

ART.T458 Advanced Machine learning

6 Advances in gradient based methods

Masamichi Shimosaka

Associate Professor
Dept. of Computer Science

This course

Foundations of machine learning (1st half)

- Interdisciplinary fields of
 - probability theory
 - statistics
 - optimization

Recent topic(s) of machine learning (2nd half)

- Deep learning!
- Deep learning for X!

Today's topic

Today's Model

- No new (prediction) models are introduced today!

Theoretical part: advances in gradient methods

- Accelerated proximal gradient methods
 - Convergence rate of proximal gradient methods
 - Acceleration, momentum, FISTA: Accelerated proximal gradient method
- Subgradient method
- Projected gradient method
- Stochastic optimization
 - From batch to stochastic way
 - vs. batch training in terms of convergence rate & computational cost
- Adaptive step size in gradient methods
 - AdaGrad, adam

Review: Proximal gradient model

Simple gradient based model for non-smooth function

- Loss function: smooth convex function
- Regularization: (non-)smooth convex function

$$\operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1 \right)$$

$\psi(\mathbf{w})$ $\Omega(\mathbf{w})$

- Minimization of upper bound of original problem sequentially

$$\operatorname{argmin}_{\mathbf{w}} \left(\psi(\mathbf{w}^{(0)}) + \nabla \psi(\mathbf{w}^{(0)})^\top (\mathbf{w} - \mathbf{w}^{(0)}) + \frac{\gamma}{2} \|\mathbf{w} - \mathbf{w}^{(0)}\|_2^2 + \Omega(\mathbf{w}) \right)$$

Upper bound of $\psi(\mathbf{w})$ around \mathbf{w}^0

$\Omega(\mathbf{w})$

- The following update is iterated until its convergence $\tilde{\psi}(\mathbf{w}|\mathbf{w}^{(0)})$

$$\mathbf{w}^{(t)} = \operatorname{prox}_{\eta_t \Omega} \left(\mathbf{w}^{(t-1)} - \eta_t \nabla \psi(\mathbf{w}^{(t-1)}) \right)$$

Gradient descent update

Proximal operation after gradient descent update

Review: Monotonic decrease property in proximal gradient methods

With proximal gradient method, the following inequality holds (a1:derive this inequality)**

$$\psi(\mathbf{w}^{(t+1)}) + \lambda \|\mathbf{w}^{(t+1)}\|_1 \leq \tilde{\psi}(\mathbf{w}^{(t+1)} | \mathbf{w}^{(t)}) \leq \tilde{\psi}(\mathbf{w}^{(t)} | \mathbf{w}^{(t)}) = \psi(\mathbf{w}^{(t)}) + \lambda \|\mathbf{w}^{(t)}\|_1$$

Objective function is monotonically decreased

- Since the function is non-negative (i.e. objective has lower bound), the parameter converges
- Actually, we can use this update for non γ -smooth function (similar to backtracking technique)

Review: Monotonic decrease property in proximal gradient methods

With proximal gradient method, the following inequality holds (a1:derive this inequality)**

$$\psi(\mathbf{w}^{(t+1)}) + \lambda \|\mathbf{w}^{(t+1)}\|_1 \leq \tilde{\psi}(\mathbf{w}^{(t+1)} | \mathbf{w}^{(t)}) \leq \tilde{\psi}(\mathbf{w}^{(t)} | \mathbf{w}^{(t)}) = \psi(\mathbf{w}^{(t)}) + \lambda \|\mathbf{w}^{(t)}\|_1$$

Objective function is monotonically decreased

- Since the function is non-negative (i.e. objective has lower bound), the parameter converges
- Actually, we can use this update for non γ -smooth function (similar to backtracking technique)

We need to consider the property of the algorithm in terms of **convergence rate**

Review: Monotonic decrease property in proximal gradient methods

With proximal gradient method, the following inequality holds (a1:derive this inequality)**

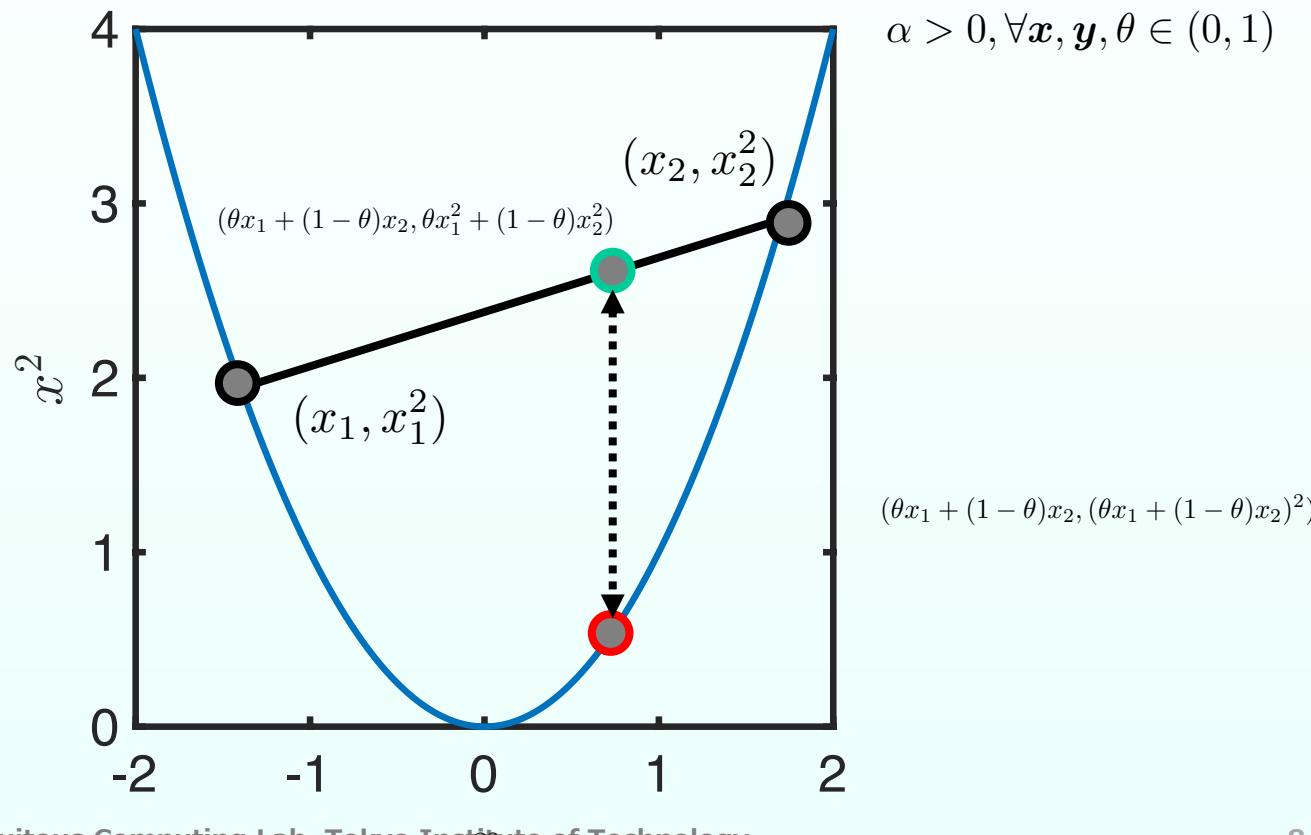
$$\psi(\mathbf{w}^{(t+1)}) + \lambda \|\mathbf{w}^{(t+1)}\|_1 \leq \tilde{\psi}(\mathbf{w}^{(t+1)} | \mathbf{w}^{(t)}) \leq \tilde{\psi}(\mathbf{w}^{(t)} | \mathbf{w}^{(t)}) = \psi(\mathbf{w}^{(t)}) + \lambda \|\mathbf{w}^{(t)}\|_1$$

Objective function is monotonically decreased

- Since the function is non-negative (i.e. objective has lower bound), the parameter converges
- Actually, we can use this update for non γ -smooth function (similar to backtracking technique)
- Convergence rate
 - Depends on types of loss function
 - Similar to smoothness, strength in convexity should be introduced.

α -strongly convex function

- Important property of loss function as well as γ -smooth function
 - c.f. Convex: $f(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) \quad \forall \mathbf{x}, \mathbf{y}, \theta \in [0, 1]$
- Definition : $\frac{\alpha}{2}\theta(1 - \theta)||\mathbf{x} - \mathbf{y}||_2^2 + f(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y})$
 - Curvature of f is not too close to zero



α -strongly convex function (cont'd)

- If function $f(\mathbf{x})$ is α -strongly convex function, $f(\mathbf{x}) - \frac{\alpha}{2} \|\mathbf{x}\|_2^2$ is convex
- If function $f(\mathbf{x})$ is α -strongly convex function, the following inequality holds, $\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}$ (i.e. $\nabla^2 f(\mathbf{x}) - \alpha \mathbf{I} \succeq \mathbf{0}$)
- Confirm whether the following functions are strongly convex or not (**a2)
 - $x \in \mathbb{R}, \exp(x)$
 - $x \in \mathbb{R}, x^2$
 - $x \in \mathbb{R}, x^4$

Review: Convergence rate

Q-linear

- Parameter can be exponentially converged

$$\exists r \in (0, 1), k' \in \mathbb{N}, \forall k \geq k', \frac{\|\mathbf{w}^{(k+1)} - \hat{\mathbf{w}}\|}{\|\mathbf{w}^{(k)} - \hat{\mathbf{w}}\|} \leq r$$

We use another metric here

- Estimating # of iteration when $J(\mathbf{w}^{(t)}) - J(\hat{\mathbf{w}}) \leq \epsilon$ holds

$$J(\mathbf{w}) = \psi(\mathbf{w}) + \Omega(\mathbf{w})$$

 
Loss Regularization
function

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$

- Ex. If $J(\mathbf{w}^{(t)}) - J(\hat{\mathbf{w}}) \leq \frac{2\gamma}{(t+1)^2} \|\hat{\mathbf{w}} - \mathbf{w}^{(0)}\|_2^2$,

$J(\mathbf{w}^{(t)}) - J(\hat{\mathbf{w}}) \leq \epsilon$ reaches after $O(1/\sqrt{\epsilon})$ iterations

Convergence rate of (proximal) gradient method

Convergence rate does not depend on regularization

- This result can be applied to standard gradient descent

$\psi(w)$	Non α -strongly convex	α -strongly convex
Non γ -smooth	$O(1/\epsilon^2)$	$O(1/\epsilon)$
γ -smooth	$O(1/\epsilon)$	$O(-\ln \epsilon)$ Corresponds to Q-linear

Acceleration of (proximal) gradient

Original of proximal gradient

$$\underline{\underline{w^{(t)} = \text{prox}_{\eta_t \Omega} \left(w^{(t-1)} - \eta_t \nabla \psi(w^{(t-1)}) \right)}}$$

Gradient descent update

Proximal operation after gradient descent update

Accelerated proximal gradient: use history of w

$$\underline{\underline{w^{(t)} = \text{prox}_{\eta_t \Omega} \left(v^{(t-1)} - \eta_t \nabla \psi(v^{(t-1)}) \right)}}$$

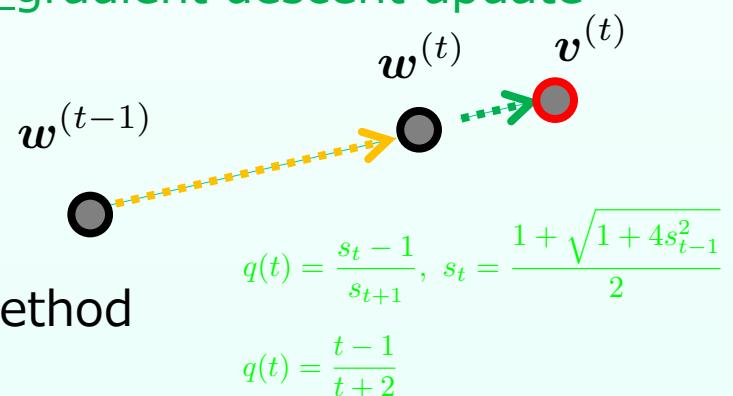
Accelerated gradient descent update

Proximal operation after accelerated gradient descent update

$$\underline{\underline{v^{(t)} = w^{(t)} + q(t)(w^{(t)} - w^{(t-1)})}}$$

Some scalar function
(carefully designed)

