Stochastic Gradient Method: Convergence

January 20, 2015

Lecture Outline

1 Auxiliary Problem Principle in the Deterministic Setting

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

Principle and Algorithm Convergence Theorem and Proof Features of APP

1 Auxiliary Problem Principle in the Deterministic Setting

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

Principle and Algorithm Convergence Theorem and Proof Features of APP

Auxiliary Problem Principle in the Deterministic Setting Principle and Algorithm

- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - APP with Constraints in the Deterministic Setting
 - APP with Constraints in the Stochastic Setting

Principle and Algorithm Convergence Theorem and Proof Features of APP

Problem under Consideration

Consider the following convex differentiable optimization problem:

 $\min_{u\in U^{\mathrm{ad}}}J(u).$

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

$$\left\langle
abla J(u^{\sharp}) \,, u - u^{\sharp} \right\rangle \geq 0 \;, \;\; \forall u \in U^{\mathrm{ad}} \;.$$

In the deterministic framework, the Auxiliary Problem Principle (APP) consists in replacing the original problem by a sequence of auxiliary problems indexed by $k \in \mathbb{N}$, and without too much perturbing the optimality conditions...

Principle and Algorithm Convergence Theorem and Proof Features of APP

APP Framework

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

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

But then the criterion is no more coercive...

Second idea: add a strongly convex term:

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

K being a real-valued differentiable function defined on \mathbb{U} and ϵ being a positive constant. At iteration *k*, given $u^{(k)} \in U^{\mathrm{ad}}$, consider the following auxiliary problem:

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

The solution $u^{(k+1)}$ of the auxiliary problem at iteration k is used to formulate a new auxiliary problem at iteration k + 1.

P. Carpentier

Master MMMEF — Cours MNOS

2014-2015 54 / 267

Principle and Algorithm Convergence Theorem and Proof Features of APP

APP Algorithm

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

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

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

Note that the optimality condition of the auxiliary problem writes

$$\left\langle
abla \mathcal{K}(u^{(k+1)}) + \epsilon
abla J(u^{(k)}) -
abla \mathcal{K}(u^{(k)}), u - u^{(k+1)} \right\rangle \geq 0 \;, \; \forall u \in U^{\mathrm{ad}} \;,$$

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

Principle and Algorithm Convergence Theorem and Proof Features of APP

Auxiliary Problem Principle in the Deterministic Setting Principle and Algorithm

• Convergence Theorem and Proof

Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions
- **3** Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - APP with Constraints in the Deterministic Setting
 - APP with Constraints in the Stochastic Setting

Principle and Algorithm Convergence Theorem and Proof Features of APP

Convergence Theorem

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, coercive on U^{ad} and differentiable, ∇J being 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 $0 < \epsilon < \frac{2b}{A}$.

Principle and Algorithm Convergence Theorem and Proof Features of APP

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 an unique solution $u^{(k+1)}$.
- **R2** The sequence $\{J(u^{(k)})\}_{k \in \mathbb{N}}$ is strictly decreasing and converges towards $J(u^{\sharp})$.
- **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.

Principle and Algorithm Convergence Theorem and Proof Features of APP

Sketch of Proof

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

The proof of the last two statements involves four steps.

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

The result holds true if \mathbb{U} is an infinite dimensional Hilbert space, and may be extended to problems with explicit constraints.

Principle and Algorithm Convergence Theorem and Proof Features of APP

Auxiliary Problem Principle in the Deterministic Setting

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

Principle and Algorithm Convergence Theorem and Proof Features of APP

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 [Cohen, 2004] 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

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

and its solution has the following closed-form expression:

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

We obtain the well-known projected gradient algorithm.

Principle and Algorithm Convergence Theorem and Proof Features of APP

Some Features of APP



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

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

the auxiliary subproblem becomes

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

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

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - APP with Constraints in the Deterministic Setting
 - APP with Constraints in the Stochastic Setting

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - APP with Constraints in the Deterministic Setting
 - APP with Constraints in the Stochastic Setting

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Standard Stochastic Gradient Method

(1)

Consider the following open-loop stochastic optimization problem:

 $\min_{u\in U^{\rm ad}}J(u)\,,$

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

The standard stochastic gradient algorithm reads as follows.

