

Generalized Stochastic Gradient Method

Lecture Outline

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Problem under Consideration

Consider the following **convex differentiable** optimization problem:

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

Let $u^{\#} \in U^{\text{ad}}$ be a solution of this problem (assume that such a solution exists). The associated optimality condition writes:

$$\langle \nabla J(u^{\#}), u - u^{\#} \rangle \geq 0, \quad \forall u \in U^{\text{ad}} .$$

In the deterministic framework, the **Auxiliary Problem Principle (APP)** consists in replacing the original problem by a sequence of auxiliary problems indexed by $k \in \mathbb{N}$, with similar optimality conditions to those of the original problem. . .

Problem under Consideration

Consider the following **convex differentiable** optimization problem:

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

Let $u^\# \in U^{\text{ad}}$ be a **solution** of this problem (assume that such a solution exists). The associated **optimality condition** writes:

$$\langle \nabla J(u^\#), u - u^\# \rangle \geq 0, \quad \forall u \in U^{\text{ad}}.$$

In the deterministic framework, the **Auxiliary Problem Principle (APP)** consists in replacing the original problem by a sequence of auxiliary problems indexed by $k \in \mathbb{N}$, with similar optimality conditions to those of the original problem...

Problem under Consideration

Consider the following **convex differentiable** optimization problem:

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

Let $u^\# \in U^{\text{ad}}$ be a **solution** of this problem (assume that such a solution exists). The associated **optimality condition** writes:

$$\langle \nabla J(u^\#), u - u^\# \rangle \geq 0, \quad \forall u \in U^{\text{ad}} .$$

In the deterministic framework, the **Auxiliary Problem Principle (APP)** consists in replacing the original problem by a **sequence of auxiliary problems** indexed by $k \in \mathbb{N}$, with similar optimality conditions to those of the original problem. . .

APP Framework

Principle: replace $J(u)$ by its **first order** approximation at $u^{(k)}$:

$$J(u) \approx J(u^{(k)}) + \left\langle \nabla J(u^{(k)}), u - u^{(k)} \right\rangle ,$$

and add a **chosen strongly convex** term to regain coercivity:

$$\frac{1}{\epsilon} \left(K(u) - K(u^{(k)}) - \left\langle \nabla K(u^{(k)}), u - u^{(k)} \right\rangle \right) ,$$

K being a function defined on \mathbb{U} and ϵ being a positive constant.

This leads to consider the following **auxiliary problem**:

$$\min_{u \in \mathbb{U}^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle ,$$

whose unique solution is denoted $u^{(k+1)}$.

APP Framework

Principle: replace $J(u)$ by its **first order** approximation at $u^{(k)}$:

$$J(u) \approx J(u^{(k)}) + \left\langle \nabla J(u^{(k)}), u - u^{(k)} \right\rangle ,$$

and add a **chosen strongly convex** term to regain coercivity:

$$\frac{1}{\epsilon} \left(K(u) - K(u^{(k)}) - \left\langle \nabla K(u^{(k)}), u - u^{(k)} \right\rangle \right) ,$$

K being a function defined on \mathbb{U} and ϵ being a positive constant.

This leads to consider the following **auxiliary problem**:

$$\min_{u \in \mathbb{U}^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle ,$$

whose unique solution is denoted $u^{(k+1)}$.

APP Framework

Principle: replace $J(u)$ by its **first order** approximation at $u^{(k)}$:

$$J(u) \approx J(u^{(k)}) + \left\langle \nabla J(u^{(k)}), u - u^{(k)} \right\rangle ,$$

and add a **chosen strongly convex** term to regain coercivity:

$$\frac{1}{\epsilon} \left(K(u) - K(u^{(k)}) - \left\langle \nabla K(u^{(k)}), u - u^{(k)} \right\rangle \right) ,$$

K being a function defined on \mathbb{U} and ϵ being a positive constant.

This leads to consider the following **auxiliary problem**:

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle ,$$

whose unique solution is denoted $u^{(k+1)}$.

APP Algorithm

- 1 Choose a core function K and a coefficient $\epsilon > 0$.
- 2 Choose $u^{(0)} \in U^{\text{ad}}$ and a tolerance $\sigma > 0$. Set $k = 0$.
- 3 Obtain the solution $u^{(k+1)}$ of the auxiliary problem

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle .$$

- 4 Set $k = k + 1$ and go to step 3 until $\|u^{(k+1)} - u^{(k)}\| < \sigma$.

Note that the optimality condition of the auxiliary problem is

$$\langle \nabla K(u^{(k+1)}) + \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u - u^{(k+1)} \rangle \geq 0, \quad \forall u \in U^{\text{ad}}$$

and coincides with the optimality condition of the initial problem in case where the sequence $\{u^{(k)}\}_{k \in \mathbb{N}}$ converges.

APP Algorithm

- 1 Choose a core function K and a coefficient $\epsilon > 0$.
- 2 Choose $u^{(0)} \in U^{\text{ad}}$ and a tolerance $\sigma > 0$. Set $k = 0$.
- 3 Obtain the solution $u^{(k+1)}$ of the auxiliary problem

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle .$$

- 4 Set $k = k + 1$ and go to step 3 until $\|u^{(k+1)} - u^{(k)}\| < \sigma$.

Note that the **optimality condition** of the auxiliary problem is

$$\left\langle \nabla K(u^{(k+1)}) + \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u - u^{(k+1)} \right\rangle \geq 0, \quad \forall u \in U^{\text{ad}},$$

and **coincides with the optimality condition of the initial problem** in case where the sequence $\{u^{(k)}\}_{k \in \mathbb{N}}$ converges.

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - **Convergence and Features**
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Convergence Theorem

(1)

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

We make the following assumptions.

H1 U^{ad} is a nonempty, closed and convex subset of an **Hilbert** space \mathbb{U} .

H2 J is a proper l.s.c. **convex** function with $\text{dom} J \cap U^{\text{ad}} \neq \emptyset$, J is coercive on U^{ad} , Gâteaux differentiable, and its gradient ∇J is Lipschitz with constant A .

H3 K is a proper l.s.c. function, **strongly convex** with modulus b and differentiable, ∇K being Lipschitz with constant B .

H4 ϵ is such that $0 < \epsilon < \frac{2b}{A}$.

Convergence Theorem

(2)

Then the following conclusions hold true.

- R1** The initial problem admits at least a solution u^\sharp , and each auxiliary problem admits a **unique** solution $u^{(k+1)}$.
- R2** The real sequence $\{J(u^{(k)})\}_{k \in \mathbb{N}}$ is **strictly decreasing** and converges towards $J^\sharp = J(u^\sharp)$.
- R3** The sequence $\{u^{(k)}\}_{k \in \mathbb{N}}$ is bounded, and every **cluster point** of this sequence is a solution of the initial problem.

Assume moreover that

H5 J is **strongly convex** with modulus α .

Then we obtain that

R4 the sequence $\{u^{(k)}\}_{k \in \mathbb{N}}$ converges towards the unique solution u^\sharp of the initial problem

Convergence Theorem

(2)

Then the following conclusions hold true.

- R1** The initial problem admits at least a solution u^\sharp , and each auxiliary problem admits a **unique** solution $u^{(k+1)}$.
- R2** The real sequence $\{J(u^{(k)})\}_{k \in \mathbb{N}}$ is **strictly decreasing** and converges towards $J^\sharp = J(u^\sharp)$.
- R3** The sequence $\{u^{(k)}\}_{k \in \mathbb{N}}$ is bounded, and every **cluster point** of this sequence is a solution of the initial problem.

Assume moreover that

- H5** J is **strongly convex** with modulus a .

Then we obtain that

- R4** the sequence $\{u^{(k)}\}_{k \in \mathbb{N}}$ converges towards the **unique** solution u^\sharp of the initial problem.

Sketch of Proof

The proof of the first statement is based on classical theorems.

The proof of the last two statements involves four steps.

- 1 Select a **Lyapunov function** Λ .
- 2 Prove that $\{\Lambda(u^{(k)})\}_{k \in \mathbb{N}}$ is a **decreasing sequence**.
Then it converges, and $\{u^{(k)}\}_{k \in \mathbb{N}}$ is a bounded sequence.
- 3 Characterize the **limit** of the sequence $\{\Lambda(u^{(k)})\}_{k \in \mathbb{N}}$.
- 4 Extract a **converging subsequence** of $\{u^{(k)}\}_{k \in \mathbb{N}}$ and characterize its limit.

*Note that the result holds true if \mathbb{U} is an **infinite dimensional Hilbert space**.*

Some Features of APP

(1)

One can take advantage of a **proper choice** of K in order to obtain many special features for the auxiliary subproblems. The reader is referred to [Carpentier et Cohen, 2017] for a detailed description of the **APP**. Two of its main properties are examined hereafter.

- APP encompasses "classical" optimization algorithms. Choosing $K(u) = \|u\|^2 / 2$, the auxiliary problem writes

$$\min_{u \in U_{\text{ad}}} \frac{1}{2} \|u\|^2 + \langle c \nabla J(u^{(k)}) - u^{(k)}, u \rangle,$$

and its solution has the following closed-form expression:

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

We obtain the well-known **projected gradient** algorithm.

Some Features of APP

(1)

One can take advantage of a **proper choice** of K in order to obtain many special features for the auxiliary subproblems. The reader is referred to [Carpentier et Cohen, 2017] for a detailed description of the **APP**. Two of its main properties are examined hereafter.

- APP encompasses “classical” optimization algorithms. Choosing $K(u) = \|u\|^2/2$, the auxiliary problem writes

$$\min_{u \in U^{\text{ad}}} \frac{1}{2} \|u\|^2 + \left\langle \epsilon \nabla J(u^{(k)}) - u^{(k)}, u \right\rangle ,$$

and its solution has the following closed-form expression:

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

We obtain the well-known **projected gradient** algorithm.

Some Features of APP

(2)

- APP allows for **decomposition**. Assume that the space \mathbb{U} is a Cartesian product of N spaces: $\mathbb{U} = \mathbb{U}_1 \times \dots \times \mathbb{U}_N$, and that $U^{\text{ad}} = U_1^{\text{ad}} \times \dots \times U_N^{\text{ad}}$, with $U_i^{\text{ad}} \subset \mathbb{U}_i$. Choosing a function K **additive** with respect to that decomposition of u , that is,

$$K(u_1, \dots, u_N) = K_1(u_1) + \dots + K_N(u_N),$$

the auxiliary subproblem becomes

$$\min_{u_1 \in U_1^{\text{ad}}, \dots, u_N \in U_N^{\text{ad}}} \sum_{i=1}^N \left(K_i(u_i) + \left\langle \epsilon \nabla_{u_i} J(u^{(k)}) - \nabla K_i(u_i^{(k)}), u_i \right\rangle \right).$$

This subproblem splits up into N independent subproblems, the i -th subproblem being

$$\min_{u_i \in U_i^{\text{ad}}} K_i(u_i) + \left\langle \epsilon \nabla_{u_i} J(u^{(k)}) - \nabla K_i(u_i^{(k)}), u_i \right\rangle.$$

Some Features of APP

(2)

- APP allows for **decomposition**. Assume that the space \mathbb{U} is a Cartesian product of N spaces: $\mathbb{U} = \mathbb{U}_1 \times \dots \times \mathbb{U}_N$, and that $U^{\text{ad}} = U_1^{\text{ad}} \times \dots \times U_N^{\text{ad}}$, with $U_i^{\text{ad}} \subset \mathbb{U}_i$. Choosing a function K **additive** with respect to that decomposition of u , that is,

$$K(u_1, \dots, u_N) = K_1(u_1) + \dots + K_N(u_N),$$

the auxiliary subproblem becomes

$$\min_{u_1 \in U_1^{\text{ad}}, \dots, u_N \in U_N^{\text{ad}}} \sum_{i=1}^N \left(K_i(u_i) + \left\langle \epsilon \nabla_{u_i} J(u^{(k)}) - \nabla K_i(u_i^{(k)}), u_i \right\rangle \right).$$

This subproblem splits up into N **independent subproblems**, the i -th subproblem being

$$\min_{u_i \in U_i^{\text{ad}}} K_i(u_i) + \left\langle \epsilon \nabla_{u_i} J(u^{(k)}) - \nabla K_i(u_i^{(k)}), u_i \right\rangle.$$

Mirror Descent — Problem settings

We consider the problem

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

where U^{ad} is a closed convex subset of a **Banach** space \mathbb{U} , J being a proper l.s.c. convex and differentiable function. We denote by

- \mathbb{U}^* the topological dual of \mathbb{U} ,
- $\langle \cdot, \cdot \rangle : \mathbb{U}^* \times \mathbb{U} \rightarrow \mathbb{R}$ the duality product associated to $(\mathbb{U}, \mathbb{U}^*)$,
- J^* the Fenchel conjugate of J : $J^*(u^*) = \sup_{u \in \mathbb{U}} \langle u^*, u \rangle - J(u)$.

