

Overview of the Stochastic Gradient Method

Lecture Outline

- 1 Stochastic Gradient Algorithm
- 2 Connexion with Stochastic Approximation
- 3 Asymptotic Efficiency and Averaging
- 4 Practical Considerations

Deterministic Constrained Optimization Problem

General optimization problem (\mathcal{P})

$$\min_{u \in U^{\text{ad}} \subset \mathbb{U}} J(u)$$

- U^{ad} closed convex subset of an Hilbert space \mathbb{U} ,
- J cost function $\mathbb{U} \rightarrow \mathbb{R}$, satisfying some assumptions
 - convexity,
 - coercivity,
 - continuity,
 - differentiability.

Extension of Problem (\mathcal{P}) — Open-Loop Case (1)

Consider Problem (\mathcal{P}) and suppose that J is the **expectation** of a function j , depending on a **random variable** \mathbf{W} defined on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and valued on $(\mathbb{W}, \mathcal{W})$:

$$J(u) = \mathbb{E}(j(u, \mathbf{W})) .$$

Then the optimization problem writes

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) .$$

Decision u is a **deterministic variable**. The available information is the probability law of \mathbf{W} (no on-line observation of \mathbf{W}), that is, an **open-loop situation**. The information structure is trivial, but...

→ **main difficulty**: calculation of the expectation.

Extension of Problem (\mathcal{P}) — Open-Loop Case (1)

Consider Problem (\mathcal{P}) and suppose that J is the **expectation** of a function j , depending on a **random variable** \mathbf{W} defined on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and valued on $(\mathbb{W}, \mathcal{W})$:

$$J(u) = \mathbb{E}(j(u, \mathbf{W})) .$$

Then the optimization problem writes

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) .$$

Decision u is a **deterministic variable**. The available information is the probability law of \mathbf{W} (no on-line observation of \mathbf{W}), that is, an **open-loop situation**. The information structure is trivial, but...

→ **main difficulty**: calculation of the expectation.

Extension of Problem (\mathcal{P}) — Open-Loop Case (1)

Consider Problem (\mathcal{P}) and suppose that J is the **expectation** of a function j , depending on a **random variable** \mathbf{W} defined on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and valued on $(\mathbb{W}, \mathcal{W})$:

$$J(u) = \mathbb{E}(j(u, \mathbf{W})) .$$

Then the optimization problem writes

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) .$$

Decision u is a **deterministic variable**. The available information is the probability law of \mathbf{W} (no on-line observation of \mathbf{W}), that is, an **open-loop** situation. The **information structure** is trivial, but. . .

→ main difficulty: calculation of the expectation.

Extension of Problem (\mathcal{P}) — Open-Loop Case (1)

Consider Problem (\mathcal{P}) and suppose that J is the **expectation** of a function j , depending on a **random variable** \mathbf{W} defined on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and valued on $(\mathbb{W}, \mathcal{W})$:

$$J(u) = \mathbb{E}(j(u, \mathbf{W})) .$$

Then the optimization problem writes

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) .$$

Decision u is a **deterministic variable**. The available information is the probability law of \mathbf{W} (no on-line observation of \mathbf{W}), that is, an **open-loop** situation. The **information structure** is trivial, but. . .

↪ **main difficulty**: **calculation** of the expectation.

Extension of Problem (\mathcal{P}) — Open-Loop Case (2)

Solution using Exact Quadrature

$$J(u) = \mathbb{E}(j(u, \mathbf{W})) \quad , \quad \nabla J(u) = \mathbb{E}(\nabla_u j(u, \mathbf{W})) .$$

Projected gradient algorithm:

$$u^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - \epsilon \nabla J(u^{(k)}) \right) .$$

Sample-Average Approximation (SAA)

Obtain a realization $(w^{(1)}, \dots, w^{(k)})$ of a k -sample of W and minimize the Monte Carlo approximation of J :

$$u^{(k)} \in \arg \min_{u \in U^{\text{ad}}} \frac{1}{k} \sum_{l=1}^k j(u, w^{(l)}) .$$

Note that $u^{(k)}$ depends on the realization $(w^{(1)}, \dots, w^{(k)})$!

Extension of Problem (\mathcal{P}) — Open-Loop Case (2)

Solution using Exact Quadrature

$$J(u) = \mathbb{E}(j(u, \mathbf{W})) \quad , \quad \nabla J(u) = \mathbb{E}(\nabla_u j(u, \mathbf{W})) .$$

Projected gradient algorithm:

$$u^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - \epsilon \nabla J(u^{(k)}) \right) .$$

Sample Average Approximation (SAA)

Obtain a realization $(w^{(1)}, \dots, w^{(k)})$ of a k -sample of \mathbf{W} and minimize the **Monte Carlo approximation** of J :

$$u^{(k)} \in \arg \min_{u \in U^{\text{ad}}} \frac{1}{k} \sum_{l=1}^k j(u, w^{(l)}) .$$

Note that $u^{(k)}$ **depends** on the realization $(w^{(1)}, \dots, w^{(k)})$!

Extension of Problem (\mathcal{P}) — Open-Loop Case

(3)

Stochastic Gradient Method

Underlying ideas:

- incorporate the realizations $(w^{(1)}, \dots, w^{(k)}, \dots)$ of samples of W **one by one** into the algorithm.
- use an **easily computable** approximation of the gradient ∇J , e.g. replace $\nabla J(u^{(k)})$ by $\nabla_u j(u^{(k)}, w^{(k+1)})$,

These considerations lead to the following algorithm:

$$u^{(k+1)} = \text{proj}_{U_{\text{ad}}} \left(u^{(k)} - c^{(k)} \nabla_u j(u^{(k)}, w^{(k+1)}) \right)$$

Iterations of the gradient algorithm are used a) to move towards the solution *and* b) to refine the Monte-Carlo sampling process.