- 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. W.
- Sompute the gradient of j and update $u^{(k+1)}$ by the formula:

$$u^{(k+1)} = \operatorname{proj}_{U^{\mathrm{ad}}} \left(u^{(k)} - \epsilon^{(k)} \nabla_{u} j(u^{(k)}, w^{(k+1)}) \right)$$

• Set k = k + 1 and go to step 2.

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Standard Stochastic Gradient Method

(2)

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

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

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

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

Robbins-Monro Theorem

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

• Reminders about the Stochastic Gradient Method

• Stochastic APP Algorithm

- Convergence Theorem and Proof
- Conclusions

3 Stochastic APP with Explicit Constraints

- Constraints in Stochastic Optimization
- APP with Constraints in the Deterministic Setting
- APP with Constraints in the Stochastic Setting

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Moving APP to the Stochastic Framework

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

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

Then, in each auxiliary problem, we replace the gradient of J by the gradient of j evaluated at sampled realizations of W. Note that the "large" (constant) step size ϵ has to be replaced by "small" (going to zero as index k goes to infinity) steps $\epsilon^{(k)}$. The k-th instance of the stochastic auxiliary problem is thus

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

 $w^{(k+1)}$ being a realization of the random variable W.

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Stochastic APP Algorithm

We thus obtain a generalized stochastic gradient algorithm.

Stochastic APP Algorithm

- Let $u^{(0)} \in U^{ad}$ 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. W.
- **(3)** Update $u^{(k+1)}$ by solving the auxiliary problem:

 $u^{(k+1)} \in \underset{u \in U^{\mathrm{ad}}}{\mathrm{arg\,min\,}} \mathcal{K}(u) + \Big\langle \epsilon^{(k)} \nabla_{u} j(u^{(k)}, w^{(k+1)}) - \nabla \mathcal{K}(u^{(k)}), u \Big\rangle.$

• Set k = k + 1 and go to step 2.

As usual, the algorithm is casted in the probabilistic framework: $\boldsymbol{U}^{(k+1)} \in \underset{u \in U^{\mathrm{ad}}}{\operatorname{arg\,min}} \mathcal{K}(u) + \left\langle \epsilon^{(k)} \nabla_{u} j(\boldsymbol{U}^{(k)}, \boldsymbol{W}^{(k+1)}) - \nabla \mathcal{K}(\boldsymbol{U}^{(k)}), u \right\rangle$.

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

P. Carpentier

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Example

With the choice

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

the auxiliary problem becomes

$$\min_{u\in U^{\mathrm{ad}}}\frac{1}{2}\|u\|^2+\left\langle\epsilon^{(k)}\nabla_u j(\boldsymbol{U}^{(k)},\boldsymbol{W}^{(k+1)})-\boldsymbol{U}^{(k)},u\right\rangle.$$

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

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

It corresponds to the standard stochastic gradient iteration.

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient MethodStochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions
- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - APP with Constraints in the Deterministic Setting
 - APP with Constraints in the Stochastic Setting

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Convergence Theorem

Make the following assumptions.

- **H1** U^{ad} is a nonempty closed convex subset of a Hilbert space \mathbb{U} .
- **H2** $j : \mathbb{U} \times \mathbb{W} \to \mathbb{R}$ is a normal integrand, and $\mathbb{E}(j(u, W))$ exists for all $u \in U^{\mathrm{ad}}$.
- **H3** $j(\cdot, w) : \mathbb{U} \to \mathbb{R}$ is a proper convex differentiable function for all $w \in \mathbb{W}$ $(j(\cdot, w)$ is l.s.c. thanks to Assumption H2).
- H4 $j(\cdot, w)$ has linearly bounded gradients (LBG):

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

- **H5** J is Lipschitz continuous and 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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

(2)

Then the following conclusions hold true.