Difficulty: In a **Banach** space, $u^{(k)} \in \mathbb{U}$ and $\nabla J(u^{(k)}) \in \mathbb{U}^*$ live in **different** spaces, so that the standard gradient formula

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

does not make sense.

Mirror Descent — Algorithm

Choose $K : \mathbb{U} \rightarrow \mathbb{R}$ strongly convex differentiable, and associated Bregman distance $D_K(u, v) = K(u) - K(v) - \langle \nabla K(v), u - v \rangle$.

Mirror Descent Algorithm

- 1 **Map** $u^{(k)} \in \mathbb{U}$ in the dual space \mathbb{U}^*

$$u^{(k)} \mapsto u^{*(k)} = \nabla K(u^{(k)}) .$$

- 2 **Perform a gradient step** in the dual space \mathbb{U}^*

$$u^{*(k)} \mapsto u^{*(k+1)} = u^{*(k)} - \epsilon \nabla J(u^{(k)}) .$$

- 3 **Map** $u^{*(k+1)} \in \mathbb{U}^*$ in the primal space \mathbb{U}

$$u^{*(k+1)} \mapsto v^{(k+1)} = \nabla K^*(u^{*(k+1)}) .$$

- 4 **Project** $v^{(k+1)} \in \mathbb{U}$ on U^{ad} for the distance D_K

$$v^{(k+1)} \mapsto u^{(k+1)} = \arg \min_{u \in U^{\text{ad}}} D_K(u, v^{(k+1)}) .$$

Mirror Descent — Link with APP

Developing the **MD algorithm** and using $\nabla K \circ \nabla K^* = I$, we have:

$$\begin{aligned}
 u^{(k+1)} &= \arg \min_{u \in U^{\text{ad}}} D_K(u, v^{(k+1)}), \\
 &= \arg \min_{u \in U^{\text{ad}}} K(u) - K(v^{(k+1)}) - \langle \nabla K(v^{(k+1)}), u - v^{(k+1)} \rangle, \\
 &= \arg \min_{u \in U^{\text{ad}}} K(u) - \langle \nabla K(v^{(k+1)}), u \rangle, \\
 &= \arg \min_{u \in U^{\text{ad}}} K(u) - \langle \nabla K(\nabla K^*(u^{*(k+1)})), u \rangle, \\
 &= \arg \min_{u \in U^{\text{ad}}} K(u) - \langle u^{*(k+1)}, u \rangle, \\
 &= \arg \min_{u \in U^{\text{ad}}} K(u) - \langle u^{*(k)} - \epsilon \nabla J(u^{(k)}), u \rangle, \\
 &= \arg \min_{u \in U^{\text{ad}}} K(u) - \langle \nabla K(u^{(k)}) - \epsilon \nabla J(u^{(k)}), u \rangle. \text{ **APP algorithm**
 \end{aligned}$$

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - **Explicit Constraints**
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Problem under Consideration

Consider the following **convex** optimization problem:

$$\min_{u \in U^{\text{ad}} \subset \mathbb{U}} J(u) \quad \text{subject to} \quad \Theta(u) \in -C \subset \mathbb{V},$$

where U^{ad} is a closed convex subset of an Hilbert space \mathbb{U} , and where C is a closed convex salient cone of another Hilbert space \mathbb{V} .

Let C^* be the dual cone^a of C . We introduce the Lagrangian L of the constrained optimization problem, defined on $U^{\text{ad}} \times C^*$:

$$L(u, p) = J(u) + \langle p, \Theta(u) \rangle.$$

Under standard convexity and continuity assumptions, and under a Constraint Qualification Condition, solving the initial problem is equivalent to determining a saddle point of the Lagrangian L .

^adefined as $C^* = \{p \in \mathbb{V}, \langle p, v \rangle \geq 0 \ \forall v \in C\}$.

Problem under Consideration

Consider the following **convex** optimization problem:

$$\min_{u \in U^{\text{ad}} \subset \mathbb{U}} J(u) \quad \text{subject to} \quad \Theta(u) \in -C \subset \mathbb{V},$$

where U^{ad} is a closed convex subset of an Hilbert space \mathbb{U} , and where C is a closed convex salient cone of another Hilbert space \mathbb{V} .

Let C^* be the dual cone^a of C . We introduce the **Lagrangian** L of the constrained optimization problem, defined on $U^{\text{ad}} \times C^*$:

$$L(u, p) = J(u) + \langle p, \Theta(u) \rangle.$$

Under standard convexity and continuity assumptions, and under a **Constraint Qualification Condition**, solving the initial problem is equivalent to determining a **saddle point** of the Lagrangian L .

^adefined as $C^* = \{p \in \mathbb{V}, \langle p, v \rangle \geq 0 \quad \forall v \in C\}$.

Uzawa and Arrow-Hurwicz Algorithms

Assuming that a saddle point of L exists, the initial problem is equivalent to the following **dual problem**:

$$\max_{\rho \in C^*} \left(\min_{u \in U^{\text{ad}}} L(u, \rho) \right).$$

This problem can be solved by using the Uzawa algorithm:

$$\begin{aligned} u^{(k+1)} &\in \arg \min_{u \in U^{\text{ad}}} J(u) + \langle \rho^{(k)}, \Theta(u) \rangle, \\ \rho^{(k+1)} &= \text{proj}_{C^*} (\rho^{(k)} + \rho \Theta(u^{(k+1)})). \end{aligned}$$

Another possibility is to use the Arrow-Hurwicz algorithm:

$$\begin{aligned} u^{(k+1)} &= \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - c(\nabla J(u^{(k)}) + (\Theta'(u^{(k)}))^T \rho^{(k)}) \right), \\ \rho^{(k+1)} &= \text{proj}_{C^*} (\rho^{(k)} + \rho \Theta(u^{(k+1)})). \end{aligned}$$

Uzawa and Arrow-Hurwicz Algorithms

Assuming that a saddle point of L exists, the initial problem is equivalent to the following **dual problem**:

$$\max_{\rho \in C^*} \left(\min_{u \in U^{\text{ad}}} L(u, \rho) \right).$$

This problem can be solved by using the **Uzawa** algorithm:

$$\begin{aligned} u^{(k+1)} &\in \arg \min_{u \in U^{\text{ad}}} J(u) + \langle \rho^{(k)}, \Theta(u) \rangle, \\ \rho^{(k+1)} &= \text{proj}_{C^*} (\rho^{(k)} + \rho \Theta(u^{(k+1)})). \end{aligned}$$

Another possibility is to use the **Arrow-Hurwicz** algorithm:

$$\begin{aligned} u^{(k+1)} &= \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - c (\nabla J(u^{(k)}) + (\Theta'(u^{(k)}))^T \rho^{(k)}) \right), \\ \rho^{(k+1)} &= \text{proj}_{C^*} (\rho^{(k)} + \rho \Theta(u^{(k+1)})). \end{aligned}$$

Uzawa and Arrow-Hurwicz Algorithms

Assuming that a saddle point of L exists, the initial problem is equivalent to the following **dual problem**:

$$\max_{\rho \in C^*} \left(\min_{u \in U^{\text{ad}}} L(u, \rho) \right).$$

This problem can be solved by using the **Uzawa** algorithm:

$$\begin{aligned} u^{(k+1)} &\in \arg \min_{u \in U^{\text{ad}}} J(u) + \langle \rho^{(k)}, \Theta(u) \rangle, \\ \rho^{(k+1)} &= \text{proj}_{C^*} (\rho^{(k)} + \rho \Theta(u^{(k+1)})). \end{aligned}$$

Another possibility is to use the **Arrow-Hurwicz** algorithm:

$$\begin{aligned} u^{(k+1)} &= \text{proj}_{U^{\text{ad}}} \left(u^{(k)} - \epsilon (\nabla J(u^{(k)}) + (\Theta'(u^{(k)}))^{\top} \rho^{(k)}) \right), \\ \rho^{(k+1)} &= \text{proj}_{C^*} (\rho^{(k)} + \rho \Theta(u^{(k+1)})). \end{aligned}$$

APP with Explicit Constraints

A “natural” extension of the APP to constrained optimization problems consists in choosing a function $K : \mathcal{U} \rightarrow \mathbb{R}$ and then replacing the resolution of the initial problem by the resolution of the following **sequence of auxiliary problems**:⁷

$$u^{(k+1)} \in \arg \min_{u \in \mathcal{U}^{\text{ad}}} K(u) + \langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \rangle + \epsilon \langle p^{(k)}, \Theta(u) \rangle,$$

$$p^{(k+1)} = \text{proj}_{\mathcal{C}^*} (p^{(k)} + \rho \Theta(u^{(k+1)})).$$

It is not difficult to show that such a framework encompasses

- the Uzawa algorithm (using $K(u) = J(u)$ and $\epsilon = 1$),
- the Arrow-Hurwicz algorithm (using $K(u) = \|u\|^2 / 2$).

Moreover, choosing an additive core K allows for decomposition in the minimization stage of the APP algorithm.

⁷Note that the term $\epsilon \langle p^{(k)}, \Theta(u) \rangle$ could be replaced by $\epsilon \langle p^{(k)}, \Theta'(u^{(k)}) \cdot u \rangle$.

APP with Explicit Constraints

A “natural” extension of the **APP** to constrained optimization problems consists in choosing a function $K : \mathcal{U} \rightarrow \mathbb{R}$ and then replacing the resolution of the initial problem by the resolution of the following **sequence of auxiliary problems**:⁷

$$u^{(k+1)} \in \arg \min_{u \in \mathcal{U}^{\text{ad}}} K(u) + \langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \rangle + \epsilon \langle p^{(k)}, \Theta(u) \rangle,$$

$$p^{(k+1)} = \text{proj}_{\mathcal{C}^*} (p^{(k)} + \rho \Theta(u^{(k+1)})).$$

It is not difficult to show that such a framework encompasses

- the **Uzawa** algorithm (using $K(u) = J(u)$ and $\epsilon = 1$),
- the **Arrow-Hurwicz** algorithm (using $K(u) = \|u\|^2 / 2$).

Moreover, choosing an additive core K allows for decomposition in the minimization stage of the APP algorithm.

⁷Note that the term $\epsilon \langle p^{(k)}, \Theta(u) \rangle$ could be replaced by $\epsilon \langle p^{(k)}, \Theta'(u^{(k)}) \cdot u \rangle$.

APP with Explicit Constraints

A “natural” extension of the APP to constrained optimization problems consists in choosing a function $K : \mathcal{U} \rightarrow \mathbb{R}$ and then replacing the resolution of the initial problem by the resolution of the following **sequence of auxiliary problems**:⁷

$$u^{(k+1)} \in \arg \min_{u \in \mathcal{U}^{\text{ad}}} K(u) + \langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \rangle + \epsilon \langle p^{(k)}, \Theta(u) \rangle,$$

$$p^{(k+1)} = \text{proj}_{\mathcal{C}^*} (p^{(k)} + \rho \Theta(u^{(k+1)})).$$

It is not difficult to show that such a framework encompasses

- the **Uzawa** algorithm (using $K(u) = J(u)$ and $\epsilon = 1$),
- the **Arrow-Hurwicz** algorithm (using $K(u) = \|u\|^2 / 2$).