Extension of Problem (\mathcal{P}) — Open-Loop Case

(3)

Stochastic Gradient Method

Underlying ideas:

- incorporate the realizations $(w^{(1)}, \dots, w^{(k)}, \dots)$ of samples of W **one by one** into the algorithm.
- use an **easily computable** approximation of the gradient ∇J , e.g. replace $\nabla J(u^{(k)})$ by $\nabla_u j(u^{(k)}, w^{(k+1)})$,

These considerations lead to the following algorithm:

$$u^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - \epsilon^{(k)} \nabla_u j(u^{(k)}, w^{(k+1)}) \right) .$$

Iterations of the gradient algorithm are used **a)** to move towards the solution *and* **b)** to refine the Monte-Carlo sampling process.

- 1 Stochastic Gradient Algorithm
- 2 Connexion with Stochastic Approximation
- 3 Asymptotic Efficiency and Averaging
- 4 Practical Considerations

Stochastic Gradient (SG) algorithm

Standard Stochastic Gradient Algorithm

$$\min_{u \in U^{\text{ad}} \subset \mathbb{U}} \mathbb{E}(j(u, \mathbf{W})) . \quad (\mathcal{P}_{\text{ol}})$$

- 1 Let $u^{(0)} \in U^{\text{ad}}$ and choose a positive real sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.
- 2 At iteration $(k + 1)$, draw a realization $w^{(k+1)}$ of the r.v. \mathbf{W} .
- 3 Compute the gradient of j and update $u^{(k+1)}$ by the formula:

$$u^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - \epsilon^{(k)} \nabla_u j(u^{(k)}, w^{(k+1)}) \right) .$$

- 4 Set $k = k + 1$ and go to step 2.

Note that $(w^{(1)}, \dots, w^{(k)}, \dots)$ is a realization of a ∞ -sample of \mathbf{W}
 \rightsquigarrow numerical implementation of the stochastic gradient method.

Stochastic Gradient (SG) algorithm

Standard Stochastic Gradient Algorithm

$$\min_{u \in U^{\text{ad}} \subset U} \mathbb{E}(j(u, \mathbf{W})) . \quad (\mathcal{P}_{\text{ol}})$$

- ① Let $u^{(0)} \in U^{\text{ad}}$ and choose a positive real sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.
- ② At iteration $(k + 1)$, draw a realization $w^{(k+1)}$ of the r.v. \mathbf{W} .
- ③ Compute the gradient of j and update $u^{(k+1)}$ by the formula:

$$u^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - \epsilon^{(k)} \nabla_u j(u^{(k)}, w^{(k+1)}) \right) .$$

- ④ Set $k = k + 1$ and go to step 2.

Note that $(w^{(1)}, \dots, w^{(k)}, \dots)$ is a **realization** of a ∞ -sample of \mathbf{W}

\rightsquigarrow **numerical implementation** of the stochastic gradient method.

Probabilistic Considerations

(1)

In order to study the convergence of the algorithm, it is necessary to cast it in the adequate **probabilistic framework**:

$$\mathbf{U}^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(\mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) \right) .$$

where $\{\mathbf{W}^{(k)}\}_{k \in \mathbb{N}}$ is a infinite-dimensional sample of \mathbf{W} .³

↔ Iterative relation involving **random variables**.

- Convergence in law.
- Convergence in probability.
- Convergence in L^p norm.
- Almost sure convergence (the “intuitive” one).

³Note that $(\Omega, \mathcal{A}, \mathbb{P})$ has to be “big enough” to support such a sample.

Probabilistic Considerations

(1)

In order to study the convergence of the algorithm, it is necessary to cast it in the adequate **probabilistic framework**:

$$\mathbf{U}^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(\mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) \right).$$

where $\{\mathbf{W}^{(k)}\}_{k \in \mathbb{N}}$ is a infinite-dimensional sample of \mathbf{W} .³

↔ Iterative relation involving **random variables**.

- Convergence in law.
- Convergence in probability.
- Convergence in L^p norm.
- **Almost sure convergence** (the “intuitive” one).

³Note that $(\Omega, \mathcal{A}, \mathbb{P})$ has to be “big enough” to support such a sample.

Probabilistic Considerations

(2)

An iteration of the algorithm is represented by the general relation:

$$\mathbf{U}^{(k+1)} = \mathcal{R}^{(k)}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) .$$

Let $\mathcal{F}^{(k)}$ be the σ -field generated by $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$.

Probabilistic Considerations

(2)

An iteration of the algorithm is represented by the general relation:

$$\mathbf{U}^{(k+1)} = \mathcal{R}^{(k)}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) .$$

Let $\mathcal{F}^{(k)}$ be the σ -field generated by $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$.

- Since $\mathbf{U}^{(k)}$ is $\mathcal{F}^{(k)}$ -mesurable for all k , we have

$$\mathbb{E}(\mathbf{U}^{(k)} \mid \mathcal{F}^{(k)}) = \mathbf{U}^{(k)} .$$

- Since $\mathbf{W}^{(k+1)}$ is independent of $\mathcal{F}^{(k)}$, we have (disintegration) that the conditional expectation of $\mathbf{U}^{(k+1)}$ w.r.t. $\mathcal{F}^{(k)}$ merely consists of a standard expectation:

$$\mathbb{E}(\mathbf{U}^{(k+1)} \mid \mathcal{F}^{(k)})(\omega) = \int_{\Omega} \mathcal{R}^{(k)}(\mathbf{U}^{(k)}(\omega), \mathbf{W}(\omega')) d\mathbb{P}(\omega') .$$

Probabilistic Considerations

(2)

An iteration of the algorithm is represented by the general relation:

$$\mathbf{U}^{(k+1)} = \mathcal{R}^{(k)}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) .$$

Let $\mathcal{F}^{(k)}$ be the σ -field generated by $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$.

