

Optimization

Part VI: Stochastic optimization

Barbara Pascal, Nelly Pustelnik

CNRS, Laboratoire de Physique de l'ENS de Lyon, Univ. Lyon 1
nelly.pustelnik@ens-lyon.fr

barbara.pascal@ens-lyon.fr

Motivations

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x)$$

Motivations

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x)$$

Stochastic optimization: sequence of *random* variables $(x_n)_{n \in \mathbb{N}}$ s.t.

$$f(x_n) \xrightarrow{\mathbb{E}} \inf_{x \in \mathcal{H}} f(x) \quad \text{and possibly} \quad x_n \xrightarrow{\mathbb{E}} \hat{x}$$

Motivations

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x)$$

Stochastic optimization: sequence of *random* variables $(x_n)_{n \in \mathbb{N}}$ s.t.

$$f(x_n) \xrightarrow{\mathbb{E}} \inf_{x \in \mathcal{H}} f(x) \quad \text{and possibly} \quad x_n \xrightarrow{\mathbb{E}} \hat{x}$$

f contains *randomness*, i.e.,

$$f(x) \equiv \mathbb{E}F(x)$$

E.g.: noisy observations

$$\{F_s(x), s = 1, \dots, S\} \text{ i.i.d.}$$

Motivations

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x)$$

Stochastic optimization: sequence of *random* variables $(x_n)_{n \in \mathbb{N}}$ s.t.

$$f(x_n) \xrightarrow{\mathbb{E}} \inf_{x \in \mathcal{H}} f(x) \quad \text{and possibly} \quad x_n \xrightarrow{\mathbb{E}} \hat{x}$$

f contains *randomness*, i.e.,

$$f(x) \equiv \mathbb{E}F(x)$$

E.g.: noisy observations

$$\{F_s(x), s = 1, \dots, S\} \text{ i.i.d.}$$

Inject randomness

- ▶ to be tractable
- ▶ to converge faster

E.g.: sum of functions $S \gg 1$

$$f(x) = \frac{1}{S} \sum_{s=1}^S F_s(x)$$

Machine learning: empirical and expected risks

Example of **binary classification**:

Machine learning: empirical and expected risks

Example of **binary classification**:

- ▶ Training set $\mathcal{T} = \{(\mathbf{u}_s, z_s), \mathbf{u}_s \in \mathbb{R}^M, z_s \in \{-1, 1\}, s = 1, \dots, S\}$

Machine learning: empirical and expected risks

Example of **binary classification**:

- ▶ Training set $\mathcal{T} = \{(\mathbf{u}_s, z_s), \mathbf{u}_s \in \mathbb{R}^M, z_s \in \{-1, 1\}, s = 1, \dots, S\}$
- ▶ Linear predictor model $z \sim \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle$, $\phi : \mathbb{R}^M \rightarrow \mathbb{R}^N$ feature map

Machine learning: empirical and expected risks

Example of **binary classification**:

- ▶ Training set $\mathcal{T} = \{(\mathbf{u}_s, z_s), \mathbf{u}_s \in \mathbb{R}^M, z_s \in \{-1, 1\}, s = 1, \dots, S\}$
- ▶ Linear predictor model $z \sim \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle$, $\phi : \mathbb{R}^M \rightarrow \mathbb{R}^N$ feature map
- ▶ Loss $F(z, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle)$ measures the prediction error

Machine learning: empirical and expected risks

Example of **binary classification**:

- ▶ Training set $\mathcal{T} = \{(\mathbf{u}_s, z_s), \mathbf{u}_s \in \mathbb{R}^M, z_s \in \{-1, 1\}, s = 1, \dots, S\}$
- ▶ Linear predictor model $z \sim \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle$, $\phi: \mathbb{R}^M \rightarrow \mathbb{R}^N$ feature map
- ▶ Loss $F(z, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle)$ measures the prediction error

Empirical risk:

error on the training set

$$\frac{1}{S} \sum_{s=1}^S F(z_s, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}_s) \rangle) = \frac{1}{S} \sum_{s=1}^S F_s(\boldsymbol{\vartheta})$$

Machine learning: empirical and expected risks

Example of **binary classification**:

- ▶ Training set $\mathcal{T} = \{(\mathbf{u}_s, z_s), \mathbf{u}_s \in \mathbb{R}^M, z_s \in \{-1, 1\}, s = 1, \dots, S\}$
- ▶ Linear predictor model $z \sim \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle$, $\phi: \mathbb{R}^M \rightarrow \mathbb{R}^N$ feature map
- ▶ Loss $F(z, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle)$ measures the prediction error