- C.f. Nesterov acceleration, momentum method



Convergence rate of accelerated (proximal) gradient method

Convergence rate does not depend on regularization

- This result can be applied to standard gradient descent

$\psi(w)$	Non α -strongly convex	α -strongly convex
Non γ -smooth	$O(1/\epsilon^2)$	$O(1/\epsilon)$
γ -smooth	$O(1/\epsilon)$  $O(1/\sqrt{\epsilon})$	$O(-\ln \epsilon)$  $O(-\ln \epsilon)$

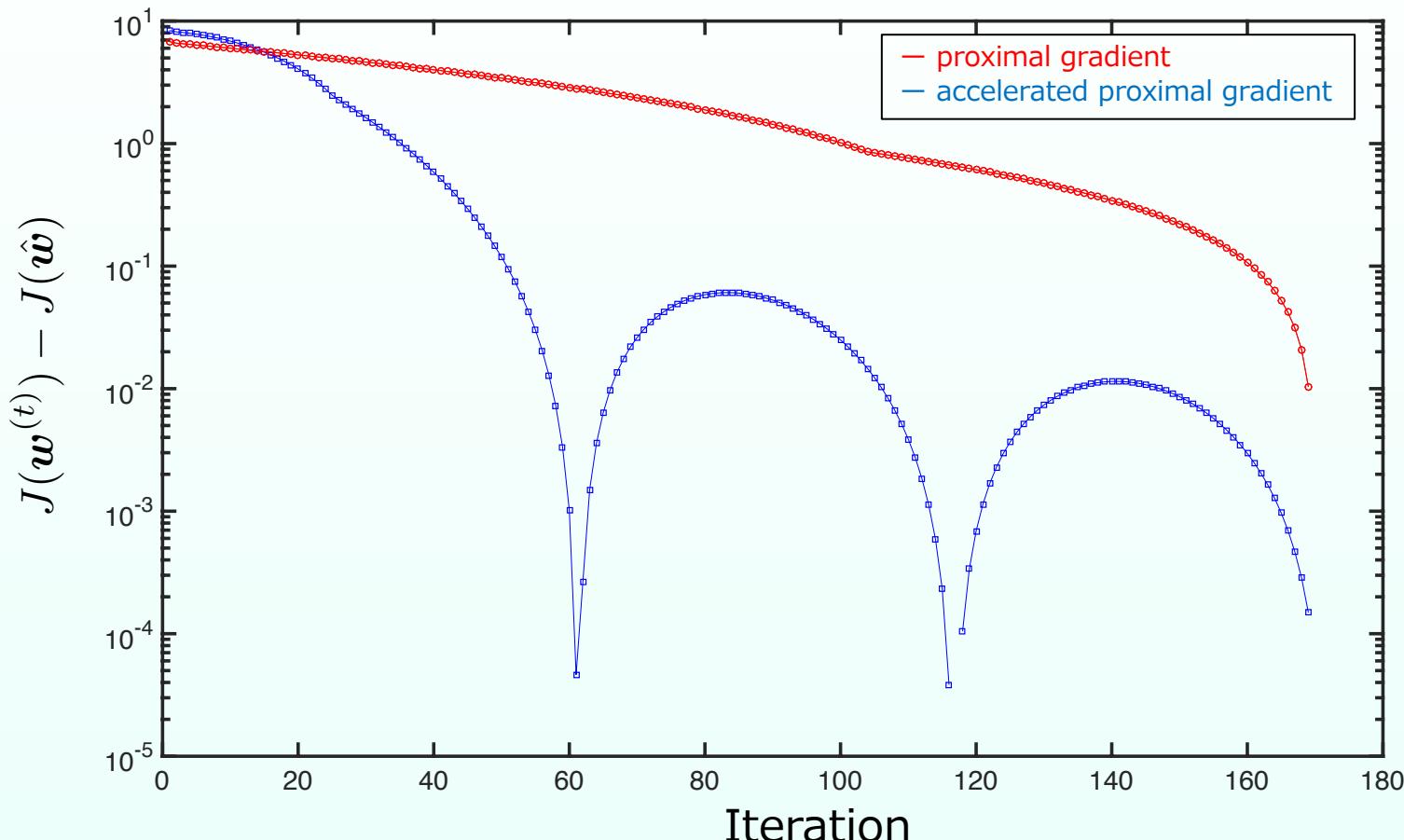
with acceleration Rate will be better with acceleration
(even if it is in the same order)

Acceleration has relation to momentum in NN community
FISTA: Fast ISTA: accelerated proximal gradient

FISTA in action with toy example

$$J(\mathbf{w}) = (\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A}(\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1$$

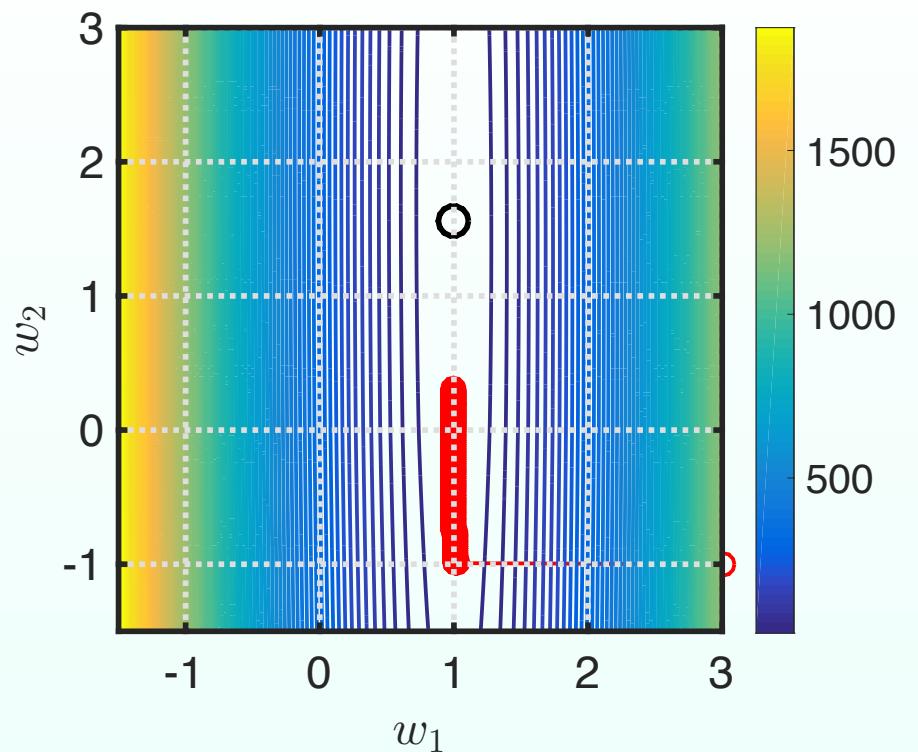
$$\mathbf{A} = \begin{pmatrix} 300 & 0.5 \\ 0.5 & 1 \end{pmatrix}, \boldsymbol{\mu}^\top = \begin{pmatrix} 1 & 2 \end{pmatrix}, \lambda = 0.89$$



FISTA in action with toy example (cont'd)

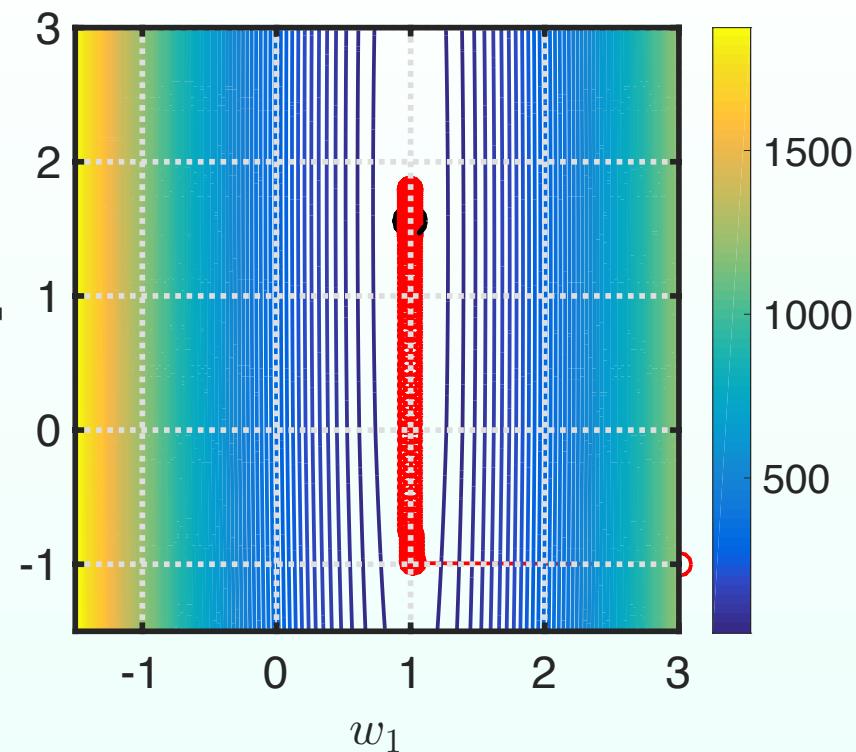
$$J(\mathbf{w}) = (\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A}(\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1$$

Proximal gradient



$$\mathbf{A} = \begin{pmatrix} 300 & 0.5 \\ 0.5 & 1 \end{pmatrix}, \boldsymbol{\mu}^\top = \begin{pmatrix} 1 & 2 \end{pmatrix}, \lambda = 0.89, \eta_t = 1/\gamma$$

Accelerated Proximal gradient



In some condition, accelerated proximal gradient gets faster convergence

Subgradient method

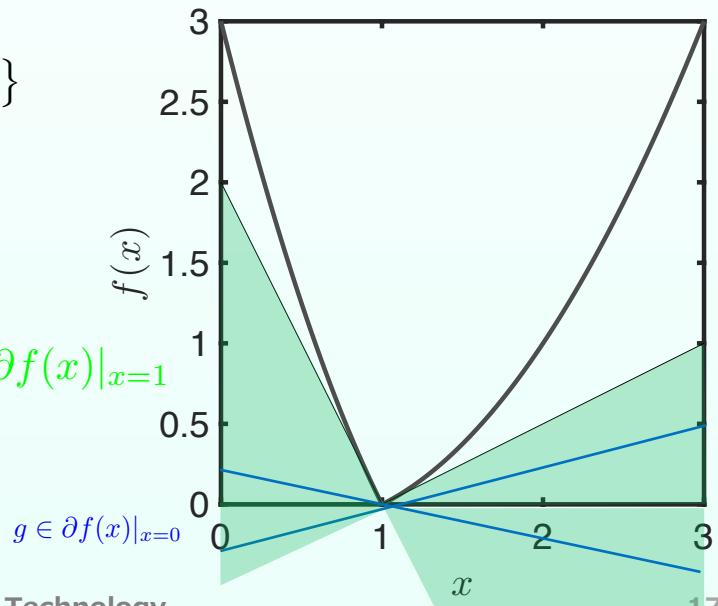
Subgradient

Generalized derivative for non-smooth function

- Ex. of non-smooth function
 - $f(\mathbf{w}) = \|\mathbf{w}\|_1$ we could not obtain derivative at $\mathbf{w} = 0$
 - $f(w) = \max(0, 1 - w)$ we could not obtain derivative at $w = 1$