- Since $\mathbf{U}^{(k)}$ is $\mathcal{F}^{(k)}$ -mesurable for all k , we have

$$\mathbb{E}(\mathbf{U}^{(k)} \mid \mathcal{F}^{(k)}) = \mathbf{U}^{(k)} .$$

- Since $\mathbf{W}^{(k+1)}$ is independent of $\mathcal{F}^{(k)}$, we have (disintegration) that the **conditional expectation** of $\mathbf{U}^{(k+1)}$ w.r.t. $\mathcal{F}^{(k)}$ merely consists of a **standard expectation**:

$$\mathbb{E}(\mathbf{U}^{(k+1)} \mid \mathcal{F}^{(k)})(\omega) = \int_{\Omega} \mathcal{R}^{(k)}(\mathbf{U}^{(k)}(\omega), \mathbf{W}(\omega')) \, d\mathbb{P}(\omega') .$$

Example: Estimation of an Expectation

(1)

Let \mathbf{W} be a real-valued random variable defined on $(\Omega, \mathcal{A}, \mathbb{P})$.
We want to compute an **estimate** of its expectation

$$\mathbb{E}(\mathbf{W}) = \int_{\Omega} \mathbf{W}(\omega) d\mathbb{P}(\omega) .$$

Monte Carlo method: obtain a k -sample $(W^{(1)}, \dots, W^{(k)})$ of \mathbf{W}
and compute the associated arithmetic mean:

$$U^{(k)} = \frac{1}{k} \sum_{l=1}^k W^{(l)} .$$

By the Strong Law of Large Numbers (SLLN), the sequence of
random variables $\{U^{(k)}\}_{k \in \mathbb{N}}$ almost surely converges to $\mathbb{E}(\mathbf{W})$.

Example: Estimation of an Expectation

(1)

Let \mathbf{W} be a real-valued random variable defined on $(\Omega, \mathcal{A}, \mathbb{P})$.
 We want to compute an **estimate** of its expectation

$$\mathbb{E}(\mathbf{W}) = \int_{\Omega} \mathbf{W}(\omega) \, d\mathbb{P}(\omega) .$$

Monte Carlo method: obtain a k -sample $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$ of \mathbf{W}
 and compute the associated arithmetic mean:

$$\mathbf{U}^{(k)} = \frac{1}{k} \sum_{l=1}^k \mathbf{W}^{(l)} .$$

By the Strong Law of Large Numbers (SLLN), the sequence of
 random variables $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ **almost surely converges** to $\mathbb{E}(\mathbf{W})$.

Example: Estimation of an Expectation

(2)

A straightforward computation leads to

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{1}{k+1} \left(\mathbf{U}^{(k)} - \mathbf{W}^{(k+1)} \right).$$

Using the notations $\epsilon^{(k)} = 1/(k+1)$ and $j(u, w) = (u - w)^2/2$, the last expression of $\mathbf{U}^{(k+1)}$ writes

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_u j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

which corresponds to the stochastic gradient algorithm applied to:⁴

$$\min_{u \in \mathbb{R}^2} \frac{1}{2} \mathbb{E}((u - W)^2).$$

⁴Recall that $\mathbb{E}(W)$ is the point which minimizes the dispersion of W .

Example: Estimation of an Expectation

(2)

A straightforward computation leads to

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{1}{k+1} \left(\mathbf{U}^{(k)} - \mathbf{W}^{(k+1)} \right) .$$

Using the notations $\epsilon^{(k)} = 1/(k+1)$ and $j(u, w) = (u - w)^2/2$, the last expression of $\mathbf{U}^{(k+1)}$ writes

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_u j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) ,$$

which corresponds to the **stochastic gradient algorithm** applied to:⁴

$$\min_{u \in \mathbb{R}} \frac{1}{2} \mathbb{E} \left((u - \mathbf{W})^2 \right) .$$

⁴Recall that $\mathbb{E}(\mathbf{W})$ is the point which minimizes the dispersion of \mathbf{W} .

Example: Estimation of an Expectation

(3)

This example makes it possible to enlighten some features of the stochastic gradient method.

Example: Estimation of an Expectation

(3)

This example makes it possible to enlighten some features of the stochastic gradient method.

- The step size $\epsilon^{(k)} = 1/(k+1)$ goes to zero as k goes to $+\infty$. Note however that $\epsilon^{(k)}$ goes to zero “not too fast”, that is,

$$\sum_{k \in \mathbb{N}} \epsilon^{(k)} = +\infty .$$

- It is reasonable to expect an almost sure convergence result for the stochastic gradient algorithm (rather than a weaker notion as convergence in distribution or convergence in probability).
- As the Central Limit Theorem (CLT) applies to this case, we may expect a similar result for the rate of convergence of the stochastic gradient algorithm.

Example: Estimation of an Expectation

(3)

This example makes it possible to enlighten some features of the stochastic gradient method.

- The step size $\epsilon^{(k)} = 1/(k+1)$ goes to zero as k goes to $+\infty$. Note however that $\epsilon^{(k)}$ goes to zero “not too fast”, that is,

$$\sum_{k \in \mathbb{N}} \epsilon^{(k)} = +\infty .$$

- It is reasonable to expect an **almost sure convergence** result for the stochastic gradient algorithm (rather than a weaker notion as convergence in distribution or convergence in probability).
- As the Central Limit Theorem (CLT) applies to this case, we may expect a similar result for the rate of convergence of the stochastic gradient algorithm.

Example: Estimation of an Expectation

(3)

This example makes it possible to enlighten some features of the stochastic gradient method.

- The step size $\epsilon^{(k)} = 1/(k+1)$ goes to zero as k goes to $+\infty$. Note however that $\epsilon^{(k)}$ goes to zero “not too fast”, that is,

$$\sum_{k \in \mathbb{N}} \epsilon^{(k)} = +\infty .$$

- It is reasonable to expect an **almost sure convergence** result for the stochastic gradient algorithm (rather than a weaker notion as convergence in distribution or convergence in probability).
- As the Central Limit Theorem (CLT) applies to this case, we may expect a similar result for the **rate of convergence** of the stochastic gradient algorithm.

- 1 Stochastic Gradient Algorithm
- 2 Connexion with Stochastic Approximation
- 3 Asymptotic Efficiency and Averaging
- 4 Practical Considerations

Stochastic Approximation (SA) Framework

A classical problem in Stochastic Approximation is to determine the **zero of a function** $h : \mathbb{U} \rightarrow \mathbb{U}$, with $\mathbb{U} = \mathbb{R}^n$, in case that the observation of $h(u)$ is **perturbed by an additive random variable** ξ .

