

Propriétés probabilistes dans les algorithmes d'optimisation sans et avec dérivées

Clément Royer - University of Wisconsin-Madison

Séminaire SPOC
Institut de Mathématiques de Bourgogne

12 avril 2017

L'aléatoire est de plus en plus présent en optimisation numérique.

Pour de multiples raisons :

- *Problèmes de grande taille* : Méthodes classiques trop coûteuses.
- *Calcul distribué* : Données stockées sur plusieurs machines/processeurs.
- *Applications* : Problèmes d'apprentissage.

L'aléatoire est de plus en plus présent en optimisation numérique.

Pour de multiples raisons :

- *Problèmes de grande taille* : Méthodes classiques trop coûteuses.
- *Calcul distribué* : Données stockées sur plusieurs machines/processeurs.
- *Applications* : Problèmes d'apprentissage.

Questions sur l'aléatoire

- Comment l'analyse des méthodes est-elle affectée ?
- Améliore-t-on les variantes déterministes ?
- Aléatoire en optimisation sans dérivées ?

Analyse de Complexité

- Etudier le **taux de convergence** d'un critère donné.
- **Borner** le comportement d'une méthode **au pire cas**.
- Avec aléatoire : **résultats en espérance/probabilité**.

Utilité de la complexité

- Quelles indications données par la complexité ?
- Quel lien avec la pratique ?
- Importance pour les **méthodes sans dérivées** ?

Trame

- 1 Introduire des aspects aléatoires dans des algorithmes sans dérivées.
- 2 Fournir des garanties théoriques (ex : complexité).
- 3 Comparer la complexité et le comportement numérique.

Trame

- 1 Introduire des aspects aléatoires dans des algorithmes sans dérivées.
- 2 Fournir des garanties théoriques (ex : complexité).
- 3 Comparer la complexité et le comportement numérique.

- Dans cet exposé : méthodes de **recherche directe**;
- Les résultats s'appliquent à d'autres méthodes, comme les **régions de confiance**.

- 1 Recherche directe déterministe
- 2 Recherche directe à base de descente probabiliste
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

- 1 Recherche directe déterministe
 - Optimisation sans dérivées
 - Recherche directe
- 2 Recherche directe à base de descente probabiliste
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

Soit le problème suivant :

$$\min_{x \in \mathbb{R}^n} f(x).$$

Hypothèses sur f

- f minorée, a priori non convexe.
- f de classe \mathcal{C}^1 , ∇f fonction lipschitzienne.

Soit le problème suivant :

$$\min_{x \in \mathbb{R}^n} f(x).$$

Hypothèses sur f

- f minorée, a priori non convexe.
- f de classe \mathcal{C}^1 , ∇f fonction lipschitzienne.

Optimisation différentiable

Depuis $x \in \mathbb{R}^n$, on peut décroître f dans la direction de $-\nabla f(x)$!

- Principe de base des méthodes de *premier ordre/gradient*.
- Objectif : converger vers un **point stationnaire d'ordre 1**:

$$\liminf_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0.$$

Le gradient de f existe mais **ne peut pas être utilisé en pratique.**

- *Code de simulation* : le gradient est trop coûteux.
- *Fonction f disponible sous forme de boîte noire* : pas de code pour la dérivée.
- *Différentiation automatique* inapplicable.

Exemples : météorologie, industrie pétrolière, médecine...

Le gradient de f existe mais **ne peut pas être utilisé en pratique.**

- *Code de simulation* : le gradient est trop coûteux.
- *Fonction f disponible sous forme de boîte noire* : pas de code pour la dérivée.
- *Différentiation automatique* inapplicable.

Exemples : météorologie, industrie pétrolière, médecine...

Mesure de performance : Nombre d'évaluations de l'objectif.

Méthodes sans dérivées déterministes

- Méthodes à modèles, comme les régions de confiance.
- Méthodes directionnelles, comme la recherche directe.

Introduction to Derivative-Free Optimization

A.R. Conn, K. Scheinberg, L.N. Vicente. (2009)

- Théorie de convergence (vers des optima locaux).
- Bornes de complexité/Vitesses de convergence activement étudiées.

Méthodes "DFO" stochastiques

- Ont pour but de trouver des **optima globaux**:
Ex) Stratégies Evolutionnaires, Algorithmes Génétiques.
- Souvent sans analogues déterministes.

- Cet exposé ne concerne pas les méthodes stochastiques;
- Nos algorithmes sont basés sur des éléments **probabilistes**.

Méthodes "DFO" avec propriétés probabilistes

- Développées à partir d'algorithmes déterministes.
- Bénéficient de **garanties théoriques grâce à cela**.
- Le côté aléatoire améliore la performance.

- 1 Recherche directe déterministe
 - Optimisation sans dérivées
 - Recherche directe
- 2 Recherche directe à base de descente probabiliste
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

- Variantes sans dérivées des méthodes de gradient.
 - Introduites vers 1960, théorie de convergence vers 1990.
 - **Simple à implémenter, fort potentiel de parallélisme.**
- **Optimization by direct search: new perspectives on some classical and modern methods.**
Kolda, Lewis and Torczon (*SIAM Review*, 2003).

1 **Initialisation:** $x_0 \in \mathbb{R}^n, \alpha_0 > 0, 0 < \theta < 1 \leq \gamma$.

2 **Pour** $k = 0, 1, 2, \dots$

- Choisir un ensemble D_k de r vecteurs.
- Si il existe $d_k \in D_k$ tel que

$$f(x_k + \alpha_k d_k) < f(x_k) - \alpha_k^2,$$

alors (*k réussie*) poser $x_{k+1} := x_k + \alpha_k d_k$ et $\alpha_{k+1} := \gamma \alpha_k$.

- Sinon (*k non réussie*) poser $x_{k+1} := x_k$ et $\alpha_{k+1} := \theta \alpha_k$.

1 **Initialisation:** $x_0 \in \mathbb{R}^n, \alpha_0 > 0, 0 < \theta < 1 \leq \gamma$.

2 **Pour** $k = 0, 1, 2, \dots$

- Choisir un ensemble D_k de r vecteurs.
- Si il existe $d_k \in D_k$ tel que

$$f(x_k + \alpha_k d_k) < f(x_k) - \alpha_k^2,$$

alors (*k réussie*) poser $x_{k+1} := x_k + \alpha_k d_k$ et $\alpha_{k+1} := \gamma \alpha_k$.