- **R1** The initial problem has a non empty set of solutions U^{\sharp} .
- **R2** Each auxiliary problem has a unique solution $U^{(k+1)}$.
- **R3** The sequence of random variables $\{J(U^{(k)})\}_{k\in\mathbb{N}}$ almost surely converges to $J^{\sharp} = \min_{u \in U^{\mathrm{ad}}} J(u)$.
- **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^{\sharp} .

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

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.

- Select a Lyapunov function. Here we choose $\Lambda(u) = K(u^{\sharp}) - K(u) - \langle \nabla K(u), u^{\sharp} - u \rangle$.
- **2** Bound from above the variation of Λ . Using assumptions and writing optimality conditions, we get: $\mathbb{E}(\Lambda(\boldsymbol{U}^{(k+1)}) | \mathcal{F}^{(k)}) \leq (1 + \alpha^{(k)})\Lambda(\boldsymbol{U}^{(k)}) + \beta^{(k)} - \epsilon^{(k)}(J(\boldsymbol{U}^{(k)}) - J(u^{\sharp}))$.
- Prove the convergence of the sequences.

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

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

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

and that

$$\sum_{k\in\mathbb{N}} oldsymbollpha^{(k)} < +\infty \quad ext{and} \quad \sum_{k\in\mathbb{N}} oldsymboleta^{(k)} < +\infty \;, \; \mathbb{P} ext{-a.s.} \;.$$

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} -a.s..

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Two Useful Lemmas

Technical Lemma

Let J be a real-valued function defined on a Hilbert space \mathbb{U} . We assume 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) $\exists \delta > 0, \forall k \in \mathbb{N}, ||u^{(k+1)} - u^{(k)}|| \le \delta \epsilon^{(k)}.$
Then the sequence $\{J(u^{(k)})\}_{k \in \mathbb{N}}$ converges to μ .

Proof of Convergence

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

(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 $u^{\sharp} \in U^{\sharp}$ be a solution of the initial problem.

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

 $\min_{u \in U^{\mathrm{ad}}} K(u) + \left\langle \epsilon^{(k)} \nabla_u j(\boldsymbol{U}^{(k)}, \boldsymbol{W}^{(k+1)}) - \nabla K(\boldsymbol{U}^{(k)}), u \right\rangle ,$

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.

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Proof of Convergence



7)

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

$$\Lambda(u) \ge \frac{b}{2} \|u - u^{\sharp}\|_{\mathbb{U}}^{2}, \qquad (a^{\sharp})$$

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

• Consider the variation of Λ during the algorithm:

$$\boldsymbol{\Delta}^{(k)} = \Lambda(\boldsymbol{U}^{(k+1)}) - \Lambda(\boldsymbol{U}^{(k)})$$

= $\underbrace{\mathcal{K}(\boldsymbol{U}^{(k)}) - \mathcal{K}(\boldsymbol{U}^{(k+1)}) - \langle \nabla \mathcal{K}(\boldsymbol{U}^{(k)}), \boldsymbol{U}^{(k)} - \boldsymbol{U}^{(k+1)} \rangle}_{T_1}$
+ $\underbrace{\langle \nabla \mathcal{K}(\boldsymbol{U}^{(k)}) - \nabla \mathcal{K}(\boldsymbol{U}^{(k+1)}), \boldsymbol{u}^{\sharp} - \boldsymbol{U}^{(k+1)} \rangle}_{T_2}$

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Proof of Convergence

(3)

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

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

• From the convexity of $j(\cdot, w)$, we have that

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

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

Using the Schwartz inequality, we obtain: $T_4 \leq \frac{\epsilon^{(k)}}{b} \| \mathbf{G}^{(k)} \|_{\mathbb{U}}^2$. The LBG assumption and the majoration (7) of Λ yield:

$$T_4 \leq rac{\epsilon^{(k)}}{b} \left(lpha \Lambda(oldsymbol{U}^{(k)}) + eta
ight) \; .$$

