

Optimisation sans dérivées

Clément Royer

Certificat Chef de Projet IA - Université Paris Dauphine-PSL

18 novembre 2021



Un tour d'horizon

- De techniques classiques en optimisation sans dérivées;
- Et de leurs liens avec l'IA.

Un tour d'horizon

- De techniques classiques en optimisation sans dérivées;
- Et de leurs liens avec l'IA.

Des révisions

- Les concepts de lundi nous seront utiles;
- Ils seront rappelés.

Un tour d'horizon

- De techniques classiques en optimisation sans dérivées;
- Et de leurs liens avec l'IA.

Des révisions

- Les concepts de lundi nous seront utiles;
- Ils seront rappelés.

En fonction du temps

- Illustrations numériques;
- Retour global.

Quand on parle d'optimisation sans dérivées, on parle de...

- 1 Derivative-free optimization (DFO);

Quand on parle d'optimisation sans dérivées, on parle de...

- ① Derivative-free optimization (DFO);
- ② Black-box optimization;

Quand on parle d'optimisation sans dérivées, on parle de...

- ① Derivative-free optimization (DFO);
- ② Black-box optimization;
- ③ Surrogate-based optimization;

Quand on parle d'optimisation sans dérivées, on parle de...

- ① Derivative-free optimization (DFO);
- ② Black-box optimization;
- ③ Surrogate-based optimization;
- ④ Response surface methodology/Design of experiments;

Quand on parle d'optimisation sans dérivées, on parle de...

- ① Derivative-free optimization (DFO);
- ② Black-box optimization;
- ③ Surrogate-based optimization;
- ④ Response surface methodology/Design of experiments;
- ⑤ Automated machine learning;

Quand on parle d'optimisation sans dérivées, on parle de...

- ① Derivative-free optimization (DFO);
- ② Black-box optimization;
- ③ Surrogate-based optimization;
- ④ Response surface methodology/Design of experiments;
- ⑤ Automated machine learning;
- ⑥ Hyperparameter tuning.

Quand on parle d'optimisation sans dérivées, on parle de...

- 1 **Derivative-free optimization (DFO);**
- 2 **Black-box optimization;**
- 3 Surrogate-based optimization;
- 4 Response surface methodology/Design of experiments;
- 5 Automated machine learning;
- 6 Hyperparameter tuning.

Nous allons parler de...

- **Tout cela** dans un même cadre;
- **Des avancées pour 1+2;**
- Du lien avec 5+6.

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
- 3 Méthodes basées sur des modèles

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
- 3 Méthodes basées sur des modèles

1 L'optimisation sans dérivées

- Exemples et définition
- Formulation du problème

2 Méthodes de recherche directe

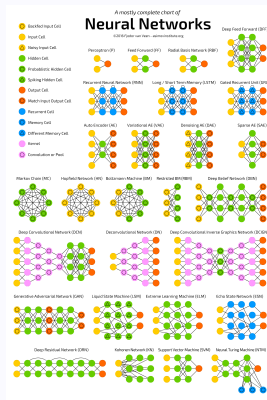
- Recherche aléatoire \Rightarrow recherche directe
- Recherche directe \Rightarrow De l'aléatoire
- Recherche directe et fonctions stochastiques

3 Méthodes basées sur des modèles

- Introduction aux méthodes basées sur des modèles
- Vers des modèles aléatoires

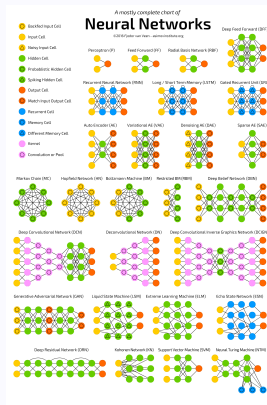
Entraînons un réseau de neurones...

- Quelle architecture ? (nombre de couches, types de couches, etc)
- Quel algorithme d'entraînement ? (Adam, RMSProp, SGD, etc)
- Quelles options pour l'algorithme (*learning rate*, etc)?



Entraînons un réseau de neurones...

- Quelle architecture ? (nombre de couches, types de couches, etc)
- Quel algorithme d'entraînement ? (Adam, RMSProp, SGD, etc)
- Quelles options pour l'algorithme (*learning rate*, etc)?



Calibration d'hyperparamètres

- Chaque test d'une configuration correspond à un nouvel entraînement (heures/jours en temps CPU + argent !);
- Énormément de choix possibles.

Calcul scientifique

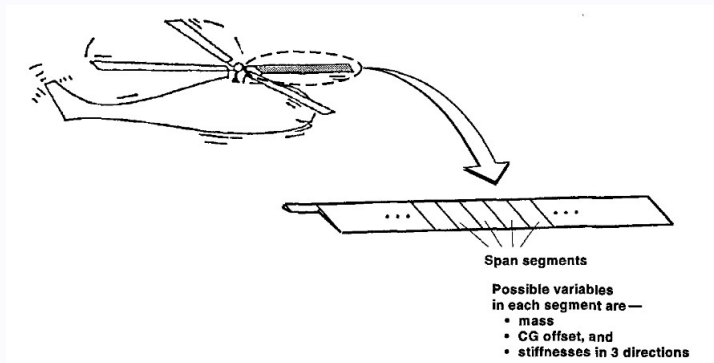
- Usage intensif de la simulation par ordinateur (CFD, CAO) dans les applications de type physique (aéronautique, automobile, météorologie);
- Beaucoup de paramètres à calibrer.



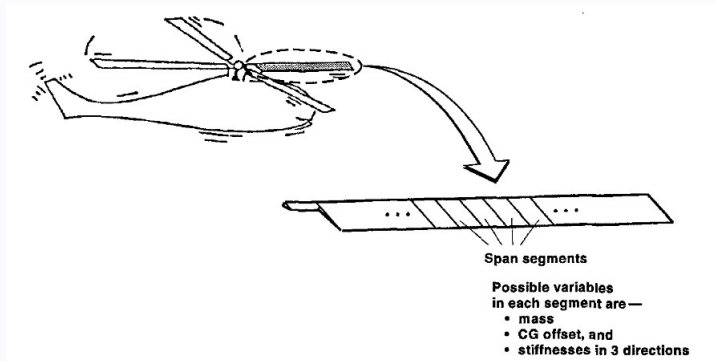
Optimisation de simulateurs

- Besoin : optimiser les paramètres de codes de simulations;
- Coût d'exécution des simulateurs élevé;
- Parfois plusieurs versions à coût variable (multifidélité).

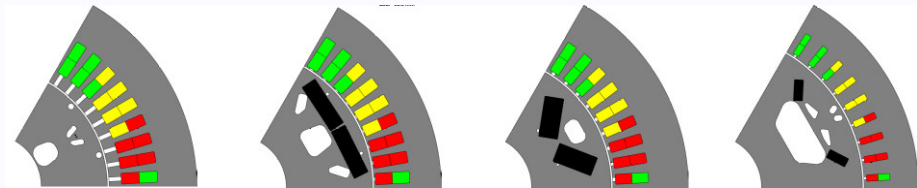
Exemple classique : Design de pale d'hélicoptère (Booker et al. 1998)



Exemple classique : Design de pale d'hélicoptère (Booker et al. 1998)



- 30 paramètres;
- 1 simulation : 2 semaines de calculs CFD;
- Échec de la simulation 60% du temps.
- Optimisation multi-disciplinaire : codes imbriqués;
- De la simulation numérique, beaucoup de calculs...
- qui peuvent échouer !



- Environ 50 paramètres (continus);
- Multiobjectif (3 objectifs), 6 fonctions de contraintes;
- La plupart des points ne sont pas réalisables !
- 1 simulation \approx 5 minutes;
- Optimisation (par algorithmes génétiques) : 3 semaines !

Une classe de problèmes en commun

Points communs : IA automatisée et optimisation de simulateurs

- Effort de calcul conséquent, basé sur la simulation;
- Choix des meilleurs paramètres non trivial;
- Préliminaire à la construction/au déploiement du système.



En termes d'optimisation

- **Le choix des meilleurs paramètres peut se formuler comme un problème d'optimisation;**
- La fonction objectif de ce problème est très coûteuse à évaluer (en temps de simulation).

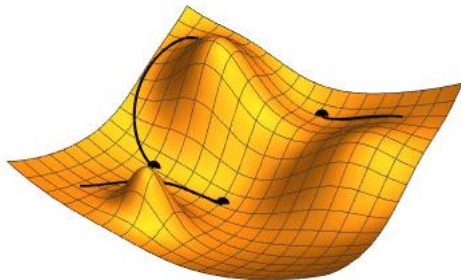
Optimisation sans dérivées

Certaines dérivées ne peuvent pas être exploitées pour optimiser.

Optimisation sans dérivées

Certaines dérivées ne peuvent pas être exploitées pour optimiser.

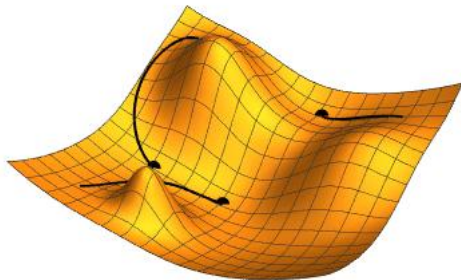
- L'optimisation est fortement basée sur les dérivées (gradient, conditions d'optimalité);
- **Certaines** : il en suffit d'une !



Optimisation sans dérivées

Certaines dérivées ne peuvent pas être exploitées pour optimiser.

- L'optimisation est fortement basée sur les dérivées (gradient, conditions d'optimalité);
- **Certaines** : il en suffit d'une !
- **ne peuvent pas** : qu'elles existent ou non !



Toujours préférable d'utiliser les dérivées si possible !

Toujours préférable d'utiliser les dérivées si possible !

- 1 Expression analytique \Rightarrow Dérivées explicites;

Toujours préférable d'utiliser les dérivées si possible !

- ① Expression analytique \Rightarrow Dérivées explicites;
- ② Approximation numérique par différences finies;

Toujours préférable d'utiliser les dérivées si possible !

- ❶ Expression analytique \Rightarrow Dérivées explicites;
- ❷ Approximation numérique par différences finies;
- ❸ Puissance de la différentiation automatique/symbolique.

