

Session 13: Stochastic gradient descent

Optimization and Computational Linear Algebra for Data Science

Final exam

- ❖ Scope: everything except today's lecture and this week's video.
- ❖ Of course, it will be a bit more focused on what we did after the midterm (PCA, linear regression, convex functions, optimization...)
- ❖ Same format as for the midterm
- ❖ “24 hours window” on Thursday December 17th.
- ❖ 1 hour 40 minutes to work + 20 minutes to scan + upload on Gradescope.
- ❖ In case you have any issue when uploading: **email me your work.**

Contents

1. Introduction: supervised learning
2. Stochastic gradient descent
3. Convergence analysis, comparison with gradient descent

Introduction

Supervised learning

Assume we have a pair of random variables (X, Y) that is distributed according to some distribution P_0 .

Example: $X = \text{MRI image of a brain}$
 $Y = \begin{cases} +1 & \text{if there is a cancer} \\ 0 & \text{otherwise.} \end{cases}$

Goal: Predict Y from X using a model

ϕ
↑
model's parameter.

Hope: $\phi(X) \simeq Y$

Supervised learning

Find θ that minimizes the risk:

$$\underline{\underline{R(\theta)}} = \mathbb{E} \left[\underset{\substack{\uparrow \\ \text{loss function, for instance } l(a,b) = (a-b)^2}}{l(\phi_\theta(x), y)} \right] \leftarrow \begin{array}{l} \text{expectation} \\ \text{wrt } (X,Y) \sim \underline{\underline{P_0}} \end{array}$$

ISSUE: we can not compute $R(\theta)$ because we do not know P_0 .

However, we may have access to samples

$(x_1, y_1), \dots, (x_N, y_N)$ $\stackrel{\text{iid}}{\sim} P_0$ "dataset"

minimize the empirical risk:

$$R_N(\theta) = \frac{1}{N} \sum_{i=1}^N l(\phi_\theta(x_i), y_i)$$

Supervised learning

$$R_N(\theta) = \frac{1}{N} \sum_{i=1}^N \underbrace{l(\varphi_\theta(x_i), y_i)}$$

Example: linear regression

$$l(a, b) = (a - b)^2$$

$$\varphi_\theta(x) = \langle \theta, x \rangle$$

$$R_N(\theta) = \frac{1}{N} \sum_{i=1}^N (\langle \theta, x_i \rangle - y_i)^2$$

$$= \frac{1}{N} \| X\theta - y \|^2 \rightarrow X = \begin{pmatrix} -x_1 & - \\ \vdots & \\ -x_N & - \end{pmatrix}$$

$f_i(\theta)$ "loss on the i^{th} sample only"

Why not using gradient descent ?

$$\underset{\text{red}}{R_N(\theta)} = \underset{\text{red}}{f(\theta)} = \frac{1}{N} \sum_{i=1}^N f_i(\theta). \quad \theta \in \mathbb{R}^n$$

Gradient descent iterations:

$$\begin{aligned} \theta_{t+1} &= \theta_t - \alpha_t \overbrace{\nabla f(\theta_t)} \\ &= \theta_t - \frac{\alpha_t}{N} \sum_{i=1}^N \underbrace{\nabla f_i(\theta_t)} \in \mathbb{R}^n \end{aligned}$$

Computing $\nabla f_i(\theta_t)$ takes $\simeq n$ operations.
Computing $\nabla f(\theta_t)$ takes $\simeq \underline{\underline{N \cdot n}}$ operations.

a lot too large if we have
of data.

Stochastic gradient descent

Stochastic gradient descent

$$f(\theta) = \frac{1}{N} \sum_{i=1}^N f_i(\theta).$$

Starting at some $\theta_0 \in \mathbb{R}^n$, perform the updates:

Pick i uniformly at random in $\{1, \dots, N\}$,

Update $\theta_{t+1} = \theta_t - \alpha_t \nabla f_i(\theta_t)$,

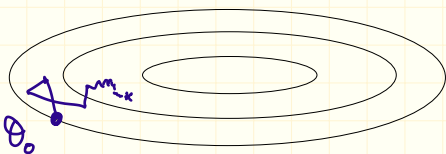
→ takes only n computations per iterations.