Empirical risk:

error on the training set

$$\frac{1}{S} \sum_{s=1}^S F(z_s, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}_s) \rangle) = \frac{1}{S} \sum_{s=1}^S F_s(\boldsymbol{\vartheta})$$

Expected risk or *generalization error*

error on the testing set

$$\mathbb{E}_{(\mathbf{u}, z)} F(z, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle) = f(\boldsymbol{\vartheta})$$

Machine learning: empirical and expected risks

Example of **binary classification**:

- ▶ Training set $\mathcal{T} = \{(\mathbf{u}_s, z_s), \mathbf{u}_s \in \mathbb{R}^M, z_s \in \{-1, 1\}, s = 1, \dots, S\}$
- ▶ Linear predictor model $z \sim \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle$, $\phi: \mathbb{R}^M \rightarrow \mathbb{R}^N$ feature map
- ▶ Loss $F(z, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle)$ measures the prediction error

Empirical risk:

error on the training set

$$\frac{1}{S} \sum_{s=1}^S F(z_s, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}_s) \rangle) = \frac{1}{S} \sum_{s=1}^S F_s(\boldsymbol{\vartheta})$$

Expected risk or *generalization error*

error on the testing set

$$\mathbb{E}_{(\mathbf{u}, z)} F(z, \langle \boldsymbol{\vartheta}, \phi(\mathbf{u}) \rangle) = f(\boldsymbol{\vartheta})$$

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

- ▶ no access to $\nabla f(x)$ (or to $u \in \partial f(x)$),

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

- ▶ no access to $\nabla f(x)$ (or to $u \in \partial f(x)$),
- ▶ but only to i.i.d. *unbiased* samples, i.e. $\mathbb{E}\nabla F_n(x) = \nabla f(x)$.

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

- ▶ no access to $\nabla f(x)$ (or to $u \in \partial f(x)$),
- ▶ but only to i.i.d. *unbiased* samples, i.e. $\mathbb{E}\nabla F_n(x) = \nabla f(x)$.

$$x_{n+1} = x_n - \gamma_n \nabla F_n(x_n), \quad \gamma_n = Cn^{-\alpha}, \quad \alpha \in [0, 1] \quad (\mathbf{SGD})$$

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

- ▶ no access to $\nabla f(x)$ (or to $u \in \partial f(x)$),
- ▶ but only to i.i.d. *unbiased* samples, i.e. $\mathbb{E}\nabla F_n(x) = \nabla f(x)$.

$$x_{n+1} = x_n - \gamma_n \nabla F_n(x_n), \quad \gamma_n = Cn^{-\alpha}, \quad \alpha \in [0, 1] \quad (\mathbf{SGD})$$

E.g.: Binary classification in Machine Learning

Online learning

$$\nabla F_n(\vartheta) = \nabla F(z_{s^{(n)}}), \langle \vartheta, \phi(\mathbf{u}_{s^{(n)}}) \rangle$$

for some $s^{(n)} \in \{1, \dots, S\}$.

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

- ▶ no access to $\nabla f(x)$ (or to $u \in \partial f(x)$),
- ▶ but only to i.i.d. *unbiased* samples, i.e. $\mathbb{E}\nabla F_n(x) = \nabla f(x)$.

$$x_{n+1} = x_n - \gamma_n \nabla F_n(x_n), \quad \gamma_n = Cn^{-\alpha}, \quad \alpha \in [0, 1] \quad \textbf{(SGD)}$$

E.g.: Binary classification in Machine Learning

Online learning

$$\nabla F_n(\vartheta) = \nabla F(z_{s^{(n)}}), \langle \vartheta, \phi(\mathbf{u}_{s^{(n)}}) \rangle$$

for some $s^{(n)} \in \{1, \dots, S\}$.

Batch learning

$$\nabla F_n(\vartheta) = \frac{1}{L} \sum_{\ell=1}^L \nabla F(z_{s_\ell^{(n)}}), \langle \vartheta, \phi(\mathbf{u}_{s_\ell^{(n)}}) \rangle$$

for L indices $s_\ell^{(n)}$.