- Sinon (*k non réussie*) poser $x_{k+1} := x_k$ et $\alpha_{k+1} := \theta \alpha_k$.

Choisir les directions de sondage

On cherche à choisir les ensembles des directions/de sondage D_k pour garantir la convergence de l'algorithme.

Choisir les directions de sondage

On cherche à choisir les ensembles des directions/de sondage D_k pour garantir la convergence de l'algorithme.

Une mesure de qualité

Pour un ensemble de vecteurs D , la mesure cosinus de D est donnée par

$$\text{cm}(D) = \min_{v \in \mathbb{R}^n \setminus \{0\}} \max_{d \in D} \frac{d^\top v}{\|d\| \|v\|}.$$

On cherche à choisir les ensembles des directions/de sondage D_k pour garantir la convergence de l'algorithme.

Une mesure de qualité

Pour un ensemble de vecteurs D , la mesure cosinus de D est donnée par

$$\text{cm}(D) = \min_{v \in \mathbb{R}^n \setminus \{0\}} \max_{d \in D} \frac{d^\top v}{\|d\| \|v\|}.$$

- Si $\text{cm}(D) > 0$, pour tout v il existe $d \in D$ tel que (d, v) est un angle aigu.
- Avec $v = -\nabla f(x) \neq 0$, D contient une direction de descente pour f en x .

Comment garantir $\text{cm}(D) > 0$?

Ensemble de générateurs positifs (PSS)

D est un PSS s'il génère \mathbb{R}^n par combinaisons linéaires à coefficients positifs ou nuls.

- D est un PSS $\Leftrightarrow \text{cm}(D) > 0$.
- Un PSS contient au moins $n + 1$ vecteurs.

Comment garantir $\text{cm}(D) > 0$?

Ensemble de générateurs positifs (PSS)

D est un PSS s'il génère \mathbb{R}^n par combinaisons linéaires à coefficients positifs ou nuls.

- D est un PSS $\Leftrightarrow \text{cm}(D) > 0$.
- Un PSS contient au moins $n + 1$ vecteurs.

Exemple

$D_{\oplus} = \{e_1, \dots, e_n, -e_1, \dots, -e_n\}$ est un PSS avec

$$\text{cm}(D_{\oplus}) = \frac{1}{\sqrt{n}}.$$

Lemma

Si l'itération k n'est pas réussie et $\text{cm}(D_k) \geq \kappa > 0$,

$$\kappa \|\nabla f(x_k)\| \leq \mathcal{O}(\alpha_k).$$

Lemma

Indépendamment de $\{D_k\}$,

$$\lim_{k \rightarrow \infty} \alpha_k = 0.$$

Lemma

Si l'itération k n'est pas réussie et $\text{cm}(D_k) \geq \kappa > 0$,

$$\kappa \|\nabla f(x_k)\| \leq \mathcal{O}(\alpha_k).$$

Lemma

Indépendamment de $\{D_k\}$,

$$\lim_{k \rightarrow \infty} \alpha_k = 0.$$

Théorème de convergence

Si $\forall k, \text{cm}(D_k) \geq \kappa$,

$$\liminf_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0.$$

Théorème de complexité

Soient $\epsilon \in (0, 1)$ et N_ϵ le nombre d'appels à f nécessaires pour satisfaire $\inf_{0 \leq l \leq k} \|\nabla f(x_l)\| < \epsilon$. On a

$$N_\epsilon \leq \mathcal{O}(r(\kappa\epsilon)^{-2}).$$

En choisissant $D_k = D_\oplus$, on a $\kappa = 1/\sqrt{n}$, $r = 2n$, la borne devient

$$N_\epsilon \leq \mathcal{O}(n^2 \epsilon^{-2}).$$

- 1 Recherche directe déterministe
- 2 Recherche directe à base de descente probabiliste
 - Descente probabiliste
 - Convergence et complexité
 - Descente probabiliste en pratique
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

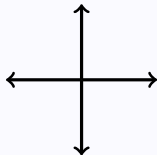
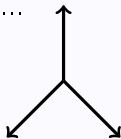
Idée (Gratton et Vicente, 2013)

Utiliser des vecteurs tirés aléatoirement et indépendamment, typiquement moins que $n + 1$!

Idée (Gratton et Vicente, 2013)

Utiliser des vecteurs tirés aléatoirement et indépendamment, typiquement moins que $n + 1$!

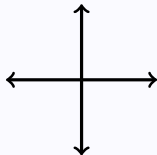
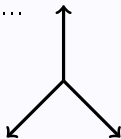
From PSS...



Idée (Gratton et Vicente, 2013)

Utiliser des vecteurs tirés aléatoirement et indépendamment, typiquement moins que $n + 1$!

From PSS...



...to random sets

Motivation numérique

- Test de convergence: $f(x_k) < f_{\text{low}} + 10^{-3} (f(x_0) - f_{\text{low}})$;
- Budget: 2000 n appels à f max.

Problème	D_{\oplus}	$Q D_{\oplus}$	$2n$	$n+1$	$n/2$	2	1
	Déterministe		Probabiliste				
arglina	3.42	16.67	10.30	6.01	3.21	1.00	–
arglinb	20.50	11.38	7.38	2.81	2.35	1.00	2.04
broydn3d	4.33	11.22	6.54	3.59	2.04	1.00	–
dqrtc	7.16	19.50	9.10	4.56	2.77	1.00	–
engval1	10.53	23.96	11.90	6.48	3.55	1.00	2.08
freuroth	56.00	1.33	1.00	1.67	1.33	1.00	4.00
integreq	16.04	18.85	12.44	6.76	3.52	1.00	–
nondquar	6.90	17.36	7.56	4.23	2.76	1.00	–
sinquad	–	2.12	1.31	1.00	1.60	1.23	–
vardim	1.00	3.30	1.80	2.40	2.30	1.80	4.30

Table: Ratio du nombre d'appels à f (moyenne sur 10 réalisations, taille $n = 40$)

Notations probabilistes

- Ensembles/Directions de sondage : $D_k = \mathfrak{D}_k(\omega)$, $d_k = \mathfrak{d}_k(\omega)$;
- Itérés : $x_k = X_k(\omega)$;
- Longueurs de pas : $\alpha_k = A_k(\omega)$.

① **Initialisation:** $x_0 \in \mathbb{R}^n$, $\alpha_0 > 0$, $0 < \theta < 1 \leq \gamma$.