Moreover, choosing an **additive** core K allows for decomposition in the minimization stage of the APP algorithm.

⁷Note that the term $\epsilon \langle p^{(k)}, \Theta(u) \rangle$ could be replaced by $\epsilon \langle p^{(k)}, \Theta'(u^{(k)}) \cdot u \rangle$.

Lagrangian Stability

The standard Lagrangian approach has the following drawback. Assume that p^\sharp is a solution of the dual problem, and consider the set $\hat{U}(p^\sharp)$ of solutions associated to the primal minimization:

$$\hat{U}(p^\sharp) = \arg \min_{u \in U^{\text{ad}}} L(u, p^\sharp).$$

Then the set U^\sharp of solutions of the initial problem may be **strictly included** in $\hat{U}(p^\sharp)$, as illustrated by the following (linear) example:

$$\min_{u \in [-1,1]} -u \quad \text{s.t.} \quad u = 0,$$

whose unique saddle point is $\{0\} \times \{1\}$ whereas $\hat{U}(1) = [-1, 1]$.

A solution $\hat{u} \in \hat{U}(p^\sharp)$ induced by the dual problem is not always a solution of the initial problem (**stability of the Lagrangian**)!

Augmented Lagrangian

A remedy to this difficulty is to use a different duality theory, based on the idea of **regularization**. This idea leads to a new Lagrangian L_c which is called the **augmented Lagrangien**. The Lagrangian L_c is defined on the set $U^{\text{ad}} \times \mathbb{V}$,⁸ and its expression, which depends on a scalar parameter $c > 0$, is given by:

$$L_c(u, p) = J(u) + \frac{1}{2c} \left(\|\text{proj}_{C^*}(p + c\Theta(u))\|_{\mathbb{V}}^2 - \|p\|_{\mathbb{V}}^2 \right).$$

The augmented Lagrangian has the two following properties.

- ① The standard Lagrangian L and the augmented Lagrangian L_c have the same set of saddle points.
- ② The augmented Lagrangian L_c is always stable:

$$U^{\text{ad}} = \arg \min_{u \in U^{\text{ad}}} L_c(u, p^{\text{ad}}).$$

⁸whereas the standard Lagrangian L is defined on $U^{\text{ad}} \times C^*$

Augmented Lagrangian

A remedy to this difficulty is to use a different duality theory, based on the idea of **regularization**. This idea leads to a new Lagrangian L_c which is called the **augmented Lagrangien**. The Lagrangian L_c is defined on the set $U^{\text{ad}} \times \mathbb{V}$,⁸ and its expression, which depends on a scalar parameter $c > 0$, is given by:

$$L_c(u, p) = J(u) + \frac{1}{2c} \left(\|\text{proj}_{C^*}(p + c\Theta(u))\|_{\mathbb{V}}^2 - \|p\|_{\mathbb{V}}^2 \right).$$

The augmented Lagrangian has the two following properties.

- 1 The standard Lagrangian L and the augmented Lagrangian L_c have the **same set of saddle points**.
- 2 The augmented Lagrangian L_c is **always stable**:

$$U^\sharp = \arg \min_{u \in U^{\text{ad}}} L_c(u, p^\sharp).$$

⁸whereas the standard Lagrangian L is defined on $U^{\text{ad}} \times C^*$

APP and Augmented Lagrangian

The solution of the initial problem can be obtained by solving the “augmented” dual problem:

$$\max_{p \in \mathbb{V}} \left(\min_{u \in U^{\text{ad}}} L_c(u, p) \right).$$

The extension of the **APP** to that dual problem consists in solving the following sequence of auxiliary problems (see [Cohen, 2004]):

$$u^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \rangle + \epsilon \langle \text{proj}_{C^*}(p^{(k)} + c\Theta(u^{(k)})), \Theta(u) \rangle,$$

$$p^{(k+1)} = \left(1 - \frac{\rho}{c}\right) p^{(k)} + \frac{\rho}{c} \text{proj}_{C^*}(p^{(k)} + c\Theta(u^{(k+1)})).$$



- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - **Stochastic Gradient Method**
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Standard Stochastic Gradient Method

(1)

Consider the following **open-loop** stochastic optimization problem:

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

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

The standard stochastic gradient algorithm is written as follows.

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

Standard Stochastic Gradient Method

(1)

Consider the following **open-loop** stochastic optimization problem:

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

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

The **standard stochastic gradient** algorithm is written as follows.

- 1 Let $u^{(0)} \in U^{\text{ad}}$ and choose a positive real sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.
- 2 At iteration k , 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.

Standard Stochastic Gradient Method

(2)

This algorithm in fact involves **random variables** on $(\Omega, \mathcal{A}, \mathbb{P})$:

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

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

Recall that a sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is called a σ -sequence if

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

Convergence Theorem (Theorem)

Under various assumptions, the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ of random variables generated by the stochastic gradient algorithm almost surely converges to u^* .

Standard Stochastic Gradient Method

(2)

This algorithm in fact involves **random variables** on $(\Omega, \mathcal{A}, \mathbb{P})$:

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

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

Recall that a sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is called a **σ -sequence** if

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

Under various assumptions, the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ of random variables generated by the stochastic gradient algorithm almost surely converges to u^* .

Standard Stochastic Gradient Method

(2)

This algorithm in fact involves **random variables** on $(\Omega, \mathcal{A}, \mathbb{P})$:

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

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

Recall that a sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is called a **σ -sequence** if

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

Robbins-Monro Theorem

Under various assumptions, the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ of random variables generated by the stochastic gradient algorithm **almost surely** converges to $\mathbf{U}^\#$.

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - **Stochastic APP Algorithm**
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Moving APP to the Stochastic Framework

In order to **mix the ideas** of both the Auxiliary Problem Principle and the Stochastic Gradient Method, we first replace the initial problem by the associated sequence of **auxiliary problems**, namely

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle.$$

Then, in each auxiliary problem, we **replace** the gradient of J by the partial gradient of j evaluated at sampled realizations of W . The k -th instance of the stochastic auxiliary problem thus writes

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{w_j} j(u^{(k)}, w^{(k+1)}) - \nabla K(u^{(k)}), u \right\rangle.$$

It should be noted that the **large step** ϵ (constant) has to be replaced by **small steps** $\epsilon^{(k)}$ (going to zero as k goes to infinity).

Moving APP to the Stochastic Framework

In order to **mix the ideas** of both the Auxiliary Problem Principle and the Stochastic Gradient Method, we first replace the initial problem by the associated sequence of **auxiliary problems**, namely

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \right\rangle .$$

Then, in each auxiliary problem, we **replace** the gradient of J by the partial gradient of j evaluated at sampled realizations of W . The k -th instance of the stochastic auxiliary problem thus writes

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{uj}(u^{(k)}, w^{(k+1)}) - \nabla K(u^{(k)}), u \right\rangle .$$

It should be noted that the **large step** ϵ (constant) has to be replaced by **small steps** $\epsilon^{(k)}$ (going to zero as k goes to infinity).

Stochastic APP Algorithm

We thus obtain a **generalized stochastic gradient** algorithm.

Stochastic APP Algorithm

- 1 Let $u^{(0)} \in U^{\text{ad}}$ and choose a positive real sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.
- 2 At iteration k , draw a realization $w^{(k+1)}$ of the r.v. \mathbf{W} .
- 3 Update $u^{(k+1)}$ by solving the auxiliary problem:
$$u^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{w_j}(u^{(k)}, w^{(k+1)}) - \nabla K(u^{(k)}), u \right\rangle.$$
- 4 Set $k = k + 1$ and go to step 2.

As usual, the algorithm is casted in its probabilistic framework:

$$U^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{w_j}(U^{(k)}, W^{(k+1)}) - \nabla K(U^{(k)}), u \right\rangle.$$

The fact that the solution $U^{(k+1)}$ of this problem corresponds to a random variable, that is, a measurable function, has to be justified.

Stochastic APP Algorithm

We thus obtain a **generalized stochastic gradient** algorithm.

Stochastic APP Algorithm

- 1 Let $u^{(0)} \in U^{\text{ad}}$ and choose a positive real sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.
- 2 At iteration k , draw a realization $w^{(k+1)}$ of the r.v. \mathbf{W} .
- 3 Update $u^{(k+1)}$ by solving the auxiliary problem:
$$u^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{wj}(u^{(k)}, w^{(k+1)}) - \nabla K(u^{(k)}), u \right\rangle.$$
- 4 Set $k = k + 1$ and go to step 2.

As usual, the algorithm is casted in its **probabilistic framework**:

$$\mathbf{U}^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{wj}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) - \nabla K(\mathbf{U}^{(k)}), u \right\rangle.$$

*The fact that the solution $\mathbf{U}^{(k+1)}$ of this problem corresponds to a random variable, that is, a **measurable function**, has to be justified.*

Example

With the **choice**

$$K(u) = \frac{1}{2} \|u\|^2 ,$$

the auxiliary problem becomes

$$\min_{u \in U^{\text{ad}}} \frac{1}{2} \|u\|^2 + \left\langle \epsilon^{(k)} \nabla_{uj}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) - \mathbf{U}^{(k)} , u \right\rangle .$$

The set of solutions of this problem (an unique solution per ω) forms an unique random variable $\mathbf{U}^{(k+1)}$, whose expression is

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

It corresponds to the **standard stochastic gradient** iteration.

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - **Convergence Theorem and Proof**
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Convergence Theorem

(1)

Make the following assumptions.

H1 U^{ad} is a nonempty closed convex subset of a Hilbert space \mathbb{U} .

H2 $j : \mathbb{U} \times \mathbb{W} \rightarrow \mathbb{R}$ is a **normal integrand**, and $\mathbb{E}(j(u, \mathbf{W}))$ exists for all $u \in U^{\text{ad}}$.

H3 $j(\cdot, w) : \mathbb{U} \rightarrow \mathbb{R}$ is a proper convex differentiable function for all $w \in \mathbb{W}$ (thanks to **H2**, $j(\cdot, w)$ is a l.s.c. function).

H4 $j(\cdot, w)$ has **linearly bounded gradients (LBG)**:

$$\exists c_1, c_2 > 0, \forall (u, w) \in U^{\text{ad}} \times \mathbb{W}, \|\nabla_{\mathbb{U}} j(u, w)\| \leq c_1 \|u\| + c_2.$$

H5 J is coercive on U^{ad} .

H6 K is a proper l.s.c. function, **strongly convex** with modulus b and differentiable.

H7 $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a **σ -sequence**.

Convergence Theorem

(2)

Then the following conclusions hold true.

- R1** The initial problem has a non empty set of solutions $U^\# \subset \mathbb{U}$.
- R2** Each auxiliary problem has a unique solution $U^{(k+1)}$ which is a **random variable**.
- R3** The sequence of random variables $\{J(U^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to $J^\# = \min_{u \in U^{\text{ad}}} J(u)$.
- R4** The sequence of random variables $\{U^{(k)}\}_{k \in \mathbb{N}}$ is almost surely bounded, and every **cluster point** of a **realization** of this sequence almost surely belongs to the optimal set $U^\#$.

At last, if J is strongly convex, the sequence $\{U^{(k)}\}_{k \in \mathbb{N}}$ almost surely converges to the unique solution $u^\#$ of the initial problem.

Convergence Theorem

(2)

Then the following conclusions hold true.

- R1** The initial problem has a non empty set of solutions $U^\# \subset \mathbb{U}$.
- R2** Each auxiliary problem has a unique solution $U^{(k+1)}$ which is a **random variable**.
- R3** The sequence of random variables $\{J(U^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to $J^\# = \min_{u \in U^{\text{ad}}} J(u)$.
- R4** The sequence of random variables $\{U^{(k)}\}_{k \in \mathbb{N}}$ is almost surely bounded, and every **cluster point** of a **realization** of this sequence almost surely belongs to the optimal set $U^\#$.

At last, if J is **strongly convex**, the sequence $\{U^{(k)}\}_{k \in \mathbb{N}}$ almost surely **converges** to the **unique solution** $u^\#$ of the initial problem.