Definition of subderivative

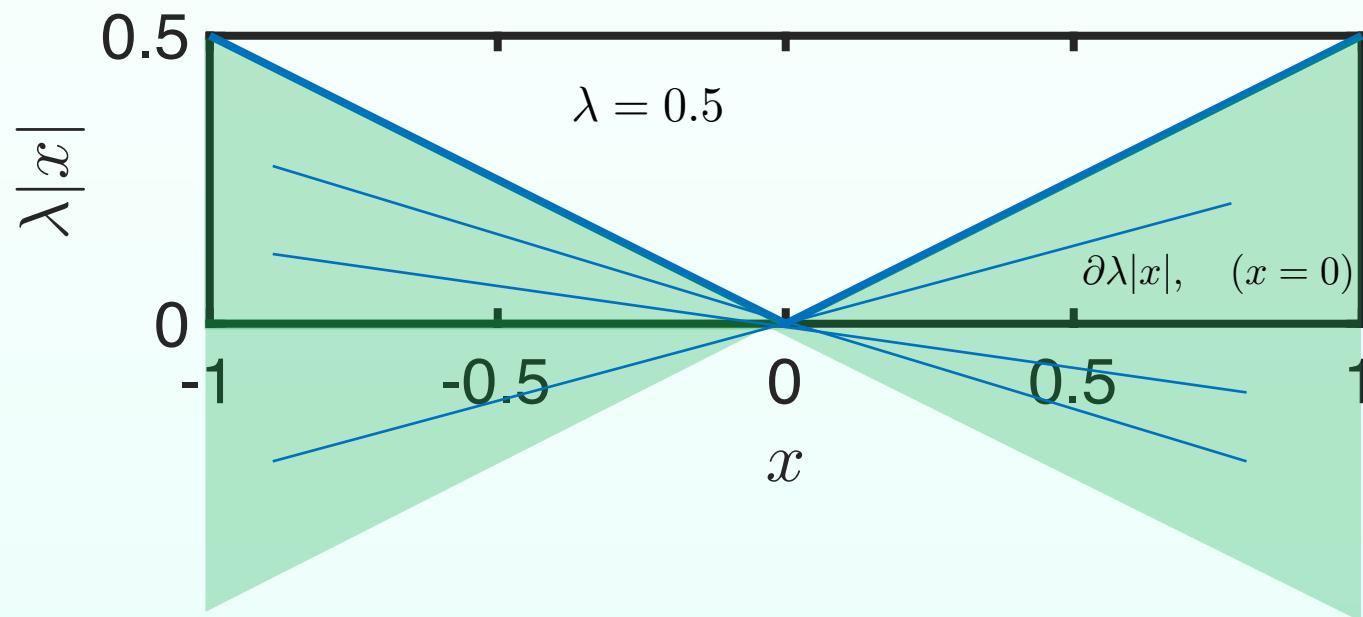
- Subderivative of convex function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is defined as
$$\partial f(\mathbf{w}) = \{\mathbf{g} \in \mathbb{R}^d | (\tilde{\mathbf{w}} - \mathbf{w})^\top \mathbf{g} + f(\mathbf{w}) \leq f(\tilde{\mathbf{w}}), \forall \tilde{\mathbf{w}} \in \mathbb{R}^d\}$$
- Subgradient $\mathbf{g} \in \partial f(\mathbf{w})$
- In convex function, this definition is equiv.
to 1st condition of convex function



Subgradient of L1 norm

L1 norm in 1d $f(x) = \lambda|x|, \lambda > 0$

$$g = \begin{cases} \lambda & \text{if } x > 0 \\ \eta\lambda & \text{if } x = 0, \eta \in [-1, 1] \\ -\lambda & \text{if } x < 0 \end{cases}$$



Necessary condition of local minima

Necessary condition of local minima with subgradient

- Even if the objective is non-smooth, we can reuse the necessary condition for smooth function
 - For continuous differentiable function $\frac{\partial \psi(\mathbf{w})}{\partial \mathbf{w}} + \frac{\partial \lambda \Omega(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{0}$
 - For non-smooth function $\mathbf{0} \in \partial \psi(\mathbf{w}) + \lambda \partial \Omega(\mathbf{w})$
- Ex. Lasso
$$\operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1 \right)$$
- Optimality condition can be efficiently obtained without Lagrangian
$$-\lambda \mathbf{1} \leq \mathbf{X}^\top (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \leq \lambda \mathbf{1}$$

Hinge-loss minimization with subgradient methods

Review: L2 regularized hinge-loss minimization (SVM)

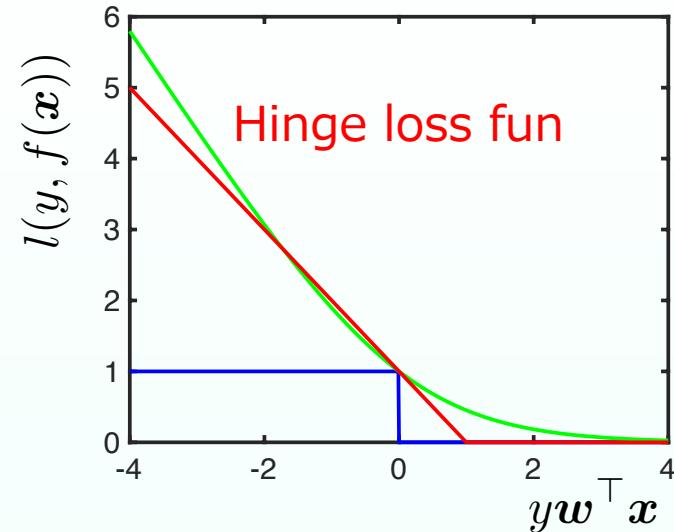
$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left(\sum_{i=1}^n \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x}_i) + \lambda \mathbf{w}^\top \mathbf{w} \right)$$

Hinge-loss: Non-smooth function at $y \mathbf{w}^\top \mathbf{x} = 1$

- $\mathbf{g} \in \partial l_i(\mathbf{w})$: Subgradient of hinge-loss can be written as

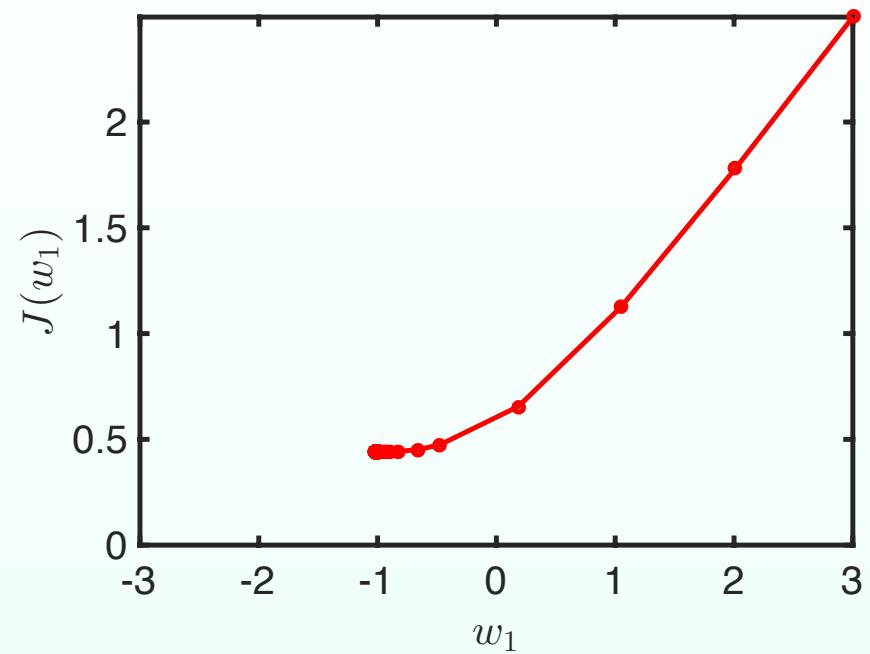
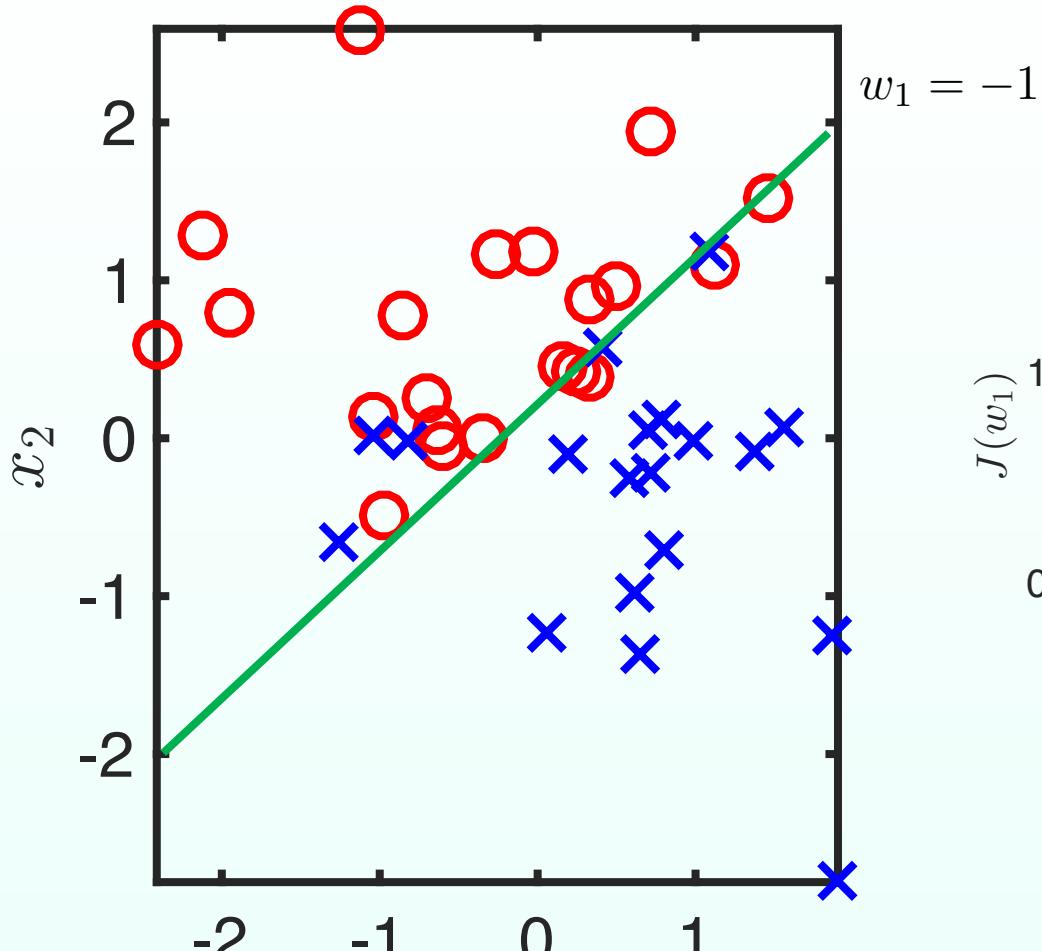
$$\mathbf{g} = \begin{cases} \boxed{ } & \text{if } y_i \mathbf{w}^\top \mathbf{x}_i > 1 \\ 1 & \text{if } y_i \mathbf{w}^\top \mathbf{x}_i = 1, \quad \eta \in [0, 1] \\ -1 & \text{if } y_i \mathbf{w}^\top \mathbf{x}_i < 1 \end{cases}$$

- We can apply gradient descent with the subgradient of hinge-loss
- The subgradient method for this model is like “**perceptron**” learning