② **Pour** $k = 0, 1, 2, \dots$

- Choisir un ensemble \mathfrak{D}_k de r vecteurs **tirés indépendamment au hasard**.
- Si il existe $\mathfrak{d}_k \in \mathfrak{D}_k$ tel que

$$f(X_k + \alpha_k \mathfrak{d}_k) < f(X_k) - \alpha_k^2,$$

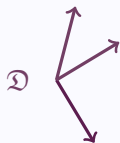
alors (*k réussie*) poser $X_{k+1} := X_k + \alpha_k \mathfrak{d}_k$ et $A_{k+1} := \gamma A_k$.

- Sinon (*k non réussie*) poser $X_{k+1} := X_k$ et $A_{k+1} := \theta A_k$.

- 1 Recherche directe déterministe
- 2 Recherche directe à base de descente probabiliste
 - Descente probabiliste
 - **Convergence et complexité**
 - Descente probabiliste en pratique
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

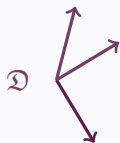
Qu'est-ce qu'un bon ensemble de sondage ?

② n'est pas un PSS...

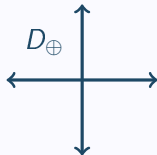


Qu'est-ce qu'un bon ensemble de sondage ?

\mathcal{D} n'est pas un PSS...

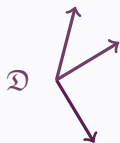


... D_{\oplus} si...

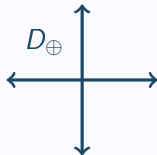


Qu'est-ce qu'un bon ensemble de sondage ?

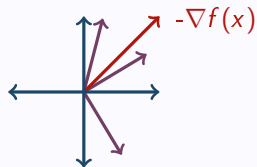
\mathcal{D} n'est pas un PSS...



... D_{\oplus} si...



... $-\nabla f(x)$ plus proche de \mathcal{D} !



Etre proche de l'opposé du gradient : un gage de qualité ?

Propriétés dans le cas déterministe

- On a requis

$$\text{cm}(D_k) = \min_{v \neq 0} \max_{d \in D_k} \frac{d^\top v}{\|d\| \|v\|} \geq \kappa.$$

- Il suffirait d'avoir

$$\text{cm}(D_k, -\nabla f(x_k)) = \max_{d \in D_k} \frac{d^\top [-\nabla f(x_k)]}{\|d\| \|\nabla f(x_k)\|} \geq \kappa.$$

- Avec de l'aléatoire, la seconde propriété peut être vraie **en probabilité**.
- Quels sont les bons **outils probabilistes** pour exprimer cela ?

Plusieurs types de résultats

Déterministe/Pour toute réalisation



Avec probabilité 1/Presque sûr



Avec une certaine probabilité

Sous-martingale

Une **sous-martingale** est une suite de variables aléatoires $\{V_k\}$ telle que $\mathbb{E}[|V_k|] < \infty$ et

$$\mathbb{E}(V_k | V_0, V_1, \dots, V_{k-1}) \geq V_{k-1}.$$

- On cherche à étudier

$$\mathbb{P}(\text{cm}(\mathcal{D}_k, -\nabla f(X_k)) \geq \kappa).$$

où X_k dépend de $\mathcal{D}_0, \dots, \mathcal{D}_{k-1}$ mais pas de \mathcal{D}_k .

- On va utiliser les probabilités conditionnelles/le conditionnement au passé.

- On cherche à étudier

$$\mathbb{P}(\text{cm}(\mathfrak{D}_k, -\nabla f(X_k)) \geq \kappa).$$

où X_k dépend de $\mathfrak{D}_0, \dots, \mathfrak{D}_{k-1}$ mais pas de \mathfrak{D}_k .

- On va utiliser les probabilités conditionnelles/le conditionnement au passé.

Propriété de descente probabiliste

Une suite d'ensembles aléatoires $\{\mathfrak{D}_k\}$ est dite à descente (ρ, κ) si:

$$\begin{aligned} \mathbb{P}(\text{cm}(\mathfrak{D}_0, -\nabla f(x_0)) \geq \kappa) &\geq \rho \\ \forall k \geq 1, \quad \mathbb{P}(\text{cm}(\mathfrak{D}_k, -\nabla f(X_k)) \geq \kappa \mid \mathfrak{D}_0, \dots, \mathfrak{D}_{k-1}) &\geq \rho, \end{aligned}$$

Lemma

Pour toute réalisation $\{\alpha_k\}$ de $\{A_k\}$, indépendamment de $\{\mathcal{D}_k\}$,

$$\lim_{k \rightarrow \infty} \alpha_k = 0.$$

Lemma

Si l'itération k n'est pas réussie,

$$\{\text{cm}(\mathcal{D}_k, -\nabla f(X_k)) \geq \kappa\} \subset \{\kappa \|\nabla f(X_k)\| \leq \mathcal{O}(A_k)\}.$$

Il s'agit de prouver que $\{\text{cm}(\mathcal{D}_k, -\nabla f(X_k)) \geq \kappa\}$ se produit suffisamment souvent.

Résultats de convergence (2)

Soit $\{\mathfrak{D}_k\}$ à descente (p, κ) et $Z_k = \mathbf{1}(\text{cm}(\mathfrak{D}_k, -\nabla f(X_k)) \geq \kappa)$.

Proposition

Soit

$$S_k = \sum_{i=0}^{k-1} [Z_i - p_0], \quad p_0 = \frac{\ln \theta}{\ln(\theta/\gamma)}.$$

- 1 $\{\liminf_k \|\nabla f(X_k)\| > 0\} \subset \{S_k \rightarrow -\infty\}$.
- 2 Si $p > p_0$, $\{S_k\}$ est une sous-martingale avec $\mathbb{P}(\limsup S_k = \infty) = 1$.

Résultats de convergence (2)

Soit $\{\mathcal{D}_k\}$ à descente (p, κ) et $Z_k = \mathbf{1}(\text{cm}(\mathcal{D}_k, -\nabla f(X_k)) \geq \kappa)$.

Proposition

Soit

$$S_k = \sum_{i=0}^{k-1} [Z_i - p_0], \quad p_0 = \frac{\ln \theta}{\ln(\theta/\gamma)}.$$

- 1 $\{\liminf_k \|\nabla f(X_k)\| > 0\} \subset \{S_k \rightarrow -\infty\}$.
- 2 Si $p > p_0$, $\{S_k\}$ est une sous-martingale avec $\mathbb{P}(\limsup S_k = \infty) = 1$.

