# A Characterization of Stochastic Mirror Descent Algorithms and Their Convergence Properties

Navid Azizan and Babak Hassibi

**Caltech**

## Motivation

- **Stochastic Mirror Descent (SMD)** is a general <u>family</u> of optimization algorithms
- *Stochastic Gradient Descent (SGD)* is a special case of SMD
- Other examples include *exponential weights algorithm*, *p-norms algorithm*, etc.
- SMD algorithms have become increasingly popular in optimization, machine learning, signal processing, control, etc.

## Problem Setup

Data: $\{(x_i, y_i) : i = 1, \ldots, n\}$
where $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$

Model: $y_i = f(x_i, w) + v_i, \quad i = 1, \ldots, n$
where $f(\cdot, \cdot)$ is a given function that represents the model class
$w \in \mathbb{R}^p$ is an unknown weight vector (parameter)
$v_i$ is the noise, which represents measurement error, modeling error, etc.

Loss Function: $l(\cdot)$ is a nonnegative differentiable loss function with $l(0) = 0$
$L_i(w) = l(y_i - f(x_i, w)) \qquad L(w) = \sum_{i=1}^{n} L_i(w)$

SGD: $w_i = w_{i-1} - \eta_i \nabla L_i(w_{i-1})$

## Minimax Optimality of SGD

Consider a linear model $f(x_i, w) = x_i^T w$, i.e., $y_i = x_i^T w + v_i$,
and the square loss $L_i(w) = \frac{1}{2}(y_i - x_i^T w)^2$.
In this case, SGD is $w_i = w_{i-1} + \eta(y_i - x_i^T w_{i-1})x_i$

> **Theorem** (Hassibi et al, NIPS '93). *For any initialization $w_0$, any sufficiently small step size $\eta$, i.e., $0 < \eta \leq \min_i \frac{1}{\|x_i\|^2}$, and any number of steps $T \geq 1$, the SGD iterates $\{w_i\}$ are the optimal solution to the following minimization problem*
> $$\min_{\{w_i\}} \max_{w, \{v_i\}} \frac{\|w - w_T\|^2 + \eta \sum_{i=1}^{T}(x_i^T w - x_i^T w_{i-1})^2}{\|w - w_0\|^2 + \eta \sum_{i=1}^{T} v_i^2},$$
> *and the optimal value is 1.*

- The ratio is the $H^\infty$ norm of the transfer operator that maps the unknown disturbances to the estimation errors
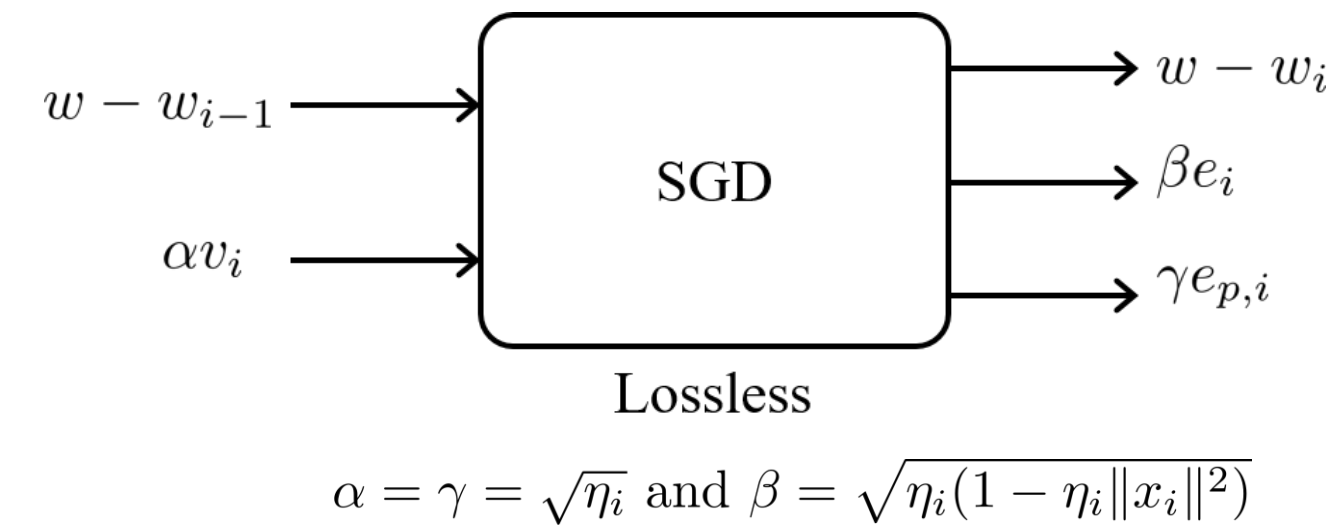- Interpretations: Robustness and Conservatism