Toy example of subgradient based SVM

Binary classification toy example with 2d input



Projected gradient

Review: constraint optimization

Constraint optimization

- Separating non-smooth function from objective into constraints

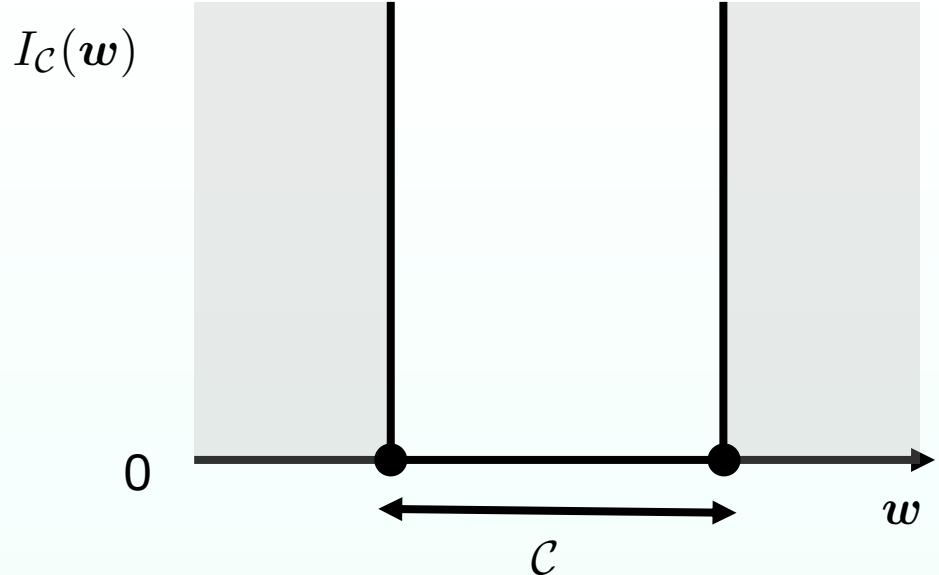
$$\begin{array}{ll}\text{minimize} & \psi(\mathbf{w}) \\ \text{subject to} & \left. \begin{array}{l} g_i(\mathbf{w}) = 0, \ i = 1, \dots, m \\ h_i(\mathbf{w}) \leq 0, \ i = 1, \dots, r \end{array} \right\} \mathbf{w} \in \mathcal{C}\end{array}$$

- We can rewrite constraint condition as region \mathcal{C}
- This problem can be written as $\text{minimize}_{\mathbf{w} \in \mathcal{C}} \psi(\mathbf{w})$
- From now, we consider (proximal) gradient method for this problem

Indicator function of Region

Indicator function for (proximal) gradient method

$$I_{\mathcal{C}}(\mathbf{w}) = \begin{cases} 0 & \mathbf{w} \in \mathcal{C} \\ \infty & \mathbf{w} \notin \mathcal{C} \end{cases}$$



The problem can be written as

$$\text{minimize}_{\mathbf{w} \in \mathcal{C}} \psi(\mathbf{w}) \leftrightarrow \text{minimize}_{\mathbf{w}} \psi(\mathbf{w}) + I_{\mathcal{C}}(\mathbf{w})$$

Applying proximal gradient with indicator function

The problem we consider now

$$\text{minimize}_{\mathbf{w}} \psi(\mathbf{w}) + I_{\mathcal{C}}(\mathbf{w})$$

Proximal gradient with indicator function

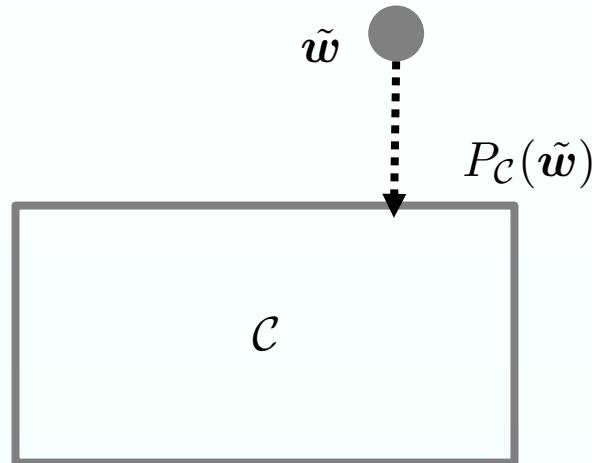
$$\mathbf{w}^{(t)} = \text{prox}_{\eta_t I_{\mathcal{C}}}(\mathbf{w}^{(t-1)} - \eta_t \nabla \psi(\mathbf{w}^{(t-1)}))$$

$$\tilde{\mathbf{w}}$$

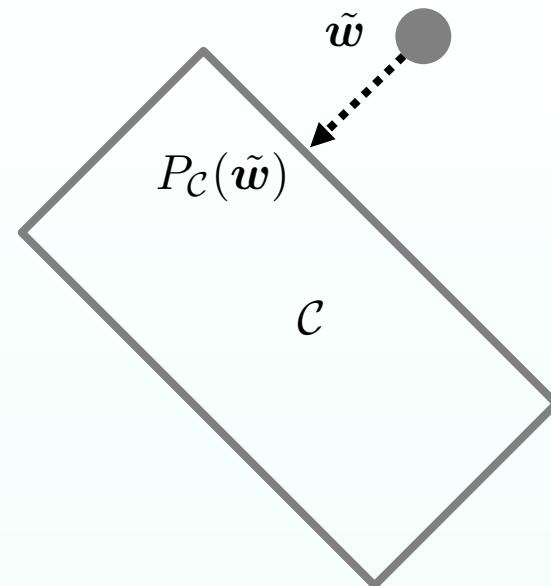
$$\begin{aligned}\text{prox}_{\eta_t I_{\mathcal{C}}}(\tilde{\mathbf{w}}) &= \operatorname*{argmin}_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{w} - \tilde{\mathbf{w}}\|_2^2 + I_{\mathcal{C}}(\mathbf{w}) \right) \\ &= \operatorname*{argmin}_{\mathbf{w} \in \mathcal{C}} (\|\mathbf{w} - \tilde{\mathbf{w}}\|_2^2) \\ &= P_{\mathcal{C}}(\tilde{\mathbf{w}})\end{aligned}$$

Toy example of projection in 2d space

Case 1



Case 2



We can analytically obtain projection step for

- Ex. Box constraint $\mathcal{C} = \{\mathbf{w} | w_i \in [r_i, R_i], \forall i = 1, \dots, d, r_i \leq R_i\}$

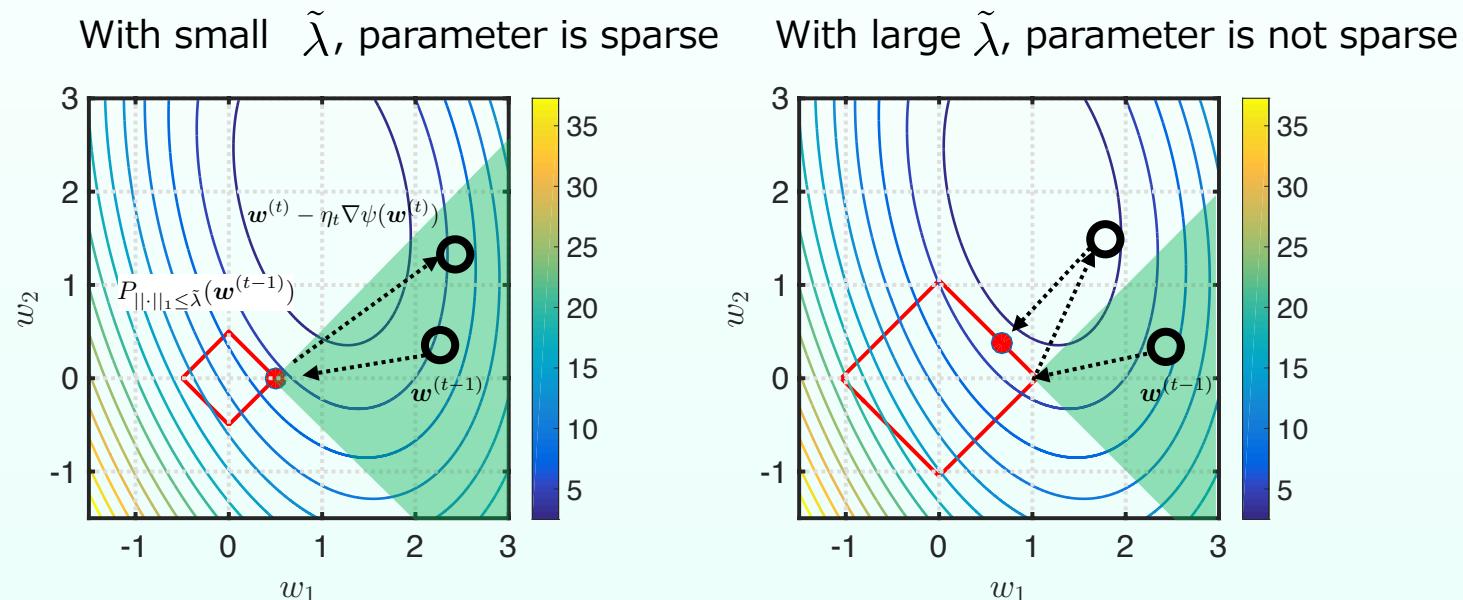
$$\{P_{\mathcal{C}}(\tilde{w})\}_i = \begin{cases} R_i & \text{if } \tilde{w}_i \geq R_i \\ \tilde{w}_i & \text{if } \tilde{w}_i \in [r_i, R_i] \\ r_i & \text{if } \tilde{w}_i < r_i \end{cases}$$

Review: Geometric interpretation of Lasso

Lasso and its equivalent problem (with L1 constraints)

$$\operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1 \right) \quad \leftrightarrow \quad \begin{array}{l} \text{minimize}_{\mathbf{w}} \quad \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \\ \text{subject to} \quad \|\mathbf{w}\|_1 \leq \tilde{\lambda} \end{array}$$

- We consider to apply projected gradient for constraint problem



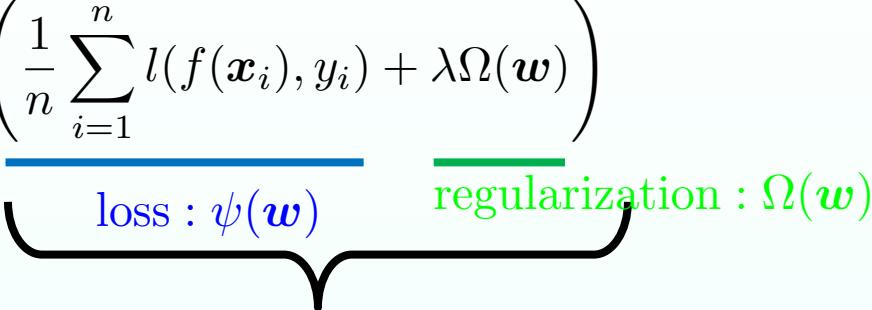
Stochastic optimization

Batch training is not always feasible in large scale machine learning

Batch training

- The optimization methods we had thought are known as batch training
- In batch training, all the training data for loss and its gradient are used

$$\hat{f}^{\text{REG}}(\mathbf{x}) = \operatorname{argmin}_f \left(\frac{1}{n} \sum_{i=1}^n l(f(\mathbf{x}_i), y_i) + \lambda \Omega(\mathbf{w}) \right)$$


 $J(\mathbf{w})$

- We always use gradient information:
• In linear logistic regression, this requires $O(dn)$ operation per iteration
- For large scale ML, we want to avoid such computational cost

Stochastic optimization

Original (batch training)

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{n} \sum_{i=1}^n l(\mathbf{x}_i, y_i, \mathbf{w}) + \lambda \Omega(\mathbf{w}) \right) = \operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{n} \sum_{i=1}^n l_i(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \right)$$

Loss per sample regularization : $\Omega(\mathbf{w})$ Loss per sample
(with simplified representation)

(Online) Stochastic optimization

- Iterate some update process $t = 1, \dots, T$ with initial guess $\mathbf{w}^{(0)}$
 - 1. $z_t \sim \mathcal{U}(z|n)$: randomly pickup index ranging from 1 to n
 - 2. update parameter $\mathbf{w}^{(t)} \leftarrow \mathbf{w}^{(t-1)}$ using $l_{z_t}(\mathbf{w}^{(t-1)})$, and $\lambda\Omega(\mathbf{w}^{(t-1)})$

Equivalent problem with (online) stochastic optimization

During iteration process, it is equivalent to solve the following problem

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{T} \sum_{t=1}^T l_{z_t}(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \right)$$