Toujours préférable d'utiliser les dérivées si possible !

- ❶ Expression analytique \Rightarrow Dérivées explicites;
- ❷ Approximation numérique par différences finies;
- ❸ Puissance de la différentiation automatique/symbolique.

Dérivées non disponibles

- ❶ Systèmes complexes \Rightarrow Risque d'erreur dans le code à la main;
- ❷ Évaluations coûteuses/bruitées \Rightarrow Problème pour les différences finies;
- ❸ Code propriétaire \Rightarrow Pas de diff. auto.

Algorithmes génétiques/évolutionnaires

- Souvent inspirés par la nature;
- Efficaces avec peu de variables et des évaluations peu coûteuses/parallélisables;
- Beaucoup d'heuristiques inspirées par la nature.

Algorithmes génétiques/évolutionnaires

- Souvent inspirés par la nature;
- Efficaces avec peu de variables et des évaluations peu coûteuses/parallélisables;
- Beaucoup d'heuristiques inspirées par la nature.

Une méthode remarquable : CMA-ES

- Maintient une matrice de covariance;
- Efficace et populaire;
- Récemment interprétée comme une méthode de gradient appliquée à de l'optimisation sur des distributions (chercheurs en IA).

1 L'optimisation sans dérivées

- Exemples et définition
- Formulation du problème

2 Méthodes de recherche directe

- Recherche aléatoire \Rightarrow recherche directe
- Recherche directe \Rightarrow De l'aléatoire
- Recherche directe et fonctions stochastiques

3 Méthodes basées sur des modèles

- Introduction aux méthodes basées sur des modèles
- Vers des modèles aléatoires

Notre formulation

$$\text{minimiser}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{s. c.} \quad \mathbf{x} \in \mathcal{F}$$

- \mathbf{x} : variables;
- f : fonction objectif;
- \mathcal{F} : ensemble admissible.

Notre formulation

$$\text{minimiser}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) \quad \text{s. c.} \quad \mathbf{x} \in \mathcal{F}$$

- \mathbf{x} : variables;
- f : fonction objectif;
- \mathcal{F} : ensemble admissible.

Cadre de travail

- f minorée sur \mathcal{F} : $f(\mathbf{x}) \geq f_{\text{low}}$;
- Evaluations de f coûteuses :
 - Entraîner un réseau de neurones;
 - Prise de sang d'un patient;
 - Ramener un véhicule en usine.

Trouver les paramètres optimaux ?

- Sans dérivées, pas de garantie d'optimalité locale...
- Optimalité globale seulement garantie sous d'autres hypothèses (convexité) ou si l'on attend indéfiniment.

Notre objectif ?

Trouver les paramètres optimaux ?

- Sans dérivées, pas de garantie d'optimalité locale...
- Optimalité globale seulement garantie sous d'autres hypothèses (convexité) ou si l'on attend indéfiniment.

Trouver une meilleure configuration ?

- Toute amélioration est bonne à prendre;
- Les configurations de départ peuvent être fixées par des experts, et difficiles à améliorer.

Notre objectif ?

Trouver les paramètres optimaux ?

- Sans dérivées, pas de garantie d'optimalité locale...
- Optimalité globale seulement garantie sous d'autres hypothèses (convexité) ou si l'on attend indéfiniment.

Trouver une meilleure configuration ?

- Toute amélioration est bonne à prendre;
- Les configurations de départ peuvent être fixées par des experts, et difficiles à améliorer.

Fournir des garanties

- Validation de principe, guide pour choisir des méthodes;
- Métrique du moment : **complexité**.

Définition

A partir

- d'un critère de convergence/d'arrêt;
- d'une précision $\epsilon > 0$;
- d'un algorithme itératif $\{\mathbf{x}_k\}_k$;

borner le **nombre d'appels de fonction** requis dans le pire des cas pour satisfaire le critère avec précision ϵ .

Définition

A partir

- d'un critère de convergence/d'arrêt;
- d'une précision $\epsilon > 0$;
- d'un algorithme itératif $\{\mathbf{x}_k\}_k$;

borner le **nombre d'appels de fonction** requis dans le pire des cas pour satisfaire le critère avec précision ϵ .

La borne (en tant que fonction de ϵ) s'appelle la **complexité au pire cas** de l'algorithme.

Définition

A partir

- d'un critère de convergence/d'arrêt;
- d'une précision $\epsilon > 0$;
- d'un algorithme itératif $\{\mathbf{x}_k\}_k$;

borner le **nombre d'appels de fonction** requis dans le pire des cas pour satisfaire le critère avec précision ϵ .

La borne (en tant que fonction de ϵ) s'appelle la **complexité au pire cas** de l'algorithme.

Exemples de critère de convergence

- Gradient de f : $\|\nabla f(\mathbf{x}_k)\|$;
- Valeur de f : $f(\mathbf{x}_k)$.

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe**
- 3 Méthodes basées sur des modèles

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
 - Recherche aléatoire \Rightarrow recherche directe
 - Recherche directe \Rightarrow De l'aléatoire
 - Recherche directe et fonctions stochastiques
- 3 Méthodes basées sur des modèles

But : Résoudre minimiser $\mathbf{x} \in \mathcal{F} f(\mathbf{x})$ avec accès à f uniquement, budget limité.

Algorithme de recherche basique

Start with: $\hat{\mathbf{x}}_0 = \mathbf{x}_0 \in \mathcal{F}$, $f = f(\mathbf{x}_0)$, $k = 0$.

- 1 Calculer \mathbf{x}_{k+1} et évaluer $f(\mathbf{x}_{k+1})$.

But : Résoudre minimiser $\mathbf{x} \in \mathcal{F} f(\mathbf{x})$ avec accès à f uniquement, budget limité.

Algorithme de recherche basique

Start with: $\hat{\mathbf{x}}_0 = \mathbf{x}_0 \in \mathcal{F}$, $f = f(\mathbf{x}_0)$, $k = 0$.

- 1 Calculer \mathbf{x}_{k+1} et évaluer $f(\mathbf{x}_{k+1})$.
- 2 Si $f(\mathbf{x}_{k+1}) < f(\hat{\mathbf{x}}_k)$ poser $\hat{\mathbf{x}}_{k+1} = \mathbf{x}_{k+1}$, sinon poser $\hat{\mathbf{x}}_{k+1} = \hat{\mathbf{x}}_k$.

But : Résoudre minimiser $\mathbf{x} \in \mathcal{F} f(\mathbf{x})$ avec accès à f uniquement, budget limité.

Algorithme de recherche basique

Start with: $\hat{\mathbf{x}}_0 = \mathbf{x}_0 \in \mathcal{F}$, $f = f(\mathbf{x}_0)$, $k = 0$.

- ➊ Calculer \mathbf{x}_{k+1} et évaluer $f(\mathbf{x}_{k+1})$.
- ➋ Si $f(\mathbf{x}_{k+1}) < f(\hat{\mathbf{x}}_k)$ poser $\hat{\mathbf{x}}_{k+1} = \mathbf{x}_{k+1}$, sinon poser $\hat{\mathbf{x}}_{k+1} = \hat{\mathbf{x}}_k$.
- ➌ Si budget dépassé terminer, sinon incrémenter k de 1.

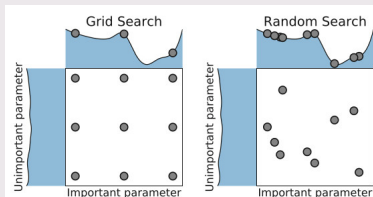
But : Résoudre minimiser $\mathbf{x} \in \mathcal{F} f(\mathbf{x})$ avec accès à f uniquement, budget limité.

Algorithme de recherche basique

Start with: $\hat{\mathbf{x}}_0 = \mathbf{x}_0 \in \mathcal{F}$, $f = f(\mathbf{x}_0)$, $k = 0$.

- 1 Calculer \mathbf{x}_{k+1} et évaluer $f(\mathbf{x}_{k+1})$.
- 2 Si $f(\mathbf{x}_{k+1}) < f(\hat{\mathbf{x}}_k)$ poser $\hat{\mathbf{x}}_{k+1} = \mathbf{x}_{k+1}$, sinon poser $\hat{\mathbf{x}}_{k+1} = \hat{\mathbf{x}}_k$.
- 3 Si budget dépassé terminer, sinon incrémenter k de 1.

- **Grille (Grid search) :** Valeurs des \mathbf{x}_k fixées a priori;
- **Aléatoire (Random search) :** Tirer \mathbf{x}_{k+1} au hasard.



Théorème

Soit $f^* = \min_{\mathbf{x} \in \mathcal{F}} f(\mathbf{x})$. Alors,

$$\mathbb{P}(f(\hat{\mathbf{x}}_K) \leq f^* + \epsilon) \geq p$$

après

$$K = \frac{\ln(p)}{\ln \left[\frac{\mu(\{\mathbf{x} \in \mathcal{F} | f(\mathbf{x}) > f^* + \epsilon\})}{\mu(\mathcal{F})} \right]}.$$

itérations.

Théorème

Soit $f^* = \min_{\mathbf{x} \in \mathcal{F}} f(\mathbf{x})$. Alors,

$$\mathbb{P}(f(\hat{\mathbf{x}}_K) \leq f^* + \epsilon) \geq p$$

après

$$K = \frac{\ln(p)}{\ln \left[\frac{\mu(\{\mathbf{x} \in \mathcal{F} | f(\mathbf{x}) > f^* + \epsilon\})}{\mu(\mathcal{F})} \right]}.$$

itérations.

- Plus : Valable pour f quelconque !
- Moins : Beaucoup d'itérations/d'évaluations de f ;
- Le budget est consommé en exploration.

Recherche exploratoire

- Variables en petite dimension;
- Algorithmes exploratoires.

