

Part VII - Stochastic optimization

1 Intro

Database $S = \{(u_i, z_i) \in \mathcal{X} \times \mathcal{Y} \mid \forall i = 1, \dots, n\}$

→ Assumption: prediction = linear function
 $P_x(u_i) = \langle x, \phi(u_i) \rangle$

$$\min_x \underbrace{\frac{1}{n} \sum_{i=1}^n l(z_i, \langle x, \phi(u_i) \rangle)}_{F: \mathbb{R}^d \rightarrow \mathbb{R}} + \lambda \Omega(x)$$

(if $\phi: \mathcal{X} \rightarrow \mathbb{R}^d$)

→ The predictor P_x is chosen to minimize the empirical risk over a parametrized set of predictors, potentially with regularization

→ F : objective function composed with empirical risk + regularization.

→ no closed form for the minimizer of F .

2 Empirical risk versus Expected risk.

→ Empirical risk = Training cost
 $\frac{1}{n} \sum_{i=1}^n l(z_i, P_x(u_i))$

← the quantity
you can observe

→ Expected risk = Testing cost
 $\mathbb{E}_{(u, z)} \{ l(z, P_x(u)) \}$

← the quantity you
really care about
"data not seen"

Goal: minimize the training objective but the error of unseen data

Remark : Two fundamental questions.

① • Compute $\hat{x} \in \text{Argmin } F(x)$

↙ optimization question

② • Analyze \hat{x} : guarantee $F(x)$ is good to minimize the expected risk.

↖ statistical question

3 * Stochastic gradient descent (SGD)

→ Standard gradient descent requires to compute the full gradient ∇F .

Limitations : can be costly
require to access the entire data

→ Alternative : compute unbiased stochastic approximation of the gradient. $g_k(x_k)$
such that :

$$\mathbb{E}\{g_k(x_{k+1}) | x_{k+1}\} = \nabla F(x_{k+1})$$

Algo : Set $(\gamma_k)_{k \geq 0}$

Set $x_0 \in \mathbb{R}^d$

$$(\forall k \geq 0) \quad x_k = x_{k-1} - \gamma_k \nabla g_k(x_{k-1})$$

↳ Batch gradient
= all dataset
is considered
For every
step of
the gradient

→ SGD for empirical risk minimization

$$F(x) = \frac{1}{n} \sum_{i=1}^n \ell(z_i, p_x(u_i)) \quad F_i(x)$$

At each iteration, choose uniformly random $i(k) = \{1, \dots, n\}$ and

$$g_k(x_k) = \nabla F_{i(k)}(x_k)$$

Remark: mini-batch variant

= Fixed number of training examples which is less than the dataset

x at each iteration, averaged gradient over a random subset of indices

$$x_{k+1} = x_k - \frac{\eta_k}{n} \sum_{i \in \mathbb{I}} \nabla F_i(x_k)$$

↳ Variance reduction of the gradient estimate
but use more gradient so increase running time

Remark: To avoid overfitting { regularization term
early stopping.

→ SGD for expected risk minimization

* Expected risk = Testing cost ← the quantity you really care about "data not seen"
 $\mathbb{E}_{(u,z)} \{ \ell(z, p_x(u)) \}$

* Stochastic approximation: we assume to observe a noisy version of the gradient = random pair of observations.

expectation of unseen data and we just see samples

$$F(x) = \mathbb{E} \{ F_i(x) \} = \mathbb{E} \{ \ell(z_i, p_x(u_i)) \}$$

expected gradient

$$\rightarrow \nabla F(x) = \mathbb{E} \{ \nabla F_i(x) \}$$

(Rq: swap $\mathbb{E}\{\cdot\}$ and ∇ unbiasedness)
↳ pass otherwise it creates dependencies.