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

Theorem (Polyak-Ruppert averaging)

Assume that

- ▶ F_n is *convex* and B -Lipschitz continuous on $\{\|x\| \leq D\}$,
- ▶ F_n are i.i.d. random functions satisfying $\mathbb{E}\nabla F_n(x) = \nabla f(x)$,
- ▶ \hat{x} , the global minimizer, is such that $\|\hat{x}\| \leq D$.

$$\text{Let } x_{n+1} = P_{\|x\| \leq D} \left(x_n - \frac{2D}{B\sqrt{n}} \nabla F_n(x_n) \right).$$

$$\text{Then for } \bar{x}_{n+1} \equiv \frac{1}{n} \sum_{k=0}^n x_k, \quad \mathbb{E}f(\bar{x}_n) - f(\hat{x}) \leq \frac{2DB}{\sqrt{n}}.$$

Strongly-convex functions

Let \mathcal{H} be a Hilbert space. Let $f: \mathcal{H} \rightarrow]-\infty, +\infty]$ and $\mu > 0$.

f is μ -strongly-convex if

$$(\forall x \in \text{dom } f)(\forall y \in \text{dom } f) \quad f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{\mu}{2} \|x - y\|^2.$$

Strongly-convex functions

Let \mathcal{H} be a Hilbert space. Let $f: \mathcal{H} \rightarrow]-\infty, +\infty]$ and $\mu > 0$.

f is μ -strongly-convex if

$$(\forall x \in \text{dom } f)(\forall y \in \text{dom } f) \quad f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{\mu}{2} \|x - y\|^2.$$

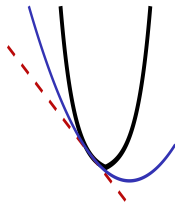


Strongly-convex functions

Let \mathcal{H} be a Hilbert space. Let $f: \mathcal{H} \rightarrow]-\infty, +\infty]$ and $\mu > 0$.

f is μ -strongly-convex if

$$(\forall x \in \text{dom } f)(\forall y \in \text{dom } f) \quad f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{\mu}{2} \|x - y\|^2.$$

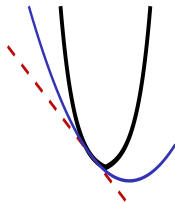


Strongly-convex functions

Let \mathcal{H} be a Hilbert space. Let $f: \mathcal{H} \rightarrow]-\infty, +\infty]$ and $\mu > 0$.

f is μ -strongly-convex if

$$(\forall x \in \text{dom } f)(\forall y \in \text{dom } f) \quad f(x) - f(y) \leq \langle \nabla f(x), x - y \rangle - \frac{\mu}{2} \|x - y\|^2.$$



Characterization of strong-convexity:

Let $f: \mathcal{H} \rightarrow]-\infty, +\infty]$ twice differentiable.

► f is μ -strongly-convex iff

$$(\forall x \in \text{dom } f), (\forall z \in \mathcal{H})$$

$$\langle z | \nabla^2 f(x) z \rangle \geq \mu \|z\|^2.$$

Stochastic (sub)gradient descent [Robbins-Monro]

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) \equiv \mathbb{E}F(x)$$

Theorem (Polyak-Ruppert averaging with **strong-convexity**)

Assume that

- ▶ F_n are i.i.d. random functions satisfying $\mathbb{E}\nabla F_n(x) = \nabla f(x)$,
- ▶ F_n is *convex* and B -Lipschitz continuous on $\{\|x\| \leq D\}$,
- ▶ f is μ -strongly-convex on $\{\|x\| \leq D\}$,
- ▶ \hat{x} , the global minimizer, is such that $\|\hat{x}\| \leq D$.

$$\text{Let } x_{n+1} = P_{\|x\| \leq D} \left(x_n - \frac{2}{\mu(n+1)} \nabla F_n(x_n) \right).$$

$$\text{Then for } \bar{x}_{n+1} \equiv \frac{2}{n(n+1)} \sum_{k=1}^n kx_{k-1}, \quad \mathbb{E}f(\bar{x}_n) - f(\hat{x}) \leq \frac{2B^2}{\mu(n+1)}.$$

Stochastic forward-backward algorithm

Let $C = \{x \in \mathcal{H} \mid \|x\| \leq D\}$, which is a *convex* subset of \mathcal{H}

Stochastic forward-backward algorithm

Let $C = \{x \in \mathcal{H} \mid \|x\| \leq D\}$, which is a *convex* subset of \mathcal{H}

Reminder

$$\text{prox}_{\iota_C} = P_C$$

Stochastic forward-backward algorithm

Let $C = \{x \in \mathcal{H} \mid \|x\| \leq D\}$, which is a *convex* subset of \mathcal{H}

Reminder $\text{prox}_{\iota_C} = P_C$

$$x_{n+1} = P_C(x_n - \gamma_n \nabla F_n(x_n))$$

Stochastic forward-backward algorithm

Let $C = \{x \in \mathcal{H} \mid \|x\| \leq D\}$, which is a *convex* subset of \mathcal{H}