Given a random process $\{\xi^{(k)}\}_{k \in \mathbb{N}}$ and a filtration $\{\mathcal{F}^{(k)}\}_{k \in \mathbb{N}}$, the **standard SA algorithm** consists in the following iteration:

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} + \epsilon^{(k)} \left(h(\mathbf{U}^{(k)}) + \xi^{(k+1)} \right),$$

Link with the stochastic gradient algorithm:

$$\begin{aligned} h(u) &= -\nabla J(u), \\ \xi^{(k+1)} &= \nabla J(\mathbf{U}^{(k)}) - \nabla_{u^i} J(\mathbf{U}^{(k)}, W^{(k+1)}). \end{aligned}$$

\leadsto Finding u^* s.t. $h(u^*) = 0$ is equivalent to solving $\nabla J(u^*) = 0$.

Stochastic Approximation (SA) Framework

A classical problem in Stochastic Approximation is to determine the **zero of a function** $h : \mathbb{U} \rightarrow \mathbb{U}$, with $\mathbb{U} = \mathbb{R}^n$, in case that the observation of $h(u)$ is **perturbed by an additive random variable** ξ .

Given a random process $\{\xi^{(k)}\}_{k \in \mathbb{N}}$ and a filtration $\{\mathcal{F}^{(k)}\}_{k \in \mathbb{N}}$, the **standard SA algorithm** consists in the following iteration:

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} + \epsilon^{(k)} \left(h(\mathbf{U}^{(k)}) + \xi^{(k+1)} \right),$$

Link with the stochastic gradient algorithm:

$$\begin{aligned} h(u) &= -\nabla J(u), \\ \xi^{(k+1)} &= \nabla J(\mathbf{U}^{(k)}) - \nabla_u j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}). \end{aligned}$$

\rightsquigarrow Finding $u^\#$ s.t. $h(u^\#) = 0$ is **equivalent** to solving $\nabla J(u^\#) = 0$.

Convergence Theorem (SA)

(1)

Assumptions

- ① The random variable $\mathbf{U}^{(0)}$ is $\mathcal{F}^{(0)}$ -mesurable.
- ② The mapping $h : \mathbb{U} \rightarrow \mathbb{U}$ is continuous, such that
 - $\exists u^\sharp \in \mathbb{R}^n, h(u^\sharp) = 0$ and $\langle h(u), u - u^\sharp \rangle < 0, \forall u \neq u^\sharp$;
 - $\exists a > 0, \forall u \in \mathbb{R}^n, \|h(u)\|^2 \leq a(1 + \|u\|^2)$.
- ③ The random variable $\xi^{(k)}$ is $\mathcal{F}^{(k)}$ -mesurable for all k , and
 - $\mathbb{E}(\xi^{(k+1)} \mid \mathcal{F}^{(k)}) = 0$,
 - $\exists d > 0, \mathbb{E}(\|\xi^{(k+1)}\|^2 \mid \mathcal{F}^{(k)}) \leq d(1 + \|\mathbf{U}^{(k)}\|^2)$.
- ④ The sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a **σ -sequence**, that is,

$$\sum_{k \in \mathbb{N}} \epsilon^{(k)} = +\infty, \quad \sum_{k \in \mathbb{N}} (\epsilon^{(k)})^2 < +\infty.$$

Convergence Theorem (SA)

(2)

Robbins-Monro Theorem

Under the previous assumptions, the sequence $\{U^{(k)}\}_{k \in \mathbb{N}}$ of random variables generated by the Stochastic Approximation algorithm **almost surely converges** to the solution $u^\#$ of $h(u) = 0$.

For a proof, see [Duflo, 1997, §1.4].

This theorem can be extended to more general situations.

- A projection operator can be added:

$$U^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(U^{(k)} + e^{(k)} (h(U^{(k)}) + \xi^{(k+1)}) \right)$$

- A “small” additional term $R^{(k+1)}$ can be added:⁵

$$U^{(k+1)} = U^{(k)} + e^{(k)} (h(U^{(k)}) + \xi^{(k+1)} + R^{(k+1)})$$

⁵for example a bias in $h(u)$, as considered in the Kiefer-Wolfowitz algorithm

Convergence Theorem (SA)

(2)

Robbins-Monro Theorem

Under the previous assumptions, the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ of random variables generated by the Stochastic Approximation algorithm **almost surely converges** to the solution $u^\#$ of $h(u) = 0$.

For a proof, see [Duflo, 1997, §1.4].

This theorem can be extended to **more general situations**.

- A projection operator can be added:

$$\mathbf{U}^{(k+1)} = \text{proj}_{U^{\text{ad}}} \left(\mathbf{U}^{(k)} + \epsilon^{(k)} (h(\mathbf{U}^{(k)}) + \xi^{(k+1)}) \right).$$

- A “small” additional term $\mathbf{R}^{(k+1)}$ can be added:⁵

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} + \epsilon^{(k)} (h(\mathbf{U}^{(k)}) + \xi^{(k+1)} + \mathbf{R}^{(k+1)}).$$

⁵for example a bias on $h(u)$, as considered in the **Kiefer-Wolfowitz** algorithm

Rate of Convergence (SA)

(1)

We recall a result about the **asymptotic normality** of the sequence $\{\mathbf{U}^{(k)}\}$ generated by the SA algorithm, that is, an estimation of its **rate of convergence**.

We first need to be more specific about the notion of σ -sequence.

Definition

A positive real sequence $\{c^{(k)}\}_{k \in \mathbb{N}}$ is a $\sigma(\alpha, \beta, \gamma)$ -sequence if it is such that

$$c^{(k)} = \frac{\alpha}{k^\gamma + \beta},$$

with $\alpha > 0$, $\beta \geq 0$ and $1/2 < \gamma \leq 1$.

A consequence of this definition is that a $\sigma(\alpha, \beta, \gamma)$ -sequence is also a σ -sequence.

Rate of Convergence (SA)

(1)

We recall a result about the **asymptotic normality** of the sequence $\{\mathbf{U}^{(k)}\}$ generated by the SA algorithm, that is, an estimation of its **rate of convergence**.

We first need to be more specific about the notion of σ -sequence.

Definition

A positive real sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a **$\sigma(\alpha, \beta, \gamma)$ -sequence** if it is such that

$$\epsilon^{(k)} = \frac{\alpha}{k^\gamma + \beta},$$

with $\alpha > 0$, $\beta \geq 0$ and $1/2 < \gamma \leq 1$.

