ICASSP, @New Orleans, LA, USA
March 5-9, Mar, 2017

Introduction

Background Now, we have many kinds of data (SNS, IoT, Digital News) [1].
 → **Machine learning and signal processing become more important !**

Challenges (demands for algorithms)

- Process **real-time streaming data**.
- Achieve **low complexity**.
- Deal with **high-dimensionality** and **sparsity** of data.

[1] McKinsey, "Big Data: The next frontier for innovation, competition, and productivity," 2011.

In Signal Processing ...

- **Squared distance cost:**
 Projection-based method
 ex) NLMS, APA, APFBS [2]
- **Change Geometry:**
 PAPA, Variable Metric [3]

In Machine Learning ...

- **Regularized Dual Averaging (RDA) type:** RDA [4]
- **Forward Backward Splitting (FBS) type:** FOBOS [5]
- **Change Geometry:** ADA GRAD [6]

[2] Murakami *et al.*, 2010, [3] Yukawa *et al.*, 2009, [4] Xiao, 2009, [5] Singer *et al.*, 2009, [6] Duchi *et al.*, 2011

In Signal Processing ...

Background Now, we have many kinds of data (SNS, IoT, Digital News) [1].
 → **Machine learning and signal processing become more important !**

Challenges (demands for algorithms)

- Process **real-time streaming data**.
- Achieve **low complexity**.
- Deal with **high-dimensionality** and **sparsity** of data.

[1] McKinsey, "Big Data: The next frontier for innovation, competition, and productivity," 2011.

In Signal Processing ...

- **Squared distance cost:**
 Projection-based method
 ex) NLMS, APA, APFBS [2]
- **Change Geometry:**
 PAPA, Variable Metric [3]

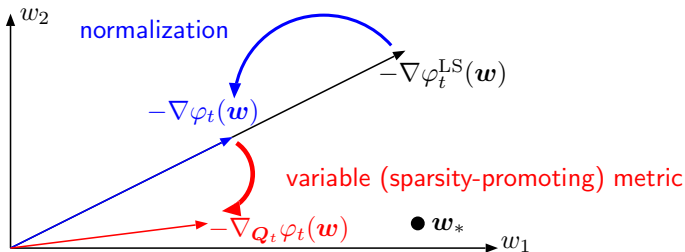
In Machine Learning ...

- **Regularized Dual Averaging (RDA) type:** RDA [4]
- **Forward Backward Splitting (FBS) type:** FOBOS [5]
- **Change Geometry:** ADA GRAD [6]

[2] Murakami *et al.*, 2010, [3] Yukawa *et al.*, 2009, [4] Xiao, 2009, [5] Singer *et al.*, 2009, [6] Duchi *et al.*, 2011

An Illustration of Projection-based Method

in case of online regression



Ordinary least square cost.

$$\varphi_t^{\text{LS}}(\mathbf{w}) := \frac{1}{2} (y_t - \mathbf{w}^\top \mathbf{x}_t)^2$$

$\mathbf{x}_t \in \mathbb{R}^n$: the input vector
 $y_t \in \mathbb{R}$: the output
 $\mathbf{w}_* \in \mathbb{R}^n$: the unknown vector

Squared distance cost with Q_t -metric.

$$\varphi_t(\mathbf{w}) = \frac{1}{2} \left(\frac{y_t - \langle \mathbf{w}, \mathbf{x}_t \rangle_{Q_t}}{\|\mathbf{x}_t\|_{Q_t}} \right)^2$$

$Q_t \in \mathbb{R}^{n \times n}$: a positive definite matrix
 Q_t -norm : $\|\mathbf{w}\|_{Q_t} := \sqrt{\langle \mathbf{w}, \mathbf{w} \rangle_{Q_t}}$,
 $\langle \mathbf{w}, \mathbf{z} \rangle_{Q_t} := \sqrt{\mathbf{w}^\top Q_t \mathbf{z}}$ for $\mathbf{w}, \mathbf{z} \in \mathbb{R}^n$

In Machine Learning ...