Update parameter with finite iteration

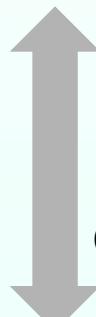
- If $T \rightarrow \infty$, the problem is equivalent to the following problem

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \left(\mathbb{E}_{z \sim \mathcal{U}(z|n)} [l_z(\mathbf{w})] + \lambda \Omega(\mathbf{w}) \right)$$

Expectation of per sample loss function

C.f. batch training

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \left(\frac{1}{n} \sum_{i=1}^n l_i(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \right)$$



Stochastic gradient descent (SGD)

Applying (proximal) GD to stochastic optimization

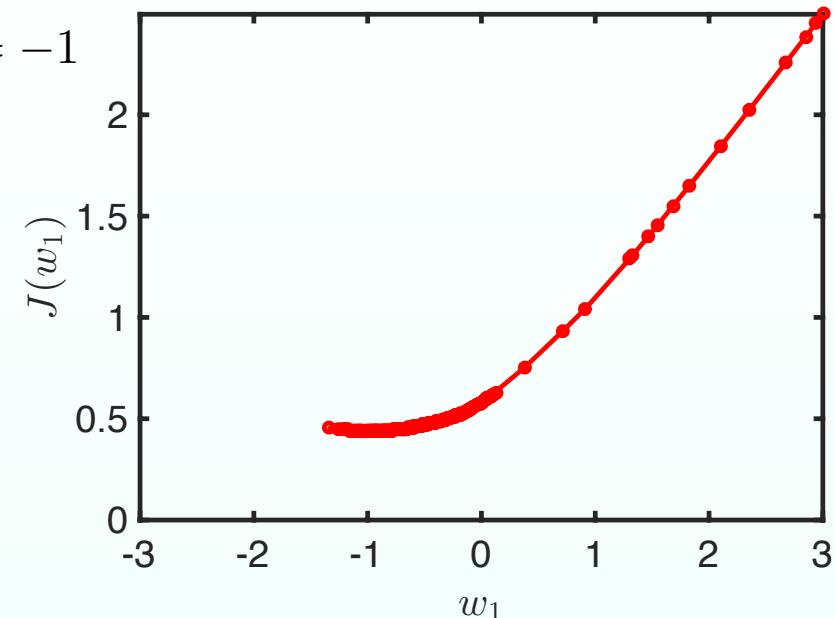
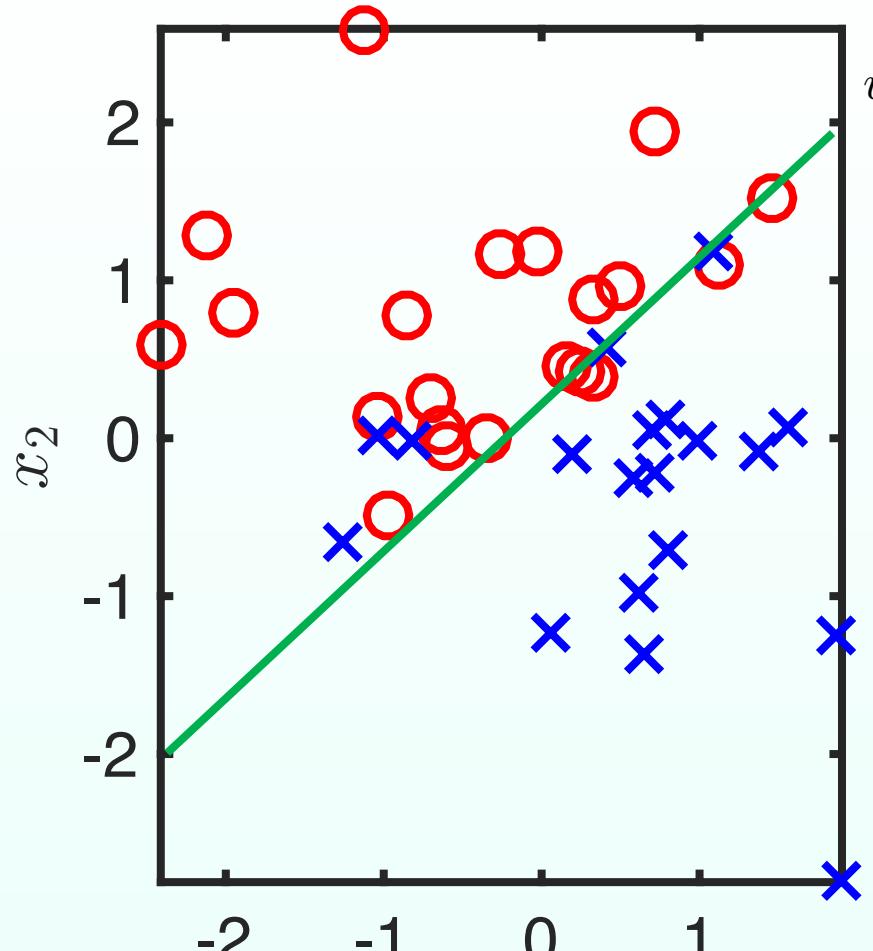
- Iterate some update process $t = 1, \dots, T$ with initial guess: $\mathbf{w}^{(0)}$
 - 1. $z_t \sim \mathcal{U}(z|n)$: randomly pickup index ranging from 1 to n
 - 2. compute gradient information from single sample: $\mathbf{g}_t \in \partial l_{z_t}(\mathbf{w}^{(t-1)})$
 - 3. update the parameter using gradient (with proximal operation)

$$\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega(\cdot)}(\mathbf{w}^{(t-1)} - \eta_t \mathbf{g}_t)$$

- Comparison with batch training
 - Gradient is made from the single sample instead of whole dataset
- Mini-batch training
 - Mix of Batch and (online) stochastic training
 - Use $m \ll n$ samples for obtaining gradient rather than single sample

Action of SGD with toy example

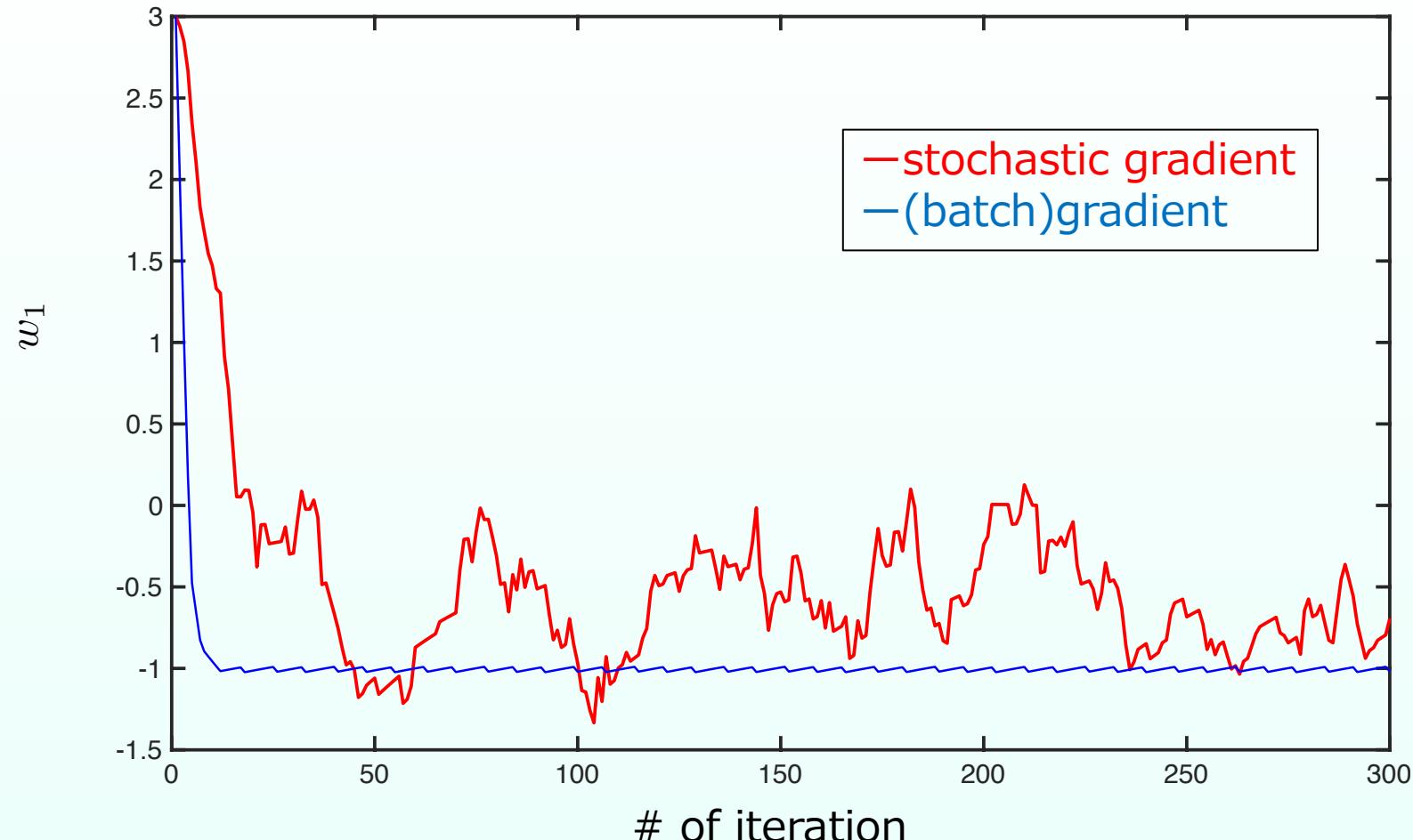
L2 regularized hinge loss minimization in 2d



Action of SGD with toy example (cont'd)

Stochastic gradient converges gradually

- With small size dataset, use batch training!



Stochastic gradient with momentum

Momentum instead of using gradient from single data

- To remove the instability due to the robustness
- Original

$$\mathbf{g}_t \in \partial l_{z_t}(\mathbf{w}^{(t-1)})$$

$$\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega(\cdot)}(\mathbf{w}^{(t-1)} - \eta_t \mathbf{g}_t)$$

- Momentum update (like acceleration)

$$\mathbf{g}_t \in \partial l_{z_t}(\mathbf{w}^{(t-1)})$$

$$\boldsymbol{\mu}_t = \beta \boldsymbol{\mu}_{t-1} + (1 - \beta) \mathbf{g}_t$$

Momentum

$$= (1 - \beta) \left(\mathbf{g}_t + \sum_{\tau} \beta^{\tau} \mathbf{g}_{t-\tau} \right), \quad \beta \in (0, 1)$$

a.k.a. Polynomial decay averaging

$$\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega(\cdot)}(\mathbf{w}^{(t-1)} - \eta_t \boldsymbol{\mu}_t)$$

Convergence rate of SGD

Gap between estimation and ground truth

$\mathcal{E}_{z_1:T \sim \mathcal{U}(z|n)} [J(\mathbf{w}^{(T)})] - J(\hat{\mathbf{w}})$ will be bounded by

Non α -strongly convex	α -strongly convex
$O\left(\frac{1}{\sqrt{T}}\right)$	$O\left(\frac{1}{T}\right)$

- Ground truth $\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w}} \left(\sum_i l_i(\mathbf{w}) + \lambda \Omega(\mathbf{w}) \right)$
 - Usually, theoretical analysis uses weighted average $\bar{\mathbf{w}}^{(T)} = \frac{1}{T+1} \sum_{\tau=0}^T \mathbf{w}^{(\tau)}$, and its variants rather than using $\mathbf{w}^{(T)}$
- C.f. Regret is also frequently used $R(T) = \sum_{t=1}^T \left(l_{z_t} \left(\mathbf{w}^{(t)} \right) + \lambda \Omega(\mathbf{w}^{(t)}) \right) - \inf_{\mathbf{w}} \sum_{t=1}^T (l_{z_t}(\mathbf{w}) + \lambda \Omega(\mathbf{w}))$