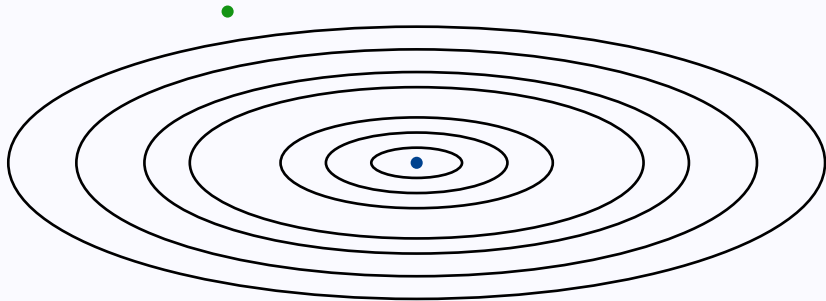
Recherche exploratoire

- Valables en petite dimension;
- Algorithmes exploratoires.

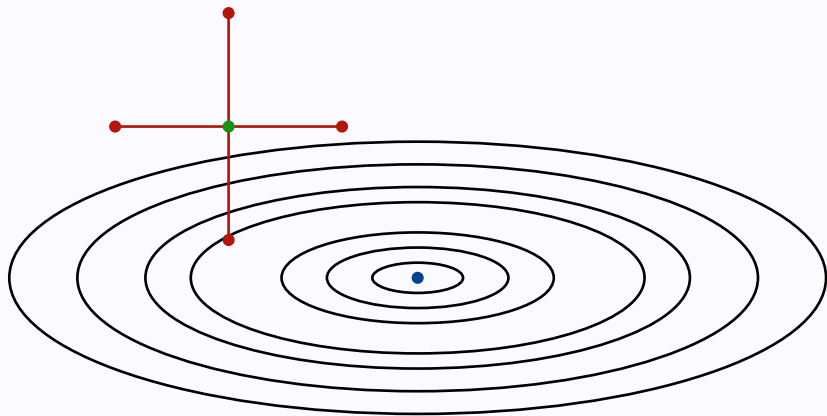
Recherche directe

- Origine : années 1960, théorie : années 1990;
- But : Aller au-delà de l'exploration pure.
- Intérêt : **simplicité, parallélisme**;
- La méthode du simplexe (Nelder-Mead, 1965) a plus de 125000 citations et est toujours la méthode sans dérivées de MATLAB!

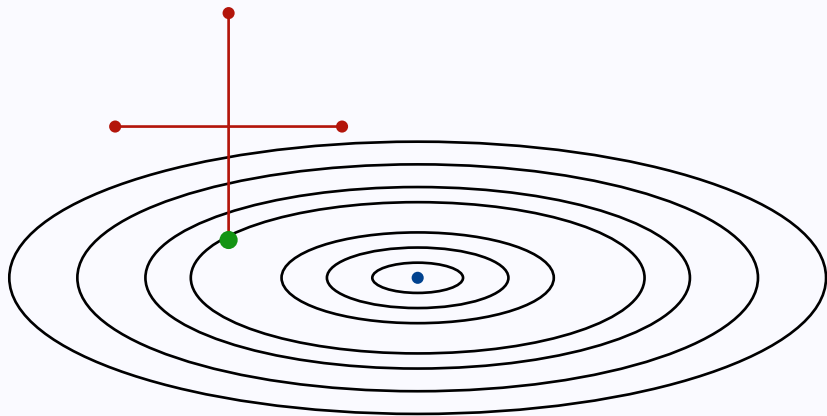
Exemple : Recherche par coordonnées



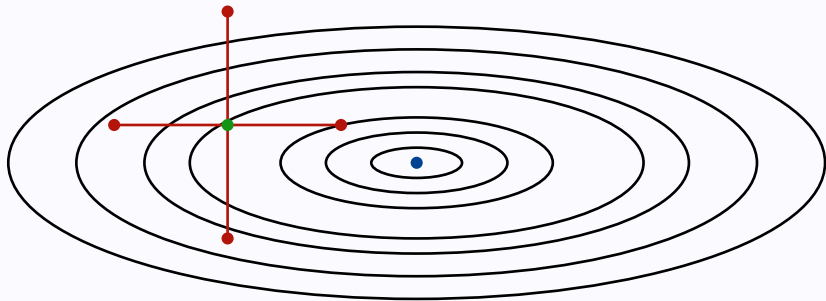
Exemple : Recherche par coordonnées



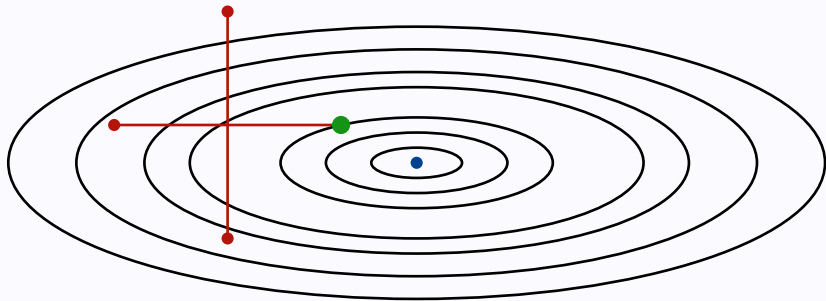
Exemple : Recherche par coordonnées



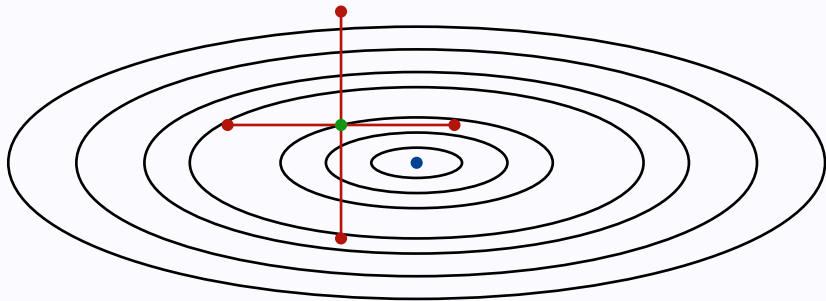
Exemple : Recherche par coordonnées



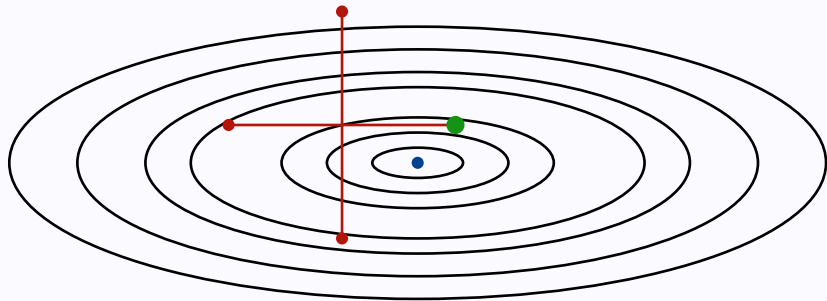
Exemple : Recherche par coordonnées



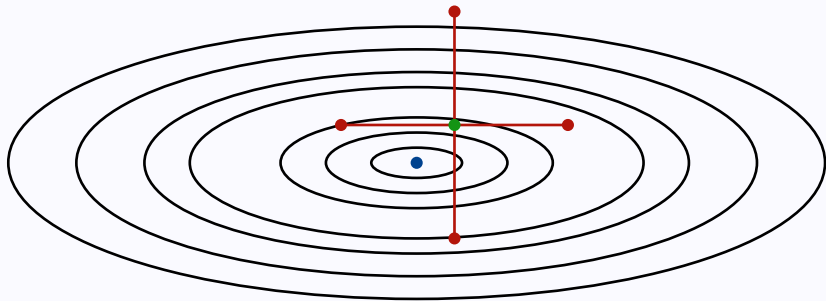
Exemple : Recherche par coordonnées



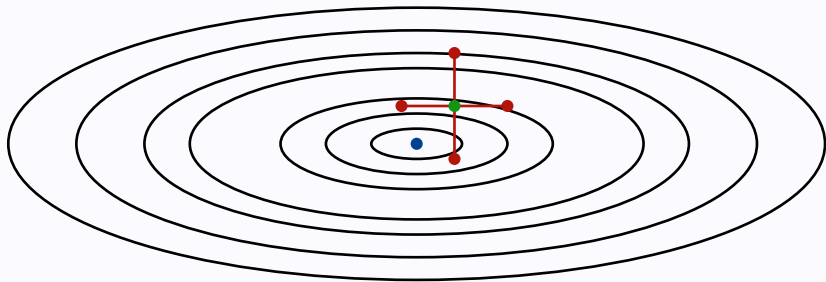
Exemple : Recherche par coordonnées



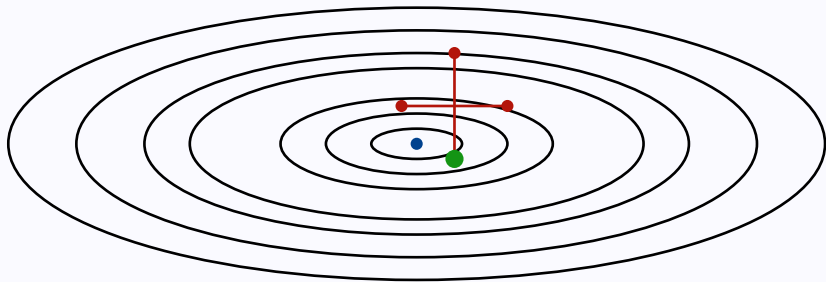
Exemple : Recherche par coordonnées



Exemple : Recherche par coordonnées



Exemple : Recherche par coordonnées



- ❶ Initialisation : $\mathbf{x}_0 \in \mathbb{R}^n, \alpha_0 > 0$.
- ❷ Pour $k = 0, 1, 2, \dots$
 - Choisir un ensemble $\mathcal{D}_k \subset \mathbb{R}^n$ de r directions.

- ① Initialisation : $\mathbf{x}_0 \in \mathbb{R}^n, \alpha_0 > 0$.
- ② Pour $k = 0, 1, 2, \dots$
 - Choisir un ensemble $\mathcal{D}_k \subset \mathbb{R}^n$ de r directions.
 - SI il existe $\mathbf{d}_k \in \mathcal{D}_k$ tel que

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

alors poser $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$ et $\alpha_{k+1} \geq \alpha_k$ (*itération réussie*).

① Initialisation : $\mathbf{x}_0 \in \mathbb{R}^n, \alpha_0 > 0$.