Sketch of Proof

The proof of the 1st statement is based on **optimization** theorems.

The proof of the 2nd statement involves **measurability** arguments.

The proof of the last two statements consists of **three steps**.

❶ **Select a Lyapunov function.**

Here we choose $\Lambda(u) = K(u^\sharp) - K(u) - \langle \nabla K(u), u^\sharp - u \rangle$.

❷ **Bound from above the variation of Λ .**

Using assumptions and writing optimality conditions, we get:

$$\mathbb{E}(\Lambda(\mathbf{U}^{(k+1)}) \mid \mathcal{F}^{(k)}) \leq (1 + \alpha^{(k)})\Lambda(\mathbf{U}^{(k)}) + \beta^{(k)} - \epsilon^{(k)}(J(\mathbf{U}^{(k)}) - J(u^\sharp)).$$

❸ **Prove the convergence of the sequences.**

Using two **technical lemmas**, we obtain that $\{\Lambda(\mathbf{U}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to a finite random variable Λ^∞ , and that $\{J(\mathbf{U}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to $J(u^\sharp)$. Using a compactness argument, it exists subsequences of $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ converging almost surely to elements belonging to the set U^\sharp .

Two Useful Lemmas

(1)

Robbins-Siegmund Theorem

Let $\{\Lambda^{(k)}\}_{k \in \mathbb{N}}$, $\{\alpha^{(k)}\}_{k \in \mathbb{N}}$, $\{\beta^{(k)}\}_{k \in \mathbb{N}}$ and $\{\eta^{(k)}\}_{k \in \mathbb{N}}$ be four **positive** sequences of real-valued random variables adapted to the filtration $\{\mathcal{F}^{(k)}\}_{k \in \mathbb{N}}$. Assume that

$$\mathbb{E}(\Lambda^{(k+1)} \mid \mathcal{F}^{(k)}) \leq (1 + \alpha^{(k)})\Lambda^{(k)} + \beta^{(k)} - \eta^{(k)}, \quad \forall k \in \mathbb{N},$$

and that

$$\sum_{k \in \mathbb{N}} \alpha^{(k)} < +\infty \quad \text{and} \quad \sum_{k \in \mathbb{N}} \beta^{(k)} < +\infty, \quad \mathbb{P}\text{-a.s. .}$$

Then, the sequence $\{\Lambda^{(k)}\}_{k \in \mathbb{N}}$ almost surely converges to a finite^a random variable Λ^∞ , and we have that $\sum_{k \in \mathbb{N}} \eta^{(k)} < +\infty, \mathbb{P}\text{-a.s.}$

^aA random variable \mathbf{X} is finite if $\mathbb{P}(\{\omega \in \Omega \mid \mathbf{X}(\omega) = +\infty\}) = 0$.

Two Useful Lemmas

(2)

Technical Lemma

Let J be a real-valued function defined on a Hilbert space \mathbb{U} , such that J is **Lipschitz** continuous with constant L .

Let $\{u^{(k)}\}_{k \in \mathbb{N}}$ be a sequence of elements of \mathbb{U} and let $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ be a sequence of positive real numbers such that

- (a) $\sum_{k \in \mathbb{N}} \epsilon^{(k)} = +\infty,$
- (b) $\exists \mu \in \mathbb{R}, \sum_{k \in \mathbb{N}} \epsilon^{(k)} |J(u^{(k)}) - \mu| < +\infty,$
- (c) $\exists \delta > 0, \forall k \in \mathbb{N}, \|u^{(k+1)} - u^{(k)}\| \leq \delta \epsilon^{(k)}.$

Then the sequence $\{J(u^{(k)})\}_{k \in \mathbb{N}}$ converges to μ .

Proof of the Convergence Theorem

(1)

The proof of the first statement is based on standard theorems in the field of convex optimization ensuring the existence of solutions in a general Hilbert space. Let $u^\# \in U^\#$ be a solution of the initial problem.

The existence of a r.v. $U^{(k+1)}$ solution of the auxiliary problem

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle c^{(k)} \nabla_{\omega} K(U^{(k)}, W^{(k+1)}) - \nabla K(U^{(k)}), u \right\rangle,$$

is a consequence of the fact that the criterion to be minimized is a normal integrand. The arg min is a closed-valued and measurable multifunction and thus at least admits a measurable selection (see [Rockafellar & Wets, 1998, Theorem 14.37] for further details).

The solution $U^{(k+1)}$ is unique because K is strongly convex.

Proof of the Convergence Theorem

(1)

The proof of the first statement is based on standard theorems in the field of convex optimization ensuring the existence of solutions in a general Hilbert space. Let $u^\# \in U^\#$ be a solution of the initial problem.

The existence of a r.v. $\mathbf{U}^{(k+1)}$ solution of the auxiliary problem

$$\min_{u \in U^{\text{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_{uj}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}) - \nabla K(\mathbf{U}^{(k)}), u \right\rangle ,$$

is a consequence of the fact that the criterion to be minimized is a **normal integrand**. The **arg min** is a closed-valued and measurable multifunction and thus at least admits a **measurable selection** (see [Rockafellar & Wets, 1998, Theorem 14.37] for further details).

The solution $\mathbf{U}^{(k+1)}$ is **unique** because K is **strongly convex**.

Proof of Convergence

(2)

- Let $\Lambda(u) = K(u^\sharp) - K(u) - \langle \nabla K(u), u^\sharp - u \rangle$. We have

$$\Lambda(u) \geq \frac{b}{2} \|u - u^\sharp\|_{\mathbb{U}}^2, \quad (5)$$

(strong convexity of K) so that Λ is **bounded from below**.

- Consider the variation of Λ from one iteration to another:

$$\begin{aligned} \Delta^{(k)} &= \Lambda(U^{(k+1)}) - \Lambda(U^{(k)}) \\ &= \underbrace{K(U^{(k)}) - K(U^{(k+1)}) - \langle \nabla K(U^{(k)}), U^{(k)} - U^{(k+1)} \rangle}_{T_1} \\ &\quad + \underbrace{\langle \nabla K(U^{(k)}) - \nabla K(U^{(k+1)}), u^\sharp - U^{(k+1)} \rangle}_{T_2} \end{aligned}$$

From the convexity of K , we have that $T_1 \leq 0$.

Proof of Convergence

(2)

- Let $\Lambda(u) = K(u^\sharp) - K(u) - \langle \nabla K(u), u^\sharp - u \rangle$. We have

$$\Lambda(u) \geq \frac{b}{2} \|u - u^\sharp\|_{\mathbb{U}}^2, \quad (5)$$

(strong convexity of K) so that Λ is **bounded from below**.

- Consider the **variation** of Λ from one iteration to another:

$$\begin{aligned} \Delta^{(k)} &= \Lambda(\mathbf{U}^{(k+1)}) - \Lambda(\mathbf{U}^{(k)}) \\ &= \underbrace{K(\mathbf{U}^{(k)}) - K(\mathbf{U}^{(k+1)}) - \langle \nabla K(\mathbf{U}^{(k)}), \mathbf{U}^{(k)} - \mathbf{U}^{(k+1)} \rangle}_{T_1} \\ &\quad + \underbrace{\langle \nabla K(\mathbf{U}^{(k)}) - \nabla K(\mathbf{U}^{(k+1)}), u^\sharp - \mathbf{U}^{(k+1)} \rangle}_{T_2}. \end{aligned}$$

From the **convexity** of K , we have that $T_1 \leq 0$.

Proof of Convergence

(3)

Let $\mathbf{G}^{(k)} = \nabla_{\mathbf{u}j}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$. From the **optimality condition** of the auxiliary problem evaluated at \mathbf{u}^\sharp , we have that

$$T_2 \leq \underbrace{\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{u}^\sharp - \mathbf{U}^{(k)} \rangle}_{T_3} + \underbrace{\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{U}^{(k)} - \mathbf{U}^{(k+1)} \rangle}_{T_4} .$$

Proof of Convergence

(3)

Let $\mathbf{G}^{(k)} = \nabla_{\mathbf{u}j}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$. From the **optimality condition** of the auxiliary problem evaluated at \mathbf{u}^\sharp , we have that

$$T_2 \leq \underbrace{\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{u}^\sharp - \mathbf{U}^{(k)} \rangle}_{T_3} + \underbrace{\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{U}^{(k)} - \mathbf{U}^{(k+1)} \rangle}_{T_4}.$$

- From the **convexity** of $j(\cdot, \mathbf{w})$, we have that

$$T_3 \leq j(\mathbf{u}^\sharp, \mathbf{W}^{(k+1)}) - j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}).$$

- The optimality condition at $\mathbf{U}^{(k)}$ and the strong convexity of K lead to: $\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{U}^{(k)} - \mathbf{U}^{(k+1)} \rangle \geq b \|\mathbf{U}^{(k+1)} - \mathbf{U}^{(k)}\|_V^2$.

Using the Schwarz inequality, we obtain: $T_4 \leq \frac{\epsilon^{(k)}}{b} \|\mathbf{G}^{(k)}\|_V^2$.

The LBG assumption and the majoration (5) of Λ yield:

$$T_4 \leq \epsilon^{(k)} (\alpha \Lambda(\mathbf{U}^{(k)}) + \beta).$$

Proof of Convergence

(3)

Let $\mathbf{G}^{(k)} = \nabla_{\mathbf{u}j}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$. From the **optimality condition** of the auxiliary problem evaluated at \mathbf{u}^\sharp , we have that

$$T_2 \leq \underbrace{\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{u}^\sharp - \mathbf{U}^{(k)} \rangle}_{T_3} + \underbrace{\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{U}^{(k)} - \mathbf{U}^{(k+1)} \rangle}_{T_4}.$$

- From the **convexity** of $j(\cdot, \mathbf{w})$, we have that

$$T_3 \leq j(\mathbf{u}^\sharp, \mathbf{W}^{(k+1)}) - j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)}).$$

- The **optimality condition** at $\mathbf{U}^{(k)}$ and the **strong convexity** of K lead to: $\epsilon^{(k)} \langle \mathbf{G}^{(k)}, \mathbf{U}^{(k)} - \mathbf{U}^{(k+1)} \rangle \geq b \|\mathbf{U}^{(k+1)} - \mathbf{U}^{(k)}\|_{\mathbf{U}}^2$.

Using the **Schwarz inequality**, we obtain: $T_4 \leq \frac{\epsilon^{(k)}}{b} \|\mathbf{G}^{(k)}\|_{\mathbf{U}}^2$.

The **LBG** assumption and the majoration (5) of Λ yield:

$$T_4 \leq \epsilon^{(k)} \left(\alpha \Lambda(\mathbf{U}^{(k)}) + \beta \right).$$

Proof of Convergence

(4)

Collecting the upper bounds obtained for T_1 , T_3 and T_4 leads to

$$\Delta^{(k)} \leq \epsilon^{(k)} (j(u^\#, \mathbf{W}^{(k+1)}) - j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})) + (\epsilon^{(k)})^2 (\alpha \Lambda(\mathbf{U}^{(k)}) + \beta).$$

Taking the conditional expectation w.r.t. the σ -field $\mathcal{F}^{(k)}$ generated by $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$, we obtain that^a

$$\mathbb{E}(\Lambda(\mathbf{U}^{(k+1)}) \mid \mathcal{F}^{(k)}) \leq (1 + \alpha^{(k)}) \Lambda(\mathbf{U}^{(k)}) + \beta^{(k)} + \epsilon^{(k)} (j(u^\#) - j(\mathbf{U}^{(k)})),$$

with $\alpha^{(k)} = \alpha(\epsilon^{(k)})^2$ and $\beta^{(k)} = \beta(\epsilon^{(k)})^2$.

^aRecall that $\mathbf{W}^{(k+1)}$ is independent of $\mathcal{F}^{(k)}$ and that $\mathbf{U}^{(k)}$ is $\mathcal{F}^{(k)}$ -measurable.

