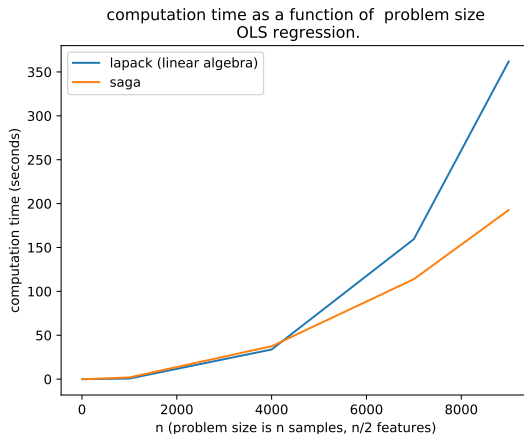


Fondamentaux théoriques du machine learning



Overview of lecture 6

Density estimation

- Motivation

- Kernel density estimation

Gradient algorithms

Gradient descent

- Convergence for least squares

- General convergence results

- Variations on gradient descent

Stochastic gradient descent

Variance reduction

Density estimation

- Motivation

- Kernel density estimation

Gradient algorithms

Gradient descent

- Convergence for least squares

- General convergence results

- Variations on gradient descent

Stochastic gradient descent

Variance reduction

Density estimation

Applications of density estimation

Kernel density estimation

Context

In machine learning, we often encounter problems in high dimension, where closed-form solutions are not available, or where even if they are available, the necessary computation time is too large.

Context

Instead, we use **iterative** algorithm such as Gradient descent (GD) or Stochastic gradient descent (SGD). SGD is the standard optimization algorithm for large-scale machine learning.

In this lecture we will study some theoretical convergence results on GD and SGD. The key properties will be convexity, strong convexity, and smoothness of the functions being optimized.

Least-squares problem

We will study ERM (Empirical risk minimization) of the OLS problem with a gradient algorithm.

- ▶ $X \in \mathbb{R}^{n,d}$
- ▶ $y \in \mathbb{R}^n$.

$$f(\theta) = \frac{1}{2n} \|X\theta - y\|_2^2 \quad (1)$$

$$\theta \leftarrow \theta - \gamma \nabla_f(\theta) \quad (2)$$

Gradient

$$f(\theta) = \frac{1}{2n} \|X\theta - y\|_2^2 \quad (3)$$

The gradient and the Hessian write :

$$\nabla_{\theta} f = \frac{1}{n} X^T (X\theta - y) \quad (4)$$

$$H = \frac{1}{n} X^T X \quad (5)$$

Minimizers

We note η^* the minimizers of f . If H is not invertible, they might be not unique, but all have the same function value $f(\eta^*)$.

All minimizers verify that

$$\nabla_{\eta^*} f = 0 \tag{6}$$

This means that

$$H\eta^* = \frac{1}{n}X^T y \tag{7}$$

Minimizers

With a Taylor expansion, we have that

$$f(\theta) - f(\eta^*) = \frac{1}{2}(\theta - \eta^*)^T H(\theta - \eta^*) \quad (8)$$

Gradient update

Exercise 1: We perform a gradient update with step size γ . t denotes the iteration number. Show that

$$\theta_t = \theta_{t-1} - \gamma H(\theta_{t-1} - \eta^*) \quad (9)$$

Gradient update

Exercise 2 : Deduce that :

$$\theta_t - \eta^* = (I - \gamma H)(\theta_{t-1} - \eta^*) \quad (10)$$

and that

$$\theta_t - \eta^* = (I - \gamma H)^t(\theta_0 - \eta^*) \quad (11)$$

Measure of performance

We can use two measures of performance of the gradient algorithm :

- ▶ Distance to minimizer :

$$\|\theta_t - \eta^*\|_2^2 = (\theta_0 - \eta^*)^T (I - \gamma H)^{2t} (\theta_0 - \eta^*) \quad (12)$$

- ▶ Convergence in function values :

$$f(\theta_t) - f(\eta^*) = \frac{1}{2} (\theta_0 - \eta^*)^T (I - \gamma H)^{2t} H (\theta_0 - \eta^*) \quad (13)$$

Distance to minimizer

If we can bound the eigenvalues of $(I - \gamma H)^{2t}$, we can bound $\|\theta_t - \eta^*\|_2^2$.