② Pour $k = 0, 1, 2, \dots$

- Choisir un ensemble $\mathcal{D}_k \subset \mathbb{R}^n$ de r directions.
- SI il existe $\mathbf{d}_k \in \mathcal{D}_k$ tel que

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

alors poser $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$ et $\alpha_{k+1} \geq \alpha_k$ (itération réussie).

- Sinon poser $\mathbf{x}_{k+1} := \mathbf{x}_k$ et $\alpha_{k+1} := 0.5\alpha_k$ (itération non réussie).

① Initialisation : $\mathbf{x}_0 \in \mathbb{R}^n, \alpha_0 > 0$.

② Pour $k = 0, 1, 2, \dots$

- Choisir un ensemble $\mathcal{D}_k \subset \mathbb{R}^n$ de r directions.
- SI il existe $\mathbf{d}_k \in \mathcal{D}_k$ tel que

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

alors poser $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$ et $\alpha_{k+1} \geq \alpha_k$ (itération réussie).

- Sinon poser $\mathbf{x}_{k+1} := \mathbf{x}_k$ et $\alpha_{k+1} := 0.5\alpha_k$ (itération non réussie).

❶ Initialisation : $\mathbf{x}_0 \in \mathbb{R}^n, \alpha_0 > 0$.

❷ Pour $k = 0, 1, 2, \dots$

- Choisir un ensemble $\mathcal{D}_k \subset \mathbb{R}^n$ de r directions.
- SI il existe $\mathbf{d}_k \in \mathcal{D}_k$ tel que

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

alors poser $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$ et $\alpha_{k+1} \geq \alpha_k$ (itération réussie).

- Sinon poser $\mathbf{x}_{k+1} := \mathbf{x}_k$ et $\alpha_{k+1} := 0.5\alpha_k$ (itération non réussie).

Aspects cruciaux

- Condition de décroissance;
- Choix de \mathcal{D}_k , valeur de r .

Une mesure de qualité

Soit $\mathcal{D} \subset \mathbb{R}^n$ et $\mathbf{v} \in \mathbb{R}^n \setminus \{0\}$, la mesure cosinus de \mathcal{D} en \mathbf{v} est

$$\text{cm}(\mathcal{D}, \mathbf{v}) = \max_{\mathbf{d} \in \mathcal{D}} \frac{\mathbf{d}^\top \mathbf{v}}{\|\mathbf{d}\| \|\mathbf{v}\|}$$

Une mesure de qualité

Soit $\mathcal{D} \subset \mathbb{R}^n$ et $\mathbf{v} \in \mathbb{R}^n \setminus \{0\}$, la mesure cosinus de \mathcal{D} en \mathbf{v} est

$$\text{cm}(\mathcal{D}, \mathbf{v}) = \max_{\mathbf{d} \in \mathcal{D}} \frac{\mathbf{d}^\top \mathbf{v}}{\|\mathbf{d}\| \|\mathbf{v}\|}$$

- Distance angulaire de l'analyse de données textuelles (word2vec);
- Ici sert à approcher tout vecteur de l'espace par un vecteur de \mathcal{D} ;

Une mesure de qualité

Soit $\mathcal{D} \subset \mathbb{R}^n$ et $\mathbf{v} \in \mathbb{R}^n \setminus \{0\}$, la mesure cosinus de \mathcal{D} en \mathbf{v} est

$$\text{cm}(\mathcal{D}, \mathbf{v}) = \max_{\mathbf{d} \in \mathcal{D}} \frac{\mathbf{d}^\top \mathbf{v}}{\|\mathbf{d}\| \|\mathbf{v}\|}$$

- Distance angulaire de l'analyse de données textuelles (word2vec);
- Ici sert à approcher tout vecteur de l'espace par un vecteur de \mathcal{D} ;
- Meilleurs ensembles (PSS) : $\text{cm}(\mathcal{D}, \mathbf{v}) > 0 \quad \forall \mathbf{v} \neq 0$ (requiert $|\mathcal{D}| \geq n + 1$).

Une mesure de qualité

Soit $\mathcal{D} \subset \mathbb{R}^n$ et $\mathbf{v} \in \mathbb{R}^n \setminus \{0\}$, la mesure cosinus de \mathcal{D} en \mathbf{v} est

$$\text{cm}(\mathcal{D}, \mathbf{v}) = \max_{\mathbf{d} \in \mathcal{D}} \frac{\mathbf{d}^\top \mathbf{v}}{\|\mathbf{d}\| \|\mathbf{v}\|}$$

- Distance angulaire de l'analyse de données textuelles (word2vec);
- Ici sert à approcher tout vecteur de l'espace par un vecteur de \mathcal{D} ;
- Meilleurs ensembles (PSS) : $\text{cm}(\mathcal{D}, \mathbf{v}) > 0 \quad \forall \mathbf{v} \neq 0$ (requiert $|\mathcal{D}| \geq n + 1$).

Une mesure de qualité

Soit $\mathcal{D} \subset \mathbb{R}^n$ et $\mathbf{v} \in \mathbb{R}^n \setminus \{0\}$, la **mesure cosinus** de \mathcal{D} en \mathbf{v} est

$$\text{cm}(\mathcal{D}, \mathbf{v}) = \max_{\mathbf{d} \in \mathcal{D}} \frac{\mathbf{d}^\top \mathbf{v}}{\|\mathbf{d}\| \|\mathbf{v}\|}$$

- Distance angulaire de l'analyse de données textuelles (word2vec);
- Ici sert à approcher tout vecteur de l'espace par un vecteur de \mathcal{D} ;
- Meilleurs ensembles (PSS) : $\text{cm}(\mathcal{D}, \mathbf{v}) > 0 \quad \forall \mathbf{v} \neq 0$ (requiert $|\mathcal{D}| \geq n + 1$).

Exemple : $\mathcal{D}_\oplus = \{\mathbf{e}_1, \dots, \mathbf{e}_n, -\mathbf{e}_1, \dots, -\mathbf{e}_n\}$

- $|\mathcal{D}_\oplus| = 2n$.
- $\forall \mathbf{v}, \quad \text{cm}(\mathcal{D}_\oplus, \mathbf{v}) \geq \frac{1}{\sqrt{n}}$.

Hypothèse : Il existe $\kappa \in (0, 1)$ tel que $\forall k, \text{cm}(\mathcal{D}_k, \mathbf{v}) \geq \kappa \forall \mathbf{v} \in \mathbb{R}^n$, et $|\mathcal{D}_k| = r$ pour tout k .

Théorème

Soit $\epsilon \in (0, 1)$ et N_ϵ le nombre d'appels à f pour satisfaire $\|\nabla f(\mathbf{x}_k)\| < \epsilon$.
Alors

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

Hypothèse : Il existe $\kappa \in (0, 1)$ tel que $\forall k, \text{cm}(\mathcal{D}_k, \mathbf{v}) \geq \kappa \forall \mathbf{v} \in \mathbb{R}^n$, et $|\mathcal{D}_k| = r$ pour tout k .

Théorème

Soit $\epsilon \in (0, 1)$ et N_ϵ le nombre d'appels à f pour satisfaire $\|\nabla f(\mathbf{x}_k)\| < \epsilon$.
Alors

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

- Avec $\mathcal{D}_k = \mathcal{D}_\oplus$, on a $\kappa = 1/\sqrt{n}$, $r = 2n$, d'où

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

- Le facteur n^2 ne peut pas être amélioré avec une méthode déterministe.

Positif

Négatif

Positif

- **Exploitation** : On se base sur le point courant pour choisir le suivant;
- **Exploration** : Mouvement contrôlé par une longueur de pas α_k .

Négatif

Positif

- **Exploitation** : On se base sur le point courant pour choisir le suivant;
- **Exploration** : Mouvement contrôlé par une longueur de pas α_k .

Négatif

- **Utilisation de PSS** : Au moins $n + 1$ évaluations par itération;
- **Dépendance en n** : Dans la complexité, mais aussi dans le nombre d'évaluations à chaque itération.

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
 - Recherche aléatoire \Rightarrow recherche directe
 - Recherche directe \Rightarrow De l'aléatoire
 - Recherche directe et fonctions stochastiques
- 3 Méthodes basées sur des modèles

Recherche déterministe

- Directions souvent fixes \Rightarrow Peu d'exploration;
- Complexité en $\mathcal{O}(n^2\epsilon^{-2})$: forte dépendance en la dimension.

Recherche déterministe

- Directions souvent fixes \Rightarrow Peu d'exploration;
- Complexité en $\mathcal{O}(n^2\epsilon^{-2})$: forte dépendance en la dimension.

Techniques aléatoires

- Recherche directe probabiliste : Utiliser des directions aléatoires \mathbf{D}_k telles que

$$\forall \mathbf{v}, \mathbb{P}(\text{cm}(\mathbf{D}_k, \mathbf{v}) \geq \kappa | \mathbf{D}_0, \dots, \mathbf{D}_{k-1}) \geq p.$$

Recherche déterministe

- Directions souvent fixes \Rightarrow Peu d'exploration;
- Complexité en $\mathcal{O}(n^2\epsilon^{-2})$: forte dépendance en la dimension.

Techniques aléatoires

- Recherche directe probabiliste : Utiliser des directions aléatoires \mathbf{D}_k telles que

$$\forall \mathbf{v}, \mathbb{P}(\text{cm}(\mathbf{D}_k, \mathbf{v}) \geq \kappa | \mathbf{D}_0, \dots, \mathbf{D}_{k-1}) \geq p.$$

- Recherche aléatoire de Nesterov : Tirer $\mathbf{u}_k \sim \mathcal{N}(0, \mathbf{I})$ et prendre

$$\frac{f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x})}{\mu} \mathbf{u} \quad \text{ou} \quad \frac{f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x} - \mu \mathbf{u})}{\mu} \mathbf{u}$$

Recherche déterministe

- Directions souvent fixes \Rightarrow Peu d'exploration;
- Complexité en $\mathcal{O}(n^2\epsilon^{-2})$: forte dépendance en la dimension.