A consequence of this definition is that a $\sigma(\alpha, \beta, \gamma)$ -sequence is also a σ -sequence.

Rate of Convergence (SA)

(2)

Assumptions

- ① h is continuously differentiable and, in a neighborhood of u^\sharp ,

$$h(u) = -H(u - u^\sharp) + O(\|u - u^\sharp\|^2),$$

where H is a symmetric positive-definite matrix.

- ② The sequence $\{\mathbb{E}(\xi^{(k+1)}(\xi^{(k+1)})^\top \mid \mathcal{F}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to a symmetric positive-definite matrix Γ .
- ③ $\exists \delta > 0$ such that $\sup_{k \in \mathbb{N}} \mathbb{E}(\|\xi^{(k+1)}\|^{2+\delta} \mid \mathcal{F}^{(k)}) < +\infty$.
- ④ The sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a $\sigma(\alpha, \beta, \gamma)$ -sequence.
- ⑤ The square matrix $(H - \lambda I)$ is positive-definite, with

$$\lambda = \begin{cases} 0 & \text{if } \gamma < 1, \\ \frac{1}{2\alpha} & \text{if } \gamma = 1. \end{cases}$$

Rate of Convergence (SA)

(3)

We retain the assumptions ensuring the almost sure convergence.

Central Limit Theorem

Under all previous assumptions, the sequence of random variables $\{(1/\sqrt{\epsilon^{(k)}})(\mathbf{U}^{(k)} - \mathbf{u}^\#)\}_{k \in \mathbb{N}}$ **converges in law** towards a centered gaussian distribution with covariance matrix Σ , that is,

$$\frac{1}{\sqrt{\epsilon^{(k)}}} (\mathbf{U}^{(k)} - \mathbf{u}^\#) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma),$$

in which Σ is the solution of the so-called **Lyapunov equation**

$$(H - \lambda I)\Sigma + \Sigma(H - \lambda I) = \Gamma.$$

For a proof, see [Duflo, 1996, Chapter 4].

Rate of Convergence (SA)

(4)

- The result is valid only for **unconstrained problems**: $U^{\text{ad}} = \mathbb{U}$.
- The result can be rephrased as

$$k^{\frac{1}{2}} \left(U^{(k)} - u^{\sharp} \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \alpha \Sigma),$$

so that β has in fact **no influence** on the convergence rate.

- The choice $\gamma = 1$ achieves the **greatest convergence rate**. We recover the rate $1/\sqrt{k}$ of a standard Monte Carlo estimator.
- If we refer back to the optimization problem $(\mathcal{P}_{\text{ol}})$, that is, $h = -\nabla J$, we notice that H is the Hessian matrix of J at u^{\sharp} :

$$H = \nabla^2 J(u^{\sharp}),$$

and that Γ is the **covariance matrix** of $\nabla_u j$ evaluated at u^{\sharp} :

$$\Gamma = \mathbb{E} \left(\nabla_u j(u^{\sharp}, W) (\nabla_u j(u^{\sharp}, W))^{\top} \right).$$

Rate of Convergence (SA)

(4)

- The result is valid only for **unconstrained problems**: $U^{\text{ad}} = \mathbb{U}$.
- The result can be rephrased as

$$k^{\frac{\gamma}{2}} \left(\mathbf{U}^{(k)} - u^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \alpha \Sigma),$$

so that β has in fact **no influence** on the convergence rate.

- The choice $\gamma = 1$ achieves the **greatest convergence rate**. We recover the rate $1/\sqrt{k}$ of a standard Monte Carlo estimator.
- If we refer back to the optimization problem $(\mathcal{P}_{\text{ol}})$, that is, $h = -\nabla J$, we notice that H is the Hessian matrix of J at $u^\#$:

$$H = \nabla^2 J(u^\#),$$

and that Γ is the **covariance matrix** of $\nabla_u j$ evaluated at $u^\#$:

$$\Gamma = \mathbb{E} \left(\nabla_u j(u^\#, W) (\nabla_u j(u^\#, W))^T \right).$$

Rate of Convergence (SA)

(4)

- The result is valid only for **unconstrained problems**: $U^{\text{ad}} = \mathbb{U}$.
- The result can be rephrased as

$$k^{\frac{\gamma}{2}} \left(\mathbf{U}^{(k)} - u^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \alpha \Sigma),$$

so that β has in fact **no influence** on the convergence rate.

- The choice $\gamma = 1$ achieves the **greatest convergence rate**. We recover the rate $1/\sqrt{k}$ of a standard **Monte Carlo estimator**.
- If we refer back to the optimization problem $(\mathcal{P}_{\text{ol}})$, that is, $h = -\nabla J$, we notice that H is the Hessian matrix of J at $u^\#$:

$$H = \nabla^2 J(u^\#),$$

and that Γ is the covariance matrix of $\nabla_u j$ evaluated at $u^\#$:

$$\Gamma = \mathbb{E} \left(\nabla_u j(u^\#, W) (\nabla_u j(u^\#, W))^T \right).$$

Rate of Convergence (SA)

(4)

- The result is valid only for **unconstrained problems**: $U^{\text{ad}} = \mathbb{U}$.
- The result can be rephrased as

$$k^{\frac{\gamma}{2}} \left(\mathbf{U}^{(k)} - u^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \alpha \Sigma),$$

so that β has in fact **no influence** on the convergence rate.