4. Convergence of SGD (Robbins-Monro algo.)

Ass: F convex and B -Lipschitz
 F admits a minimizer \hat{x} s.t. $\|x_0 - \hat{x}\|_2 \leq D$
 Unbiased gradient $\mathbb{E}\{g_k(x_{k-1}) | x_{k-1}\} = \nabla F(x_{k-1})$
 Bounded gradient $\|g_k(x_{k-1})\|_2^2 \leq B^2 \quad \forall k$ almost everywhere.

Set $\gamma_k = \frac{D}{B\sqrt{k}}$, SGD of F satisfies

$$\mathbb{E}\{F(\bar{x}_k) - F(x_0)\} \leq DB \frac{2 + \log(L)}{\sqrt{k}}$$

where $\bar{x}_k = \left(\sum_{s=1}^k \gamma_s x_{s-1}\right) / \sum_{s=1}^k \gamma_s$.

Proof: $\mathbb{E}\{\|x_k - \hat{x}\|_2^2\} = \mathbb{E}\{\|x_{k-1} - \gamma_k g_k(x_{k-1}) - \hat{x}\|_2^2\}$
 $= \mathbb{E}\{\|x_{k-1} - \hat{x}\|_2^2\} + \gamma_k^2 \mathbb{E}\{\|g_k(x_{k-1})\|_2^2\}$
 $- 2\gamma_k \mathbb{E}\{g_k(x_{k-1})^\top (x_{k-1} - \hat{x})\}$

$$\mathbb{E}\{g_k(x_{k-1})^\top (x_{k-1} - \hat{x})\} = \mathbb{E}\{\mathbb{E}\{g_k(x_{k-1})^\top (x_{k-1} - \hat{x}) | x_{k-1}\}\}$$

$$= \mathbb{E}\{\mathbb{E}\{g_k(x_{k-1})^\top | x_{k-1}\}^\top (x_{k-1} - \hat{x})\}$$

$$= \mathbb{E}\{\nabla F(x_{k-1})^\top (x_{k-1} - \hat{x})\}$$

$$\mathbb{E}\{\|x_k - \hat{x}\|_2^2\} \leq \mathbb{E}\{\|x_{k-1} - \hat{x}\|_2^2\} - 2\gamma_k \mathbb{E}\{\nabla F(x_{k-1})^\top (x_{k-1} - \hat{x})\} + \gamma_k^2 B^2$$

$$\gamma_k \mathbb{E}\{\nabla F(x_{k-1})^\top (x_{k-1} - \hat{x})\} \leq \frac{1}{2} \mathbb{E}\{\|x_{k-1} - \hat{x}\|_2^2\} - \frac{1}{2} \mathbb{E}\{\|x_k - \hat{x}\|_2^2\} + \frac{1}{2} \gamma_k^2 B^2$$

By convexity analysis $F(x_{n-1}) - F(\bar{x}) \leq \nabla F(x_{n-1})^T (x_{n-1} - \bar{x})$

• F convex $\Rightarrow \cdot \leq -$

$$\wedge f(x) + (1-\alpha)f(y) \geq f(\alpha x + (1-\alpha)y)$$

$$f(x) - f(y) \geq \frac{f(\alpha x + (1-\alpha)y) - f(y)}{\alpha}$$

$\lim_{\alpha \rightarrow 0, \alpha \neq 0}$

$$f(x) - f(y) \geq \langle \nabla f(y), x - y \rangle$$

$$f(y) - f(x) \leq \langle \nabla f(y), y - x \rangle$$

• F convex $\Leftarrow \cdot \leq -$

$$(x, y) \in \text{dom } F \times \text{dom } F$$

$$\alpha \in [0, 1] \quad \alpha x + (1-\alpha)y \in \text{dom } F$$

$$(1) \quad f(x) \geq f(\alpha x + (1-\alpha)y) + \langle \nabla f(\alpha x + (1-\alpha)y), \alpha x + (1-\alpha)y - x \rangle$$