$\nabla f_i(\theta_t)$ can be seen as a noisy version of $\nabla f(\theta_t)$

SGD: $\theta_{t+1} = \underbrace{\theta_t - \alpha_t \nabla f(\theta_t)}_{\text{Gradient descent}} + \underbrace{\alpha_t \text{"Noise"}}_{\text{noise term.}}$

Tradeoffs in SGD

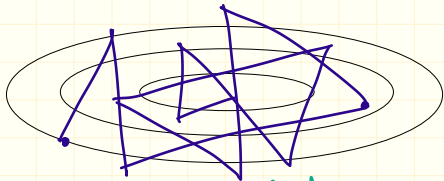
Rapidly decaying step sizes



⊕ Kills the noise

⊖ We may "never forget" the initial condition"

Slowly decaying step sizes



⊕ We move fast.

⊖ We may oscillate too much.

It is usually recommended that

- $\alpha_t \xrightarrow{t \rightarrow \infty} 0$
- $\sum_{t=0}^{+\infty} \alpha_t = +\infty$

"Kill the noise"

"Forget the initial condition"

- $\alpha_t = (0,9)^t \xrightarrow[t \rightarrow \infty]{} 0$

$$\sum_{t=0}^{+\infty} \alpha_t = \frac{1}{1-0,9} \quad \text{finite.}$$

- $\alpha_t = \frac{1}{t} \xrightarrow[t \rightarrow \infty]{} 0$

$$\sum_{t=1}^{+\infty} \frac{1}{t} = +\infty$$

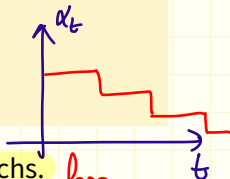
SGD in practice

Mini-batch stochastic gradient descent:

← reduces the noise

Pick a mini-batch i_1, \dots, i_k in $\{1, \dots, N\}$,

$$\text{Update } \theta_{t+1} = \theta_t - \frac{\alpha_t}{k} \sum_{m=1}^k \nabla f_{i_m}(\theta_t),$$



- ❖ Decrease the step size after a fixed number of epochs.
 - ❖ Use momentum + “adaptive gradient”: Adagrad, RMSprop, Adedelta, Adam, Adamax, Nadam...
- loop over the data set*

Excellent reference:

<https://arxiv.org/pdf/1609.04747.pdf>

Convergence analysis

Convergence rates

- if the f_i are convex and L -smooth: SGD with $\alpha_t = 1/\sqrt{t}$ achieves an error $\leq C/\sqrt{t}$.
- if the f_i are μ -strongly convex and L -smooth: SGD with $\alpha_t = 1/(\mu t)$ achieves an error $\leq C/t$.

Compare with gradient descent:

- L smooth $\rightarrow \leq \frac{\text{Constant}}{t}$
- $\left. \begin{array}{l} L\text{-smooth} \\ + \\ \mu\text{-strongly conv} \end{array} \right\} \rightarrow \leq \text{Constant} \cdot \left(1 - \frac{\mu}{L}\right)^t \leq e^{-\frac{\mu}{L}t}$

GD vs SGD

Gradient descent

Time per step $n \cdot N$
nb. of param. \rightarrow nb. of samples

Error after t steps

$$\leq c \cdot e^{-\rho t} \frac{N}{L}$$

(strongly convex case)