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Proof of Convergence



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

$$\boldsymbol{\Delta}^{(k)} \leq \epsilon^{(k)} \left(j(\boldsymbol{u}^{\sharp}, \boldsymbol{W}^{(k+1)}) - j(\boldsymbol{U}^{(k)}, \boldsymbol{W}^{(k+1)}) \right) + \frac{\left(\epsilon^{(k)}\right)^2}{b} \left(\alpha \Lambda(\boldsymbol{U}^{(k)}) + \beta \right) \,.$$

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

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

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

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

Reminder. We have also obtained the two following inequalities:

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

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Proof of Convergence

From the Robbins-Siegmund theorem, {Λ(U^(k))}_{k∈ℕ} almost surely converges to a finite random variable Λ[∞] and we have

$$\sum_{k=0}^{+\infty} \epsilon^{(k)} ig(J(oldsymbol{U}^{(k)}) - J(u^{\sharp}) ig) < +\infty \ , \ \ \mathbb{P} ext{-a.s.} \ .$$

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

We deduce that both sequences {U^(k)}_{k∈ℕ} and {G^(k)}_{k∈ℕ} are a.s. bounded, so that the same holds true for the sequence {1/(ϵ^(k)) || U^(k+1) - U^(k) ||_U}_{k∈ℕ}. This makes it possible to use the second technical lemma and claim that {J(U^(k))}_{k∈ℕ} almost surely converges to J(u[#]).

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

Proof of Convergence



Pick some ω ∈ Ω₀. The sequence of realizations {u^(k)}_{k∈ℕ} of {U^(k)}_{k∈ℕ} associated with ω is bounded, and u^(k) ∈ U^{ad}. By a compactness argument,⁶ it exists a convergent subsequence {u^{(Φ(k))}}_{k∈ℕ}, with limit ū. Using the lower semi-continuity of function J, we have that

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

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

⁶A 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]).

P. Carpentier

2014-2015 82 / 267

Proof of Convergence

We ultimately consider the case when J is strongly convex with modulus a. Then the initial problem has a unique solution u^{\sharp} . Thanks to the strong convexity property of J, we have

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

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

The proof is complete.

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof

Conclusions

- 3 Stochastic APP with Explicit Constraints
 - Constraints in Stochastic Optimization
 - APP with Constraints in the Deterministic Setting
 - APP with Constraints in the Stochastic Setting

Reminders about the Stochastic Gradient Method Stochastic APP Algorithm Convergence Theorem and Proof Conclusions

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 proved under 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 the σ -sequence $\{\epsilon^{(k)}\}_{k\in\mathbb{N}}$.

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP
- 2 Auxiliary Problem Principle in the Stochastic Setting
 - Reminders about the Stochastic Gradient Method
 - Stochastic APP Algorithm
 - Convergence Theorem and Proof
 - Conclusions

3 Stochastic APP with Explicit Constraints

- Constraints in Stochastic Optimization
- APP with Constraints in the Deterministic Setting
- APP with Constraints in the Stochastic Setting

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions

3 Stochastic APP with Explicit Constraints

Constraints in Stochastic Optimization

• APP with Constraints in the Deterministic Setting

• APP with Constraints in the Stochastic Setting

A Constrained Stochastic Optimization Problem

Recall that in our stochastic optimization setting, the probability space is denoted $(\Omega, \mathcal{A}, \mathbb{P})$, and W is a random variable defined on the space $(\mathbb{W}, \mathcal{W})$.

- We are interested in the case where the criterion J is defined as J(u) = E(j(u, W)), with j : U × W → R. This is the standard framework when studying open-loop stochastic optimization problems.
- We will hereafter consider only constraints ⊖ which are of a deterministic nature: ⊖ : U → V.
- The problem we deal with has thus the following expression:

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

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Nature of Constraints in Stochastic Optimization

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

- A constraint may be deterministic: $\Theta(u) \in -C$.
- A constraint may be formulated in the almost sure sense:
 θ(u, W) ∈ −C ℙ-a.s.. It is generally used to express hard constraints (physical laws, ...).
- Another (more realistic) way is to formulate stochastic constraints in probability: P(θ(u, W) ∈ −C) ≥ π, which means that the constraints can sometimes be violated.
- Another possibility is to have a constraint in expectation: E(θ(u, W)) ∈ −C. Although usually non intuitive, such a formulation proves useful in some specific problems, as it will be illustrated in the lecture next week.