## Proof: The Conservation Law of SGD

Define "innovations" and "predicted error" as
$e_i := y_i - x_i^T w_{i-1}$ and $e_{p,i} := x_i^T w - x_i^T w_{i-1}$

**Conservation of Uncertainty**

For each step of SGD:



$\alpha = \gamma = \sqrt{\eta_i}$ and $\beta = \sqrt{\eta_i(1 - \eta_i\|x_i\|^2)}$

> **Lemma.** *For any noise values $\{v_i\}$, any true parameter $w$, and any step-size sequence $\{\eta_i\}$, the following relation holds for the SGD iterates $\{w_i\}$*
> $$\|w - w_{i-1}\|^2 + \eta_i v_i^2 = \|w - w_i\|^2 + \eta_i(1 - \eta_i\|x_i\|^2)e_i^2 + \eta_i e_{p,i}^2, \quad \forall i \geq 1.$$

## Implications for Overparameterized Models

Set of solutions: $\mathcal{W} = \{w \mid y_i = x_i^T w, \ i = 1, \ldots, n\}$

**Convergence and Implicit Regularization:**

> For $\eta < \min_i \frac{2}{\|x_i\|^2}$, the SGD iterates converge to a solution $w_\infty \in \mathcal{W}$. Further
> $$w_\infty = \arg\min_{w \in \mathcal{W}} \|w - w_0\|$$

In particular, if initialized at zero, SGD converges to the minimum $l_2$ norm solution $\boxed{w_\infty = \arg\min_{w \in \mathcal{W}} \|w\|}$
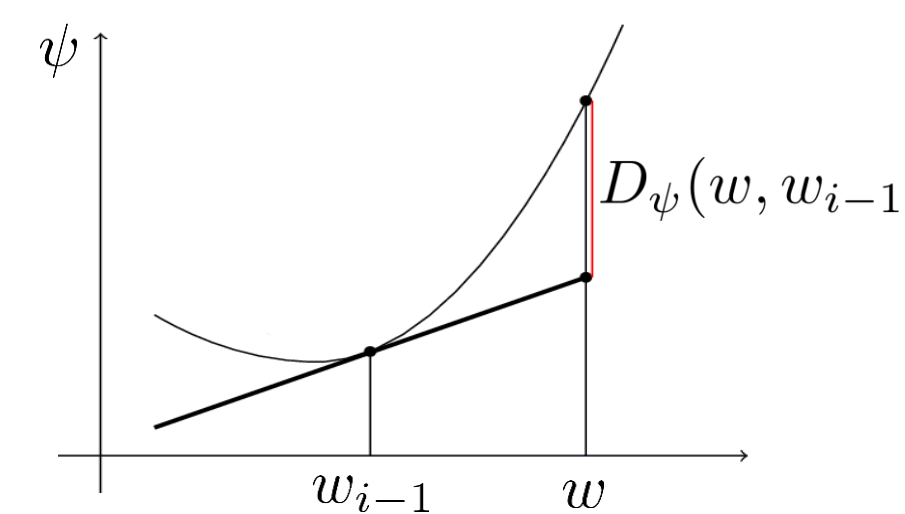
This is called implicit regularization

> What if we want a different regularizer?