$$(2) \quad f(y) \geq f(\alpha x + (1-\alpha)y) + \langle \nabla f(\alpha x + (1-\alpha)y), \alpha(y-x) \rangle$$

$$(1) + (2) \quad \alpha f(x) + (1-\alpha)f(y) \geq f(\alpha x + (1-\alpha)y)$$

$$\delta_n \mathbb{E} \{ F(x_{n-1}) - F(\bar{x}) \} \leq \frac{1}{2} \mathbb{E} \{ \|x_{n-1} - \bar{x}\|^2 \} + \frac{1}{2} \mathbb{E} \{ \|x_n - \bar{x}\|^2 \} + \frac{1}{2} \delta_n^2 B^2$$

$$\frac{1}{\sum_{s=1}^t \gamma_s} \sum_{s=1}^t \gamma_s \mathbb{E} \{ F(x_{s-1}) - F(\bar{x}) \} \leq \frac{1}{2 \sum_{s=1}^t \gamma_s} \mathbb{E} \{ \|x_0 - \bar{x}\|^2 \} + \frac{\sum_{s=1}^t \gamma_s^2 B^2}{2 \sum_{s=1}^t \gamma_s}$$

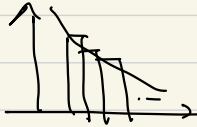
larger than $\min_{0 \leq s \leq t-1} F(x_s) - F(\bar{x})$

and than $F(\bar{x}_k) - F(\bar{x})$ with $\bar{x}_k = \frac{\sum_{s=1}^k \gamma_s x_{s-1}}{\sum_{s=1}^k \gamma_s}$
(Jensen Inequality)

goes to 0 if $\frac{1}{\sum_{s=1}^t \gamma_s} \rightarrow 0$
 $\delta t \rightarrow 0$

Recall : Series Integral comparisons

$$\bullet \sum_{s=1}^t \frac{1}{\sqrt{s}} \geq \int_0^t \frac{ds}{\sqrt{s+3}} = \left[2\sqrt{s+1} \right]_0^t = 2\sqrt{t+1} - 2 \geq \frac{1}{2}\sqrt{t}$$



$$\bullet \sum_{s=1}^t \frac{1}{s} \leq 1 + \sum_{s=2}^t \frac{1}{s} \leq 1 + \int_1^t \frac{ds}{s} = 1 + \log(t)$$

Choose $\sigma_s = \frac{\sigma}{\sqrt{s}}$ $\sigma > 0$

$$\begin{aligned} \mathbb{E} \{ F(\bar{x}_n) - F(\bar{x}) \} &\leq \frac{1}{2 \sum_{s=1}^n \frac{\sigma}{\sqrt{s}}} \mathbb{E} \{ \|x_0 - \bar{x}\|^2 \} + \frac{\sum_{s=1}^n \frac{\sigma^2}{s} B^2}{2 \sum_{s=1}^n \frac{\sigma}{\sqrt{s}}} \\ &\leq \frac{1}{\sigma \sqrt{n}} \mathbb{E} \{ \|x_0 - \bar{x}\|^2 \} + \frac{\sigma B^2 \sum_{s=1}^n \frac{1}{s}}{\sqrt{n}} \\ &\leq \frac{1}{\sqrt{n}} \left(\frac{D^2}{\sigma} + \sigma B^2 (1 + \log(n)) \right) (1 + \log(n)) \end{aligned}$$

Setting $\sigma = \frac{D}{B}$ \Rightarrow Result proved.

Remarks

- Similar proofs for subgradient descent as relying on $F(u) \leq F(v) + \eta(u-v)$ with $\eta \in \partial F(u)$
- Work with projected gradient descent as using non-expansivity of P_C .
(First steps of the proof).
- convergence of the gradient descent

$$F(x_n) - F(\bar{x}) \leq \frac{L}{2k} \|x_0 - \bar{x}\|^2$$

→ faster than stochastic.