Background Now, we have many kinds of data (SNS, IoT, Digital News) [1].
 → **Machine learning and signal processing become more important !**

Challenges (demands for algorithms)

- Process **real-time streaming data**.
- Achieve **low complexity**.
- Deal with **high-dimensionality** and **sparsity** of data.

[1] McKinsey, "Big Data: The next frontier for innovation, competition, and productivity," 2011.

In Signal Processing ...

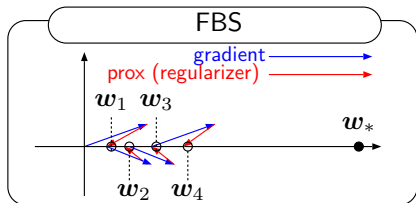
- **Squared distance cost:**
 Projection-based method
 ex) NLMS, APA, APFBS [2]
- **Change Geometry:**
 PAPA, Variable Metric [3]

In Machine Learning ...

- **Regularized Dual Averaging (RDA) type:** RDA [4]
- **Forward Backward Splitting (FBS) type:** FOBOS [5]
- **Change Geometry:** ADA GRAD [6]

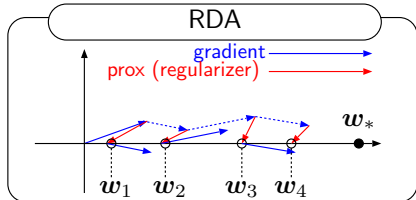
[2] Murakami *et al.*, 2010, [3] Yukawa *et al.*, 2009, [4] Xiao, 2009, [5] Singer *et al.*, 2009, [6] Duchi *et al.*, 2011

FBS type versus RDA type



- **FBS:** The effects of the proximity operator accumulate over the iterations.

Tradeoff between the sparsity level and the estimation accuracy.



- **RDA:** FREE from the accumulation.

A high level of sparsity comes with high estimation accuracy !

Abstract of This Work

Motivation

- 1 **RDA**: Sparse solution.
- 2 **Squared distance cost**: Stable adaptation.
- 3 **Variable metric**: Promoting sparsity to improve the performance.

Combination of 3 properties.

Sparsity-promoting and stable learning !

Proposed Algorithm

Projection-based Dual Averaging (PDA)

- Features: **RDA** with **Squared distance cost**, employing **variable metric**.
- Show efficacy by numerical examples (simulated data, real data).

Preliminaries

Problem Setting : an online regularized optimization

$$\min_{\mathbf{w} \in \mathbb{R}^n} \mathbb{E} [\varphi_t(\mathbf{w})] + \psi_t(\mathbf{w}), \quad t \in \mathbb{N} \quad (1)$$

ψ_t, φ_t : a possibly nonsmooth function

\mathbf{w} : supposed to be sparse or compressible

Basic Stochastic Optimization Methods

■ SGD

$$\mathbf{w}_t := \mathbf{w}_{t-1} - \eta \nabla \varphi_t(\mathbf{w}_{t-1}), \quad \eta > 0 \quad (2)$$

■ Dual Averaging [1]

$$\mathbf{w}_t := \arg \min_{\mathbf{w} \in \mathbb{R}^n} \left(\left\langle \frac{\sum_{i=1}^t \nabla \varphi_{i-1}(\mathbf{w}_{i-1})}{t}, \mathbf{w} \right\rangle + \mu_t h(\mathbf{w}) \right), \quad \mu_t = \mathcal{O} \left(\frac{1}{\sqrt{t}} \right) \quad (3)$$

$h(\mathbf{w})$: the prox-function

[1] Y. Nesterov, "Primal-dual subgradient methods for convex problems", Mathematical Programming, 2009.