- The choice $\gamma = 1$ achieves the **greatest convergence rate**. We recover the rate $1/\sqrt{k}$ of a standard **Monte Carlo estimator**.
- If we refer back to the **optimization problem** (\mathcal{P}_{ol}), that is, $h = -\nabla J$, we notice that H is the **Hessian matrix** of J at $u^\#$:

$$H = \nabla^2 J(u^\#),$$

and that Γ is the **covariance matrix** of $\nabla_u j$ evaluated at $u^\#$:

$$\Gamma = \mathbb{E} \left(\nabla_u j(u^\#, \mathbf{W}) (\nabla_u j(u^\#, \mathbf{W}))^\top \right).$$

- 1 Stochastic Gradient Algorithm
- 2 Connexion with Stochastic Approximation
- 3 Asymptotic Efficiency and Averaging**
- 4 Practical Considerations

Stochastic Newton Algorithm

(1)

Here, the step sizes $\epsilon^{(k)}$ are built using the optimal choice $\gamma = 1$. The scalar gain α is replaced by a matrix gain A , where A is a symmetric positive-definite matrix. The SG algorithm becomes

$$U^{(k+1)} = U^{(k)} - \frac{1}{k+\beta} A \nabla_{U^k} J(U^{(k)}, W^{(k+1)}),$$

which in the Stochastic Approximation setting writes

$$U^{(k+1)} = U^{(k)} + \frac{1}{k+\beta} \left(A h(U^{(k)}) + A \xi^{(k+1)} \right).$$

The Central Limit Theorem is thus available, and we have

$$\sqrt{k} \left(U^{(k)} - U^* \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_A),$$

where Σ_A is the unique solution of

$$\left(AH - \frac{I}{2} \right) \Sigma_A + \Sigma_A \left(HA - \frac{I}{2} \right) = A \Gamma A.$$

Stochastic Newton Algorithm

(1)

Here, the step sizes $\epsilon^{(k)}$ are built using the optimal choice $\gamma = 1$. The *scalar* gain α is replaced by a **matrix** gain A , where A is a symmetric positive-definite matrix. The **SG** algorithm becomes

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{1}{k + \beta} A \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

which in the Stochastic Approximation setting writes

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} + \frac{1}{k + \beta} \left(A h(\mathbf{U}^{(k)}) + A \xi^{(k+1)} \right).$$

The Central Limit Theorem is thus available, and we have

$$\sqrt{k} \left(\mathbf{U}^{(k)} - \mathbf{U}^{\#} \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_A),$$

where Σ_A is the unique solution of

$$\left(AH - \frac{I}{2} \right) \Sigma_A + \Sigma_A \left(HA - \frac{I}{2} \right) = A \Gamma A.$$

Stochastic Newton Algorithm

(1)

Here, the step sizes $\epsilon^{(k)}$ are built using the optimal choice $\gamma = 1$. The *scalar* gain α is replaced by a **matrix** gain A , where A is a symmetric positive-definite matrix. The **SG** algorithm becomes

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{1}{k + \beta} A \nabla_u j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

which in the Stochastic Approximation setting writes

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} + \frac{1}{k + \beta} \left(A h(\mathbf{U}^{(k)}) + A \xi^{(k+1)} \right).$$

The **Central Limit Theorem** is thus available, and we have

$$\sqrt{k} \left(\mathbf{U}^{(k)} - \mathbf{u}^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \Sigma_A),$$

where Σ_A is the unique solution of

$$\left(AH - \frac{I}{2} \right) \Sigma_A + \Sigma_A \left(HA - \frac{I}{2} \right) = A \Gamma A.$$

Stochastic Newton Algorithm

(2)

Let \mathcal{C}_H be the set of symmetric positive-definite matrices A , such that $AH - I/2$ is a positive-definite matrix.

Theorem

The choice $A^\# = H^{-1}$ for the matrix A minimizes the asymptotic covariance matrix Σ_A over the set \mathcal{C}_H . The expression of the minimal asymptotic covariance matrix is

$$\Sigma_{A^\#} = H^{-1}\Gamma H^{-1}.$$

Sketch of proof. Rewrite the covariance matrix Σ_A as $\Delta_A + H^{-1}\Gamma H^{-1}$. Then the matrix Δ_A satisfies a Lyapunov equation, whose solution is thus semi-definite positive, hence the result. \square

Stochastic Newton Algorithm

(3)

Definition

A stochastic gradient algorithm is **Newton-efficient** if the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ it generates has the same asymptotic convergence rate as the optimal Newton algorithm, namely

$$\sqrt{k} \left(\mathbf{U}^{(k)} - u^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, H^{-1} \Gamma H^{-1}) .$$

Note that the improvement is on the covariance matrix of the Gaussian distribution. The rate of convergence remains $1/\sqrt{k}$.

Question. How to obtain an **implementable** Newton-efficient stochastic gradient algorithm?

Stochastic Gradient Algorithm with Averaging

(1)

[Polyak, 1992] proposed to implement a Newton-efficient algorithm by incorporating a new **averaging stage** in the standard algorithm.

Assuming that the admissible set U^{ad} is equal to the space U , the standard stochastic gradient algorithm iteration

$$U^{(k+1)} = U^{(k)} - \epsilon^{(k)} \nabla_u J(U^{(k)}, W^{(k+1)}),$$

is replaced by

$$U^{(k+1)} = U^{(k)} - \epsilon^{(k)} \nabla_u J(U^{(k)}, W^{(k+1)}),$$

$$U_M^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} U^{(l)}.$$

Note that an equivalent recursive form for the averaging stage is

$$U_M^{(k+1)} = U_M^{(k)} + \frac{1}{k+1} (U^{(k+1)} - U_M^{(k)}).$$

Stochastic Gradient Algorithm with Averaging

(1)

[Polyak, 1992] proposed to implement a Newton-efficient algorithm by incorporating a new **averaging stage** in the standard algorithm. Assuming that the admissible set U^{ad} is equal to the space \mathbb{U} , the standard stochastic gradient algorithm iteration

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

is replaced by

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

$$\mathbf{U}_M^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} \mathbf{U}^{(l)}.$$

Note that an equivalent recursive form for the averaging stage is

$$\mathbf{U}_M^{(k+1)} = \mathbf{U}_M^{(k)} + \frac{1}{k+1} (\mathbf{U}^{(k+1)} - \mathbf{U}_M^{(k)}).$$

Stochastic Gradient Algorithm with Averaging

(1)

[Polyak, 1992] proposed to implement a Newton-efficient algorithm by incorporating a new **averaging stage** in the standard algorithm. Assuming that the admissible set U^{ad} is equal to the space \mathbb{U} , the standard stochastic gradient algorithm iteration

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

is replaced by

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

$$\mathbf{U}_{\text{M}}^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} \mathbf{U}^{(l)}.$$