## Stochastic Mirror Descent (SMD)

- A general family of optimization algorithms that includes stochastic gradient descent
- For any strictly convex and differentiable potential $\psi$, the SMD update rule is $\boxed{w_i = \arg\min_w \ \eta w^T \nabla L_i(w_{i-1}) + D_\psi(w, w_{i-1})}$

where $D_\psi(w, w_{i-1}) = \psi(w) - \psi(w_{i-1}) - \nabla\psi(w_{i-1})^T(w - w_{i-1})$ is the Bregman divergence w.r.t. $\psi$



- Equivalently, the SMD update can be expressed as
$$\boxed{\nabla\psi(w_i) = \nabla\psi(w_{i-1}) - \eta_i \nabla L_i(w_{i-1})}$$
- For SGD $\psi(w) = \frac{1}{2}\|w\|^2$

## Minimax Optimality of SMD

> **Theorem.** *Consider any (nonlinear) model $f$, any differentiable loss $l$ with property $l(0) = l'(0) = 0$, and any initialization $w_0$. For sufficiently small sequence of step sizes $\{\eta_i\}$, i.e., one for which $\psi(w) - \eta_i L_i(w)$ is convex for all $i$, and for any number of steps $T \geq 1$, the SMD iterates $\{w_i\}$, w.r.t. any strictly convex potential $\psi$, are the optimal solution to the following minimization problem*
> $$\min_{\{w_i\}} \max_{w, \{v_i\}} \frac{D_\psi(w, w_T) + \sum_{i=1}^{T} \eta_i D_{L_i}(w, w_{i-1})}{D_\psi(w, w_0) + \sum_{i=1}^{T} \eta_i l(v_i)},$$
> *and the optimal value is 1.*

- Generalizes several results, e.g. (SGD/square loss/linear model) [Hassibi et al '93] and (p-norms/square loss/linear model) [Kivinen et al '06]

- Proof by the conservation law of SMD:



> **Lemma.** *For any model $f(\cdot, \cdot)$, any differentiable loss $l(\cdot)$, any parameter $w$ and noise values $\{v_i\}$ that satisfy $y_i = f(x_i, w) + v_i$ for $i = 1, \ldots, n$, and any step-size sequence $\{\eta_i\}$, the following relation holds for the SMD iterates*
> $$D_\psi(w, w_{i-1}) + \eta_i l(v_i) = D_\psi(w, w_i) + E_i(w_i, w_{i-1}) + \eta_i D_{L_i}(w, w_{i-1}),$$
> *for all $i \geq 1$, where $E_i(w_i, w_{i-1}) := D_\psi(w_i, w_{i-1}) - \eta_i D_{L_i}(w_i, w_{i-1}) + \eta_i L_i(w_{i-1})$.*

## Implicit Regularization of SMD

> **Proposition.** *If $l$ is differentiable and convex and has a unique root at 0, $\psi$ is strictly convex, and positive sequence $\{\eta_i\}$ is such that $\psi - \eta_i L_i$ is convex for all $i$, then for any initialization $w_0$, the SMD iterates converge to*
> $$w_\infty = \arg\min_{w \in \mathcal{W}} D_\psi(w, w_0).$$

In particular, if we initialize SMD with $w_0 = \arg\min_{w \in \mathbb{R}^m} \psi(w)$, it converges to $\boxed{w_\infty = \arg\min_{w \in \mathcal{W}} \psi(w)}$

i.e., the minimum-potential solution. This is another implicit regularization.