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions

3 Stochastic APP with Explicit Constraints

Constraints in Stochastic Optimization

• APP with Constraints in the Deterministic Setting

• APP with Constraints in the Stochastic Setting

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Problem under Consideration

Consider the following convex optimization problem:

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

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*.

addefined as $C^{\star} = \{ p \in \mathbb{V}, \langle p, v \rangle \ge 0 \ \forall v \in C \}.$

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Uzawa and Arrow-Hurwicz Algorithms

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

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

This problem can be solved by using the Uzawa algorithm:

$$u^{(k+1)} \in \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} J(u) + \left\langle p^{(k)}, \Theta(u) \right\rangle,$$
$$p^{(k+1)} = \operatorname{proj}_{C^{\star}} \left(p^{(k)} + \rho \, \Theta(u^{(k+1)}) \right).$$

Another possibility is to use the Arrow-Hurwicz algorithm:

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

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

APP with Explicit Constraints

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

$$u^{(k+1)} = \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} \mathcal{K}(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla \mathcal{K}(u^{(k)}), u \right\rangle + \epsilon \left\langle p^{(k)}, \Theta(u) \right\rangle,$$
$$p^{(k+1)} = \operatorname{proj}_{\mathcal{C}^{\star}} \left(p^{(k)} + \rho \Theta(u^{(k+1)}) \right).$$

It is not difficult to show that this APP 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 $\epsilon \langle p^{(k)}, \Theta(u) \rangle$ may be replaced by $\epsilon \langle p^{(k)}, \Theta'(u^{(k)}).u \rangle$.

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

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 $\widehat{U}(p^{\sharp})$ of solutions associated to the primal minimization:

 $\widehat{U}(p^{\sharp}) = \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} L(u, p^{\sharp}) \ .$

Then the set U^{\sharp} of solutions of the initial problem may be strictly included in $\widehat{U}(p^{\sharp})$, 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 $\widehat{U}(1) = [-1, 1]$.

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)!

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

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 $U^{ad} \times \mathbb{V}$,⁸ and its expression, which depends on a scalar parameter c > 0, is given by:

$$L_c(u,p) = J(u) + \frac{1}{2c} \Big(\| \operatorname{proj}_{C^{\star}}(p + c\Theta(u)) \|^2 - \|p\|^2 \Big).$$

The augmented Lagrangian has the two following properties.

- 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:

$$U^{\sharp} = \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} L_c(u, p^{\sharp}) \ .$$

⁸whereas the standard Lagrangian L is defined on $U^{\mathrm{ad}} imes C^{\star}$

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

APP and Augmented Lagrangian

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

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

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

$$u^{(k+1)} = \underset{u \in U^{\mathrm{ad}}}{\operatorname{arg\,min}} \mathcal{K}(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla \mathcal{K}(u^{(k)}), u \right\rangle \\ + \epsilon \left\langle \operatorname{proj}_{C^{\star}} \left(p^{(k)} + c \Theta(u^{(k)}) \right), \Theta(u) \right\rangle,$$

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

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

1 Auxiliary Problem Principle in the Deterministic Setting

- Principle and Algorithm
- Convergence Theorem and Proof
- Features of APP

2 Auxiliary Problem Principle in the Stochastic Setting

- Reminders about the Stochastic Gradient Method
- Stochastic APP Algorithm
- Convergence Theorem and Proof
- Conclusions

3 Stochastic APP with Explicit Constraints

- Constraints in Stochastic Optimization
- APP with Constraints in the Deterministic Setting
- APP with Constraints in the Stochastic Setting

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

On the Agenda

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

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

A Useful Tool in Stochastic Approximation

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

$$\boldsymbol{U}^{(k+1)} = \boldsymbol{U}^{(k)} + \epsilon^{(k)} \Big(h(\boldsymbol{U}^{(k)}) + \boldsymbol{\xi}^{(k+1)} \Big)$$

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

u=h(u),

(see [Kushner & Clark, 1978]). A useful corollary is the following.