Techniques aléatoires

- Recherche directe probabiliste : Utiliser des directions aléatoires \mathbf{D}_k telles que

$$\forall \mathbf{v}, \mathbb{P}(\text{cm}(\mathbf{D}_k, \mathbf{v}) \geq \kappa | \mathbf{D}_0, \dots, \mathbf{D}_{k-1}) \geq p.$$

- Recherche aléatoire de Nesterov : Tirer $\mathbf{u}_k \sim \mathcal{N}(0, \mathbf{I})$ et prendre

$$\frac{f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x})}{\mu} \mathbf{u} \quad \text{ou} \quad \frac{f(\mathbf{x} + \mu \mathbf{u}) - f(\mathbf{x} - \mu \mathbf{u})}{\mu} \mathbf{u}$$

- Dans les deux cas : Complexité $\mathcal{O}(n\epsilon^{-2})!$

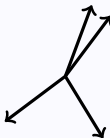
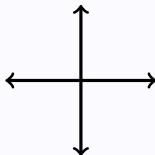
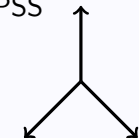
Idée

Utiliser des directions aléatoires en [recherche directe](#).

Idée

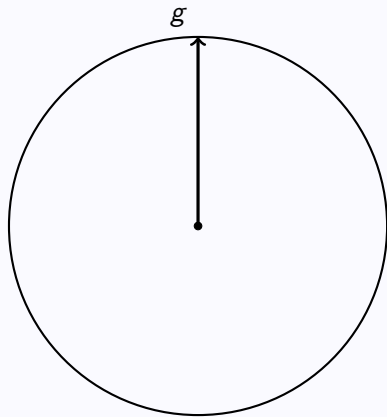
Utiliser des directions aléatoires en [recherche directe](#).

Avant : PSS

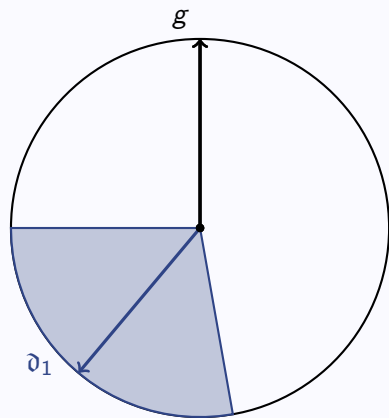


Après : aléatoire

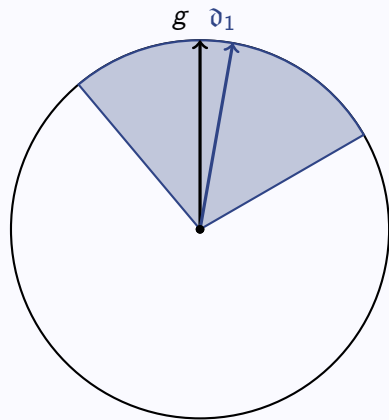
En pratique : deux directions (pas une !)



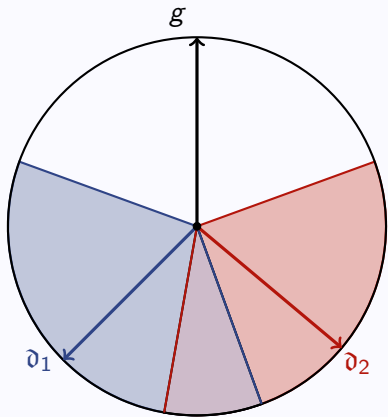
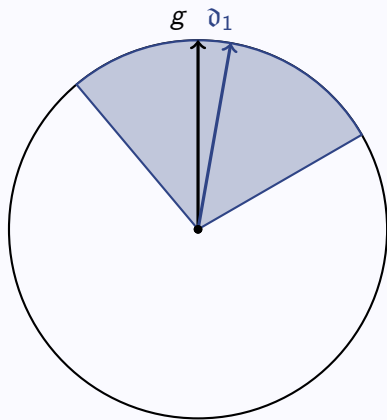
En pratique : deux directions (pas une !)



En pratique : deux directions (pas une !)

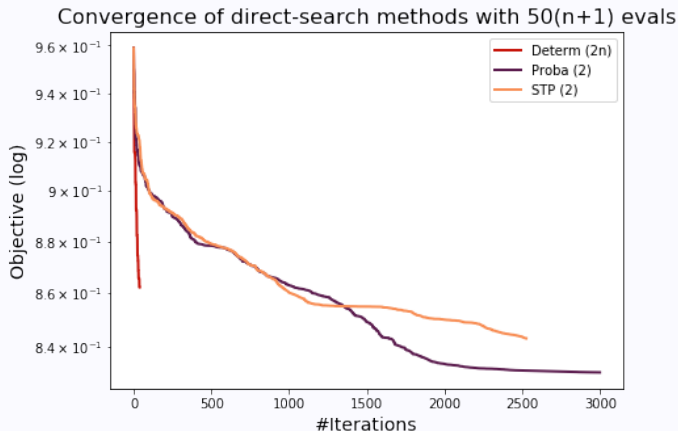


En pratique : deux directions (pas une !)



Et ça marche !

- Régression linéaire robuste avec perte de Tukey;
- Données synthétiques : $n = 100$, budget de $50(n + 1)$ évals.



- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
 - Recherche aléatoire \Rightarrow recherche directe
 - Recherche directe \Rightarrow De l'aléatoire
 - Recherche directe et fonctions stochastiques
- 3 Méthodes basées sur des modèles

$$\text{minimiser}_{\mathbf{x} \in \mathcal{F}} f(\mathbf{x})$$

Nouvelles hypothèses

- f seulement disponible via un oracle stochastique $\tilde{f}(\mathbf{x}; \xi)$;
- Le vecteur ξ est une quantité aléatoire dans un ensemble Ξ ;
- Typique : $\tilde{f}(\cdot; \xi)$ convexe en \mathbf{x} pour toute réalisation de ξ .
- Minimum de f atteint en \mathbf{x}_* .

minimiser $_{\mathbf{x} \in \mathcal{F}} f(\mathbf{x})$

Nouvelles hypothèses

- f seulement disponible via un oracle stochastique $\tilde{f}(\mathbf{x}; \xi)$;
 - Le vecteur ξ est une quantité aléatoire dans un ensemble Ξ ;
 - Typique : $\tilde{f}(\cdot; \xi)$ convexe en \mathbf{x} pour toute réalisation de ξ .
 - Minimum de f atteint en \mathbf{x}_* .
-
- Oracle stochastic \leftrightarrow “Bandit feedback”;
 - Liens modernes avec la littérature en IA sur les bandits et l’optimisation en ligne.

Cadre des bandits à plusieurs bras

- Ensemble de bras $\{1, \dots, A\}$;
- À l'itération k , on joue le bras \mathbf{x}_k , ξ_k est tirée, on obtient $f(\mathbf{x}_k; \xi_k)$;
- **Regret cumulé espéré :**

$$\mathbb{E} \left[\sum_{k=0}^{K-1} f(\mathbf{x}_k; \xi_k) \right] - Kf(\mathbf{x}_*),$$

Bandits avec infinités de bras (Auer, 2002)

- 1 On joue \mathbf{x}_k , ξ_k généré;
- 2 On observe $f(\mathbf{x}_k; \xi_k)$.

But (complexité) : $f(\bar{\mathbf{x}}_K) - f(\mathbf{x}_*) \leq \epsilon$, $\bar{\mathbf{x}}_K = \frac{1}{K} \sum_{k=0}^{K-1} \mathbf{x}_k$.

Méthodes à un point

- Tirer $\mathbf{u}_k \sim \mathcal{U}(\mathbb{S}^{n-1})$ et utiliser

$$\frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu} \mathbf{u}_k \quad \text{ou} \quad \frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k^+) - \tilde{f}(\mathbf{x}_k - \mu \mathbf{u}_k; \boldsymbol{\xi}_k^-)}{\mu} \mathbf{u}_k$$

Méthodes à un point

- Tirer $\mathbf{u}_k \sim \mathcal{U}(\mathbb{S}^{n-1})$ et utiliser

$$\frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu} \mathbf{u}_k \quad \text{ou} \quad \frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k^+) - \tilde{f}(\mathbf{x}_k - \mu \mathbf{u}_k; \boldsymbol{\xi}_k^-)}{\mu} \mathbf{u}_k$$

- Complexité : $\mathcal{O}(n\epsilon^{-3})$ pour problèmes convexes.

Méthodes à un point

- Tirer $\mathbf{u}_k \sim \mathcal{U}(\mathbb{S}^{n-1})$ et utiliser

$$\frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu} \mathbf{u}_k \quad \text{ou} \quad \frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k^+) - \tilde{f}(\mathbf{x}_k - \mu \mathbf{u}_k; \boldsymbol{\xi}_k^-)}{\mu} \mathbf{u}_k$$

- Complexité : $\mathcal{O}(n\epsilon^{-3})$ pour problèmes convexes.
- Bonne dépendance en n , moins en ϵ .

Méthodes à un point

- Tirer $\mathbf{u}_k \sim \mathcal{U}(\mathbb{S}^{n-1})$ et utiliser

$$\frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu} \mathbf{u}_k \quad \text{ou} \quad \frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k^+) - \tilde{f}(\mathbf{x}_k - \mu \mathbf{u}_k; \boldsymbol{\xi}_k^-)}{\mu} \mathbf{u}_k$$

- Complexité : $\mathcal{O}(n\epsilon^{-3})$ pour problèmes convexes.
- Bonne dépendance en n , moins en ϵ .

Méthodes à un point

- Tirer $\mathbf{u}_k \sim \mathcal{U}(\mathbb{S}^{n-1})$ et utiliser

$$\frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu} \mathbf{u}_k \quad \text{ou} \quad \frac{\tilde{f}(\mathbf{x}_k + \mu \mathbf{u}_k; \boldsymbol{\xi}_k^+) - \tilde{f}(\mathbf{x}_k - \mu \mathbf{u}_k; \boldsymbol{\xi}_k^-)}{\mu} \mathbf{u}_k$$

- Complexité : $\mathcal{O}(n\epsilon^{-3})$ pour problèmes convexes.
- Bonne dépendance en n , moins en ϵ .