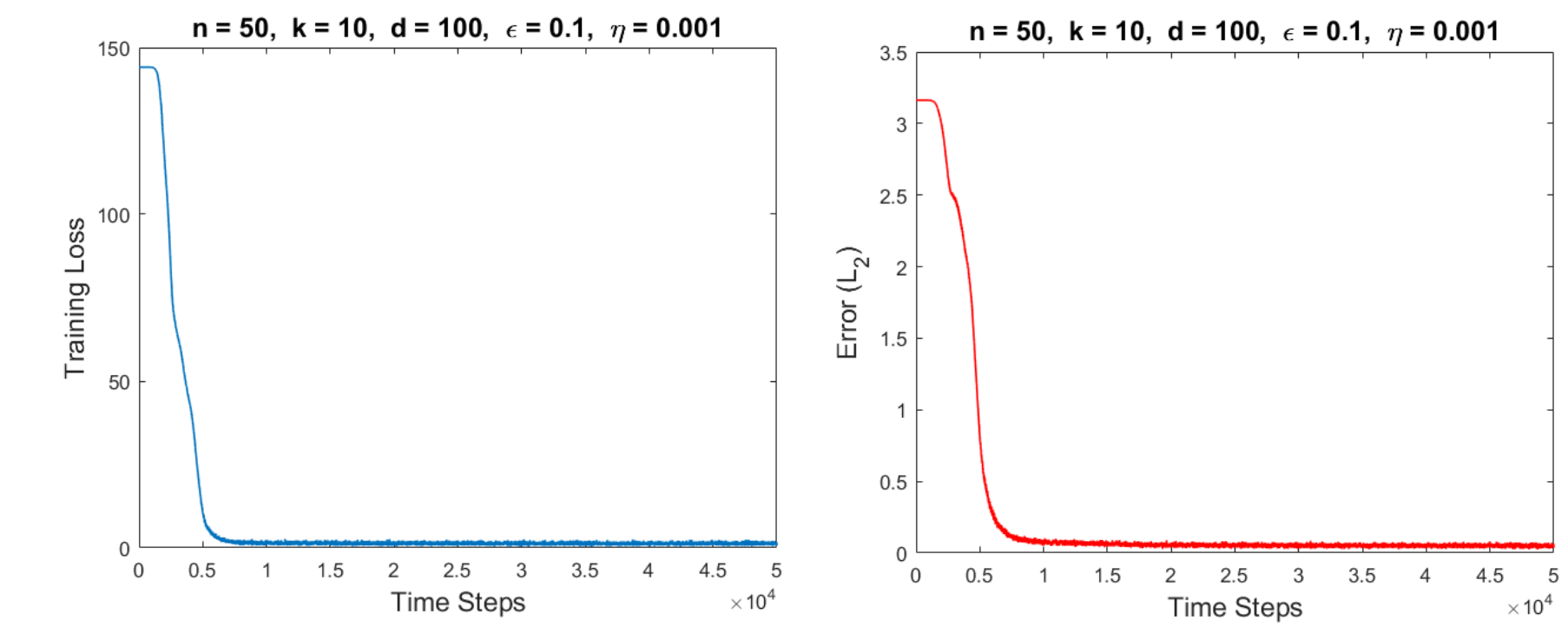> One can choose the potential function of SMD for any desired regularization

## Example: Compressed Sensing via SMD

recovering a sparse signal

$$\min_w \quad \|w\|_1$$
$$\text{s.t.} \quad x_i^T w = y_i, \ i = 1, \ldots, n.$$

$\psi(w) = \|w\|_1$ is not differentiable but we can use $\psi(w) = \|w\|_{1+\epsilon}^{1+\epsilon}$:



SMD w. $\psi(w) = \|w\|_{1+\epsilon}$ recovers the sparse solution!

## Stochastic Convergence in Underparameterized models

- Under-parameterized (online streaming) linear regression
- Vanishing step size
- Classical result:

> **Proposition.** *Consider $y_i = x_i^T w + v_i, i \geq 1$, where $\mathbb{E}[v_i] = 0$, $\mathbb{E}[v_i v_j] = \sigma^2 \delta_{ij}$, and the $x_i$ are persistently exciting. For any step size sequence $\{\eta_i\}$ such that $\sum_{i=1}^{\infty} \eta_i = \infty, \sum_{i=1}^{\infty} \eta_i^2 < \infty$, the SMD iterates with respect to any strongly convex potential $\psi(\cdot)$, converge to $w$ in the mean-square sense.*

- Direct and elementary proof using the conservation law of SMD
- Avoids ergodic averaging or appealing to stochastic differential equations

## + New Experimental Results

- SMD with different potential functions ran on MNIST
- The problem is **non-linear**



6 initial points x 4 different mirrors = 24 points on the manifold
Bregman divergences between the final and initial points, in 4 different norms

SMD converges to the point with smallest Bregman divergence from the initial point