Let $\{\boldsymbol{U}^{(k)}\}_{k\in\mathbb{N}}$ be the sequence generated by the Stochastic Approximation algorithm, and assume that $\exists u^{\sharp} \in \mathbb{U}$, such that $\mathbb{P}\left(\lim_{k \to +\infty} \boldsymbol{U}^{(k)} = u^{\sharp}\right) > 0$. Then u^{\sharp} 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:

 $\min_{u \in U^{\mathrm{ad}}} \mathbb{E} \left(j(u, \boldsymbol{W})
ight)$ subject to $\Theta(u) \in -C$,

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 $j(u, w^{(k+1)})$.⁹ We thus obtain a tentative Stochastic Uzawa Algorithm:

$$\begin{split} \boldsymbol{U}^{(k+1)} &= \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} j(u, \boldsymbol{W}^{(k+1)}) + \left\langle \boldsymbol{P}^{(k)}, \boldsymbol{\Theta}(u) \right\rangle, \\ \boldsymbol{P}^{(k+1)} &= \operatorname{proj}_{C^{\star}} \left(\boldsymbol{P}^{(k)} + \rho^{(k)} \boldsymbol{\Theta}(\boldsymbol{U}^{(k+1)}) \right). \end{split}$$

Question: what about the convergence of this algorithm?

⁹Note that we have replaced here the evaluation of J by the one of j, whereas ∇J is replaced by $\nabla_u j$ in the stochastic gradient method...

P. Carpentier

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Stochastic Uzawa Algorithm Counter-Example

Consider a constrained stochastic optimization problem with:

- $\mathbb{U} = \mathbb{R}^2$ and $U^{\mathrm{ad}} = \mathbb{U}$,
- $\mathbb{V} = \mathbb{R}$ and $C = \{0\}$ (equality constraint),
- $\mathbb{W} = \mathbb{R}^4$ and $\boldsymbol{W} = (\boldsymbol{A}_1, \boldsymbol{A}_2, \boldsymbol{B}_1, \boldsymbol{B}_2)$,
- $j(u, w) = \frac{1}{2}(a_1u_1^2 + a_2u_2^2) + (b_1u_1 + b_2u_2),$
- $\Theta(u) = \theta_1 u_1 + \theta_2 u_2.$

The optimality conditions (KKT) of this problem write:

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

so that the value of the optimal multiplier is:

$$p^{\sharp} = -rac{rac{\mathbb{E}\left(oldsymbol{B}_{1}
ight)}{\mathbb{E}\left(oldsymbol{A}_{1}
ight)} heta_{1}+rac{\mathbb{E}\left(oldsymbol{B}_{2}
ight)}{\mathbb{E}\left(oldsymbol{A}_{2}
ight)} heta_{2}}{rac{ heta_{1}^{2}}{\mathbb{E}\left(oldsymbol{A}_{1}
ight)}+rac{ heta_{2}^{2}}{\mathbb{E}\left(oldsymbol{A}_{2}
ight)}}\;.$$

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Stochastic Uzawa Algorithm Counter-Example

Apply our Uzawa algorithm. The minimization stage leads to:

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

and the update of the multiplier writes:

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

We thus obtain

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

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

$$\overline{p} = -rac{\mathbb{E}\left(rac{m{B_1}}{m{A_1}}
ight) heta_1 + \mathbb{E}\left(rac{m{B_2}}{m{A_2}}
ight) heta_2}{rac{ heta_1^2}{\mathbb{E}(m{A_1})} + rac{ heta_2^2}{\mathbb{E}(m{A_2})}} \,.$$

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Stochastic Uzawa Algorithm Counter-Example

As soon as the random variables A_i and B_i are non independent, we have



The stochastic Uzawa algorithm does not solve the problem!