Méthodes deux/multi-pas

- Hypothèse : Le même $\boldsymbol{\xi}$ permet de faire plusieurs évaluations.
- Tirer $\mathbf{u}_k \sim \mathcal{U}(\mathbb{S}^{n-1})$ et prendre

$$\frac{\tilde{f}(\mathbf{x}_k + \mu_k \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu_k} \mathbf{u}_k \quad \text{or} \quad \frac{\tilde{f}(\mathbf{x}_k + \mu_k \mathbf{u}_k; \boldsymbol{\xi}_k) - \tilde{f}(\mathbf{x}_k - \mu_k \mathbf{u}_k; \boldsymbol{\xi}_k)}{\mu_k} \mathbf{u}_k$$

- Complexité : $\mathcal{O}(n\epsilon^{-2})$ pour problèmes convexes.

Quelques idées-clés

- Exploration via décroissance de fonction (toujours un meilleur point;
- Échantillonnage mais pas de construction de dérivées;
- Les méthodes classiques rajoutent maintenant de l'aléatoire pour des raisons de **performance** en gardant des garanties.

Liens avec l'IA

- Géométriques : distance angulaire;
- Bandit/Optimisation en ligne : nouveaux outils stochastiques.

NOMAD/HyperNOMAD: <https://github.com/bbopt/HyperNOMAD>

- Dédié au départ aux problèmes physiques (HydroQuébec);
- C++/Matlab/Python;
- Gère de nombreuses difficultés non abordées ici :
 - Variables catégorielles, entières;
 - Contraintes plus ou moins relâchables.
- HyperNOMAD (2019) : Extension appliquée pour optimiser architectures et hyperparamètres de réseaux de neurones.

Méthode à base de bandits

- Hyperband (Jamieson et al 2016), BOHB (Falkner et al 2018): approches de bandits + optimisation bayésienne (cf ci-après);
- Garanties supérieures à la recherche aléatoire, applicables à un grand nombre de variables.

- ❶ Quelle est la différence entre une recherche sur grille (*grid search*) et une recherche aléatoire (*random search*) ?
- ❷ Comment améliorer la dépendance en la dimension d'une recherche directe standard ?
- ❸ Sur quel concept d'IA se reposent les méthodes de "feedback" ?

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
- 3 Méthodes basées sur des modèles

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
- 3 Méthodes basées sur des modèles
 - Introduction aux méthodes basées sur des modèles
 - Vers des modèles aléatoires

Résumé de ce qui précède

- Les méthodes de recherche directe explorent...
- ...et certaines **exploitent** localement.
- Utilisation de nouveaux points (notamment pour les méthodes aléatoires), pas de ré-utilisation d'information antérieure.

Résumé de ce qui précède

- Les méthodes de recherche directe explorent...
- ...et certaines **exploitent** localement.
- Utilisation de nouveaux points (notamment pour les méthodes aléatoires), pas de ré-utilisation d'information antérieure.

DFO basée sur des modèles

- Utilise des évaluations **passées** de la fonction pour en construire un modèle;
- Ré-utilise des points, beaucoup moins coûteux que des différences finies.

Algorithme : Régions de confiance sans dérivées

- But : minimiser $\mathbf{x} \in \mathbb{R}^n$ $f(\mathbf{x})$;
- Évaluations de f coûteuses.

Algorithme : Régions de confiance sans dérivées

- But : minimiser $f(\mathbf{x})$ sur $\mathbf{x} \in \mathbb{R}^n$;
- Évaluations de f coûteuses.

Entrées : $\mathbf{x}_0 \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\delta_0 > 0$.

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

- Calculer un modèle $\mathbf{s} \mapsto m_k(\mathbf{x}_k + \mathbf{s})$ de f en \mathbf{x}_k ;
- Calculer $\mathbf{s}_k \approx \operatorname{argmin}_{\|\mathbf{s}\| \leq \delta_k} m_k(\mathbf{x}_k + \mathbf{s})$;

Algorithme : Régions de confiance sans dérivées

- But : minimiser $\mathbf{x} \in \mathbb{R}^n f(\mathbf{x})$;
- Évaluations de f coûteuses.

Entrées : $\mathbf{x}_0 \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\delta_0 > 0$.

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

- Calculer un modèle $\mathbf{s} \mapsto m_k(\mathbf{x}_k + \mathbf{s})$ de f en \mathbf{x}_k ;
- Calculer $\mathbf{s}_k \approx \operatorname{argmin}_{\|\mathbf{s}\| \leq \delta_k} m_k(\mathbf{x}_k + \mathbf{s})$;
- Évaluer $\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(\mathbf{x}_k) - m_k(\mathbf{x}_k + \mathbf{s}_k)}$.

Algorithme : Régions de confiance sans dérivées

- But : minimiser $\mathbf{x} \in \mathbb{R}^n f(\mathbf{x})$;
- Évaluations de f coûteuses.

Entrées : $\mathbf{x}_0 \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\delta_0 > 0$.

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

- Calculer un modèle $\mathbf{s} \mapsto m_k(\mathbf{x}_k + \mathbf{s})$ de f en \mathbf{x}_k ;
- Calculer $\mathbf{s}_k \approx \operatorname{argmin}_{\|\mathbf{s}\| \leq \delta_k} m_k(\mathbf{x}_k + \mathbf{s})$;
- Évaluer $\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(\mathbf{x}_k) - m_k(\mathbf{x}_k + \mathbf{s}_k)}$.
- Si $\rho_k \geq \eta$, poser $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ et $\delta_{k+1} \geq \delta_k$.

Algorithme : Régions de confiance sans dérivées

- But : minimiser $\mathbf{x} \in \mathbb{R}^n$ $f(\mathbf{x})$;
- Évaluations de f coûteuses.

Entrées : $\mathbf{x}_0 \in \mathbb{R}^n$, $\eta \in (0, 1)$, $\delta_0 > 0$.

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

- Calculer un modèle $\mathbf{s} \mapsto m_k(\mathbf{x}_k + \mathbf{s})$ de f en \mathbf{x}_k ;
- Calculer $\mathbf{s}_k \approx \operatorname{argmin}_{\|\mathbf{s}\| \leq \delta_k} m_k(\mathbf{x}_k + \mathbf{s})$;
- Évaluer $\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{m_k(\mathbf{x}_k) - m_k(\mathbf{x}_k + \mathbf{s}_k)}$.
- Si $\rho_k \geq \eta$, poser $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$ et $\delta_{k+1} \geq \delta_k$.
- Sinon, poser $\mathbf{x}_{k+1} = \mathbf{x}_k$ et $\delta_{k+1} = \delta_k/2$.

Qualité du modèle

But : Approcher une fonction dérivable f par un modèle m .

- Bornes sur l'erreur d'approximation;
- Garanties locales, dans une **région de confiance**.

Qualité du modèle

But : Approcher une fonction dérivable f par un modèle m .

- Bornes sur l'erreur d'approximation;
- Garanties locales, dans une **région de confiance**.

Modèles pleinement linéaires

Le modèle m est κ -pleinement linéaire pour f en (\mathbf{x}, δ) si pour tout $\mathbf{y} \in B(\mathbf{x}, \delta)$,

$$\begin{aligned} |m(\mathbf{y}) - f(\mathbf{y})| &\leq \kappa\delta^2 \\ \|\nabla m(\mathbf{y}) - \nabla f(\mathbf{y})\| &\leq \kappa\delta. \end{aligned}$$

Construire des modèles pleinement linéaires

- \mathcal{P}_n^d : polynômes de degré d sur \mathbb{R}^n , $\dim \mathcal{P}_n^d = q + 1$;
- $\Phi = \{\phi_0(\cdot), \dots, \phi_q(\cdot)\}$: base de \mathcal{P}_n^d ;
- $\mathcal{Y} = \{\mathbf{y}^0, \dots, \mathbf{y}^p\}$: ensemble de $p + 1$ points de \mathbb{R}^n ;

Construire des modèles pleinement linéaires

- \mathcal{P}_n^d : polynômes de degré d sur \mathbb{R}^n , $\dim \mathcal{P}_n^d = q + 1$;
- $\Phi = \{\phi_0(\cdot), \dots, \phi_q(\cdot)\}$: base de \mathcal{P}_n^d ;
- $\mathcal{Y} = \{\mathbf{y}^0, \dots, \mathbf{y}^p\}$: ensemble de $p + 1$ points de \mathbb{R}^n ;
- **But** : modèle $m(\mathbf{x}) = \sum_{i=0}^q \alpha_i \phi_i(\mathbf{x})$ tels que

$$\forall j = 0, \dots, p, \quad m(\mathbf{y}^j) \approx f(\mathbf{y}^j).$$

- Reformulé comme $M(\Phi, \mathcal{Y})\alpha \approx f(\mathcal{Y})$, avec

$$M(\Phi, \mathcal{Y}) = \begin{bmatrix} \phi_0(\mathbf{y}^0) & \cdots & \phi_q(\mathbf{y}^0) \\ \vdots & \vdots & \vdots \\ \phi_0(\mathbf{y}^p) & \cdots & \phi_q(\mathbf{y}^p) \end{bmatrix}, \quad f(\mathcal{Y}) = \begin{bmatrix} f(\mathbf{y}^0) \\ \vdots \\ f(\mathbf{y}^p) \end{bmatrix}.$$

Régression polynomiale

Calculer α^* solution ed

$$\text{minimiser}_{\alpha \in \mathbb{R}^{q+1}} \|M(\Phi, \mathcal{Y})\alpha - f(\mathcal{Y})\|^2.$$

et prendre $m(\mathbf{x}) = \sum_{i=0}^q \alpha_i^* \phi_i(\mathbf{x})$.

Construire des modèles pleinement linéaires (2)

Régression polynomiale

Calculer α^* solution de

$$\text{minimiser}_{\alpha \in \mathbb{R}^{q+1}} \|M(\Phi, \mathcal{Y})\alpha - f(\mathcal{Y})\|^2.$$

et prendre $m(\mathbf{x}) = \sum_{i=0}^q \alpha_i^* \phi_i(\mathbf{x})$.

Résultat

Si $\mathcal{Y} \subset B(\mathbf{y}^0, \delta)$ est équilibré, alors m est pleinement linéaire $B(\mathbf{y}^0, \delta)$.