Théorème : convergence presque sûre

Si $\{\mathcal{D}_k\}$ est à descente (p, κ) avec $p > p_0$, on a

$$\mathbb{P}\left(\liminf_{k \rightarrow \infty} \|\nabla f(X_k)\| = 0\right) = 1.$$

Idée intuitive

Soient $G_k = \nabla f(X_k)$ et $Z_k = \mathbf{1}(\text{cm}(\mathfrak{D}_k, -G_k) \geq \kappa)$.

- Si $Z_k = 1$ et k réussie, on a $\kappa \|G_k\| < \mathcal{O}(A_k)$...

Idée intuitive

Soient $G_k = \nabla f(X_k)$ et $Z_k = \mathbf{1}(\text{cm}(\mathfrak{D}_k, -G_k) \geq \kappa)$.

- Si $Z_k = 1$ et k réussie, on a $\kappa \|G_k\| < \mathcal{O}(A_k)$...
- ... A_k tend vers 0...
- ...donc si $\inf_{0 \leq l \leq k} \|G_l\|$ est grand, $\sum_{l=0}^k Z_l$ doit être faible.

Idée intuitive

Soient $G_k = \nabla f(X_k)$ et $Z_k = \mathbf{1}(\text{cm}(\mathfrak{D}_k, -G_k) \geq \kappa)$.

- Si $Z_k = 1$ et k réussie, on a $\kappa \|G_k\| < \mathcal{O}(A_k)$...
- ... A_k tend vers 0...
- ...donc si $\inf_{0 \leq l \leq k} \|G_l\|$ est grand, $\sum_{l=0}^k Z_l$ doit être faible.

Une borne utile

Pour chaque réalisation de l'algorithme,

$$\sum_{l=0}^k z_l \leq \mathcal{O}\left(\frac{1}{\kappa^2 \|\tilde{g}_k\|^2}\right) + p_0 k,$$

où $\|\tilde{g}_k\| = \inf_{0 \leq l \leq k} \|g_l\|$.

Rappel : $Z_l = \mathbf{1}(\text{cm}(\mathfrak{D}_l, -\nabla f(X_l)) \geq \kappa)$.

Argument d'inclusion

$$\left\{ \inf_{0 \leq l \leq k} \|\nabla f(X_k)\| \geq \epsilon \right\} \subset \left\{ \sum_{l=0}^k Z_l \leq \lambda k \right\}$$

avec $\lambda = \mathcal{O}\left(\frac{1}{k \kappa^2 \epsilon^{-2}}\right) + p_0$.

Borne de type Chernoff

Pour tout $\lambda \in (0, p)$,

$$\mathbb{P}\left(\sum_{l=0}^{k-1} Z_l \leq \lambda k\right) \leq \exp\left[-\frac{(p-\lambda)^2}{2p}k\right].$$

Théorème : Complexité probabiliste

Soient $\{\mathfrak{D}_k\}$ à descente (ρ, κ) , $\epsilon \in (0, 1)$ et N_ϵ le nombre d'appels à f nécessaires pour obtenir $\inf_{0 \leq l \leq k} \|\nabla f(X_l)\| \leq \epsilon$. Alors

$$\mathbb{P} \left(N_\epsilon \leq \mathcal{O} \left(\frac{r(\kappa\epsilon)^{-2}}{\rho - \rho_0} \right) \right) \geq 1 - \exp \left(-\mathcal{O} \left(\frac{\rho - \rho_0}{\rho} (\kappa\epsilon)^{-2} \right) \right).$$

Théorème : Complexité probabiliste

Soient $\{\mathfrak{D}_k\}$ à descente (p, κ) , $\epsilon \in (0, 1)$ et N_ϵ le nombre d'appels à f nécessaires pour obtenir $\inf_{0 \leq l \leq k} \|\nabla f(X_l)\| \leq \epsilon$. Alors

$$\mathbb{P} \left(N_\epsilon \leq \mathcal{O} \left(\frac{r (\kappa \epsilon)^{-2}}{p - p_0} \right) \right) \geq 1 - \exp \left(-\mathcal{O} \left(\frac{p - p_0}{p} (\kappa \epsilon)^{-2} \right) \right).$$

- Déterministe : $\mathcal{O}(n^2 \epsilon^{-2})$.
- Probabiliste : $\mathcal{O}(r n \epsilon^{-2})$ en probabilité
 $\Rightarrow \mathcal{O}(n \epsilon^{-2})$ lorsque $r = 2$!
- Serait-on meilleur (en probabilité) avec moins de directions ?

- 1 Recherche directe déterministe
- 2 Recherche directe à base de descente probabiliste
 - Descente probabiliste
 - Convergence et complexité
 - Descente probabiliste en pratique
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

Construire une suite à descente (p, κ)

On cherche à avoir

$$p > p_0 = \frac{\ln(\theta)}{\ln(\theta/\gamma)}$$

avec le plus petit $r = |\mathfrak{D}_k|$ possible.

Une technique : génération uniforme sur la sphère unité

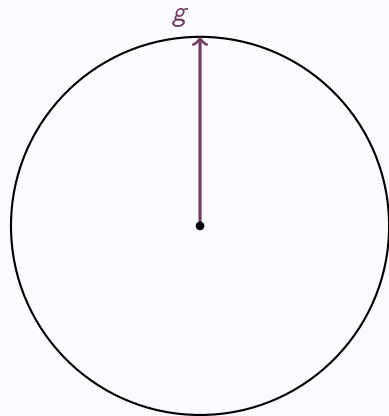
Si

$$r > \log_2 \left(1 - \frac{\ln \theta}{\ln \gamma} \right),$$

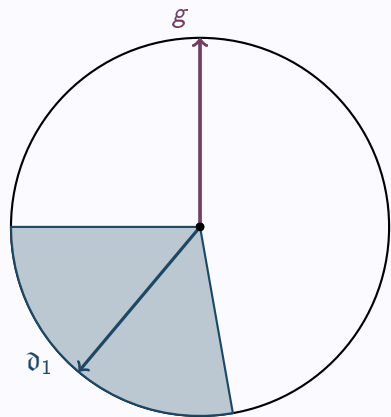
il existe (p, τ) indépendants de n tels que la suite \mathfrak{D}_k est à descente $(p, \tau/\sqrt{n})$ et $p > p_0$.