Reminder. We have also obtained the two following inequalities:

$$\Lambda(\mathbf{U}^{(k)}) \geq \frac{b}{2} \|\mathbf{U}^{(k)} - u^\#\|_0^2 \quad \text{and} \quad \|\mathbf{U}^{(k)} - \mathbf{U}^{(k+1)}\|_0 \leq \frac{\epsilon^{(k)}}{b} \|G^{(k)}\|_0.$$

Proof of Convergence

(4)

Collecting the upper bounds obtained for T_1 , T_3 and T_4 leads to $\Delta^{(k)} \leq \epsilon^{(k)}(j(u^\#, \mathbf{W}^{(k+1)}) - j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})) + (\epsilon^{(k)})^2(\alpha\Lambda(\mathbf{U}^{(k)}) + \beta)$.

Taking the **conditional expectation** w.r.t. the σ -field $\mathcal{F}^{(k)}$ generated by $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$, we obtain that^a

$$\mathbb{E}(\Lambda(\mathbf{U}^{(k+1)}) \mid \mathcal{F}^{(k)}) \leq (1 + \alpha^{(k)})\Lambda(\mathbf{U}^{(k)}) + \beta^{(k)} + \epsilon^{(k)}(J(u^\#) - J(\mathbf{U}^{(k)})),$$

with $\alpha^{(k)} = \alpha(\epsilon^{(k)})^2$ and $\beta^{(k)} = \beta(\epsilon^{(k)})^2$.

^aRecall that $\mathbf{W}^{(k+1)}$ is independent of $\mathcal{F}^{(k)}$ and that $\mathbf{U}^{(k)}$ is $\mathcal{F}^{(k)}$ -measurable.

Reminder. We have also obtained the two following inequalities:

$$\Lambda(\mathbf{U}^{(k)}) \geq \frac{b}{2} \|\mathbf{U}^{(k)} - u^\#\|_0^2 \quad \text{and} \quad \|\mathbf{U}^{(k)} - \mathbf{U}^{(k+1)}\|_0 \leq \frac{\epsilon^{(k)}}{b} \|\mathbf{G}^{(k)}\|_0.$$

Proof of Convergence

(4)

Collecting the upper bounds obtained for T_1 , T_3 and T_4 leads to $\Delta^{(k)} \leq \epsilon^{(k)}(j(u^\#, \mathbf{W}^{(k+1)}) - j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})) + (\epsilon^{(k)})^2(\alpha\Lambda(\mathbf{U}^{(k)}) + \beta)$.

Taking the **conditional expectation** w.r.t. the σ -field $\mathcal{F}^{(k)}$ generated by $(\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)})$, we obtain that^a

$$\mathbb{E}(\Lambda(\mathbf{U}^{(k+1)}) \mid \mathcal{F}^{(k)}) \leq (1 + \alpha^{(k)})\Lambda(\mathbf{U}^{(k)}) + \beta^{(k)} + \epsilon^{(k)}(J(u^\#) - J(\mathbf{U}^{(k)})),$$

with $\alpha^{(k)} = \alpha(\epsilon^{(k)})^2$ and $\beta^{(k)} = \beta(\epsilon^{(k)})^2$.

^aRecall that $\mathbf{W}^{(k+1)}$ is independent of $\mathcal{F}^{(k)}$ and that $\mathbf{U}^{(k)}$ is $\mathcal{F}^{(k)}$ -measurable.

Reminder. We have also obtained the two following inequalities:

$$\Lambda(\mathbf{U}^{(k)}) \geq \frac{b}{2} \|\mathbf{U}^{(k)} - u^\#\|_{\mathbb{U}}^2 \quad \text{and} \quad \|\mathbf{U}^{(k)} - \mathbf{U}^{(k+1)}\|_{\mathbb{U}} \leq \frac{\epsilon^{(k)}}{b} \|\mathbf{G}^{(k)}\|_{\mathbb{U}}.$$

Proof of Convergence

(5)

- From the **Robbins-Siegmund** theorem, $\{\Lambda(\mathbf{U}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to a finite random variable Λ^∞ and we have

$$\sum_{k=0}^{+\infty} \epsilon^{(k)} (J(\mathbf{U}^{(k)}) - J(u^\#)) < +\infty, \quad \mathbb{P}\text{-a.s. .}$$

Let Ω_0 denote the subset of Ω such that the two almost sure properties in the point above are fulfilled ($\mathbb{P}(\Omega_0) = 1$).

- We deduce that both sequences $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ and $\{G^{(k)}\}_{k \in \mathbb{N}}$ are a.s. bounded, so that the same holds true for the sequence $\left\{\frac{1}{\sqrt{k}} \|\mathbf{U}^{(k+1)} - \mathbf{U}^{(k)}\|_U\right\}_{k \in \mathbb{N}}$. This makes it possible to use the second technical lemma⁹ and claim that $\{J(\mathbf{U}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to $J(u^\#)$.

⁹Assumptions also imply that J is Lipschitz on each bounded subset of U^{nd} .

Proof of Convergence

(5)

- From the **Robbins-Siegmund** theorem, $\{\Lambda(\mathbf{U}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to a finite random variable Λ^∞ and we have

$$\sum_{k=0}^{+\infty} \epsilon^{(k)} (J(\mathbf{U}^{(k)}) - J(u^\#)) < +\infty, \quad \mathbb{P}\text{-a.s. .}$$

Let Ω_0 denote the subset of Ω such that the two almost sure properties in the point above are fulfilled ($\mathbb{P}(\Omega_0) = 1$).

- We deduce that both sequences $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ and $\{\mathbf{G}^{(k)}\}_{k \in \mathbb{N}}$ are **a.s. bounded**, so that the same holds true for the sequence $\left\{ \frac{1}{\epsilon^{(k)}} \|\mathbf{U}^{(k+1)} - \mathbf{U}^{(k)}\|_{\mathbb{U}} \right\}_{k \in \mathbb{N}}$. This makes it possible to use the second **technical lemma**⁹ and claim that $\{J(\mathbf{U}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to $J(u^\#)$.

⁹Assumptions also imply that J is Lipschitz on each bounded subset of U^{ad} .

Proof of Convergence

(6)

- Pick some $\omega \in \Omega_0$. The sequence of realizations $\{u^{(k)}\}_{k \in \mathbb{N}}$ of $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ associated with ω is bounded, and $u^{(k)} \in U^{\text{ad}}$. By a **compactness argument**,¹⁰ it exists a convergent subsequence $\{u^{(\Phi(k))}\}_{k \in \mathbb{N}}$, with limit \bar{u} . Using the lower semi-continuity of function J , we have that

$$J(\bar{u}) \leq \liminf_{k \rightarrow +\infty} J(u^{(\Phi(k))}) = J(u^\#).$$

Since $\bar{u} \in U^{\text{ad}}$, we deduce that $\bar{u} \in U^\#$.

¹⁰A subset $U^{\text{ad}} \subset \mathbb{U}$ is **compact** if it is closed and bounded, provided that \mathbb{U} is a **finite-dimensional** Hilbert space. If U^{ad} is an **infinite-dimensional** Hilbert space, such a property remains true only in the **weak topology**. If U^{ad} is closed in the strong topology and is **convex**, then it is also closed in the weak topology, and hence compact if bounded. In the same vein, the l.s.c. property of J is preserved in the weak topology if J is **convex** (see [Ekeland & Temam, 1999]).

Proof of Convergence

(7)

We ultimately consider the case when J is **strongly convex** with modulus a . Then the initial problem has a **unique** solution u^\sharp (standard optimization argument).

Thanks to the strong convexity property of J , we have

$$\begin{aligned} J(\mathbf{U}^{(k)}) - J(u^\sharp) &\geq \langle \nabla J(u^\sharp), \mathbf{U}^{(k)} - u^\sharp \rangle + \frac{a}{2} \|\mathbf{U}^{(k)} - u^\sharp\|_{\mathbb{U}}^2 \\ &\geq \frac{a}{2} \|\mathbf{U}^{(k)} - u^\sharp\|_{\mathbb{U}}^2. \end{aligned}$$

Since $J(\mathbf{U}^{(k)})$ **almost surely converges** to $J(u^\sharp)$, we deduce that $\|\mathbf{U}^{(k)} - u^\sharp\|_{\mathbb{U}}$ **almost surely converges** to zero.

The proof is complete.

Conclusions

The **stochastic APP algorithm** encompasses the stochastic gradient algorithm (obtained using $K(u) = \|u\|^2/2$), as well as the so-called matrix-gain algorithm (K being in this case $K(u) = \langle u, Au \rangle / 2$ and A being a positive definite matrix).

From a theoretical point of view, the convergence theorem has been proven under rather **natural assumptions**. As a matter of fact, the convexity and differentiability assumptions are standard in the framework of convex optimization. Note that, even if an explicit convexity property is not required in the **Robbins-Monro theorem**, another assumption playing a very similar role is used.

As far as **decomposition** is concerned, the stochastic APP algorithm opens this possibility as a way to solve large stochastic optimization problems. Of course, the convergence remains “slow” because it is driven by a σ -sequence, namely $\{c^{(k)}\}_{k \in \mathbb{N}}$.

Conclusions

The **stochastic APP algorithm** encompasses the stochastic gradient algorithm (obtained using $K(u) = \|u\|^2/2$), as well as the so-called matrix-gain algorithm (K being in this case $K(u) = \langle u, Au \rangle / 2$ and A being a positive definite matrix).

From a theoretical point of view, the convergence theorem has been proven under rather **natural assumptions**. As a matter of fact, the convexity and differentiability assumptions are standard in the framework of convex optimization. Note that, even if an explicit convexity property is not required in the **Robbins-Monro theorem**, another assumption playing a very similar role is used.

As far as **decomposition** is concerned, the stochastic APP algorithm opens this possibility as a way to solve large stochastic optimization problems. Of course, the convergence remains "slow" because it is driven by a σ -sequence, namely $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.

Conclusions

The **stochastic APP algorithm** encompasses the stochastic gradient algorithm (obtained using $K(u) = \|u\|^2/2$), as well as the so-called matrix-gain algorithm (K being in this case $K(u) = \langle u, Au \rangle / 2$ and A being a positive definite matrix).

From a theoretical point of view, the convergence theorem has been proven under rather **natural assumptions**. As a matter of fact, the convexity and differentiability assumptions are standard in the framework of convex optimization. Note that, even if an explicit convexity property is not required in the **Robbins-Monro theorem**, another assumption playing a very similar role is used.

As far as **decomposition** is concerned, the stochastic APP algorithm opens this possibility as a way to solve large stochastic optimization problems. Of course, the convergence remains “slow” because it is driven by a σ -sequence, namely $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$.

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Nature of Constraints in Stochastic Optimization

Constraints in stochastic optimization arise in different ways, have different meanings and require various mathematical treatments.

Nature of Constraints in Stochastic Optimization

Constraints in stochastic optimization arise in different ways, have different meanings and require various mathematical treatments.

- A constraint may be **deterministic**: $\Theta(u) \in -C$.
- A constraint may be formulated in the **almost sure sense**, that is $\theta(u, W) \in -C$ P-a.s.. This formulation is generally used to express “hard constraints” (physical laws, ...).
- Another (more realistic) way is to formulate **stochastic constraints in probability**: $\mathbb{P}(\theta(u, W) \in -C) \geq \pi$, which means that the constraints can “sometimes” be violated.
- Another possibility is to have a constraint in **expectation**: $\mathbb{E}(\theta(u, W)) \in -C$. Although usually **non intuitive**, such a formulation proves useful in some specific problems, as it will be briefly illustrated at the end of the lecture.

Nature of Constraints in Stochastic Optimization

Constraints in stochastic optimization arise in different ways, have different meanings and require various mathematical treatments.

- A constraint may be **deterministic**: $\Theta(u) \in -C$.
- A constraint may be formulated in the **almost sure sense**, that is $\theta(u, \mathbf{W}) \in -C$ \mathbb{P} -a.s.. This formulation is generally used to express “hard constraints” (physical laws, ...).
- Another (more realistic) way is to formulate stochastic constraints in probability: $\mathbb{P}(\theta(u, \mathbf{W}) \in -C) \geq \pi$, which means that the constraints can “sometimes” be violated.
- Another possibility is to have a constraint in expectation: $\mathbb{E}(\theta(u, \mathbf{W})) \in -C$. Although usually non intuitive, such a formulation proves useful in some specific problems, as it will be briefly illustrated at the end of the lecture.