Projection-based Methods

Cost Function of Projection-based Methods

- The Q_t -metric distance cost (normalized MSE in regression case).

$$\varphi_t(\mathbf{w}) := \frac{1}{2} d_{Q_t}^2(\mathbf{w}, C_t) \quad (4)$$

$$d_{Q_t}(\mathbf{w}, C_t) := \min_{z \in C_t} \|\mathbf{w} - z\|_{Q_t} \quad (5)$$

$Q_t \in \mathbb{R}^{n \times n}$: a positive definite matrix

$C_t \subset \mathbb{R}^n$: a closed convex set

Q_t -norm for $\mathbf{w}, \mathbf{z} \in \mathbb{R}^n$: $\|\mathbf{w}\|_{Q_t} := \sqrt{\langle \mathbf{w}, \mathbf{w} \rangle_{Q_t}}, \langle \mathbf{w}, \mathbf{z} \rangle_{Q_t} := \sqrt{\mathbf{w}^\top Q_t \mathbf{z}}$

Gradient Calculation

- The Q_t -gradient of φ_t

$$\mathbf{g}_t := \nabla_{Q_t} \varphi_t(\mathbf{w}_{t-1}) = \mathbf{w}_{t-1} - P_{C_t}^{Q_t}(\mathbf{w}_{t-1}), \quad \mathbf{w}_{t-1} \in \mathbb{R}^n. \quad (6)$$

- The Q_t -projection onto C_t

$$P_{C_t}^{Q_t}(\mathbf{w}) := \arg \min_{z \in C_t} \|\mathbf{w} - z\|_{Q_t}. \quad (7)$$

Proposed Algorithm

Projection-based Dual Averaging (PDA)

$$\begin{aligned}
 \mathbf{w}_t &:= \arg \min_{\mathbf{w} \in \mathbb{R}^n} \left(\langle \mathbf{s}_t, \mathbf{w} \rangle_{\mathbf{Q}_t} + \frac{\|\mathbf{w}\|_{\mathbf{Q}_t}^2}{2\eta} + \psi_t(\mathbf{w}) \right) \\
 &= \arg \min_{\mathbf{w} \in \mathbb{R}^n} \left(\eta\psi_t(\mathbf{w}) + \frac{1}{2} \|\mathbf{w} + \eta\mathbf{s}_t\|_{\mathbf{Q}_t}^2 \right) \\
 &= \text{prox}_{\eta\psi_t}^{\mathbf{Q}_t}(-\eta\mathbf{s}_t), \quad \mathbf{s}_t = \sum_{i=1}^t \nabla_{\mathbf{Q}_t} \varphi_{i-1}(\mathbf{w}_{i-1}), \quad \eta \in [0, 2] \quad (8)
 \end{aligned}$$

■ The proximity operator

$$\text{prox}_{\eta\psi_t}^{\mathbf{Q}_t}(\mathbf{w}) := \arg \min_{\mathbf{z} \in \mathbb{R}^n} \left(\eta\psi_t(\mathbf{w}) + \frac{1}{2} \|\mathbf{w} - \mathbf{z}\|_{\mathbf{Q}_t}^2 \right), \quad \forall \mathbf{w} \in \mathbb{R}^n. \quad (9)$$

Application to Online Regression

Problem Setting : Online Regression

$$y_t := \mathbf{w}_*^T \mathbf{x}_t + \nu_t \in \mathbb{R}$$

- $\mathbf{x}_t \in \mathbb{R}^n$: the input vector
- $y_t \in \mathbb{R}$: the output
- $\mathbf{w}_* \in \mathbb{R}^n$: the unknown vector
- $\nu_t \in \mathbb{R}$: the additive white noise

Definition of the Projection

■ The linear variety

$$C_t := \arg \min_{\mathbf{w} \in \mathbb{R}^n} \|\mathbf{X}_t^T \mathbf{w} - \mathbf{y}_t\|_{I_n}. \quad (10)$$

$$\mathbf{X}_t := [\mathbf{x}_t \cdots \mathbf{x}_{t-r+1}] \in \mathbb{R}^{n \times r}$$

$$\mathbf{y}_t := [y_t, \dots, y_{t-r+1}]^T \in \mathbb{R}^r$$

■ The projection onto C_t

$$P_{C_t}^{\mathbf{Q}_t}(\mathbf{w}_{t-1}) := \mathbf{w}_{t-1} - \mathbf{Q}_t^{-1} \mathbf{X}_t^\dagger (\mathbf{X}_t^T \mathbf{w}_{t-1} - \mathbf{y}_t). \quad (11)$$