Reminder

$$\text{prox}_{\iota_C} = P_C$$

$$x_{n+1} = P_C(x_n - \gamma_n \nabla F_n(x_n)) \quad \Leftrightarrow \quad \underbrace{x_{n+1} = \text{prox}_{\gamma_n \iota_C}(x_n - \gamma_n \nabla F_n(x_n))}_{\text{forward-backward iteration}}$$

Stochastic forward-backward algorithm

Let $C = \{x \in \mathcal{H} \mid \|x\| \leq D\}$, which is a *convex* subset of \mathcal{H}

Reminder $\text{prox}_{\iota_C} = P_C$

$$x_{n+1} = P_C(x_n - \gamma_n \nabla F_n(x_n)) \quad \Leftrightarrow \quad \underbrace{x_{n+1} = \text{prox}_{\gamma_n \iota_C}(x_n - \gamma_n \nabla F_n(x_n))}_{\text{forward-backward iteration}}$$

General setting: $F_s, g \in \Gamma_0(\mathcal{H})$, F_s differentiable with β -Lipschitz gradient

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) = F(x) + g(x) = \frac{1}{S} \sum_{s=1}^S F_s(x) + g(x).$$

Stochastic forward-backward algorithm

General setting: $F_s, g \in \Gamma_0(\mathcal{H})$, F_s differentiable with β -Lipschitz gradient

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) = F(x) + g(x) = \frac{1}{S} \sum_{s=1}^S F_s(x) + g(x).$$

Stochastic forward-backward algorithm

General setting: $F_s, g \in \Gamma_0(\mathcal{H})$, F_s differentiable with β -Lipschitz gradient

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) = F(x) + g(x) = \frac{1}{S} \sum_{s=1}^S F_s(x) + g(x).$$

Stochastic Proximal Gradient (SPG)

Let $\gamma_n = Cn^{-\alpha}$

Stochastic forward-backward algorithm

General setting: $F_s, g \in \Gamma_0(\mathcal{H})$, F_s differentiable with β -Lipschitz gradient

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) = F(x) + g(x) = \frac{1}{S} \sum_{s=1}^S F_s(x) + g(x).$$

Stochastic Proximal Gradient (SPG)

Let $\gamma_n = Cn^{-\alpha}$ and $(s^{(n)})_{n \in \mathbb{N}}$ i.i.d. random indexes s. t. $\mathbb{P}[s^{(n)} = s] = 1/S$.

Stochastic forward-backward algorithm

General setting: $F_s, g \in \Gamma_0(\mathcal{H})$, F_s differentiable with β -Lipschitz gradient

$$\hat{x} \in \underset{x \in \mathcal{H}}{\text{Argmin}} f(x) = F(x) + g(x) = \frac{1}{S} \sum_{s=1}^S F_s(x) + g(x).$$

Stochastic Proximal Gradient (SPG)

Let $\gamma_n = Cn^{-\alpha}$ and $(s^{(n)})_{n \in \mathbb{N}}$ i.i.d. random indexes s. t. $\mathbb{P}[s^{(n)} = s] = 1/S$.

$$x_{n+1} = \text{prox}_{\gamma_n g}(x_n - \gamma_n \nabla F_{s^{(n)}}(x_n))$$

Stochastic forward-backward algorithm

Stochastic Proximal Gradient (SPG) $F(x) = 1/S \sum_{s=1}^S F_s(x)$

Let $\gamma_n = Cn^{-\alpha}$ and $(s^{(n)})_{n \in \mathbb{N}}$ i.i.d. random indexes s. t. $\mathbb{P}[s^{(n)} = s] = 1/S$.

$$x_{n+1} = \text{prox}_{\gamma_n g}(x_n - \gamma_n \nabla F_{s^{(n)}}(x_n))$$

Theorem

Let $F_s, g \in \Gamma_0(\mathcal{H})$, $\hat{x} \in \text{Argmin}_{x \in \mathcal{H}} F(x) + g(x)$

Assume that

- ▶ F_s differentiable with β -Lipschitz gradient
- ▶ F or g strongly-convex,
- ▶ $\exists \sigma > 0, \exists \eta_n > 0$ s.t. $\|\nabla F_s(x) - \nabla F(x)\|^2 \leq \sigma^2 (1 + \eta_n \|\nabla F(x)\|^2)$,
- ▶ $\exists \varepsilon > 0$ s.t. $\gamma_n \leq (1 - \varepsilon)\beta^{-1} (1 + 2\sigma^2\eta_n)^{-1}$.

Then $\mathbb{E}\|x_n - \hat{x}\|^2 = \mathcal{O}(n^{-\alpha})$ (for $\alpha = 1$, need well-chosen C).