Nature of Constraints in Stochastic Optimization

Constraints in stochastic optimization arise in different ways, have different meanings and require various mathematical treatments.

- A constraint may be **deterministic**: $\Theta(u) \in -C$.
- A constraint may be formulated in the **almost sure sense**, that is $\theta(u, \mathbf{W}) \in -C$ \mathbb{P} -a.s.. This formulation is generally used to express “hard constraints” (physical laws, ...).
- Another (more realistic) way is to formulate stochastic constraints **in probability**: $\mathbb{P}(\theta(u, \mathbf{W}) \in -C) \geq \pi$, which means that the constraints can “sometimes” be violated.
- Another possibility is to have a constraint in **expectation**: $\mathbb{E}(\theta(u, \mathbf{W})) \in -C$. Although usually non intuitive, such a formulation proves useful in some specific problems, as it will be briefly illustrated at the end of the lecture.

Nature of Constraints in Stochastic Optimization

Constraints in stochastic optimization arise in different ways, have different meanings and require various mathematical treatments.

- A constraint may be **deterministic**: $\Theta(u) \in -C$.
- A constraint may be formulated in the **almost sure sense**, that is $\theta(u, \mathbf{W}) \in -C$ \mathbb{P} -a.s.. This formulation is generally used to express “hard constraints” (physical laws, ...).
- Another (more realistic) way is to formulate stochastic constraints **in probability**: $\mathbb{P}(\theta(u, \mathbf{W}) \in -C) \geq \pi$, which means that the constraints can “sometimes” be violated.
- Another possibility is to have a constraint **in expectation**: $\mathbb{E}(\theta(u, \mathbf{W})) \in -C$. Although usually **non intuitive**, such a formulation proves useful in some specific problems, as it will be briefly illustrated at the end of the lecture.

Constrained Stochastic Optimization Problem under Study

We consider the following stochastic optimization setting.

- The probability space is denoted $(\Omega, \mathcal{A}, \mathbb{P})$, and W is a random variable valued on the space (\bar{W}, \mathcal{W}) .
- We are interested in the case where the criterion J is defined as $J(u) = \mathbb{E}(j(u, W))$, with $j : U \times W \rightarrow \bar{\mathbb{R}}$. This is the standard framework when studying open-loop stochastic optimization problems.
- We will hereafter consider only constraints Θ which are of a deterministic nature: $\Theta : U \rightarrow \mathcal{V}$.

The problem we deal with has the following expression:

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, W)) \quad \text{subject to} \quad \Theta(u) \in -C.$$

Constrained Stochastic Optimization Problem under Study

We consider the following stochastic optimization setting.

- The probability space is denoted $(\Omega, \mathcal{A}, \mathbb{P})$, and \mathbf{W} is a random variable valued on the space $(\mathbb{W}, \mathcal{W})$.
- We are interested in the case where the criterion J is defined as $J(u) = \mathbb{E}(j(u, \mathbf{W}))$, with $j : \mathbb{U} \times \mathbb{W} \rightarrow \overline{\mathbb{R}}$. This is the standard framework when studying **open-loop** stochastic optimization problems.
- We will hereafter consider only constraints Θ which are of a **deterministic** nature: $\Theta : \mathbb{U} \rightarrow \mathbb{V}$.

The problem we deal with has the following expression:

$$\min_{u \in \mathbb{U}^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \Theta(u) \in -C$$

Constrained Stochastic Optimization Problem under Study

We consider the following stochastic optimization setting.

- The probability space is denoted $(\Omega, \mathcal{A}, \mathbb{P})$, and \mathbf{W} is a random variable valued on the space $(\mathbb{W}, \mathcal{W})$.
- We are interested in the case where the criterion J is defined as $J(u) = \mathbb{E}(j(u, \mathbf{W}))$, with $j : \mathbb{U} \times \mathbb{W} \rightarrow \overline{\mathbb{R}}$. This is the standard framework when studying **open-loop** stochastic optimization problems.
- We will hereafter consider only constraints Θ which are of a **deterministic** nature: $\Theta : \mathbb{U} \rightarrow \mathbb{V}$.

The problem we deal with has the following expression:

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \Theta(u) \in -C.$$

On the Agenda

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \Theta(u) \in -C.$$

We want to solve this open-loop problem using **duality** methods. . .

- Extension of the Uzawa Algorithm.
- Stochastic APP Algorithm with Constraints.
- Stochastic APP and Augmented Lagrangian.
- What happens if $\Theta(u) = \mathbb{E}(\theta(u, \mathbf{W}))$?

On the Agenda

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \Theta(u) \in -C .$$

We want to solve this open-loop problem using **duality** methods. . .

- **Extension of the Uzawa Algorithm.**
- **Stochastic APP Algorithm with Constraints.**
- **Stochastic APP and Augmented Lagrangian.**
- **What happens if $\Theta(u) = \mathbb{E}(\theta(u, \mathbf{W}))$?**

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

A Useful Tool in Stochastic Approximation

In the context of **Stochastic Approximation**, strong connections exist between the convergence of the standard SA algorithm:

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

and the behavior of the **ordinary differential equation** (ODE) associated to this algorithm:

$$\dot{\mathbf{u}} = h(\mathbf{u}),$$

(see [Kushner & Clark, 1978]). A useful property is given below.

Let $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ be the sequence generated by the **SA** algorithm. We assume that

$$\exists \mathbf{u}^\# \in \mathbb{U}, \text{ such that } \mathbb{P} \left(\lim_{k \rightarrow +\infty} \mathbf{U}^{(k)} = \mathbf{u}^\# \right) > 0.$$

Then $\mathbf{u}^\#$ is a **stable equilibrium point** of the associated **ODE**.

Extension of the Uzawa Algorithm to the Stochastic Case

Our first attempt for solving the stochastic constrained problem:

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \Theta(u) \in -C,$$

is to propose an extension of the **Uzawa** algorithm. More precisely, during the minimization stage w.r.t. u , we propose to replace the expectation $J(u)$ by the value $j(u, w^{(k+1)})$.¹¹ We thus obtain a (naive) first tentative Stochastic Uzawa Algorithm:

$$U^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} j(u, W^{(k+1)}) + \langle P^{(k)}, \Theta(u) \rangle,$$

$$P^{(k+1)} = \text{proj}_C(P^{(k)} + \rho^{(k)} \Theta(U^{(k+1)})).$$

Question: what about the convergence of this algorithm?

¹¹Note that we replace here the evaluation of J by the one of j , whereas ∇J was replaced by $\nabla_u j$ in the standard stochastic gradient method...

Extension of the Uzawa Algorithm to the Stochastic Case

Our first attempt for solving the stochastic constrained problem:

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \Theta(u) \in -C,$$

is to propose an extension of the **Uzawa** algorithm. More precisely, during the minimization stage w.r.t. u , we propose to replace the expectation $J(u)$ by the value $j(u, w^{(k+1)})$.¹¹ We thus obtain a (naive) first **tentative Stochastic Uzawa Algorithm**:

$$\begin{aligned} \mathbf{U}^{(k+1)} &\in \arg \min_{u \in U^{\text{ad}}} j(u, \mathbf{W}^{(k+1)}) + \langle \mathbf{P}^{(k)}, \Theta(u) \rangle, \\ \mathbf{P}^{(k+1)} &= \text{proj}_{C^*}(\mathbf{P}^{(k)} + \rho^{(k)} \Theta(\mathbf{U}^{(k+1)})). \end{aligned}$$

Question: what about the **convergence** of this algorithm?

¹¹Note that we replace here the evaluation of J by the one of j , whereas ∇J was replaced by ∇_{uj} in the standard stochastic gradient method...

Stochastic Uzawa Algorithm Counter-Example

(1)

Consider a constrained stochastic optimization problem with:

- $\mathbf{U} = \mathbb{R}^2$ and $\mathbf{U}^{\text{ad}} = \mathbf{U}$,
- $\mathbf{V} = \mathbb{R}$ and $\mathbf{C} = \{0\}$ (equality constraint),
- $\mathbf{W} = \mathbb{R}^4$ and $\mathbf{W} = (\mathbf{A}_1, \mathbf{A}_2, \mathbf{B}_1, \mathbf{B}_2)$,
- $j(u, w) = \frac{1}{2}(a_1 u_1^2 + a_2 u_2^2) + (b_1 u_1 + b_2 u_2)$,
- $\Theta(u) = \theta_1 u_1 + \theta_2 u_2$.

The optimality conditions (KKT) of this problem are:

$$\mathbb{E}(\mathbf{A}_1)u_1^2 + \mathbb{E}(\mathbf{B}_1) + \theta_1 p^2 = 0, \quad \mathbb{E}(\mathbf{A}_2)u_2^2 + \mathbb{E}(\mathbf{B}_2) + \theta_2 p^2 = 0, \quad \theta_1 u_1^2 + \theta_2 u_2^2 = 0,$$

so that the value of the optimal multiplier is:

$$p^2 = \frac{\frac{\mathbb{E}(\mathbf{B}_1)}{\mathbb{E}(\mathbf{A}_1)} \theta_1 + \frac{\mathbb{E}(\mathbf{B}_2)}{\mathbb{E}(\mathbf{A}_2)} \theta_2}{\frac{\theta_1^2}{\mathbb{E}(\mathbf{A}_1)} + \frac{\theta_2^2}{\mathbb{E}(\mathbf{A}_2)}}.$$

Stochastic Uzawa Algorithm Counter-Example

(1)

Consider a constrained stochastic optimization problem with:

- $\mathbf{U} = \mathbb{R}^2$ and $\mathbf{U}^{\text{ad}} = \mathbf{U}$,
- $\mathbf{V} = \mathbb{R}$ and $\mathbf{C} = \{0\}$ (equality constraint),
- $\mathbf{W} = \mathbb{R}^4$ and $\mathbf{W} = (\mathbf{A}_1, \mathbf{A}_2, \mathbf{B}_1, \mathbf{B}_2)$,
- $j(u, w) = \frac{1}{2}(a_1 u_1^2 + a_2 u_2^2) + (b_1 u_1 + b_2 u_2)$,
- $\Theta(u) = \theta_1 u_1 + \theta_2 u_2$.