Note that an equivalent recursive form for the averaging stage is

$$\mathbf{U}_{\text{M}}^{(k+1)} = \mathbf{U}_{\text{M}}^{(k)} + \frac{1}{k+1} (\mathbf{U}^{(k+1)} - \mathbf{U}_{\text{M}}^{(k)}).$$

Stochastic Gradient Algorithm with Averaging

(1)

[Polyak, 1992] proposed to implement a Newton-efficient algorithm by incorporating a new **averaging stage** in the standard algorithm. Assuming that the admissible set U^{ad} is equal to the space \mathbb{U} , the standard stochastic gradient algorithm iteration

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

is replaced by

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} \nabla_{\mathbf{u}} j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}),$$

$$\mathbf{U}_{\text{M}}^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} \mathbf{U}^{(l)}.$$

Note that an **equivalent recursive form** for the averaging stage is

$$\mathbf{U}_{\text{M}}^{(k+1)} = \mathbf{U}_{\text{M}}^{(k)} + \frac{1}{k+1} (\mathbf{U}^{(k+1)} - \mathbf{U}_{\text{M}}^{(k)}).$$

Stochastic Gradient Algorithm with Averaging

(2)

Theorem

Under the additional assumption that the $\sigma(\alpha, \beta, \gamma)$ -sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is such that $\gamma < 1$, the averaged stochastic gradient algorithm is Newton-efficient:

$$\sqrt{k} \left(\mathbf{U}_M^{(k)} - u^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, H^{-1} \Gamma H^{-1}) .$$

For a proof, see [Duflo, 1996, Chapter 4].

According to the standard theorem, the convergence rate achieved by the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ with $\gamma < 1$ is smaller than $1/\sqrt{k}$ and hence not optimal. The “nice” convergence properties are obtained regarding the averaged sequence $\{\mathbf{U}_M^{(k)}\}_{k \in \mathbb{N}}$.

Stochastic Gradient Algorithm with Averaging

(2)

Theorem

Under the additional assumption that the $\sigma(\alpha, \beta, \gamma)$ -sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is such that $\gamma < 1$, the averaged stochastic gradient algorithm is Newton-efficient:

$$\sqrt{k} \left(\mathbf{U}_M^{(k)} - u^\# \right) \xrightarrow{\mathcal{D}} \mathcal{N}(0, H^{-1} \Gamma H^{-1}).$$

For a proof, see [Duflo, 1996, Chapter 4].

According to the standard theorem, the convergence rate achieved by the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ with $\gamma < 1$ is smaller than $1/\sqrt{k}$ and hence not optimal. The “nice” convergence properties are obtained regarding the **averaged sequence** $\{\mathbf{U}_M^{(k)}\}_{k \in \mathbb{N}}$.

- 1 Stochastic Gradient Algorithm
- 2 Connexion with Stochastic Approximation
- 3 Asymptotic Efficiency and Averaging
- 4 Practical Considerations

A Toy Problem

Let us consider the following optimization problem:

$$\min_{u \in \mathbb{R}^{10}} \mathbb{E} \left(\frac{1}{2} u^\top B u + \mathbf{W}^\top u \right),$$

B being a symmetric positive definite matrix, and \mathbf{W} being a \mathbb{R}^{10} -valued Gaussian random variable $\mathcal{N}(m, \Gamma)$.

The optimal solution of this problem is $u^* = -B^{-1}m$.

It can be estimated either by Monte Carlo

$$\hat{U}^{(k+1)} = -\frac{1}{k+1} \sum_{l=1}^{k+1} B^{-1} \mathbf{W}^{(l)},$$

or by the standard stochastic gradient algorithm

$$U^{(k+1)} = U^{(k)} - e^{(k)} (B U^{(k)} + \mathbf{W}^{(k+1)})$$

A Toy Problem

Let us consider the following optimization problem:

$$\min_{u \in \mathbb{R}^{10}} \mathbb{E} \left(\frac{1}{2} u^\top B u + \mathbf{W}^\top u \right),$$

B being a symmetric positive definite matrix, and \mathbf{W} being a \mathbb{R}^{10} -valued Gaussian random variable $\mathcal{N}(m, \Gamma)$.

The **optimal solution** of this problem is $u^\# = -B^{-1}m$.

It can be estimated either by Monte Carlo

$$\hat{u}^{(k+1)} = -\frac{1}{k+1} \sum_{l=1}^{k+1} B^{-1} \mathbf{W}^{(l)},$$

or by the standard stochastic gradient algorithm

$$u^{(k+1)} = u^{(k)} - \epsilon^{(k)} (B u^{(k)} + \mathbf{W}^{(k+1)})$$

A Toy Problem

Let us consider the following optimization problem:

$$\min_{u \in \mathbb{R}^{10}} \mathbb{E} \left(\frac{1}{2} u^\top B u + \mathbf{W}^\top u \right),$$

B being a symmetric positive definite matrix, and \mathbf{W} being a \mathbb{R}^{10} -valued Gaussian random variable $\mathcal{N}(m, \Gamma)$.

The **optimal solution** of this problem is $u^\# = -B^{-1}m$.

It can be estimated either by **Monte Carlo**

$$\hat{\mathbf{U}}^{(k+1)} = -\frac{1}{k+1} \sum_{l=1}^{k+1} B^{-1} \mathbf{W}^{(l)},$$

or by the standard **stochastic gradient** algorithm

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \epsilon^{(k)} (B \mathbf{U}^{(k)} + \mathbf{W}^{(k+1)}).$$

Tuning the Standard Algorithm

(1)

Let $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ be a $\sigma(\alpha, \beta, \gamma)$ -sequence, that is, $\epsilon^{(k)} = \frac{\alpha}{k\gamma + \beta}$.

Tuning the Standard Algorithm

(1)

Let $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ be a $\sigma(\alpha, \beta, \gamma)$ -sequence, that is, $\epsilon^{(k)} = \frac{\alpha}{k\gamma + \beta}$.