Si $\gamma = \theta^{-1} = 2$, choisir $r \geq 2$ suffit à garantir $p > \frac{1}{2}$.

Deux directions uniformes suffisent, pas une

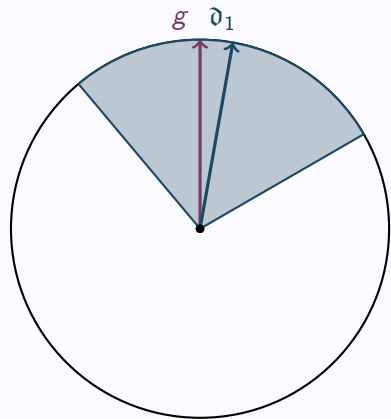


Deux directions uniformes suffisent, pas une



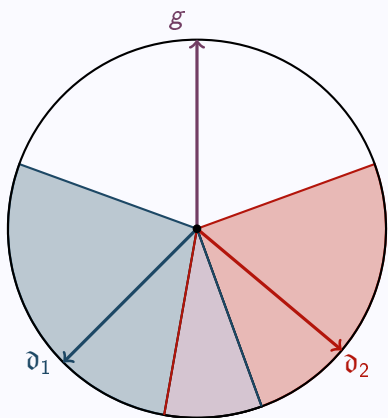
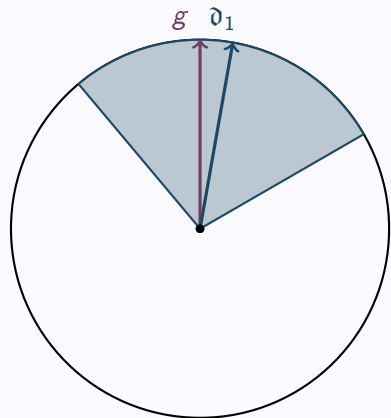
$$d_1 \sim \mathcal{U}(S^1) \Rightarrow \forall \kappa \in (0, 1), \quad \mathbb{P}(\text{cm}(d_1, g) = d_1^\top g \geq \kappa) < 1/2.$$

Deux directions uniformes suffisent, pas une



$$\vartheta_1 \sim \mathcal{U}(\mathbb{S}^1) \Rightarrow \forall \kappa \in (0, 1), \quad \mathbb{P} \left(\text{cm}(\vartheta_1, g) = \vartheta_1^\top g \geq \kappa \right) < 1/2.$$

Deux directions uniformes suffisent, pas une



$$\vartheta_1 \sim \mathcal{U}(\mathbb{S}^1) \Rightarrow \forall \kappa \in (0, 1), \quad \mathbb{P}(\text{cm}(\vartheta_1, g) = \vartheta_1^\top g \geq \kappa) < 1/2.$$

$$\vartheta_1, \vartheta_2 \sim \mathcal{U}(\mathbb{S}^1) \Rightarrow \exists \kappa^* \in (0, 1), \quad \mathbb{P}(\text{cm}(\{\vartheta_1, \vartheta_2\}, g) \geq \kappa^*) > 1/2.$$

- 1 Recherche directe déterministe
- 2 Recherche directe à base de descente probabiliste
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

Contraintes linéaires d'égalité

$$\begin{cases} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & Ax = b. \end{cases}$$

- Equivalent au problème **sans contraintes** $\min_{\tilde{x} \in \mathbb{R}^{n-m}} f(x_0 + W\tilde{x})$ où $W \in \mathbb{R}^{n \times (n-m)}$ est une base pour $\text{null}(A)$ et $Ax_0 = b$.
- Les résultats déterministes et probabilistes restent valables !

Contraintes linéaires d'égalité

$$\begin{cases} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & Ax = b. \end{cases}$$

- Equivalent au problème **sans contraintes** $\min_{\tilde{x} \in \mathbb{R}^{n-m}} f(x_0 + W\tilde{x})$ où $W \in \mathbb{R}^{n \times (n-m)}$ est une base pour $\text{null}(A)$ et $Ax_0 = b$.
- Les résultats déterministes et probabilistes restent valables !

Contraintes d'intervalle

$$\begin{cases} \min_{x \in \mathbb{R}^n} & f(x) \\ \text{s.t.} & l \leq x \leq u. \end{cases}$$

- **En pratique (déterministe)** : Utiliser $D_{\oplus} = \{e_1, \dots, e_n, -e_1, \dots, -e_n\}$ permet de converger et de se déplacer **parallèlement aux contraintes**.

1 **Initialisation:** $x_0 \in \mathbb{R}^n, \alpha_0 > 0, 0 < \theta < 1 \leq \gamma$.

2 **Pour** $k = 0, 1, 2, \dots$

- Choisir un ensemble D_k d'au plus r vecteurs.
- Si il existe $d_k \in D_k$ tel que $x_k + \alpha_k d_k$ **est admissible** and

$$f(x_k + \alpha_k d_k) < f(x_k) - \alpha_k^2,$$

alors (*k réussie*) poser $x_{k+1} := x_k + \alpha_k d_k$ et $\alpha_{k+1} := \gamma \alpha_k$.

- Sinon (*k non réussie*) poser $x_{k+1} := x_k$ et $\alpha_{k+1} := \theta \alpha_k$.

- Domaine admissible : $\mathcal{F} = \{l \leq x \leq u\}$.

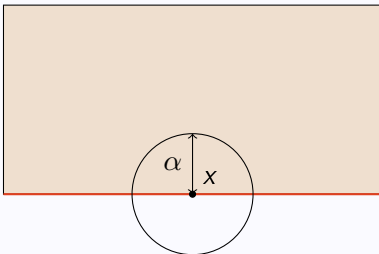
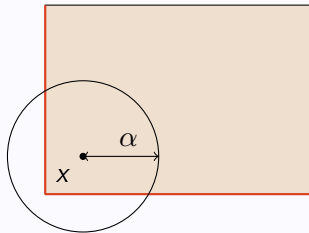
Contraintes proches