Log-error after τ units of time

$$-\rho \frac{\tau}{N \cdot n} + \text{cte.}$$

Stochastic gradient descent

Time per step n

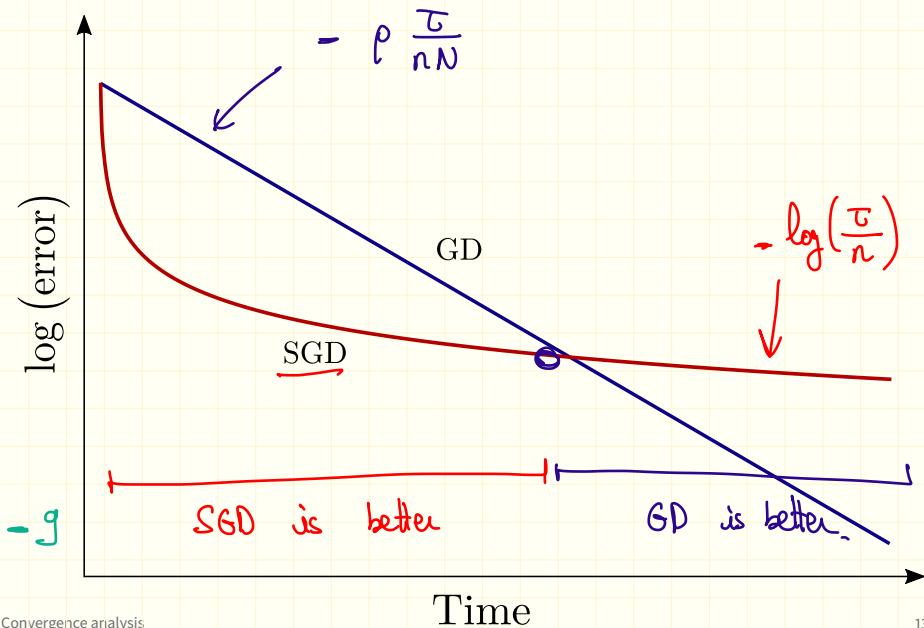
Error after t steps

$$\leq \frac{c}{t}$$

Log-error after τ units of time

$$-\log\left(\frac{\tau}{n}\right) + \text{cte.}$$

GD vs SGD



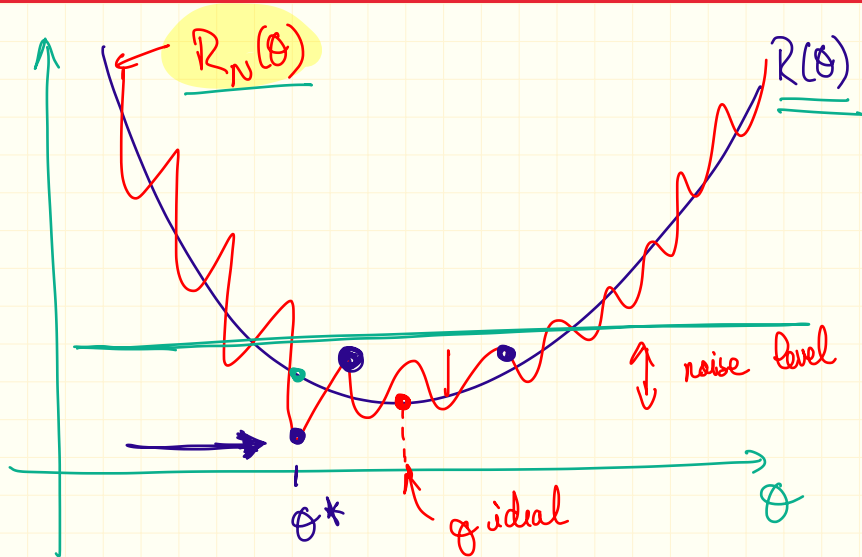
GD vs SGD: who wins ?

- ❖ If one is looking for a very small optimization error $f(\theta_t) - \min f$, then gradient descent wins.
- ❖ If one has a limited time budget and does not need a very small $f(\theta_t) - \min f$, then stochastic gradient descent wins.

IN machine learning, it makes no sense
to optimize "below the noise level"

→ SGD is a popular choice.

Questions?



Questions?

$$(u_t) : \begin{cases} u_{t+2}^{x^2} = \underbrace{u_{t+1}}_x - 2 \underbrace{u_t}_1 \\ u_0 = 1 \\ u_1 = 1 \end{cases}$$

Solve: $x^2 = x - 2$

$x^2 - x + 2 = 0$: $\Delta = 1 - 8 = -7$

solutions: $x_+ = \frac{1 + i\sqrt{7}}{2}$ $x_- = \frac{1 - i\sqrt{7}}{2}$

Questions?

Then

$$u_t = a (X_+)^t + b (X_-)^t$$

constants that needs to be determined using the initial conditions.

$$|u_t| \leq C \cdot (|X_+|^t + |X_-|^t) \leq \underline{C'} \sqrt{2}^t$$

$$|X_+| = \left| \frac{1 + i\sqrt{7}}{2} \right| = \sqrt{\frac{1}{4} + \frac{7}{4}} = \sqrt{2}$$

$$|X_-| = \sqrt{2}$$

Questions?

$$\begin{pmatrix} u_{t+2} \\ u_{t+1} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & -2 \\ 1 & 0 \end{pmatrix}}_M \cdot \begin{pmatrix} u_{t+1} \\ u_t \end{pmatrix}$$

$$u_{t+2} = u_{t+1} - 2u_t$$

$$0 = \det(M - \lambda \text{Id}) = \det \begin{pmatrix} 1-\lambda & -2 \\ 1 & -\lambda \end{pmatrix}$$
$$= -\lambda(1-\lambda) + 2 = \boxed{\lambda^2 - \lambda + 2} = 0$$

$$|u_{t+2}| \leq \sqrt{u_{t+2}^2 + u_{t+1}^2} \leq \underbrace{2C}_{\geq 0} \cdot \left(\lambda_1^{t+2} + \lambda_2^{t+2} \right)$$