■ The Moore-Penrose pseudo-inverse \mathbf{X}_t^\dagger

$$\mathbf{X}_t (\mathbf{X}_t^T \mathbf{Q}_t^{-1} \mathbf{X}_t + \delta \mathbf{I}_n)^{-1}, \quad \delta > 0.$$

Design of Metric Q_t & Regularizer ψ_t

Design of metric Q_t (following [Yukawa *et al.* 2010])

$$Q_t := \frac{\alpha}{n} I_n + \frac{1-\alpha}{S_t} \tilde{Q}_t^{-1}, \quad \alpha \in [0, 1]. \quad (12)$$

- $\tilde{Q}_t := \text{diag}(|w_{t-1,1}|, \dots, |w_{t-1,n}|) + \epsilon I_n$ for some $\epsilon > 0$
- $S_t := \sum_{i=1}^n (|w_{t-1,i}| + \epsilon)^{-1}$ to reduce $\text{tr}(I_n/n) = \text{tr}(\tilde{Q}_t^{-1}/S_t) = 1$

The Regularization Term ψ_t

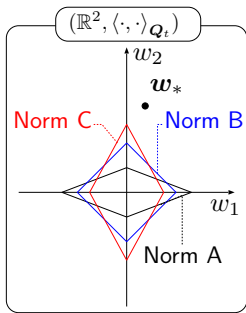
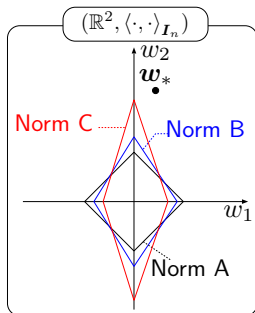
$$\psi_t(\mathbf{w}) = \lambda \|\mathbf{w}\|_{Q_t^2, 1} := \lambda \sum_{i=1}^n q_{t,i}^2 |w_i|, \quad \lambda > 0. \quad (13)$$

- $Q_t := \text{diag}(q_{t,1}, \dots, q_{t,n}) \in \mathbb{R}^{n \times n}$
- **The proximity operator** for the ψ_t in (13)

$$\text{prox}_{\eta \psi_t}^{Q_t}(\mathbf{w}) = \sum_{i=1}^n \mathbf{e}_i \text{sgn}(w_i) [|w_i| - q_{t,i} \lambda \eta]_+. \quad (14)$$

- $\{\mathbf{e}_i\}_{i=1}^n$: the standard basis of \mathbb{R}^n

The Hilbert spaces $(\mathbb{R}^2, \langle \cdot, \cdot \rangle_{I_n})$ and $(\mathbb{R}^2, \langle \cdot, \cdot \rangle_{Q_t})$



* **Norm A**

$$\|w\|_{I_n,1} := \sum_{i=1}^n |w_i|$$

* **Norm B**

$$\|w\|_{Q_t^1,1} := \sum_{i=1}^n q_{t,i} |w_i|$$

* **Norm C (we employ)**

$$\|w\|_{Q_t^2,1} := \sum_{i=1}^n q_{t,i}^2 |w_i|$$

- **Norm A** gives a **fat** unit ball in $(\mathbb{R}^2, \langle \cdot, \cdot \rangle_{Q_t})$:

The proximity operator shrinks the large component more than the small one.

Undesirable bias.

- **Norm C** gives a **tall** unit ball in $(\mathbb{R}^2, \langle \cdot, \cdot \rangle_{Q_t})$:

Shrink the small components more.

Reduce the bias !