Les ensembles

$$I_u(x, \alpha) = \{i : |u_i - [x]_i| \leq \alpha\}$$

$$I_l(x, \alpha) = \{i : |l_i - [x]_i| \leq \alpha\}$$

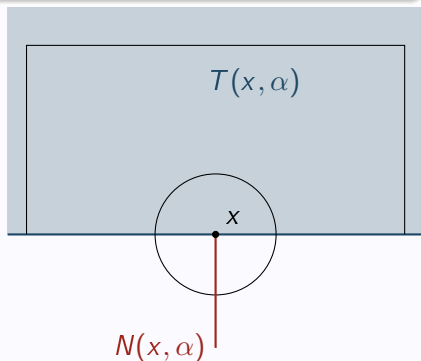
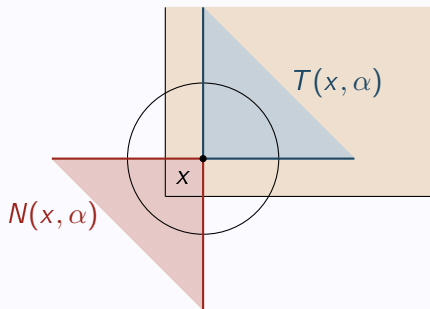
définissent les **contraintes proches** en $x \in \mathcal{F}$ pour $\alpha > 0$.



- **Cône normal approché** $N(x, \alpha)$: généré positivement par

$$\{e_i\}_{i \in I_u(x, \alpha)} \cup \{-e_i\}_{i \in I_l(x, \alpha)}.$$

- **Cône tangent** $T(x, \alpha)$: cône polaire de $N(x, \alpha)$.



- Rappel : la mesure cosinus permet d'identifier les directions de descente

$$\text{cm}(D, -\nabla f(x)) = \max_{d \in D} \frac{d^\top [-\nabla f(x)]}{\|d\| \|-\nabla f(x)\|}.$$

Descente réalisable

D est un ensemble à descente κ -admissible pour $T(x, \alpha)$ si $D \subset T(x, \alpha)$ et

$$\text{cm}_{T(x, \alpha)}(D, -\nabla f(x)) = \max_{d \in D} \frac{d^\top [-\nabla f(x)]}{\|d\| \|P_{T(x, \alpha)}[-\nabla f(x)]\|} \geq \kappa.$$

- Rappel : la mesure cosinus permet d'identifier les directions de descente

$$\text{cm}(D, -\nabla f(x)) = \max_{d \in D} \frac{d^\top [-\nabla f(x)]}{\|d\| \|-\nabla f(x)\|}.$$

Descente réalisable

D est un ensemble à descente κ -admissible pour $T(x, \alpha)$ si $D \subset T(x, \alpha)$ et

$$\text{cm}_{T(x, \alpha)}(D, -\nabla f(x)) = \max_{d \in D} \frac{d^\top [-\nabla f(x)]}{\|d\| \|P_{T(x, \alpha)}[-\nabla f(x)]\|} \geq \kappa.$$

- Avec des ensembles à descente κ -admissible : convergence + bornes sur la complexité (analyse similaire au cas sans contraintes).
- $D_{\oplus} \cap T(x, \alpha)$ est toujours à descente $\frac{1}{\sqrt{n}}$ -admissible.

Définition

Une suite $\{\mathfrak{D}_k\}$ est à descente (probabiliste) (p, κ) -admissible si

$$\begin{aligned} \mathbb{P}(\text{cm}_{T_0}(\mathfrak{D}_0, -\nabla f(x_0)) \geq \kappa) &\geq p \\ \forall k \geq 1, \quad \mathbb{P}(\text{cm}_{T_k}(\mathfrak{D}_k, -\nabla f(X_k)) \geq \kappa \mid \mathfrak{D}_0, \dots, \mathfrak{D}_{k-1}) &\geq p. \end{aligned}$$

avec $T_k = T(X_k, A_k)$.

Garanties théoriques

Si $\{\mathfrak{D}_k\}$ est à descente (p, κ) -admissible avec $p > p_0$,

- **Convergence presque sûre vers un point stationnaire;**
- **Borne de complexité en probabilité :**

$$\mathbb{P}\left(N_\epsilon \leq \mathcal{O}\left(\frac{r(\kappa\epsilon)^{-2}}{p - p_0}\right)\right) \geq 1 - \exp\left(-\mathcal{O}\left(\frac{p - p_0}{p}(\kappa\epsilon)^{-2}\right)\right).$$

Principales difficultés

- Définir des ensembles à descente probabiliste admissible.
- Estimer r et κ .
- Utiliser **moins de directions** que dans le cas déterministe ?

Principales difficultés

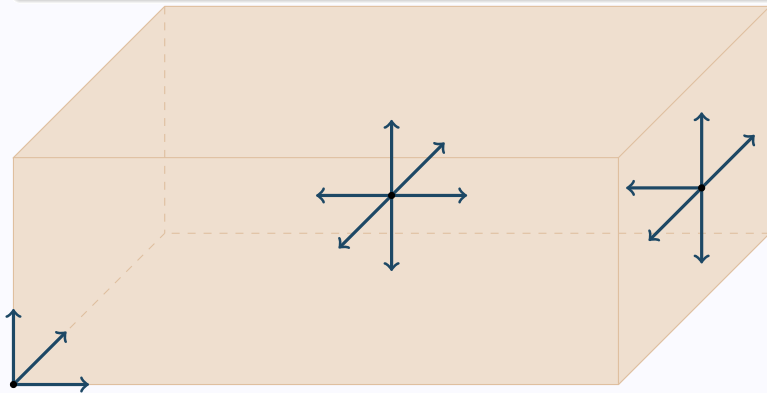
- Définir des ensembles à descente probabiliste admissible.
- Estimer r et κ .
- Utiliser **moins de directions** que dans le cas déterministe ?

Nos techniques

- Basées sur les générateurs du cône tangent (le choix déterministe);
- Directions aléatoires mais **admissibles**;
- **Au pire aussi coûteuses que la stratégie déterministe.**

Echantillonnage aléatoire parmi les générateurs

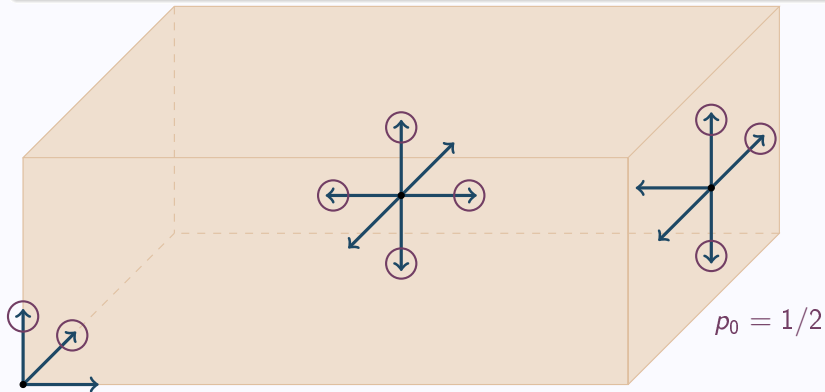
- 1 Calcul d'un ensemble déterministe V_k générant T_k ;



Premier choix de directions

Echantillonnage aléatoire parmi les générateurs

- 1 Calcul d'un ensemble **déterministe** V_k générant T_k ;
- 2 Tirer **au hasard** $\mathcal{D}_k \subset V_k$ de taille $> |V_k|p_0$;
- 3 $\{\mathcal{D}_k\}$ est à descente (p, κ) -admissible avec $p > p_0$.