Ex)

- Interpolation/régression linéaire avec $p = q = n$ (**technique de base d'analyse de données**);
- $\mathcal{Y} = \{\mathbf{y}^0, \mathbf{y}^1, \dots, \mathbf{y}^n\}$ sommets d'un simplexe de \mathbb{R}^n .

Hypothèse

À chaque itération k ,

- m_k est κ -pleinement linéaire with $\kappa > 0$;
- m_k construit avec au plus r évaluations de f dans $B(\mathbf{x}_k, \delta_k)$.

Hypothèse

À chaque itération k ,

- m_k est κ -pleinement linéaire with $\kappa > 0$;
- m_k construit avec au plus r évaluations de f dans $B(\mathbf{x}_k, \delta_k)$.

Résultat

On obtient $\|\nabla f(\mathbf{x}_k)\| < \epsilon$ en au plus

- $\mathcal{O}(\kappa^2 \epsilon^{-2})$ itérations;
- $\mathcal{O}(r \kappa^2 \epsilon^{-2})$ évaluations.

Ex) *Interpolation/régression linéaire* : $r = \mathcal{O}(n)$, $\kappa = \mathcal{O}(\sqrt{n})$

\Rightarrow Complexité en $\mathcal{O}(n^2 \epsilon^{-2})$.

Qualité du modèle

But : Approcher une fonction **deux fois dérivable** f par un modèle m .

- Bornes sur l'erreur d'approximation;
- Garanties locales, dans une **région de confiance**.

Qualité du modèle

But : Approcher une fonction **deux fois dérivable** f par un modèle m .

- Bornes sur l'erreur d'approximation;
- Garanties locales, dans une **région de confiance**.

Modèles pleinement quadratiques

Le modèle m_k est κ -**pleinement quadratique** pour f sur (\mathbf{x}_k, δ_k) si pour tout $\mathbf{y} \in B(\mathbf{x}_k, \delta_k)$,

$$\begin{aligned} |m_k(\mathbf{y}) - f(\mathbf{y})| &\leq \kappa \delta_k^3 \\ \|\nabla m_k(\mathbf{y}) - \nabla f(\mathbf{y})\| &\leq \kappa \delta_k^2 \\ \|\nabla^2 m_k(\mathbf{y}) - \nabla^2 f(\mathbf{y})\| &\leq \kappa \delta_k. \end{aligned}$$

Construire des modèles pleinement quadratiques en pratique

Choix classiques

Modèles m_k construits à partir de valeurs de f en $\mathcal{Y}_k = \{\mathbf{x}_k, \mathbf{y}^1, \dots, \mathbf{y}^r\}$:

- Interpolation/Régression;
- Radial basis functions (RBF, ou noyaux gaussiens).

Point clé

Bonne géométrie de $\mathcal{Y}_k \Rightarrow m_k$ pleinement quadratique sur $B(\mathbf{x}_k, \delta_k)$.

Ex) Interpolation quadratique avec $r = \mathcal{O}(n^2)$ échantillons

$$\mathcal{Y}_k = \left\{ \mathbf{x}_k, \{\mathbf{x}_k \pm \delta_k \mathbf{e}_i\}_{i=1}^n, \{\mathbf{x}_k + \delta_k \frac{\mathbf{e}_i + \mathbf{e}_j}{2}\}_{1 \leq i < j \leq n} \right\}.$$

Construire des modèles pleinement quadratiques en pratique

Choix classiques

Modèles m_k construits à partir de valeurs de f en $\mathcal{Y}_k = \{\mathbf{x}_k, \mathbf{y}^1, \dots, \mathbf{y}^r\}$:

- Interpolation/Régression;
- Radial basis functions (RBF, ou noyaux gaussiens).

Point clé

Bonne géométrie de $\mathcal{Y}_k \Rightarrow m_k$ pleinement quadratique sur $B(\mathbf{x}_k, \delta_k)$.

Ex) Interpolation quadratique avec $r = \mathcal{O}(n^2)$ échantillons

$$\mathcal{Y}_k = \left\{ \mathbf{x}_k, \{\mathbf{x}_k \pm \delta_k \mathbf{e}_i\}_{i=1}^n, \{\mathbf{x}_k + \delta_k \frac{\mathbf{e}_i + \mathbf{e}_j}{2}\}_{1 \leq i < j \leq n} \right\}.$$

Dans la pratique

- On ré-utilise autant de points que possible!
- On contrôle la géométrie et on la corrige si besoin:

Construire des modèles pleinement quadratiques en pratique

Choix classiques

Modèles m_k construits à partir de valeurs de f en $\mathcal{Y}_k = \{\mathbf{x}_k, \mathbf{y}^1, \dots, \mathbf{y}^r\}$:

- Interpolation/Régression;
- Radial basis functions (RBF, ou noyaux gaussiens).

Point clé

Bonne géométrie de $\mathcal{Y}_k \Rightarrow m_k$ pleinement quadratique sur $B(\mathbf{x}_k, \delta_k)$.

Ex) Interpolation quadratique avec $r = \mathcal{O}(n^2)$ échantillons

$$\mathcal{Y}_k = \left\{ \mathbf{x}_k, \{\mathbf{x}_k \pm \delta_k \mathbf{e}_i\}_{i=1}^n, \{\mathbf{x}_k + \delta_k \frac{\mathbf{e}_i + \mathbf{e}_j}{2}\}_{1 \leq i < j \leq n} \right\}.$$

Dans la pratique

- On ré-utilise autant de points que possible!
- On contrôle la géométrie et on la corrige si besoin:
 - Peut demander r nouvelles valeurs;

Construire des modèles pleinement quadratiques en pratique

Choix classiques

Modèles m_k construits à partir de valeurs de f en $\mathcal{Y}_k = \{\mathbf{x}_k, \mathbf{y}^1, \dots, \mathbf{y}^r\}$:

- Interpolation/Régression;
- Radial basis functions (RBF, ou noyaux gaussiens).

Point clé

Bonne géométrie de $\mathcal{Y}_k \Rightarrow m_k$ pleinement quadratique sur $B(\mathbf{x}_k, \delta_k)$.

Ex) Interpolation quadratique avec $r = \mathcal{O}(n^2)$ échantillons

$$\mathcal{Y}_k = \left\{ \mathbf{x}_k, \{\mathbf{x}_k \pm \delta_k \mathbf{e}_i\}_{i=1}^n, \{\mathbf{x}_k + \delta_k \frac{\mathbf{e}_i + \mathbf{e}_j}{2}\}_{1 \leq i < j \leq n} \right\}.$$

Dans la pratique

- On ré-utilise autant de points que possible!
- On contrôle la géométrie et on la corrige si besoin:
 - Peut demander r nouvelles valeurs;
 - En pratique, bien plus économe.

Hypothèses

À chaque itération k ,

- m_k est κ -pleinement quadratique, $\kappa > 0$;
- m_k construit avec au plus r nouvelles évaluations of f .

Hypothèses

À chaque itération k ,

- m_k est κ -pleinement quadratique, $\kappa > 0$;
- m_k construit avec au plus r nouvelles évaluations of f .

Résultat

On obtient $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon_g$ et $\nabla^2 f(\mathbf{x}_k) \succeq -\epsilon_g^{1/2} \mathbf{I}_n$ en au plus

- $\mathcal{O}(\kappa^3 \epsilon_g^{-3/2})$ itérations;
- $\mathcal{O}(r \kappa^3 \epsilon_g^{-3/2})$ évaluations.

Hypothèses

À chaque itération k ,

- m_k est κ -pleinement quadratique, $\kappa > 0$;
- m_k construit avec au plus r nouvelles évaluations of f .

Résultat

On obtient $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon_g$ et $\nabla^2 f(\mathbf{x}_k) \succeq -\epsilon_g^{1/2} \mathbf{I}_n$ en au plus

- $\mathcal{O}(\kappa^3 \epsilon_g^{-3/2})$ itérations;
- $\mathcal{O}(r \kappa^3 \epsilon_g^{-3/2})$ évaluations.

Exemple : Interpolation/régression

- $r = \mathcal{O}(n^2)$, $\kappa = \mathcal{O}(n)$;
- Complexité en $\mathcal{O}(n^5 \epsilon_g^{-3/2})$.

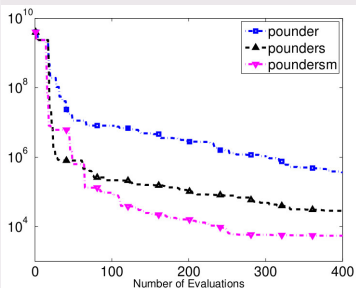
Remarque : Utiliser la structure

- Nous avons traité la fonction en “boîte noire”;
- Souvent on peut calculer les dérivées d'une partie de la fonction;
- Globalement, connaître la structure peut aider!

Remarque : Utiliser la structure

- Nous avons traité la fonction en “boîte noire”;
- Souvent on peut calculer les dérivées d'une partie de la fonction;
- Globalement, connaître la structure peut aider!

Ex) minimiser $\mathbf{x} \in \mathbb{R}^n$ $f(\mathbf{x}) = \frac{1}{2} \|r(\mathbf{x})\|^2$ $r : \mathbb{R}^n \rightarrow \mathbb{R}^m$.



- Dérivées de r
inaccessibles mais on a :

$$\nabla f(\mathbf{x}) = \mathbf{J}_r(\mathbf{x})^T r(\mathbf{x}).$$

- On construit ainsi des modèles plus précis.

- 1 L'optimisation sans dérivées
- 2 Méthodes de recherche directe
- 3 Méthodes basées sur des modèles
 - Introduction aux méthodes basées sur des modèles
 - Vers des modèles aléatoires

Avec des modèles linéaires

- En pratique, ré-utiliser des points/prendre moins de $n + 1$ points marche;
- En théorie, il faut $\mathcal{O}(n)$ points pour être pleinement linéaires!

Avec des modèles linéaires

- En pratique, ré-utiliser des points/prendre moins de $n + 1$ points marche;
- En théorie, il faut $\mathcal{O}(n)$ points pour être pleinement linéaires!

Idée

Supposer que les modèles sont pleinement linéaires **en probabilité**.

- Processus aléatoire;
- Analysé via des arguments statistiques (martingales).