Relation to Prior Work

SGD type: SGD, **NLMS**, **APA**, **PAPA**

Sparsity-promoting: FBS (Forward Backward Splitting) type
 FOBOS [1], ADA GRAD-FBS [2], **APFBS** [3]

Dual Averaging type: Dual Averaging [4]

Sparsity-promoting: RDA (Regularized Dual Averaging) type
 RDA [5], ADA GRAD-RDA [2]

*bold : projection-based method

[1] Singer *et al.*, 2009 [2] Duchi *et al.*, 2011 [3] Murakami *et al.*, 2010 [4] Nesterov, 2009 [5] Xiao, 2009

	Ordinal Cost Function	Projection-based
FBS type	FOBOS, ADA GRAD-FBS	APFBS
RDA type	RDA, ADA GRAD-RDA	

Relation to Prior Work

SGD type: SGD, **NLMS**, **APA**, **PAPA**

Sparsity-promoting: FBS (Forward Backward Splitting) type
 FOBOS [1], ADAGRAD-FBS [2], **APFBS** [3]

Dual Averaging type: Dual Averaging [4]

Sparsity-promoting: RDA (Regularized Dual Averaging) type
 RDA [5], ADAGRAD-RDA [2], **PDA**

*bold : projection-based method

[1] Singer *et al.*, 2009 [2] Duchi *et al.*, 2011 [3] Murakami *et al.*, 2010 [4] Nesterov, 2009 [5] Xiao, 2009

	Ordinal Cost Function	Projection-based
FBS type	FOBOS, ADAGRAD-FBS	APFBS
RDA type	RDA, ADAGRAD-RDA	PDA

Numerical Example

■ Experiment:

1 Sparse-System Estimation (Simulated Data)

Model $y_t := \mathbf{w}_*^T \mathbf{x}_t + \nu_t \in \mathbb{R}$

- The proportion of the zero components of $\mathbf{w}_* \in \mathbb{R}^{1000}$ is 90%.
- The noise ν_t is zero-mean i.i.d. Gaussian with variance 0.01.

2 Echo Cancellation (Real Data)

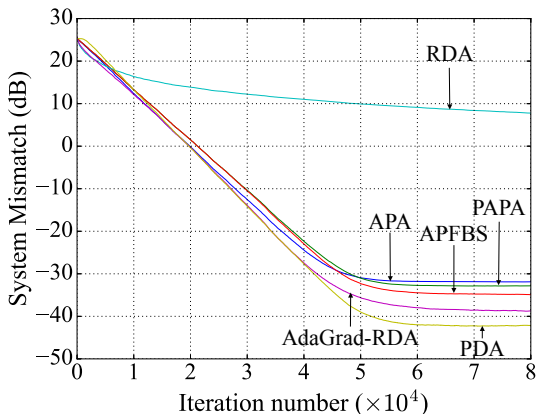
- The sampling frequency of speech signal and echo path is 8000 Hz.
- The learning is stopped whenever the amplitude of input is below 10^{-4} .
- The noise is zero-mean i.i.d. Gaussian with the signal noise ratio 20 dB.

■ Compared Algorithms: APA, PAPA, APFBS, RDA, ADA GRAD-RDA

	APA	PAPA	APFBS	RDA	ADA GRAD-RDA	PDA
Cost	$\varphi_t(\mathbf{w})$	$\varphi_t(\mathbf{w})$	$\varphi_t(\mathbf{w})$	$\varphi_t^{\text{LS}}(\mathbf{w})$	$\varphi_t^{\text{LS}}(\mathbf{w})$	$\varphi_t(\mathbf{w})$
Reg	-	-	$\lambda \ \mathbf{w}\ _{Q_t^2, 1}$	$\lambda \ \mathbf{w}\ _{I_n, 1}$	$\lambda \ \mathbf{w}\ _{I_n, 1}$	$\lambda \ \mathbf{w}\ _{Q_t^2, 1}$

Experiment 1: Sparse -System Estimation

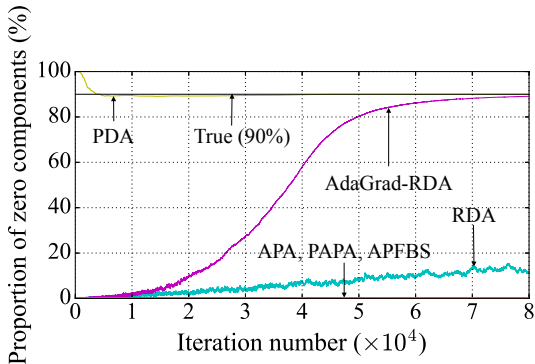
System Mismatch



- $\|w_* - w_t\|_{I_n}^2 / \|w_*\|_{I_n}^2$, w_* is true parameter, w_t is estimation at t .
- **PDA shows the best performance.**