The optimality conditions (KKT) of this problem are:

$$\mathbb{E}(\mathbf{A}_1) u_1^\sharp + \mathbb{E}(\mathbf{B}_1) + \theta_1 p^\sharp = 0, \quad \mathbb{E}(\mathbf{A}_2) u_2^\sharp + \mathbb{E}(\mathbf{B}_2) + \theta_2 p^\sharp = 0, \quad \theta_1 u_1^\sharp + \theta_2 u_2^\sharp = 0,$$

so that the value of the optimal multiplier is:

$$p^\sharp = - \frac{\frac{\mathbb{E}(\mathbf{B}_1)}{\mathbb{E}(\mathbf{A}_1)} \theta_1 + \frac{\mathbb{E}(\mathbf{B}_2)}{\mathbb{E}(\mathbf{A}_2)} \theta_2}{\frac{\theta_1^2}{\mathbb{E}(\mathbf{A}_1)} + \frac{\theta_2^2}{\mathbb{E}(\mathbf{A}_2)}}.$$

Stochastic Uzawa Algorithm Counter-Example

(2)

Apply the **tentative Uzawa algorithm**. The minimization in u gives:

- $A_1^{(k+1)} U_1^{(k+1)} + B_1^{(k+1)} + \theta_1 P^{(k)} = 0,$
- $A_2^{(k+1)} U_2^{(k+1)} + B_2^{(k+1)} + \theta_2 P^{(k)} = 0,$

and the update of the multiplier writes:

$$P^{(k+1)} = P^{(k)} + \rho^{(k)} (\theta_1 U_1^{(k+1)} + \theta_2 U_2^{(k+1)}).$$

We thus obtain

$$P^{(k+1)} = P^{(k)} - \rho^{(k)} \left(\left(\frac{\theta_1^2}{A_1^{(k+1)}} + \frac{\theta_2^2}{A_2^{(k+1)}} \right) P^{(k)} + \left(\theta_1 \frac{B_1^{(k+1)}}{A_1^{(k+1)}} + \theta_2 \frac{B_2^{(k+1)}}{A_2^{(k+1)}} \right) \right).$$

$\{P^{(k)}\}_{k \in \mathbb{N}}$ can only converge to a stable equilibrium point of the associated differential equation (ODE argument), that is,

$$\bar{p} = - \frac{\mathbb{E} \left(\frac{B_1}{A_1} \right) \theta_1 + \mathbb{E} \left(\frac{B_2}{A_2} \right) \theta_2}{\mathbb{E} \left(\frac{1}{A_1} \right) \theta_1^2 + \mathbb{E} \left(\frac{1}{A_2} \right) \theta_2^2}.$$

Stochastic Uzawa Algorithm Counter-Example

(2)

Apply the **tentative Uzawa algorithm**. The minimization in u gives:

- $\mathbf{A}_1^{(k+1)} \mathbf{U}_1^{(k+1)} + \mathbf{B}_1^{(k+1)} + \theta_1 \mathbf{P}^{(k)} = 0,$
- $\mathbf{A}_2^{(k+1)} \mathbf{U}_2^{(k+1)} + \mathbf{B}_2^{(k+1)} + \theta_2 \mathbf{P}^{(k)} = 0,$

and the update of the multiplier writes:

- $\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} + \rho^{(k)} (\theta_1 \mathbf{U}_1^{(k+1)} + \theta_2 \mathbf{U}_2^{(k+1)}).$

We thus obtain

$$\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} - \rho^{(k)} \left(\left(\frac{\theta_1^2}{\mathbf{A}_1^{(k+1)}} + \frac{\theta_2^2}{\mathbf{A}_2^{(k+1)}} \right) \mathbf{P}^{(k)} + \left(\theta_1 \frac{\mathbf{B}_1^{(k+1)}}{\mathbf{A}_1^{(k+1)}} + \theta_2 \frac{\mathbf{B}_2^{(k+1)}}{\mathbf{A}_2^{(k+1)}} \right) \right).$$

$\{\mathbf{P}^{(k)}\}_{k \in \mathbb{N}}$ can only converge to a stable equilibrium point of the associated differential equation (ODE argument), that is,

$$\bar{\mathbf{p}} = - \frac{\mathbb{E} \left(\frac{\mathbf{B}_1}{\mathbf{A}_1} \right) \theta_1 + \mathbb{E} \left(\frac{\mathbf{B}_2}{\mathbf{A}_2} \right) \theta_2}{\mathbb{E} \left(\frac{1}{\mathbf{A}_1} \right) \theta_1^2 + \mathbb{E} \left(\frac{1}{\mathbf{A}_2} \right) \theta_2^2}.$$

Stochastic Uzawa Algorithm Counter-Example

(2)

Apply the **tentative Uzawa algorithm**. The minimization in u gives:

- $\mathbf{A}_1^{(k+1)} \mathbf{U}_1^{(k+1)} + \mathbf{B}_1^{(k+1)} + \theta_1 \mathbf{P}^{(k)} = 0,$
- $\mathbf{A}_2^{(k+1)} \mathbf{U}_2^{(k+1)} + \mathbf{B}_2^{(k+1)} + \theta_2 \mathbf{P}^{(k)} = 0,$

and the update of the multiplier writes:

- $\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} + \rho^{(k)} (\theta_1 \mathbf{U}_1^{(k+1)} + \theta_2 \mathbf{U}_2^{(k+1)}).$

We thus obtain

$$\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} - \rho^{(k)} \left(\left(\frac{\theta_1^2}{\mathbf{A}_1^{(k+1)}} + \frac{\theta_2^2}{\mathbf{A}_2^{(k+1)}} \right) \mathbf{P}^{(k)} + \left(\theta_1 \frac{\mathbf{B}_1^{(k+1)}}{\mathbf{A}_1^{(k+1)}} + \theta_2 \frac{\mathbf{B}_2^{(k+1)}}{\mathbf{A}_2^{(k+1)}} \right) \right).$$

$\{\mathbf{P}^{(k)}\}_{k \in \mathbb{N}}$ can only converge to a stable equilibrium point of the associated differential equation (ODE argument), that is,

$$\bar{\mathbf{p}} = - \frac{\mathbb{E} \left(\frac{\mathbf{B}_1}{\mathbf{A}_1} \right) \theta_1 + \mathbb{E} \left(\frac{\mathbf{B}_2}{\mathbf{A}_2} \right) \theta_2}{\mathbb{E} \left(\frac{1}{\mathbf{A}_1} \right) \theta_1^2 + \mathbb{E} \left(\frac{1}{\mathbf{A}_2} \right) \theta_2^2}.$$

Stochastic Uzawa Algorithm Counter-Example

(2)

Apply the **tentative Uzawa algorithm**. The minimization in u gives:

- $\mathbf{A}_1^{(k+1)} \mathbf{U}_1^{(k+1)} + \mathbf{B}_1^{(k+1)} + \theta_1 \mathbf{P}^{(k)} = 0,$
- $\mathbf{A}_2^{(k+1)} \mathbf{U}_2^{(k+1)} + \mathbf{B}_2^{(k+1)} + \theta_2 \mathbf{P}^{(k)} = 0,$

and the update of the multiplier writes:

- $\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} + \rho^{(k)} (\theta_1 \mathbf{U}_1^{(k+1)} + \theta_2 \mathbf{U}_2^{(k+1)}).$

We thus obtain

$$\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} - \rho^{(k)} \left(\left(\frac{\theta_1^2}{\mathbf{A}_1^{(k+1)}} + \frac{\theta_2^2}{\mathbf{A}_2^{(k+1)}} \right) \mathbf{P}^{(k)} + \left(\theta_1 \frac{\mathbf{B}_1^{(k+1)}}{\mathbf{A}_1^{(k+1)}} + \theta_2 \frac{\mathbf{B}_2^{(k+1)}}{\mathbf{A}_2^{(k+1)}} \right) \right).$$

$\{\mathbf{P}^{(k)}\}_{k \in \mathbb{N}}$ can only converge to a stable equilibrium point of the associated differential equation (**ODE argument**), that is,

$$\bar{\mathbf{p}} = - \frac{\mathbb{E} \left(\frac{\mathbf{B}_1}{\mathbf{A}_1} \right) \theta_1 + \mathbb{E} \left(\frac{\mathbf{B}_2}{\mathbf{A}_2} \right) \theta_2}{\mathbb{E} \left(\frac{1}{\mathbf{A}_1} \right) \theta_1^2 + \mathbb{E} \left(\frac{1}{\mathbf{A}_2} \right) \theta_2^2}.$$

Stochastic Uzawa Algorithm Counter-Example

(3)

In general, $\mathbb{E}(\mathbf{X})/\mathbb{E}(\mathbf{Y})$ and $\mathbb{E}(\mathbf{X}/\mathbf{Y})$ are **different**, so that

$$p^\# = -\frac{\frac{\mathbb{E}(\mathbf{B}_1)}{\mathbb{E}(\mathbf{A}_1)}\theta_1 + \frac{\mathbb{E}(\mathbf{B}_2)}{\mathbb{E}(\mathbf{A}_2)}\theta_2}{\frac{\theta_1^2}{\mathbb{E}(\mathbf{A}_1)} + \frac{\theta_2^2}{\mathbb{E}(\mathbf{A}_2)}} \neq -\frac{\mathbb{E}\left(\frac{\mathbf{B}_1}{\mathbf{A}_1}\right)\theta_1 + \mathbb{E}\left(\frac{\mathbf{B}_2}{\mathbf{A}_2}\right)\theta_2}{\mathbb{E}\left(\frac{1}{\mathbf{A}_1}\right)\theta_1^2 + \mathbb{E}\left(\frac{1}{\mathbf{A}_2}\right)\theta_2^2} = \bar{p}.$$

This stochastic Uzawa algorithm does not solve the problem!

Feature. The standard stochastic gradient produces a Monte Carlo effect on the iterates $U^{(k)}$ by means of the coefficients $\lambda^{(k)}$. In the proposed Uzawa algorithm, $U^{(k)}$ is obtained by a procedure which does not incorporate such an effect, hence the failure.

Stochastic Uzawa Algorithm Counter-Example

(3)

In general, $\mathbb{E}(\mathbf{X})/\mathbb{E}(\mathbf{Y})$ and $\mathbb{E}(\mathbf{X}/\mathbf{Y})$ are **different**, so that

$$p^\# = -\frac{\frac{\mathbb{E}(B_1)}{\mathbb{E}(A_1)}\theta_1 + \frac{\mathbb{E}(B_2)}{\mathbb{E}(A_2)}\theta_2}{\frac{\theta_1^2}{\mathbb{E}(A_1)} + \frac{\theta_2^2}{\mathbb{E}(A_2)}} \neq -\frac{\mathbb{E}\left(\frac{B_1}{A_1}\right)\theta_1 + \mathbb{E}\left(\frac{B_2}{A_2}\right)\theta_2}{\mathbb{E}\left(\frac{1}{A_1}\right)\theta_1^2 + \mathbb{E}\left(\frac{1}{A_2}\right)\theta_2^2} = \bar{p}.$$

This stochastic Uzawa algorithm does not solve the problem!

Feature. The standard stochastic gradient produces a Monte Carlo effect on the iterates $U^{(k)}$ by means of the coefficients $\Delta^{(k)}$. In the proposed Uzawa algorithm, $U^{(k)}$ is obtained by a procedure which does not incorporate such an effect, hence the failure.

Stochastic Uzawa Algorithm Counter-Example

(3)

In general, $\mathbb{E}(\mathbf{X})/\mathbb{E}(\mathbf{Y})$ and $\mathbb{E}(\mathbf{X}/\mathbf{Y})$ are **different**, so that

$$p^\# = -\frac{\frac{\mathbb{E}(\mathbf{B}_1)}{\mathbb{E}(\mathbf{A}_1)}\theta_1 + \frac{\mathbb{E}(\mathbf{B}_2)}{\mathbb{E}(\mathbf{A}_2)}\theta_2}{\frac{\theta_1^2}{\mathbb{E}(\mathbf{A}_1)} + \frac{\theta_2^2}{\mathbb{E}(\mathbf{A}_2)}} \neq -\frac{\mathbb{E}\left(\frac{\mathbf{B}_1}{\mathbf{A}_1}\right)\theta_1 + \mathbb{E}\left(\frac{\mathbf{B}_2}{\mathbf{A}_2}\right)\theta_2}{\mathbb{E}\left(\frac{1}{\mathbf{A}_1}\right)\theta_1^2 + \mathbb{E}\left(\frac{1}{\mathbf{A}_2}\right)\theta_2^2} = \bar{p}.$$

This stochastic Uzawa algorithm does not solve the problem!

Feature. The standard stochastic gradient produces a **Monte Carlo** effect on the iterates $\mathbf{U}^{(k)}$ by means of the coefficients $\epsilon^{(k)}$. In the proposed Uzawa algorithm, $\mathbf{U}^{(k)}$ is obtained by a procedure which does not incorporate such an effect, hence the **failure**.

Stochastic APP Algorithm with Constraints

Consider the **APP** algorithm in the **deterministic** setting:

$$u^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \rangle + \epsilon \langle p^{(k)}, \Theta(u) \rangle,$$

$$p^{(k+1)} = \text{proj}_{C^*} (p^{(k)} + \rho \Theta(u^{(k+1)})).$$

The extension to the **stochastic** case is obtained in a canonical way by replacing in the above minimization stage the gradient of J by the partial gradient of j w.r.t. u , evaluated at a sample $W^{(k+1)}$ of W . Using the notation $G^{(k)} = \nabla_u j(U^{(k)}, W^{(k+1)})$, we obtain:

Stochastic APP Algorithm in the Constrained Case

$$U^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \langle c^{(k)} G^{(k)} - \nabla K(U^{(k)}), u \rangle + c^{(k)} \langle P^{(k)}, \Theta(u) \rangle,$$

$$P^{(k+1)} = \text{proj}_{C^*} (P^{(k)} + c^{(k)} \Theta(U^{(k+1)})).$$

Note that it **never** leads to the Uzawa algorithm because $c^{(k)} \rightarrow 0$.