Remark. The standard stochastic gradient produces an averaging effect on the iterates $U^{(k)}$ by means of the coefficients $\epsilon^{(k)}$. In the Uzawa algorithm, $U^{(k)}$ is obtained by a minimization procedure which does not incorporate any averaging effect, hence the failure.

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Stochastic APP Algorithm with Constraints

Consider the APP algorithm in the deterministic setting:

 $u^{(k+1)} = \underset{u \in U^{\mathrm{ad}}}{\mathrm{arg\,min}} \mathcal{K}(u) + \left\langle \epsilon \nabla J(u^{(k)}) - \nabla \mathcal{K}(u^{(k)}), u \right\rangle + \epsilon \left\langle p^{(k)}, \Theta(u) \right\rangle,$

 $p^{(k+1)} = \operatorname{proj}_{C^{\star}} \left(p^{(k)} + \rho \Theta(u^{(k+1)}) \right)$.

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

 $\begin{aligned} \boldsymbol{U}^{(k+1)} &= \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} \boldsymbol{K}(u) + \left\langle \epsilon^{(k)} \boldsymbol{G}^{(k)} - \nabla \boldsymbol{K}(\boldsymbol{U}^{(k)}), u \right\rangle + \epsilon^{(k)} \left\langle \boldsymbol{P}^{(k)}, \Theta(u) \right\rangle, \\ \boldsymbol{P}^{(k+1)} &= \operatorname{proj}_{\mathcal{C}^{\star}} \left(\boldsymbol{P}^{(k)} + \epsilon^{(k)} \Theta(\boldsymbol{U}^{(k+1)}) \right). \end{aligned}$

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

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Convergence Theorem

Make the following assumptions.

- H1 U^{ad} is a nonempty closed convex subset of a Hilbert space \mathbb{U} , and C is a closed convex salient cone of a Hilbert space \mathbb{V} .
- **H2** $j : \mathbb{U} \times \mathbb{W} \to \mathbb{R}$ is a normal integrand, and $\mathbb{E}(j(u, W))$ exists for all $u \in U^{\mathrm{ad}}$.
- **H3** $j(\cdot, w) : \mathbb{U} \to \mathbb{R}$ is a proper convex l.s.c. differentiable function with linearly bounded gradients (LBG), for all $w \in \mathbb{W}$.
- H4 J is strictly convex, Lipschitz and coercive on U^{ad} .
- **H5** Θ is *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.

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Convergence Theorem

Then the following conclusions hold true.

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

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Sketch of Proof

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

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

O Select a Lyapunov function.

 $\Lambda(u,p) = K(u^{\sharp}) - K(u) - \langle \nabla K(u), u^{\sharp} - u \rangle + \left\| p - p^{\sharp} \right\|^{2} / 2.$

Bound from above the variation of Λ. Using assumptions and writing optimality conditions, we get:

 $\mathbb{E}\left(\Lambda(\boldsymbol{U}^{(k+1)}, \boldsymbol{P}^{(k+1)}) \mid \mathcal{F}^{(k)}\right) \leq (1 + \alpha^{(k)})\Lambda(\boldsymbol{U}^{(k)}, \boldsymbol{P}^{(k)}) \\ + \beta^{(k)} - \epsilon^{(k)} \left(L(\boldsymbol{U}^{(k)}, \boldsymbol{p}^{\sharp}) - L(\boldsymbol{u}^{\sharp}, \boldsymbol{p}^{\sharp})\right) \,.$

Ore a convergence of the sequences.

Using the two lemmas, we obtain that $\{\Lambda(\boldsymbol{U}^{(k)}, \boldsymbol{P}^{(k)})\}_{k \in \mathbb{N}}$ almost surely converges to a finite random variable Λ^{∞} , and that $\{L(\boldsymbol{U}^{(k)}, p^{\sharp})\}_{k \in \mathbb{N}}$ almost surely converges to $L(u^{\sharp}, p^{\sharp})$. By a compactness argument and uniqueness of u^{\sharp} , the sequence $\{\boldsymbol{U}^{(k)}\}_{k \in \mathbb{N}}$ almost surely converges to u^{\sharp} .