Experiment 1: Sparse -System Estimation

Proportion of the Zero Components of the Estimated Coefficient Vector

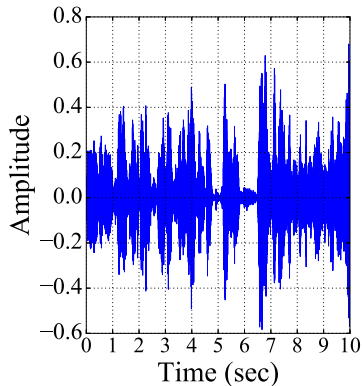


Algorithms	APA	PAPA	APFBS	RDA	ADAGRAD-RDA	PDA
Proportion	0%	0%	0%	11.6%	89%	90%

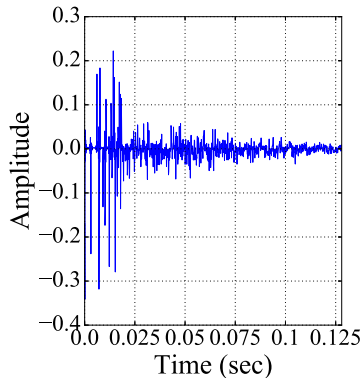
- **PDA achieve accurate sparsity.**

Experiment 2: Echo Cancellation

Amplitudes of Speech Signal and Echo Path



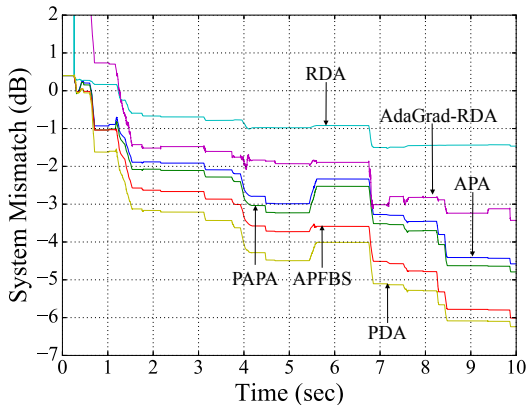
(a) Speech signal



(b) Echo path

Experiment 2: Echo Cancellation

System Mismatch



- PDA shows the best performance.

Conclusion

Conclusion

- We proposed the **projection-based dual averaging (PDA)** algorithm.
 - **projection-based**: Input-vector normalization and the sparsity-seeking variable-metric
 - **RDA**: Better sparsity-seeking.
- An application of PDA to an **online regression** problem was presented.
 - The numerical examples demonstrated the **better sparsity-seeking** and **learning properties**.

Future Work

- Application to machine learning problems (classification).
- Self-tuning method for λ and α .

Parameters for Experiments

Table: Parameters for sparse-system estimation.

Algorithms	η	λ	α	r	δ	ϵ
APA	0.16	-	-	1	10^{-5}	-
PAPA	0.14	-	0.8	1	10^{-5}	10^{-5}
APFBS	0.14	10^{-3}	0.8	1	10^{-5}	10^{-5}
RDA	0.01	10^{-3}	-	-	-	-
ADAGRAD	0.17	10^{-3}	-	-	-	-
PDA	0.13	3×10^3	0.8	1	10^{-5}	10^{-5}

Table: Parameters for echo cancellation.

Algorithms	η	λ	α	r	δ	ϵ
APA	0.3	-	-	2	10^{-15}	-
PAPA	0.3	-	0.2	2	10^{-15}	10^{-15}
APFBS	0.2	10^{-2}	0.3	2	10^{-15}	10^{-15}
RDA	1	10^{-4}	-	-	-	-
ADAGRAD	0.3	10^{-4}	-	-	-	-
PDA	0.2	25.5	0.3	2	10^{-15}	10^{-15}