Convergence rate of SGD comparison with batch training

- We assume alpha strongly convex function here
- We omit dimensionality here (since this has no effects in comparison)

	Batch	(online) stochastic
Per iteration	$O(n)$	$O(1)$
# of iteration to reach the gap is bounded by ϵ	$O(-\ln \epsilon)$	$O(\epsilon^{-1})$
Total cost	$O(-n \ln \epsilon)$	$O(\epsilon^{-1})$
Generalization error bounded by $O(1/n)$.	$O(n \ln n)$	$O(n)$

We assume loss is
a-strongly convex function

- We need additional iteration to reach the accuracy in stochastic case, but we do not need heavy computational cost for loss calculation in each iteration
- **The total cost can be reduced to $O(n)$ from $O(n \ln n)$**

Adaptive step size estimation in (stochastic) optimization

Review: (stochastic)(proximal) gradient method

Now, you can use “gradient” method for large scale ML

- Thanks to stochastic (sub-)gradient
- With Non-smooth (sparse regularization) term thanks to prox. ops

$$\underline{\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega}(\mathbf{w}^{(t-1)} - \eta_t \mathbf{g}_t)}$$

$$\mathbf{g}_t \in \partial l_{z_t}(\mathbf{w}^{(t-1)})$$

$$z_t \sim \mathcal{U}(z|n)$$

- As for step size, we assume it is fix (w.r.t. γ smooth) or back-tracking

Review: (stochastic)(proximal) gradient method

Now, you can use “gradient” method for large scale ML

- Thanks to stochastic (sub-)gradient
- With Non-smooth (sparse regularization) term thanks to prox. ops

$$\underline{\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega}(\mathbf{w}^{(t-1)} - \eta_t \mathbf{g}_t)}$$

$$\mathbf{g}_t \in \partial l_{z_t}(\mathbf{w}^{(t-1)})$$

$$z_t \sim \mathcal{U}(z|n)$$

- As for step size, we assume it is fix (w.r.t. γ smooth) or back-tracking

Is it really reasonable for any (convex) optimization?

Review: (stochastic)(proximal) gradient method

Now, you can use “gradient” method for large scale ML

- Thanks to stochastic (sub-)gradient
- With Non-smooth (sparse regularization) term thanks to prox. ops

$$\underline{\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega}(\mathbf{w}^{(t-1)} - \eta_t \mathbf{g}_t)}$$

$$\mathbf{g}_t \in \partial l_{z_t}(\mathbf{w}^{(t-1)})$$

$$z_t \sim \mathcal{U}(z|n)$$

- As for step size, we assume it is fix (w.r.t. γ smooth) or back-tracking

Is it really reasonable for any (convex) optimization?

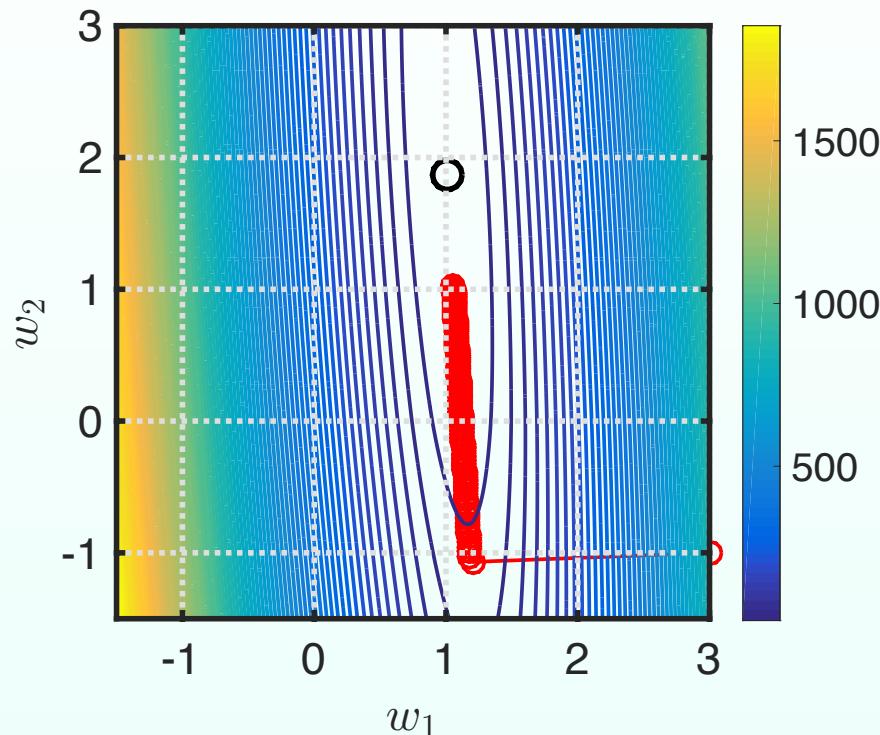
Actually, $\eta_t = \gamma^{-1}$ is not so efficient, but we could not control it automatically

GD in action with toy example

$$\underset{\mathbf{w}}{\operatorname{argmin}} \left((\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A} (\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1 \right)$$

$$\mathbf{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = (1 \ 2), \lambda = 0.89, \eta_t = 1/\gamma$$

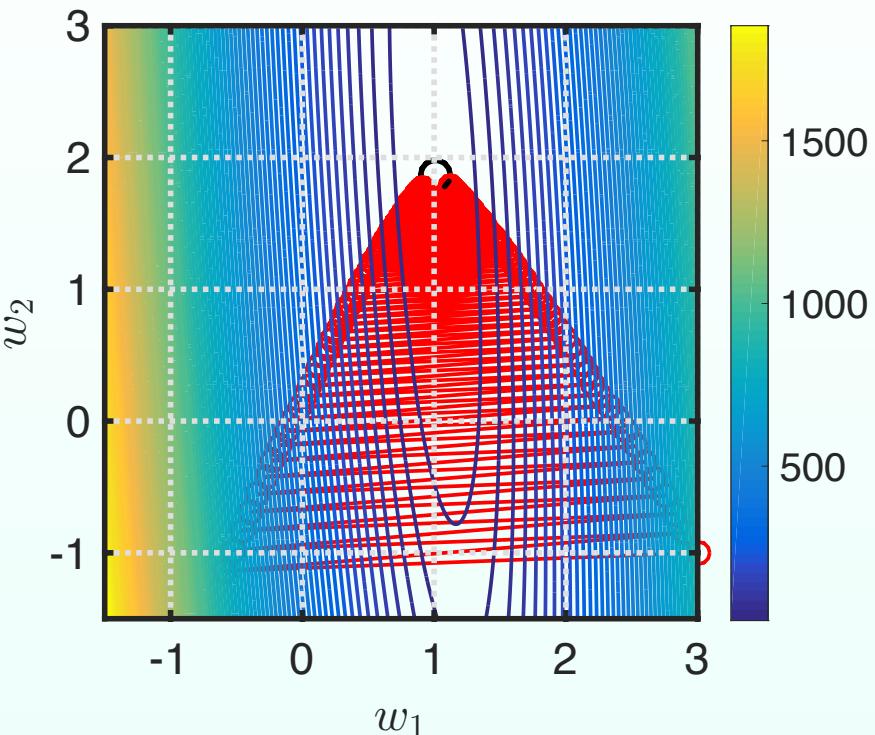
Too slow to reach minima



$$\gamma = \max \operatorname{eig}(2\mathbf{A})$$

$$\mathbf{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = (1 \ 2), \lambda = 0.89, \eta_t = 2/\gamma$$

Too vibrated to reach minima

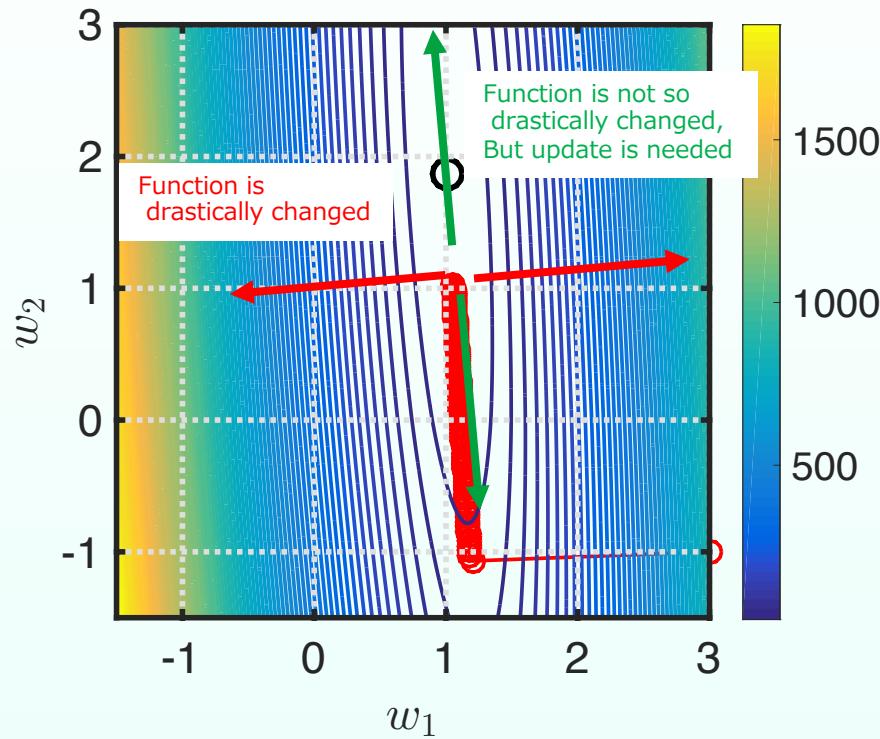


GD in action with toy example

$$\operatorname{argmin}_{\mathbf{w}} ((\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A} (\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1)$$

$$\mathbf{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = (1 \ 2), \lambda = 0.89, \eta_t = 1/\gamma$$

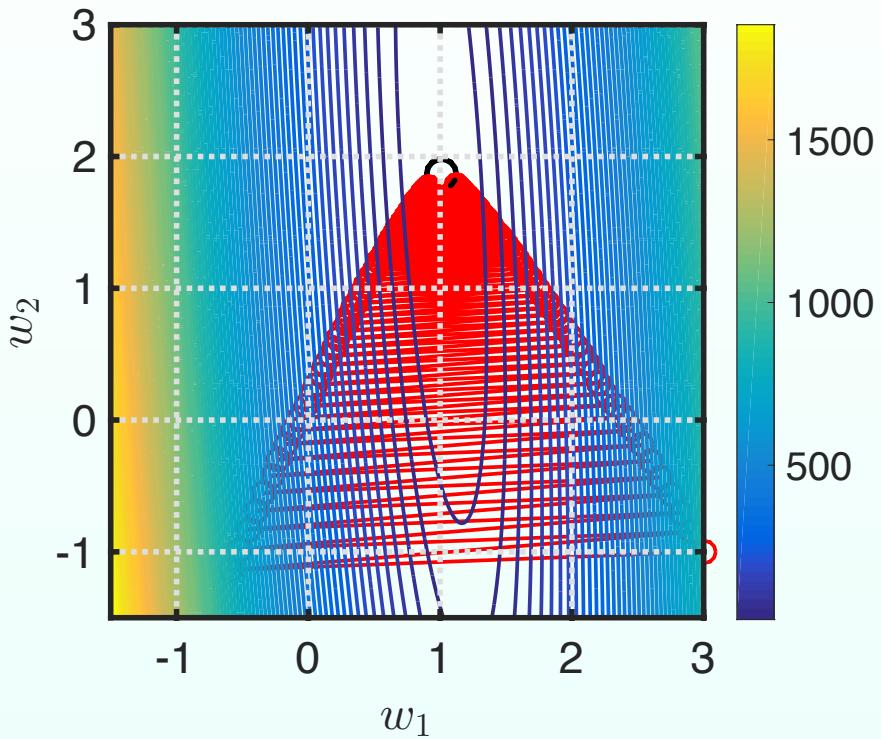
Too slow to reach minima



$$\gamma = \max \operatorname{eig}(2\mathbf{A})$$

$$\mathbf{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = (1 \ 2), \lambda = 0.89, \eta_t = 2/\gamma$$