Exercise 3 : We note λ_i the eigenvalues of H . What are the eigenvalues of $(I - \gamma H)^{2t}$?

Bounding eigenvalues

We introduce the **condition number** $\kappa = \frac{K}{\mu}$.

- ▶ L is the largest eigenvalue of H .
- ▶ μ is the largest eigenvalue of H .

All eigenvalues of $(I - \gamma H)^{2t}$ have a magnitude that is smaller than

$$\left(\max_{\lambda \in [\mu, L]} |1 - \gamma \lambda| \right)^{2t} \quad (14)$$

Bounding eigenvalues

We introduce the **condition number** $\kappa = \frac{K}{\mu}$.

- ▶ L is the largest eigenvalue of H .
- ▶ μ is the largest eigenvalue of H .

All eigenvalues of $(I - \gamma H)^{2t}$ have a magnitude that is smaller than

$$\left(\max_{\lambda \in [\mu, L]} |1 - \gamma \lambda| \right)^{2t} \quad (15)$$

Hence, we want to find γ such that $\max_{\lambda \in [\mu, L]} |1 - \gamma \lambda|$ is **minimum** (or at least small).

Bounding eigenvalues

Exercise 4: Find γ such that

$$\max_{\lambda \in [\mu, L]} |1 - \gamma\lambda| \leq \left(1 - \frac{\mu}{L}\right) = \left(1 - \frac{1}{\kappa}\right) \quad (16)$$

Exponential convergence

With $\gamma = \frac{1}{L}$, we obtain an exponential convergence

$$\|\theta_t - \eta^*\|_2^2 \leq \left(1 - \frac{1}{\kappa}\right)^{2t} \|\theta_0 - \eta^*\|_2^2 \quad (17)$$

Approximation error

We have that

$$\left(1 - \frac{1}{\kappa}\right)^{2t} \leq \exp\left(-\frac{1}{\kappa}\right)^{2t} = \exp\left(-\frac{2t}{\kappa}\right) \quad (18)$$

Exercise 5 : What number of iterations is sufficient in order to have a relative reduction of $\|\theta_t - \eta^*\|_2^2$ of ϵ ?

Approximation error

$$\begin{aligned}\epsilon &\geq \exp\left(-\frac{2t}{\kappa}\right) \\ \Leftrightarrow -\log(\epsilon) &\leq \frac{2t}{\kappa} \\ \Leftrightarrow \frac{\kappa}{2} \log\left(\frac{1}{\epsilon}\right) &\leq t\end{aligned}\tag{19}$$

Large condition number

If $\kappa = +\infty$ ($\mu = 0$), we do not have a convergence guarantee.

Convergence in function values

If H is not invertible, we only have a convergence rate in $\mathcal{O}(\frac{1}{t})$.

Generalization

We will now study more general functions.

Smoothness

Definition

Smoothness

A differentiable function f with real values is said L -smooth if and only if

$$\forall x, y \in \mathbb{R}^d, |f(y) - f(x) - \nabla_x f(y - x)| \leq \frac{L}{2} \|y - x\|^2$$

Smoothness

Lemma

f is L -smooth if and only if it has L -Lipshitz continuous gradients.

Smoothness of least-squares

Exercise 6 : **Smoothness** : Consider

$$\begin{aligned}R_n(\theta) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2 \\ &= \frac{1}{n} \|Y - X\theta\|_2^2\end{aligned}$$

Is $R_n(\theta)$ smooth ?

Non smooth optimization

Lipshitz

Non smooth optimization

Exercise 7 : Smoothness : Is f has L -lipshitz continuous gradients, is f necessary Lipshitz-continuous ?

Smooth, strongly convex functions

Theorem

Convergence of GD for a strongly convex function

Let $f : \mathbb{R}^d \Rightarrow \mathbb{R}$ be a μ -strongly convex unction with L -Lipshitz continuous gradients. Let x^ be the global minimum of f (which we know exists since f is strongly convex), $x_0 \in \mathbb{R}$, $T \in \mathbb{N}$.*

With constant step size $\gamma_t = \frac{1}{L}$, we have

$$f(x_t) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^t (f(x_0) - f(x^*))$$

Optimamity

Nesterov acceleration

page 113

Line search

Conjugate gradient

Algorithmic complexity

Stochastic gradient descent

Convergence result

With averaging

Algorithmic complexity

SAGA algorithm