- The **best convergence rate** is reached for $\gamma = 1$.
- The coefficient α influences the **asymptotic behavior**. The covariance matrix $\alpha \Sigma$ grows as α goes to $+\infty$,⁶ but using too small values for α may generate very small gradient steps. The choice of α corresponds to a **trade-off** between stability and precision.
- Ultimately, the coefficient β makes it possible to regulate the **transient behavior** of the algorithm. During the first iterations, the coefficient $\epsilon^{(k)}$ is approximately equal to α/β . If α/β is too small, the transient phase may be slow. On the contrary, taking a too large ratio may lead to a **numerical burst**.

⁶remember that Σ depends on α

Tuning the Standard Algorithm

(1)

Let $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ be a $\sigma(\alpha, \beta, \gamma)$ -sequence, that is, $\epsilon^{(k)} = \frac{\alpha}{k\gamma + \beta}$.

- The **best convergence rate** is reached for $\gamma = 1$.
- The coefficient α influences the **asymptotic behavior**.
 The covariance matrix $\alpha\Sigma$ grows as α goes to $+\infty$,⁶ but using too small values for α may generate very small gradient steps. The choice of α corresponds to a **trade-off** between stability and precision.
- Ultimately, the coefficient β makes it possible to regulate the **transient behavior** of the algorithm. During the first iterations, the coefficient $\epsilon^{(k)}$ is approximately equal to α/β . If α/β is too small, the transient phase may be slow. On the contrary, taking a too large ratio may lead to a **numerical burst**.

⁶remember that Σ depends on $\alpha \dots$

Tuning the Standard Algorithm

(1)

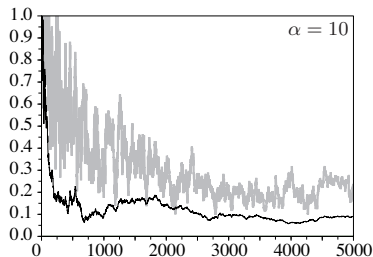
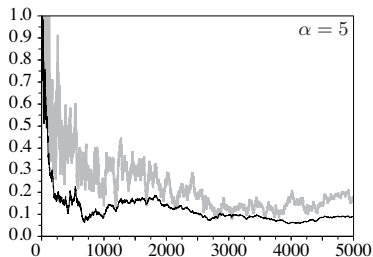
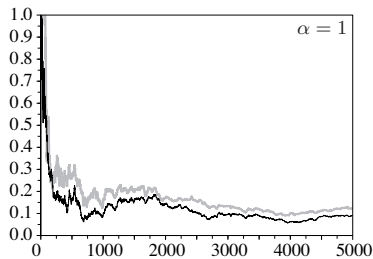
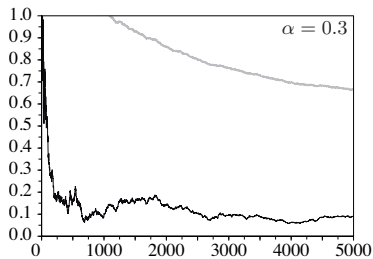
Let $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ be a $\sigma(\alpha, \beta, \gamma)$ -sequence, that is, $\epsilon^{(k)} = \frac{\alpha}{k\gamma + \beta}$.

- The **best convergence rate** is reached for $\gamma = 1$.
- The coefficient α influences the **asymptotic behavior**.
 The covariance matrix $\alpha\Sigma$ grows as α goes to $+\infty$,⁶ but using too small values for α may generate very small gradient steps. The choice of α corresponds to a **trade-off** between stability and precision.
- Ultimately, the coefficient β makes it possible to regulate the **transient behavior** of the algorithm. During the first iterations, the coefficient $\epsilon^{(k)}$ is approximately equal to α/β . If α/β is too small, the transient phase may be slow. On the contrary, taking a too large ratio may lead to a **numerical burst**.

⁶remember that Σ depends on $\alpha \dots$

Tuning the Standard Algorithm ($\alpha/\beta = 0.1$)

(2)



Tuning the Averaged Algorithm

(1)

Here, $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a $\sigma(\alpha, \beta, \gamma)$ -sequence, with $1/2 < \gamma < 1$.

The **averaged stochastic gradient algorithm** writes on our example

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{\alpha}{k^\gamma + \beta} (\mathbf{B}\mathbf{U}^{(k)} + \mathbf{W}^{(k+1)}) \quad , \quad \mathbf{U}_M^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} \mathbf{U}^{(l)} .$$

Tuning the Averaged Algorithm

(1)

Here, $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a $\sigma(\alpha, \beta, \gamma)$ -sequence, with $1/2 < \gamma < 1$.

The **averaged stochastic gradient algorithm** writes on our example

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{\alpha}{k^\gamma + \beta} (B\mathbf{U}^{(k)} + \mathbf{W}^{(k+1)}) \quad , \quad \mathbf{U}_M^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} \mathbf{U}^{(l)} .$$

- The value $\gamma = 2/3$ is considered as a **good choice**.
- The tuning of parameters α and β is much easier than for the standard algorithm. Indeed, the problem of “too small” step sizes arising from a bad choice of α is not so critical because the term $k^{-\gamma}$ goes down more slowly towards zero. Of course, the ratio α/β must be chosen in such a way that numerical bursts do not occur during the first iterations of the algorithm.

Tuning the Averaged Algorithm

(1)

Here, $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a $\sigma(\alpha, \beta, \gamma)$ -sequence, with $1/2 < \gamma < 1$.

The **averaged stochastic gradient algorithm** writes on our example

$$\mathbf{U}^{(k+1)} = \mathbf{U}^{(k)} - \frac{\alpha}{k^\gamma + \beta} (B\mathbf{U}^{(k)} + \mathbf{W}^{(k+1)}) \quad , \quad \mathbf{U}_M^{(k+1)} = \frac{1}{k+1} \sum_{l=1}^{k+1} \mathbf{U}^{(l)} .$$

- The value $\gamma = 2/3$ is considered as a **good choice**.
- The tuning of parameters α and β is **much easier** than for the standard algorithm. Indeed, the problem of “too small” step sizes arising from a bad choice of α is not so critical because the term $k^{-\gamma}$ goes down more slowly towards zero. Of course, the ratio α/β must be chosen in such a way that numerical bursts do not occur during the first iterations of the algorithm.

Tuning the Averaged Algorithm ($\alpha/\beta = 0.1$)

(2)