Convergence of SGD in the strongly convex case

Ass: F convex, B Lipschitz, $G = F + \frac{\mu}{2} \|\cdot\|^2$ admits a unique minimizer \bar{x} .

, unbiased gradient

, bounded gradient.

Set $\delta_k = \frac{1}{\mu k}$, the iterates $(x_k)_{k \geq 0}$ of SGD satisfy.

$$\mathbb{E}\{F(\bar{x}_k) - F(\bar{x})\} \leq \frac{2B^2}{\mu k} (1 + \log k).$$

Results Bach & Moulines 2011

→ dedicated to stochastic gradient descent with learning rate $\gamma_n = C n^{-\alpha}$

→ strongly smooth case.

$\alpha = O(k^{-1})$ for $\alpha = 1$ without averaging

$\alpha = O(k^{-1})$ for $\alpha \in [\frac{1}{2}, 1]$ with averaging

α Robust to the choice of C .

5. Stochastic Proximal Gradient [Rosasco-Villa-Vu 2014]

Ass $\therefore (\delta_n)_{n \in \mathbb{N}}$ positive sequence

$(d_n)_{n \in \mathbb{N}}$ a sequence in $[0, 1]$

$(G_n)_{n \in \mathbb{N}}$ be a \mathbb{H} -valued random process s.t $\mathbb{E}\{\|G_n\|^2\} < +\infty$

Fix x_n a \mathbb{H} -valued integrable vector with $\mathbb{E}\{\|x_n\|^2\} < +\infty$

and $\forall n \in \mathbb{N}$

$$\begin{cases} z_n = x_n - \delta_n G_n \\ y_n = \text{prox}_{\delta_n R}(z_n) \\ w_{n+1} = (1-d_n) x_n + d_n y_n \end{cases}$$

• Si $G_n = \nabla F(w_n)$, we get FBS. [Combettes, Wajs 2004]

f , convex, diff, β -Lipschitz continuous gradient. $h \in \Gamma_0(\mathcal{X})$

Ass: $\mathbb{E}\{G_k | \mathcal{F}_k\} = \nabla F(x_k)$ ← unbiased estimate of the gradient

weaker condition than

$\exists B > 0$ and $\alpha_k > 0$ such that $\mathbb{E}\{\|G_k - \nabla F(x_k)\|^2 | \mathcal{F}_k\}$

$$\|g_k(x_{k-1})\|_2^2 \leq B^2$$

$$\leq B^2(1 + \alpha_k \|\nabla F(x_k)\|^2)$$

$$\exists \varepsilon > 0 \text{ s.t. } \forall k > 0 \quad 0 < \gamma_k < \frac{1 - \varepsilon}{\beta(1 + 2\sigma^2\alpha_k)}$$

Under technical assumptions.

For any solution \hat{x} of the problem $\min f(x) + h(x)$, set
 $(\forall k \in \mathbb{N}^*) \quad \chi_k^2 = \frac{1}{k} \gamma_k^2 (1 + 2\alpha_k \|\nabla f(\hat{x})\|^2)$

$$\text{Assume that } \begin{cases} \sum_{k \in \mathbb{N}^*} \frac{1}{k} \gamma_k = +\infty \\ \sum_{k \in \mathbb{N}^*} \chi_k^2 < +\infty \end{cases}$$

⋮

$$\mathbb{E}\{\|x_n - \hat{x}\|^2\} = \begin{cases} O(k^{-\theta}) & \text{if } \theta \in]0, 1[\\ O(k^{-c}) + O(k^{-1}) & \text{if } \theta = 1 \end{cases}$$

→ Extension de Bach et Moulines 2011.

- Dropout : deactivate neurons output i.e components of $x = (x_j)_{1 \leq j \leq d}$.

It is done randomly.

Each neuron being possibly deactivated during one learning iteration, this forces each unit to correctly learn independently from the others.

It may help to accelerate the learning.