Rappel : Le modèle m_k est κ -pleinement linéaire en (\mathbf{x}_k, δ_k) si pour tout $\mathbf{y} \in \mathcal{B}(\mathbf{x}_k, \delta_k)$:

$$|m_k(\mathbf{y}) - f(\mathbf{y})| \leq \kappa \delta_k^2, \quad \|\nabla m_k(\mathbf{y}) - \nabla f(\mathbf{y})\| \leq \kappa \delta_k$$

Rappel : Le modèle m_k est κ -pleinement linéaire en (\mathbf{x}_k, δ_k) si pour tout $\mathbf{y} \in \mathcal{B}(\mathbf{x}_k, \delta_k)$:

$$|m_k(\mathbf{y}) - f(\mathbf{y})| \leq \kappa \delta_k^2, \quad \|\nabla m_k(\mathbf{y}) - \nabla f(\mathbf{y})\| \leq \kappa \delta_k$$

Modèle probabiliste

Une **suite aléatoire** de modèles $\{m_k\}$ est **(p, κ) -pleinement linéaire** si

$$\begin{aligned} \mathbb{P}(m_0 \text{ } \kappa\text{-plein. lin.}) &\geq p \\ \forall k \geq 1, \quad \mathbb{P}(m_k \text{ } \kappa\text{-plein. lin.} \mid m_0, \dots, m_{k-1}) &\geq p, \end{aligned}$$

Hypothèses

Pour chaque itération k ,

- $\{m_k\}_k$ est (p, κ) -pleinement linéaire avec $p > 1/2$;
- m_k construit avec r nouveaux appels à f .

Hypothèses

Pour chaque itération k ,

- $\{m_k\}_k$ est (p, κ) -pleinement linéaire avec $p > 1/2$;
- m_k construit avec r nouveaux appels à f .

Théorème

Soient $\epsilon \in (0, 1)$ et N_ϵ le nombre d'évaluations nécessaires pour obtenir $\|\nabla f(\mathbf{x}_k)\| \leq \epsilon$. Alors, $N_\epsilon = \mathcal{O}(r\kappa^2\epsilon^{-2})$ avec forte probabilité.

Sous-échantillonnage

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}), \nabla f_i \text{ disponible mais pas } \nabla f$$

- Choisir $S \subset \{1, \dots, N\}$ aléatoirement;
- Poser $m(\mathbf{x} + \mathbf{s}) = f(\mathbf{x}) + \frac{1}{|S|} \sum_{i \in S} \nabla f_i(\mathbf{x})^\top \mathbf{s}$;
- Avec $|S| = \mathcal{O}(\delta^{-2})$, le modèle est pleinement linéaire en probabilité.

Sous-échantillonnage

$$f(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}), \nabla f_i \text{ disponible mais pas } \nabla f$$

- Choisir $S \subset \{1, \dots, N\}$ aléatoirement;
- Poser $m(\mathbf{x} + \mathbf{x}) = f(\mathbf{x}) + \frac{1}{|S|} \sum_{i \in S} \nabla f_i(\mathbf{x})^\top \mathbf{s}$;
- Avec $|S| = \mathcal{O}(\delta^{-2})$, le modèle est pleinement linéaire en probabilité.

Modèles pleinement quadratiques en proba.

- Déterministe : $\mathcal{O}(n^2)$ évaluations;
- Si la matrice hessienne est creuse, vrai (en proba.) avec $\mathcal{O}(n(\log n)^4)$ évaluations via des techniques de **compression de signal**.

Paradigme de l'optimisation bayésienne

- ➊ À chaque itération, calculer une *distribution a posteriori* relativement aux évaluations connues;
- ➋ Trouver un nouveau point en maximisant une *fonction d'acquisition*;
- ➌ Répéter.

Paradigme de l'optimisation bayésienne

- 1 À chaque itération, calculer une *distribution a posteriori* relativement aux évaluations connues;
- 2 Trouver un nouveau point en maximisant une *fonction d'acquisition*;
- 3 Répéter.

- Populaire chez les *data scientists*:
- Les modèles sont basés sur des processus gaussiens et donc des fonctions RBF :

$$m_k(\mathbf{x}_k + \mathbf{s}) = \sum_{i=1}^{|\mathcal{Y}|} \exp(-\|\mathbf{y}_i - \mathbf{s}\|^2)$$

⇒ Ce sont des modèles pleinement linéaires !

Construire des modèles

- Utilise l'historique de l'algorithme;
- Implémentations efficaces meilleures qu'estimer les dérivées directement;
- Il faut exploiter de la structure s'il en existe.

Popularité des modèles

- Méta-modèles/*Surrogates* en calcul scientifique;
- Méta-modèles/processus gaussiens en optimisation bayésienne.

Codes de M. J. D. Powell

- En FORTRAN 77, mais toujours l'un des meilleurs codes disponibles;
- PDFO (<https://www.pdf0.net/>), interfaces Python et MATLAB;
- Récemment utilisé en IA (apprentissage adverse).

Autres codes

- POUNDERS : Moindres carrés sans dérivées (structure);
- ORBIT : Régions de confiance, modèle RBF;
- Packages Python/R pour l'optimisation de *surrogates*/ l'optimisation bayésienne.

- ① En quoi les méthodes basées sur des modèles exploitent-elles plus l'information de l'objectif ?
- ② Expliquer pourquoi les bornes de complexité sont très pessimistes pour les méthodes basées sur des modèles.
- ③ Donner un exemple de structure de problème qui permet d'utiliser certaines dérivées.

Optimisation sans dérivées/DFO

- Pas de dérivées dans l'algorithme (ou pas toutes);
- Différents noms, même but;
- Tout dépend de l'évaluation.

Optimisation sans dérivées/DFO

- Pas de dérivées dans l'algorithme (ou pas toutes);
- Différents noms, même but;
- Tout dépend de l'évaluation.

DFO et IA

- Les techniques DFO “randomisées” se basent sur des outils de l'IA;
- Le chemin vers l'IA automatique passe par des approches formalisées (comme en DFO).

La grande dimension

- Challenge constant;
- Beaucoup d'efforts fournis en IA.

La grande dimension

- Challenge constant;
- Beaucoup d'efforts fournis en IA.

Comprendre les algorithmes “randomisés”

- Rapport *RASC* du DOE (États-Unis);
- Approche systémique requise, nombreux points d'amélioration;
- Clé pour la grande dimension !

Apprendre un modèle des données

- Cadre d'apprentissage classique;
- Souci : peu de données;
- Utilisation de la structure du problème.

Apprendre un modèle des données

- Cadre d'apprentissage classique;
- Souci : peu de données;
- Utilisation de la structure du problème.

Fournir des outils mathématiques

- Statistique en grande dimension;
- Points de vue typiquement IA (ex : bandits).

Apprendre un modèle des données

- Cadre d'apprentissage classique;
- Souci : peu de données;
- Utilisation de la structure du problème.

Fournir des outils mathématiques

- Statistique en grande dimension;
- Points de vue typiquement IA (ex : bandits).

Rejoint les besoins de l'IA automatisée !

- A. R. Conn, K. Scheinberg, L. N. Vicente, *Introduction to derivative-free optimization*, SIAM, 2009.
 - C. Audet, W. Hare, *Derivative-Free and Blackbox Optimization*, Springer, 2017.
 - J. Larson, M. Menickelly and S. M. Wild, *Derivative-free optimization methods*, Acta Numerica, 2019.
-
- M. Feurer and F. Hutter, *Hyperparameter Optimization*, in *Automated Machine Learning*, Springer, 2019.
 - A. Flaxman, A. T. Kalai and H. B. McMahan, *Online convex optimization in the bandit setting: Gradient descent without a gradient*, Symposium on Discrete Algorithms (SODA), 2005.
 - P. I. Frazier, *Bayesian Optimization*, Tutorials in Operations Research, 2018.

- P. Auer, *Using Confidence Bounds for Exploitation-Exploration Trade-offs*, Journal of Machine Learning Research, 2002.
- A. Bandeira, K. Scheinberg and L. N. Vicente, *Convergence of trust-region methods based on probabilistic models*, SIAM Journal on Optimization, 2014.
- A.J. Booker, J.E. Dennis Jr., P.D. Frank, D.W. Moore and D.B. Serafini, *Managing surrogate objectives to optimize a helicopter rotor design – further experiments*, AIAA/ISMOO Symposium on Multidisciplinary Analysis and Optimization, 1998.
- A.R. Conn, K. Scheinberg and L.N. Vicente, *Geometry of interpolation sets in derivative free optimization*, Mathematical Programming, 2018.
- S. Gratton, C. W. Royer, L. N. Vicente and Z. Zhang, *Direct search based on probabilistic descent*, SIAM Journal on Optimization, 2015.
- S. Gratton, C. W. Royer, L. N. Vicente and Z. Zhang, *Complexity and global rates of trust-region methods based on probabilistic models*, IMA Journal of Numerical Analysis, 2018.
- K. Jamieson et al., *Hyperband: A Novel Bandit-Based Approach to Hyperparameter Optimization*, Journal of Machine Learning Research, 2018.
- D. Lakhmiri et al., *HyperNOMAD*: <https://github.com/bbopt/HyperNOMAD>
- Yu. Nesterov and V. Spokoiny, *Random gradient-free minimization of convex functions*, Foundations of Computational Mathematics, 2017.
- T. M. Ragonneau and Z. Zhang, *PDFO*: <https://www.pdf0.net/>
- K. Scheinberg and Ph. L. Toint, *Self-correcting geometry in model-based algorithms for derivative-free unconstrained optimization*, SIAM Journal on Optimization, 2010.

Merci de votre attention !

`clement.royer@dauphine.psl.eu`

Crédits image

- <https://towardsdatascience.com/>
- <https://commons.wikimedia.org/>
- A. J. Booker, J. E. Dennis Jr., P. D. Frank, D. B. Serafini and V. Torczon, *A rigorous framework for optimization of expensive functions by surrogates*, Structural Optim., 1999.
- M. Feurer and F. Hutter, *Hyperparameter Optimization*, in *Automated Machine Learning*, Springer, 2019.
- D. Gaudrie (Stellantis).
- S. M. Wild, *Beyond the Black Box in Derivative-Free and Simulation-Based Optimization*, SIAM Annual Meeting, 2016.