Too vibrated to reach minima



We could not use info on ratio on 1st and 2nd eigen values
One approach is to use momentum (averaged of gradients)

Review: Newton method

Use 2nd order information for the optimization

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \mathbf{H}(\mathbf{w}^{(t-1)})^{-1} \mathbf{g}(\mathbf{w}^{(t-1)})$$

- For batch training
- For smooth function optimization
- Let's consider toy problem $\underset{\mathbf{w}}{\operatorname{argmin}} ((\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A} (\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1)$
- If $\lambda = 0$, $\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \eta_t (\mathbf{w}^{(t-1)} - \boldsymbol{\mu})$ The param is directly moved towards minima
- 2nd order info cancel the variation of scale of gradients
 - Like whitening with Mahalanobis distance
- But we could not use Newton method for non-smooth, stochastic case

Newton-like proximal operation

Proximal operation with gradient information

- Use quadratic function that covers objective around $\mathbf{w}^{(t-1)}$

$$\mathbf{w}^{(t)} = \text{prox}_{\eta_t \lambda \Omega}(\mathbf{w}^{(t-1)} - \eta_t \mathbf{g}_t)$$

$$= \underset{\mathbf{w}}{\text{argmin}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{g}_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_2^2 \right)$$

Common in two methods

Use Euclid distance

Approximating loss function with 2nd Taylor expansion

- Newton like approach

$$\mathbf{w}^{(t)} = \underset{\mathbf{w}}{\text{argmin}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{g}_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t-1)})^\top \hat{\mathbf{H}}(\mathbf{w}^{(t-1)}) (\mathbf{w} - \mathbf{w}^{(t-1)}) \right)$$

Common in two methods

Use curvature of the objective

- It looks nice only if we can ignore computational cost

$$\hat{\mathbf{H}}(\mathbf{w}^{(t-1)}) \approx \nabla^2 \left(\sum_i l_i(\mathbf{w}^{(t-1)}) \right)$$

Obtaining pseudo Hessian in stochastic optimization

We consider to obtain Hessian like matrix H_t from sequence of gradient information

$$\mathbf{w}^{(t)} = \operatorname{argmin}_{\mathbf{w}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{g}_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}^{(t-1)})^\top \hat{\mathbf{H}}(\mathbf{w}^{(t-1)}) (\mathbf{w} - \mathbf{w}^{(t-1)}) \right)$$

Infeasible to obtain Hessian directly

$$\mathbf{w}^{(t)} = \operatorname{argmin}_{\mathbf{w}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{g}_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2\eta_0} \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{\mathbf{H}_t}^2 \right)$$

$$\eta_t = \eta_0$$

$$\|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{\mathbf{H}_t}^2 = (\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{H}_t (\mathbf{w} - \mathbf{w}^{(t-1)})$$

- A generalization form of proximal gradient method
- If we assume, $\mathbf{H}_t = \mathbf{I}, \eta_t = \gamma^{-1}$,
the algorithm is equivalent to standard proximal gradient method

Design policy of pseudo Hessian

Use BFGS formula (quasi Newton)?

- It is not suitable to obtain pseudo Hessian in stochastic optimization

Design policy for creating pseudo Hessian

- Computationally efficient (calculation cost + memory storage)

Designing diagonal matrix from gradient history

$$H_t = G_t^{1/2} + \delta I$$

$$\{G_t\}_{i,j} = \begin{cases} 0 & \text{if } i \neq j \\ \sum_{\tau=1}^t \{g_\tau\}_i^2 & \text{otherwise} \end{cases}$$

- The gradient is normalized w.r.t. the size of gradients
- The algorithm is known as AdaGrad

AdaGrad as adaptive proximal gradient

AdaGrad

- Sequential stochastic optimization problem with proximal gradients

$$\mathbf{w}^{(t)} = \operatorname{argmin}_{\mathbf{w}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{g}_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2\eta_0} \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{\mathbf{H}_t}^2 \right)$$

- Use Hessian like information

$$\mathbf{H}_t = \mathbf{G}_t^{1/2} + \delta \mathbf{I}$$

$$\{\mathbf{G}_t\}_{i,j} = \begin{cases} 0 & \text{if } i \neq j \\ \sum_{\tau=1}^t \{\mathbf{g}_\tau\}_i^2 & \text{otherwise} \end{cases}$$

- Since \mathbf{H}_t is diagonal, the parameter can be analytically solved by proximal operation. (operation cost is proportional to d)
 - Derive update formula
for L1 case (with soft-thresholding), and square of L2 regularization(**a3)

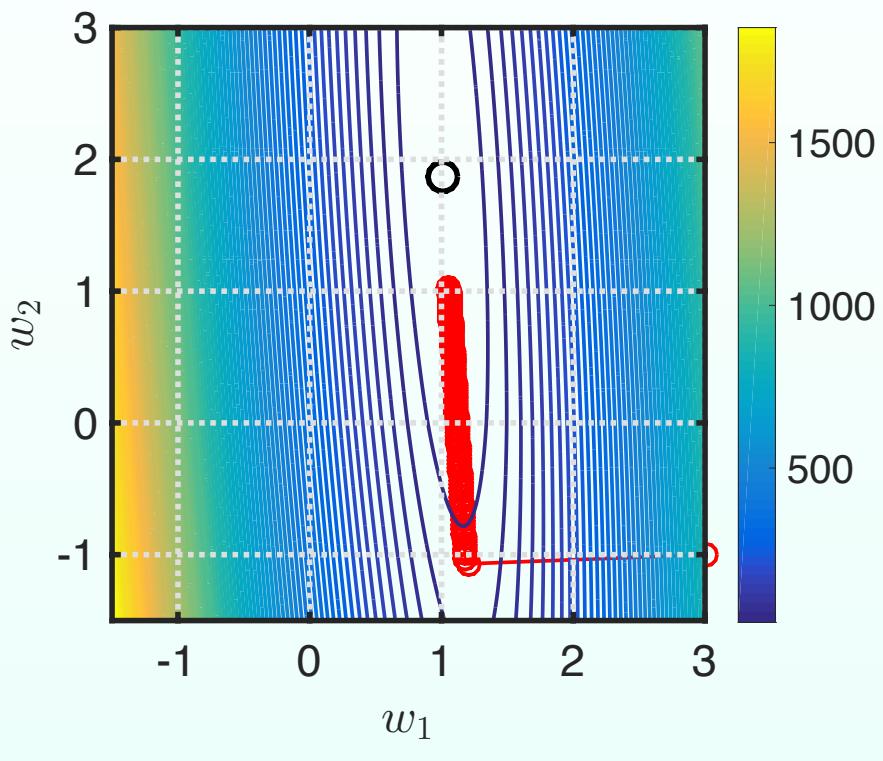
$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_1$$

$$\Omega(\mathbf{w}) = \|\mathbf{w}\|_2^2$$

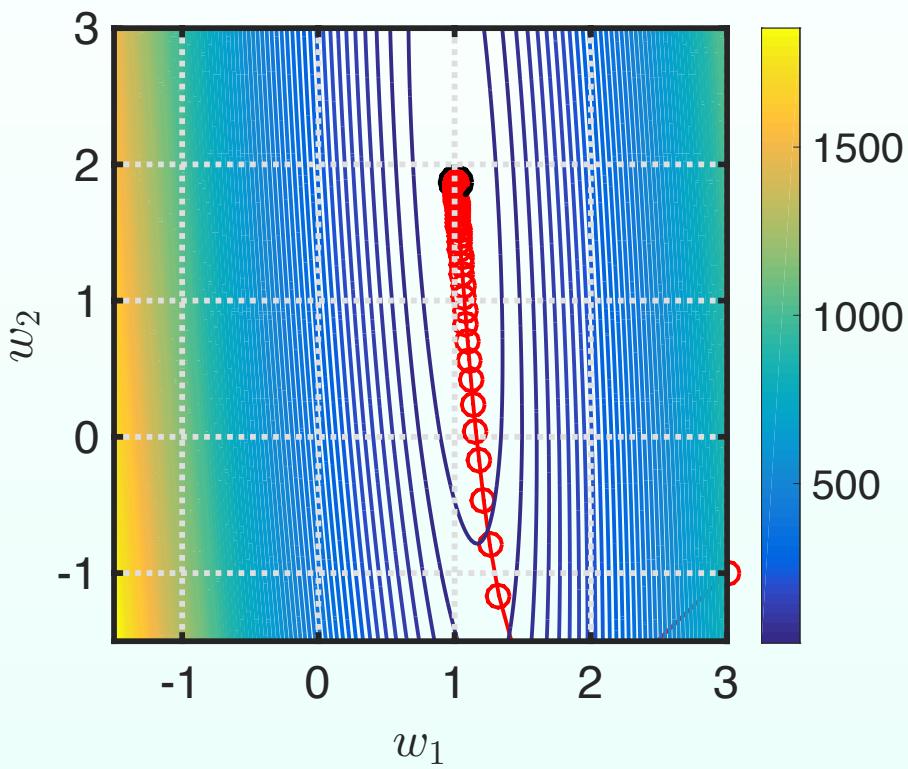
AdaGrad in action

$$\operatorname{argmin}_{\mathbf{w}} ((\mathbf{w} - \boldsymbol{\mu})^\top \mathbf{A} (\mathbf{w} - \boldsymbol{\mu}) + \lambda \|\mathbf{w}\|_1) \quad \mathbf{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = \begin{pmatrix} 1 & 2 \end{pmatrix}, \lambda = 0.89$$

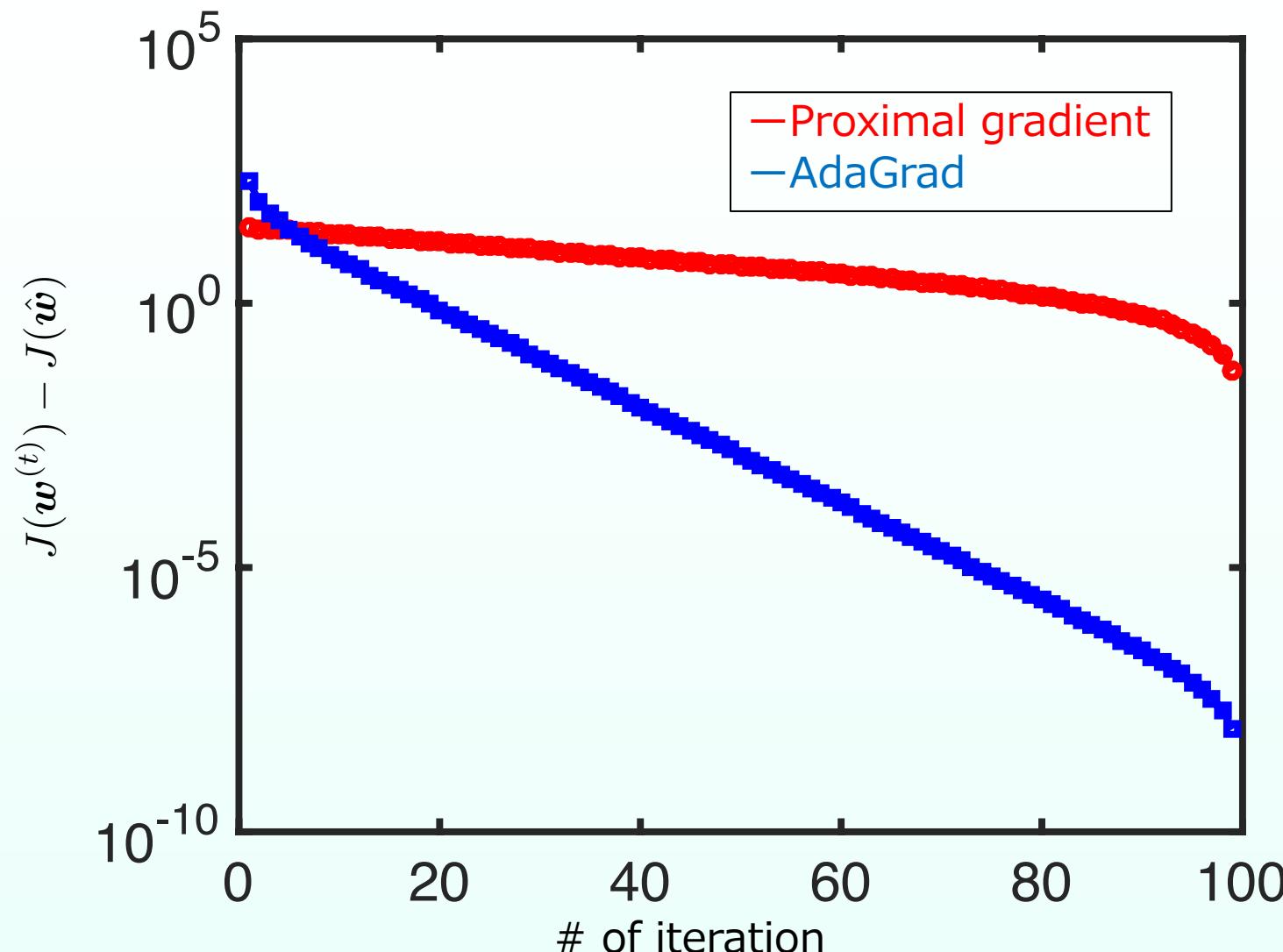
Proximal GD $\eta_t = \gamma^{-1}$



AdaGrad $\eta_t = 500\gamma^{-1}$ $\delta = 0.02$



AdaGrad in action (cont'd)