Principe

- Cas sans contraintes : besoin de peu de directions;
- Idem avec contraintes d'égalité : problème sans contraintes dans le noyau de A ;
- Bénéfique d'exploiter les sous-espaces non contraints ?

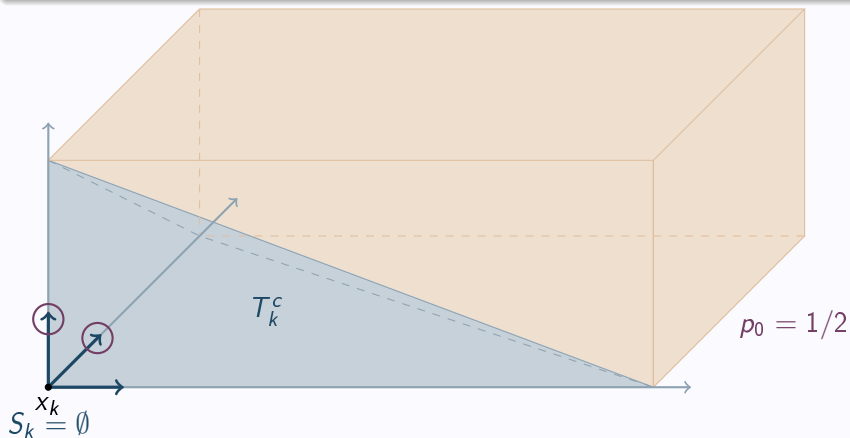
Lemme

Soit S_k un sous-espace inclus dans un cône T_k . Alors on a $T_k = S_k + T_k^c$, où T_k^c est un cône inclus dans S_k^\perp .

Second choix de directions

Deux types de vecteurs

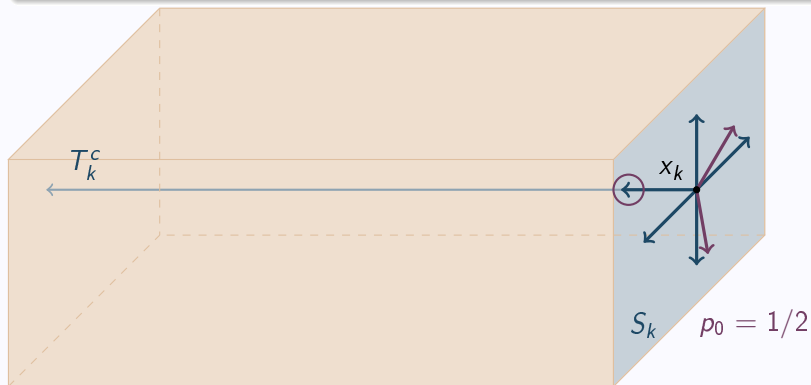
- Sous-espace S_k : Directions aléatoires;
- Complément T_k^c : Sous-ensemble aléatoire des générateurs.



Second choix de directions

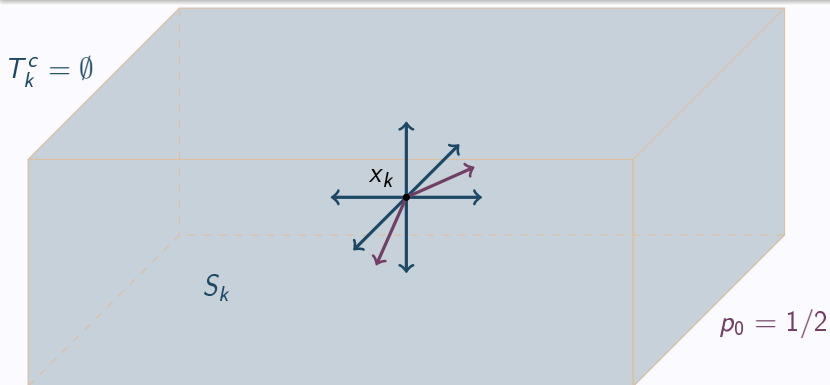
Deux types de vecteurs

- Sous-espace S_k : Directions aléatoires;
- Complément T_k^c : Sous-ensemble aléatoire des générateurs.



Deux types de vecteurs

- Sous-espace S_k : Directions aléatoires;
- Complément T_k^c : Sous-ensemble aléatoire des générateurs.



- Borne générale : $\mathcal{O}(r\kappa^{-2}\epsilon^{-2})$.

Comparaison - Egalités linéaires

Méthode	r	κ	Borne
Déterm.	$2(n - m)$	$\frac{1}{\sqrt{n-m}}$	$\mathcal{O}((n - m)^2\epsilon^{-2})$
Proba. 1	$\mathcal{O}(2(n - m)p_0)$	$\frac{1}{\sqrt{n-m}}$	$\mathcal{O}((n - m)^2\epsilon^{-2})$
Proba. 2 (subspace)	$\mathcal{O}(1)$	$\frac{\tau}{\sqrt{n-m}}$	$\mathcal{O}((n - m)\epsilon^{-2})$

Comparaison - Contraintes d'intervalle sur $n_b < n$ variables seulement

Méthode	r	κ	Borne
Déterm.	$2n$	$\frac{1}{\sqrt{n}}$	$\mathcal{O}(n^2\epsilon^{-2})$
Proba. 1	$\mathcal{O}(2np_0)$	$\frac{1}{\sqrt{n}}$	$\mathcal{O}(n^2\epsilon^{-2})$
Proba. 2 (subspace)	$\mathcal{O}(1) + \mathcal{O}(n_b p_0)$	$\frac{1}{\sqrt{n}}$	$\mathcal{O}(n n_b \epsilon^{-2})$

- Comparaison avec le solveur patternsearch de MATLAB.