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

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. The extension is obtained 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:

Regularized Stochastic APP algorithm with Constraints

$$\boldsymbol{U}^{(k+1)} = \underset{u \in U^{\mathrm{ad}}}{\mathrm{arg\,min}} \, \boldsymbol{K}(u) + \left\langle \epsilon^{(k)} \boldsymbol{G}^{(k)} - \nabla \boldsymbol{K}(\boldsymbol{U}^{(k)}) , u \right\rangle \\ + \epsilon^{(k)} \left\langle \operatorname{proj}_{\mathcal{C}^{\star}} \left(\boldsymbol{P}^{(k)} + c \Theta(\boldsymbol{U}^{(k)}) \right) , \Theta(u) \right\rangle$$

$$\boldsymbol{P}^{(k+1)} = \left(1 - \frac{\epsilon^{(k)}}{c}\right) \boldsymbol{P}^{(k)} + \frac{\epsilon^{(k)}}{c} \operatorname{proj}_{C^{\star}} \left(\boldsymbol{P}^{(k)} + c\Theta(\boldsymbol{U}^{(k+1)})\right).$$

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Convergence Theorem

Make the following assumptions.

- H1 U^{ad} is a nonempty closed convex subset of a Hilbert space \mathbb{U} , and *C* is a closed convex cone of another Hilbert space \mathbb{V} .
- **H2** $j : \mathbb{U} \times \mathbb{W} \to \mathbb{R}$ is a normal integrand, and $\mathbb{E}(j(u, \boldsymbol{W}))$ exists for all $u \in U^{\mathrm{ad}}$.
- **H3** $j(\cdot, w) : \mathbb{U} \to \mathbb{R}$ is a proper convex l.s.c. differentiable function with linearly bounded gradients (LBG), for all $w \in \mathbb{W}$.
- **H4** J is Lipschitz continuous, coercive on U^{ad} .
- **H5** Θ is *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.

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Convergence Theorem

Then the following conclusions hold true.

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

Sketch of Proof

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

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

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

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

and a substantial amount of technicalities...

Stochastic APP and Constraints in Expectation

We finally aim at solving the following stochastic optimization problem, in which the constraint corresponds to an expectation:

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

In the spirit of the stochastic gradient method, we use values of θ evaluated in realizations of \boldsymbol{W} rather than expected values of Θ . With $\boldsymbol{G}^{(k)} = \nabla_{\boldsymbol{u}} j(\boldsymbol{U}^{(k)}, \boldsymbol{W}^{(k+1)})$ and $\vartheta^{(k)} = \theta'_{\boldsymbol{u}}(\boldsymbol{U}^{(k)}, \boldsymbol{W}^{(k+1)})$, the extension of the APP method is:

Stochastic APP Algorithm with Expected Constraints

$$\begin{split} \boldsymbol{U}^{(k+1)} &\in \operatorname*{arg\,min}_{u \in U^{\mathrm{ad}}} \boldsymbol{K}(u) + \left\langle \boldsymbol{\epsilon}^{(k)} \boldsymbol{G}^{(k)} - \nabla \boldsymbol{K}(\boldsymbol{U}^{(k)}) \,, u \right\rangle \\ &+ \boldsymbol{\epsilon}^{(k)} \left\langle \boldsymbol{P}^{(k)} \,, \boldsymbol{\vartheta}^{(k)} \cdot u \right\rangle \,, \\ \boldsymbol{P}^{(k+1)} &= \operatorname{proj}_{\mathcal{C}^{\star}} \left(\boldsymbol{P}^{(k)} + \boldsymbol{\rho}^{(k)} \, \boldsymbol{\theta}(\boldsymbol{U}^{(k+1)}, \boldsymbol{W}^{(k+1)}) \right) \,. \end{split}$$

Constraints in Stochastic Optimization APP with Constraints in the Deterministic Setting APP with Constraints in the Stochastic Setting

Convergence Theorem and Proof

Long and intricate... See the lecture notes.