Advances in adaptive gradient methods

Adaptive gradient methods are actively explored in ML

- AdaGrad (2011)
- RMS Prop (2012)
- Ada Delta (2012)
- Adam (2015)
- Nadam (2016)
- ...

Researchers want to create efficient optimization

- without hand tuning parameters (such as step size tuning)
- within computationally efficient manner

Adam: Adaptive moment estimation

Adam is similar to AdaGrad except that

- Adam uses momentum μ_t rather than gradient
- Adam uses momentum V_t rather than sum of grad sq

$$\text{AdaGrad} \quad \mathbf{w}^{(t)} = \operatorname{argmin}_{\mathbf{w}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mathbf{g}_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2\eta_0} \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{\mathbf{H}_t}^2 \right)$$

$$\text{Adam} \quad \mathbf{w}^{(t)} = \operatorname{argmin}_{\mathbf{w}} \left((\mathbf{w} - \mathbf{w}^{(t-1)})^\top \mu_t + \lambda \Omega(\mathbf{w}) + \frac{1}{2\eta_0} \|\mathbf{w} - \mathbf{w}^{(t-1)}\|_{V_t}^2 \right)$$

- Disclaimer
The above formulation is just a variant of Adam.
Though Adam paper does not present any issues related to proximal gradient.
I would like to emphasize that proximal gradient could be used in adaptive gradient methods.

Please refer the paper if you want to have original behavior.

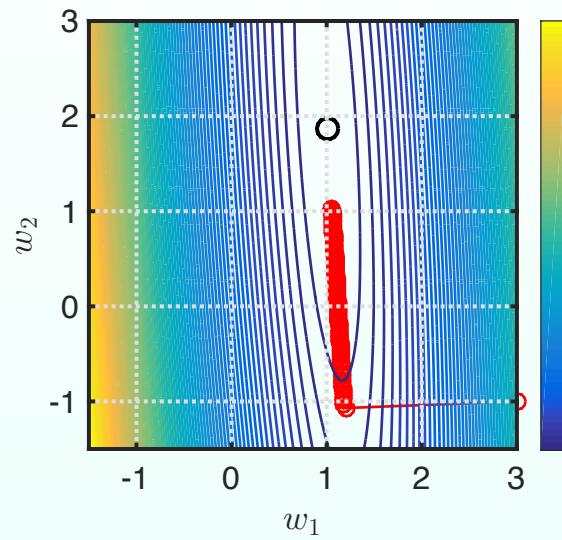
The paper provides theoretical guarantee

- Regret bound for convex function

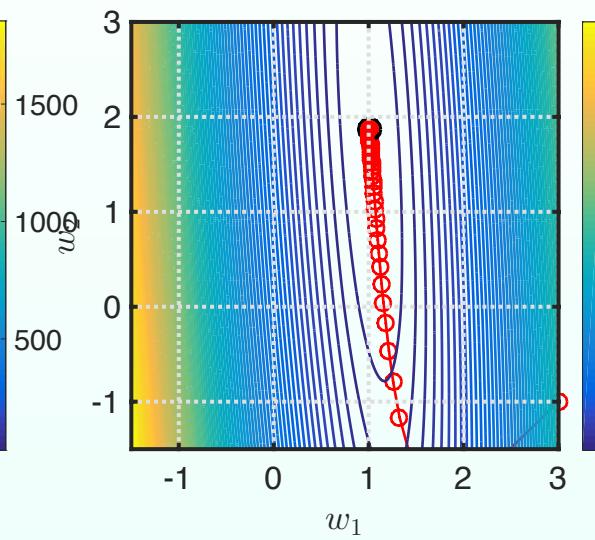
Adam in action

$$\underset{\boldsymbol{w}}{\operatorname{argmin}} \left((\boldsymbol{w} - \boldsymbol{\mu})^\top \boldsymbol{A} (\boldsymbol{w} - \boldsymbol{\mu}) + \lambda \|\boldsymbol{w}\|_1 \right) \quad \boldsymbol{A} = \begin{pmatrix} 250 & 15 \\ 15 & 4 \end{pmatrix}, \boldsymbol{\mu}^\top = \begin{pmatrix} 1 & 2 \end{pmatrix}, \lambda = 0.89$$

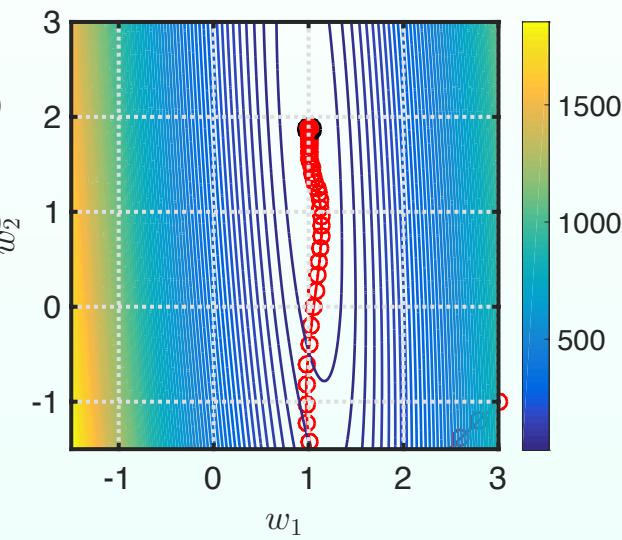
SGD



AdaGrad



Adam-like



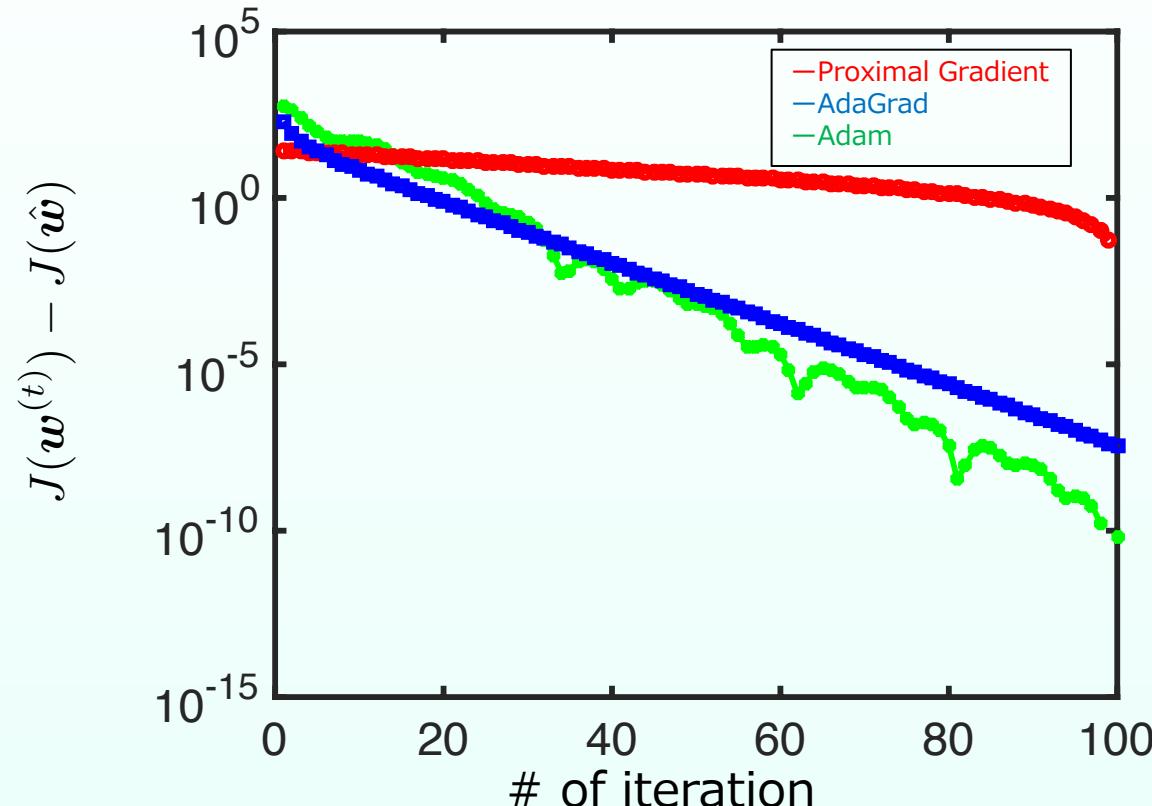
Adam in action(cont'd)

Adam, AdaGrad (vs SGD)

- Gets much better convergence with some condition $\sigma_1 \gg \sigma_2, \dots$

Adam vs AdaGrad

- I could not give clear decision from this experiment…



Review of this lecture

Theoretical part: advances in gradient methods

- Accelerated proximal gradient methods
 - Convergence rate of proximal gradient methods
 - Acceleration, momentum, FISTA: Accelerated proximal gradient method
- Subgradient method
- Projected gradient method
- Stochastic optimization
 - From batch to stochastic way
 - vs. batch training in terms of convergence rate & computational cost
- Adaptive step size in gradient methods
 - AdaGrad, Adam

Bridge to deep learning

The optimization methods discussed today

- 1. Stochastic (vs batch)
- 2. Adaptive step-size (vs fixed rate or exact line search)
- 3. Momentum (vs standard gradient)
- 4. Proximal operation (vs smooth function minimization)
- verified by the convex programming theory with convergence rate
- Also used in deep learning community
 - The technique can be (practically) applied to non-convex programming
 - Researcher in ML tackle the theoretical analysis of non-convex programming

The issues for deep learning

- How to obtain **gradient** in deep layered model?
 - With **gradient**, you can use the optimization, such as AdaGrad, Adam!

Assignment “Shimosaka-6”

Solve the following problems & submit report via OCW-i

Format: pdf (scanned hand written report / tex)

Problem:

- Derive, or confirm the equations by your hand (**a1-3)

Preview of final assignment

Evaluation by me: 50pt in total

- Report in each lecture (I will not check the quality ☹)
 - 3pts x 6 times
- Final report (I will check the quality!)
 - Details will be announced soon
 - 30 pts in total.
 - You will be asked to select a couple of problems from 8–10 problems.
 - If you have questions about the final report,
please send email to **titech-ml-course@miubiq.cs.titech.ac.jp**
(DO NOT send email to this on any issues related to Okazaki-sensei's lecture)
- 2pts by your constructive criticism to my lecture (in final report)