Quatre méthodes

Nom	Directions dans $T(x_k, \alpha_k) = T_k = S_k + T_k^c$	Garanties
dspfd-0	$D_{\oplus} \cap T_k$, ordre aléatoire	Déterm.
dspfd-1	Tirage aléatoire dans $D_{\oplus} \cap T_k$	Proba.
dspfd-2	Vecteurs dans S_k /tirage dans $D_{\oplus} \cap T_k^c$	Proba.
matlab	$D_{\oplus} \cap T(x_k, t\alpha_k)$, $t \in (0, 1)$	Déterm.

Profils de performance

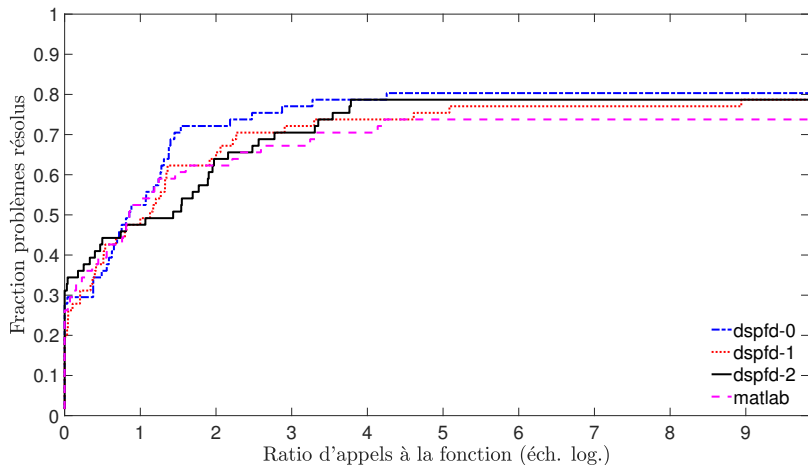
- Critère : # d'appels à f (budget de $2000n$) pour satisfaire

$$f(x_k) - f_{best} < 10^{-3}(f(x_0) - f_{best}).$$

- Problèmes de la bibliothèque CUTEst.

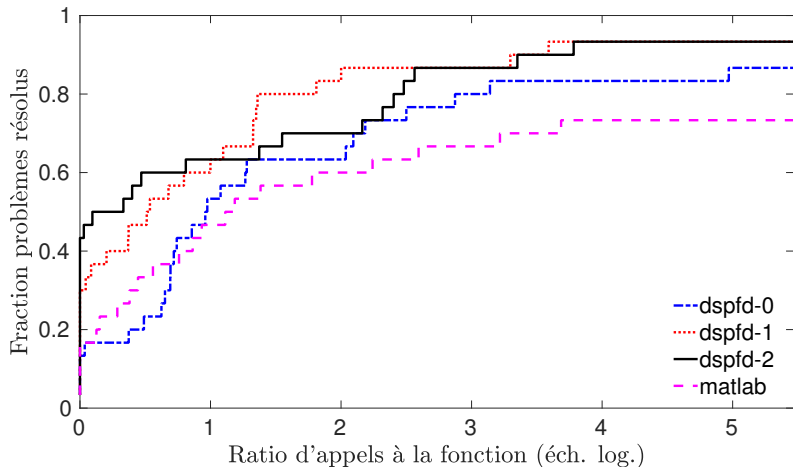
Profils avec contraintes d'intervalles

- 63 problèmes avec contraintes d'intervalles, de petites tailles : $2 \leq n \leq 20$.



Profils avec contraintes d'intervalles (2)

- 31 problèmes avec contraintes d'intervalles, **de tailles plus grandes** : $20 \leq n \leq 52$.



- 1 Recherche directe déterministe
- 2 Recherche directe à base de descente probabiliste
- 3 Extension aux problèmes avec contraintes linéaires
- 4 Propriétés probabilistes en optimisation avec dérivées

$$\min_{x \in \mathbb{R}^n} f(x)$$

- f de classe \mathcal{C}^2 ,
- f en général non convexe.

Du premier ordre au second dans les algorithmes

- L'accès à la matrice Hessienne est coûteux...
- ...tout comme l'algèbre linéaire associée:
 - Vecteurs propres;
 - Systèmes linéaires.

Peut-on s'aider de propriétés probabilistes ?

Analyse probabiliste de méthodes d'algèbre linéaire

Avec un point initial aléatoire...

- La puissance itérée trouve un vecteur propre à ϵ près en $\mathcal{O}(\epsilon^{-1})$ itérations;
- La méthode de Lanczos trouve un vecteur propre à ϵ près en $\mathcal{O}(\epsilon^{-1/2})$ itérations.

...avec une forte probabilité.

Utilité

- En lien avec les méthodes de premier ordre;
- Pour problèmes convexes et non convexes.

Points clés

- Echapper aux points selles;
- Détecter les **valeur propres négatives** de la matrice Hessienne;
- Utiliser des **directions de courbure négative**;

Les meilleures méthodes en terme de complexité garantissent la convergence à l'ordre deux !

- Revisiter les méthodes efficaces en pratique;
- Via leur analyse de complexité;
- Incorporer des progrès récents.

Blocs de base

- Méthodes de Newton-Krylov (ex : gradient conjugué):
Basées sur des produits matrice/vecteur.
- Recherches linéaires;
- Analyse probabiliste.

Problèmes sans contraintes

- Convergence via propriétés probabilistes.
- Moins d'évaluations en théorie **et en pratique**.

Direct search based on probabilistic descent. Gratton, Royer, Vicente and Zhang, *SIAM J. Optim.*, 2015.

Problèmes avec contraintes linéaires

- Descente probabiliste admissible.
- Se ramener à des sous-espaces "non contraints".
- Méthode **efficace en pratique**.

Direct search based on probabilistic feasible descent for bound and linearly constrained problems. Gratton, Royer, Vicente and Zhang, *Submitted*, 2017.

- Cas sans dérivées
 - Contraintes non linéaires;
 - Parallélisation.
- Contexte général
 - Etude de complexité;
 - **Courbure négative.**

- Cas sans dérivées
 - Contraintes non linéaires;
 - Parallélisation.
- Contexte général
 - Etude de complexité;
 - **Courbure négative.**

Merci de votre attention !

`croyer2@wisc.edu`