



Tutorial 5 : Stochastic gradient descent and theoretical properties

Exercise 1 (Optimization problems satisfy the hypothesis). Prove that the stochastic gradient descent of the following cases satisfy our hypothesis on g_k , i.e., $\mathbb{E}[g_k] = \nabla \mathcal{L}$ and $\text{Var}[g_k] \leq \sigma^2$ for some fixed constant $\sigma > 0$.

1. Gradient perturbed by Gaussian noise, i.e., $g_k(\theta) = \nabla \mathcal{L}(\theta) + \epsilon_k$ and $\epsilon_k \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$.
2. We consider the optimization problem of a linear regression problem of the form:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \underbrace{\log(1 + \exp(-y_i x_i^\top \theta))}_{\ell_i(\theta)}. \quad (1)$$

and $g_i(x) = \nabla \ell_i(x)$, $i \sim \mathcal{U}(\{1, \dots, n\})$ (i is sampled uniformly from the set $\{1, \dots, n\}$).

Exercise 2 (Reduce the variant in stochastic gradient descent). We saw that the variance of $\sigma^2 = \text{Var}[\|g_k\|^2]$ plays an important role in the analysis of stochastic gradient descent. The smaller is σ , the better the bound is. In reality, practitioners employ many techniques to reduce σ . Minibatching is one of them. Indeed, consider an optimization problem of the form:

$$\text{Minimize}_{\theta \in \mathbb{R}^d} \quad \mathcal{L}(\theta) = \sum_{i=1}^n \ell_i(\theta).$$

Instead of taking $g_k = \nabla \ell_i(\theta_k)$, minibatching takes:

$$\tilde{g}_k^2 = \frac{1}{|S|} \sum_{i \in S} \nabla \ell_i(\theta_k).$$

where $S \subseteq \{1, \dots, n\}$ is uniformly sampled from $\{1, \dots, n\}$ with replacement. Prove that:

$$\text{Var}[\|\tilde{g}_k\|^2] = \frac{1}{|S|} \text{Var}[\|g_k\|^2].$$

Exercise 3 (Stochastic gradient descent with convex and L -smooth functions). Consider the stochastic gradient descent when optimizing a convex and L -smooth function f . Assume that the stochastic gradient g_k satisfy the hypothesis as in the lecture. We want to prove that if the step size $\alpha_k = \alpha < \frac{1}{2L}$, then:

$$\mathbb{E}[f(\bar{x}^T) - \inf f] \leq \frac{\|x_0 - x^*\|^2}{T\alpha} + \sigma^2 \alpha.$$

where $\bar{x}^T = \frac{1}{T} \sum_{i=0}^{T-1} x_i$. Follow these steps:

1. Prove that:

$$\mathbb{E}[\|x_{k+1} - x^*\|^2 | x_k] = \|x_k - x^*\|^2 - 2\alpha \nabla f(x_k)^\top (x_k - x^*) + \alpha^2 (\|\nabla f(x_k)\|^2 + \sigma^2)$$

2. Prove that:

$$\mathbb{E}[\|x_{k+1} - x^*\|^2 | x_k] \leq \|x_k - x^*\|^2 + 2\alpha(1 - L\alpha)(f(x_k) - f(x^*)) + \alpha^2 \sigma^2$$

3. Conclude.

4. Optimizing α to get the best bound for the timestep K .

Exercise 4 (Choice of step-size). Consider the same setting as the previous exercise. However, this time, we will choose α_k differently and $\alpha_k < \frac{1}{2L}, \forall k \in \mathbb{N}$. In that case, the previous result becomes:

$$\mathbb{E}[f(\bar{x}^T) - \inf f] \leq \frac{\|x_0 - x^*\|^2}{\sum_i \alpha_i} + \sigma^2 \frac{\sum_i \alpha_k^2}{\sum_i \alpha_i}$$

where $\bar{x}^T = \frac{\sum_{i=0}^{T-1} x_i \alpha_i}{\sum_i \alpha_i}$ (weighted average).

1. Using the same schema, reprove the above theoretical guarantee.

2. The Robbins-Monro step-sizes are sequences of $\{\alpha_k\}_{k \in \mathbb{N}}$ satisfies:

(a) $\sum_k \alpha_k^2 < \infty$.

(b) $\sum_k \alpha_k$ is not bounded.

Prove that this choice allows $\mathbb{E}[f(\bar{x}^T) - \inf f] \rightarrow 0$.

3. For which value of β that $\alpha_k = k^{-\beta}$ satisfies Robbins Monro criteria.

4. Which values of β yields the best asymptotic convergence rate.

Exercise 5 (Subgaussian distributions). In mathematical analysis of stochastic optimization algorithms, many works assume that the noise ϵ_k belongs to certain families of distribution. Subgaussian, a generalization of Gaussian distribution is among the most popular. In this exercise, we take a look at this family. Prove the following equivalent definition of a subgaussian distribution: Let X be a random variable:

1. There exists $K_1 > 0$ such that:

$$\mathbb{P}(|X| \geq t) \leq 2 \exp\left(-\frac{t^2}{K_1^2}\right), \forall t \geq 0.$$

2. There exists $K_2 > 0$ such that:

$$\mathbb{E}[|X|^p]^{\frac{1}{p}} \leq K_2 \sqrt{p}, \forall p > 1.$$

3. There exists $K_3 > 0$ such that:

$$\mathbb{E} \exp\left(\frac{X^2}{K_3}\right) \leq 2.$$