Stochastic APP Algorithm with Constraints

Consider the **APP** algorithm in the **deterministic** setting:

$$u^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \langle \epsilon \nabla J(u^{(k)}) - \nabla K(u^{(k)}), u \rangle + \epsilon \langle p^{(k)}, \Theta(u) \rangle,$$

$$p^{(k+1)} = \text{proj}_{C^*} (p^{(k)} + \rho \Theta(u^{(k+1)})).$$

The extension to the **stochastic** case is obtained in a canonical way by **replacing** in the above minimization stage the gradient of J by the partial gradient of j w.r.t. u , evaluated at a sample $\mathbf{W}^{(k+1)}$ of \mathbf{W} . Using the notation $\mathbf{G}^{(k)} = \nabla_u j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$, we obtain:

Stochastic APP Algorithm in the Constrained Case

$$\mathbf{U}^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \langle \epsilon^{(k)} \mathbf{G}^{(k)} - \nabla K(\mathbf{U}^{(k)}), u \rangle + \epsilon^{(k)} \langle \mathbf{P}^{(k)}, \Theta(u) \rangle,$$

$$\mathbf{P}^{(k+1)} = \text{proj}_{C^*} (\mathbf{P}^{(k)} + \epsilon^{(k)} \Theta(\mathbf{U}^{(k+1)})).$$

Note that it **never** leads to the Uzawa algorithm because $\epsilon^{(k)} \rightarrow 0$.

Convergence Theorem

(1)

We make the following assumptions.

- H1** U^{ad} is a nonempty closed convex subset of a Hilbert space \mathbb{U} , and \mathcal{C} is a closed convex salient cone of a Hilbert space \mathbb{V} .
- H2** $j : \mathbb{U} \times \mathbb{W} \rightarrow \mathbb{R}$ is a **normal integrand**, and $\mathbb{E}(j(u, \mathbf{W}))$ exists for all $u \in U^{\text{ad}}$.
- H3** $j(\cdot, w) : \mathbb{U} \rightarrow \mathbb{R}$ is a proper convex differentiable function with linearly bounded gradients (**LBG**), for all $w \in \mathbb{W}$.
- H4** J is **strictly convex** and coercive on U^{ad} .
- H5** Θ is \mathcal{C} -convex, Lipschitz with constant L_Θ .
- H6** A **constraint qualification condition** holds true.
- H7** K is a proper l.s.c. function, **strongly convex** with modulus b and differentiable.
- H8** The sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a σ -sequence.

Convergence Theorem

(2)

Then the following conclusions hold true.

- R1** The initial constrained problem has a **non empty set** of saddle points $\{u^\#\} \times P^\#$.
- R2** Each auxiliary problem has a **unique solution** $U^{(k+1)}$.
- R3** The sequence of random variables $\{L(U^{(k)}, p^\#)\}_{k \in \mathbb{N}}$ **almost surely converges** to $L(u^\#, p^\#)$ for all $p^\# \in P^\#$.
- R4** The sequences of random variables $\{U^{(k)}\}_{k \in \mathbb{N}}$ and $\{P^{(k)}\}_{k \in \mathbb{N}}$ are **almost surely bounded**.
- R5** The sequence of random variables $\{U^{(k)}\}_{k \in \mathbb{N}}$ **almost surely converges** to $u^\#$.

Sketch of Proof

The proof of the first two statements is based on standard theorems.

The proof of the last two statements consists of three steps.

❶ **Select a Lyapunov function.**

$$\Lambda(u, p) = K(u^\#) - K(u) - \langle \nabla K(u), u^\# - u \rangle + \|p - p^\#\|^2 / 2.$$

❷ **Bound from above the variation of Λ .**

Using assumptions and writing optimality conditions, we get:

$$\begin{aligned} \mathbb{E}(\Lambda(\mathbf{U}^{(k+1)}, \mathbf{P}^{(k+1)}) \mid \mathcal{F}^{(k)}) &\leq (1 + \alpha^{(k)})\Lambda(\mathbf{U}^{(k)}, \mathbf{P}^{(k)}) \\ &\quad + \beta^{(k)} - \epsilon^{(k)}(L(\mathbf{U}^{(k)}, p^\#) - L(u^\#, p^\#)). \end{aligned}$$

❸ **Prove the convergence of the sequences.**

Using the two lemmas, we obtain that $\{\Lambda(\mathbf{U}^{(k)}, \mathbf{P}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to a finite random variable Λ^∞ , and that $\{L(\mathbf{U}^{(k)}, p^\#)\}_{k \in \mathbb{N}}$ almost surely converges to $L(u^\#, p^\#)$.

By a compactness argument and uniqueness of $u^\#$, the sequence $\{\mathbf{U}^{(k)}\}_{k \in \mathbb{N}}$ almost surely converges to $u^\#$.

- 1 Auxiliary Problem Principle in the Deterministic Setting
 - Principle and Algorithm
 - Convergence and Features
 - Explicit Constraints
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - Stochastic APP with Deterministic Constraints
 - Extensions of Stochastic APP with Constraints

Stochastic APP and Augmented Lagrangian

In order to deal with **non stable** problems, we extend the use of the **Augmented Lagrangian** to the stochastic framework by replacing the gradient of J by the partial gradient of j w.r.t. u .

Using the notation $\mathbf{G}^{(k)} = \nabla_u j(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$, we obtain the following augmented Lagrangian stochastic algorithm:

Regularized Stochastic APP algorithm with Constraints

$$\mathbf{U}^{(k+1)} \in \arg \min_{u \in \mathcal{U}^{\text{ad}}} K(u) + \langle \epsilon^{(k)} \mathbf{G}^{(k)} - \nabla K(\mathbf{U}^{(k)}), u \rangle + \epsilon^{(k)} \langle \text{proj}_{\mathcal{C}^*}(\mathbf{P}^{(k)} + c\Theta(\mathbf{U}^{(k)})), \Theta(u) \rangle,$$

$$\mathbf{P}^{(k+1)} = \left(1 - \frac{\epsilon^{(k)}}{c}\right) \mathbf{P}^{(k)} + \frac{\epsilon^{(k)}}{c} \text{proj}_{\mathcal{C}^*}(\mathbf{P}^{(k)} + c\Theta(\mathbf{U}^{(k+1)})).$$



Convergence Theorem

(1)

We make the following assumptions.

- H1** U^{ad} is a nonempty closed convex subset of a Hilbert space \mathbb{U} , and \mathcal{C} is a closed convex cone of another Hilbert space \mathbb{V} .
- H2** $j : \mathbb{U} \times \mathbb{W} \rightarrow \mathbb{R}$ is a **normal integrand**, and $\mathbb{E}(j(u, \mathbf{W}))$ exists for all $u \in U^{\text{ad}}$.
- H3** $j(\cdot, w) : \mathbb{U} \rightarrow \mathbb{R}$ is a proper convex differentiable function with linearly bounded gradients (**LBG**), for all $w \in \mathbb{W}$.
- H4** J is coercive on U^{ad} .
- H5** Θ is \mathcal{C} -convex, Lipschitz with constant L_Θ .
- H6** A **constraint qualification condition** holds true.
- H7** K is a proper l.s.c. function, **strongly convex** with modulus b and differentiable.
- H8** The sequence $\{\epsilon^{(k)}\}_{k \in \mathbb{N}}$ is a σ -sequence.

Convergence Theorem

(2)

Then the following conclusions hold true.

- R1** The initial constrained problem has a **non empty set** of saddle points $U^\# \times P^\#$.
- R2** Each auxiliary problem has an **unique solution** $U^{(k+1)}$.
- R3** The sequence of r.v. $\{L_c(u^\#, P^{(k)}) - L_c(U^{(k)}, p^\#)\}_{k \in \mathbb{N}}$ **almost surely converges** to zero for all saddle point $(u^\#, p^\#)$.
- R4** The sequences of random variables $\{U^{(k)}\}_{k \in \mathbb{N}}$ and $\{P^{(k)}\}_{k \in \mathbb{N}}$ are **almost surely bounded**.
- R5** Each **cluster point** of a **realization** of the sequence $\{U^{(k)}\}_{k \in \mathbb{N}}$ (resp. $\{P^{(k)}\}_{k \in \mathbb{N}}$) **almost surely converges** to an element of $U^\#$ (resp. $P^\#$).

Sketch of Proof

The proof of the first two statements is based on standard theorems.

The proof of the last two statements follows the usual scheme, with

$$\Lambda(u, p) = K(u^\#) - K(u) - \langle \nabla K(u), u^\# - u \rangle + \frac{1}{2} \|p - p^\#\|^2,$$

and a substantial amount of technicalities. . .

See the lecture notes about stochastic gradient (in French) on the Web site.

Stochastic APP and Constraints in Expectation

(1)

We finally aim at solving stochastic optimization problems in which the **deterministic** constraint in fact corresponds to the **expectation** of a stochastic constraint.

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \mathbb{E}(\theta(u, \mathbf{W})) \in -C .$$

Such a framework proves useful when dealing with **constraints in probability**:

$$\mathbb{P}(\theta(u, \mathbf{W}) \in -C) \geq \pi \quad \iff \quad \mathbb{E}(\mathbf{1}_{\{\theta(u, \mathbf{W}) \in -C\}}) \geq \pi .$$

In the spirit of the stochastic gradient method, we use values of θ evaluated at realizations of \mathbf{W} rather than **expected** values of Θ .

Stochastic APP and Constraints in Expectation

(1)

We finally aim at solving stochastic optimization problems in which the **deterministic** constraint in fact corresponds to the **expectation** of a stochastic constraint.

$$\min_{u \in U^{\text{ad}}} \mathbb{E}(j(u, \mathbf{W})) \quad \text{subject to} \quad \mathbb{E}(\theta(u, \mathbf{W})) \in -C .$$

Such a framework proves useful when dealing with **constraints in probability**:

$$\mathbb{P}(\theta(u, \mathbf{W}) \in -C) \geq \pi \quad \iff \quad \mathbb{E}(\mathbf{1}_{\{\theta(u, \mathbf{W}) \in -C\}}) \geq \pi .$$

In the spirit of the stochastic gradient method, we use values of θ evaluated at realizations of \mathbf{W} rather than **expected** values of Θ .

Stochastic APP and Constraints in Expectation

(2)

We use the **APP** framework in the stochastic setting, and we **replace** both the gradients ∇J and $\nabla \Theta$ by the partial gradients ∇_{uj} and $\nabla_u \theta$ evaluated at realizations of \mathbf{W} .

With the notations

- $\mathbf{G}^{(k)} = \nabla_{uj}(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$,
- $\vartheta^{(k)} = \theta'_u(\mathbf{U}^{(k)}, \mathbf{W}^{(k+1)})$,

the **extension** of the APP method to expected constraints is:

Stochastic APP Algorithm with Expected Constraints

$$\begin{aligned} \mathbf{U}^{(k+1)} \in \arg \min_{u \in U^{\text{ad}}} K(u) + \langle \epsilon^{(k)} \mathbf{G}^{(k)} - \nabla K(\mathbf{U}^{(k)}), u \rangle \\ + \epsilon^{(k)} \langle \mathbf{P}^{(k)}, \vartheta^{(k)} \cdot u \rangle, \\ \mathbf{P}^{(k+1)} = \text{proj}_{C^*} \left(\mathbf{P}^{(k)} + \rho^{(k)} \theta(\mathbf{U}^{(k+1)}, \mathbf{W}^{(k+1)}) \right). \end{aligned}$$

Convergence Theorem and Proof

Long and intricate. . .

See the lecture notes (in French).