

Optimization for Machine Learning

November 7, 2022

Program:

Large-scale and distributed optimization

↳ Automatic differentiation

↳ Coordinate descent

↳ Consensus optimization

Motivation:

$$\text{minimize}_{w \in \mathbb{R}^d} f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$$

Approach: Apply a stochastic gradient method because $n \gg 1$ and computing $\nabla f(w)$ is expensive

3 questions:

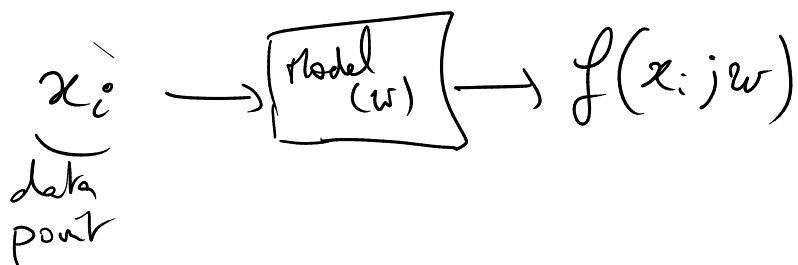
1) How do we compute $\nabla f(w)$ or $\nabla f_i(w)$ for very complex models (like neural nets)?

2) What kind of methods can we use when $d \gg 1$?

3) How do we deal with data that is distributed / on several machines?
| decentralized |

① Automatic differentiation

Complex ML model:



Q) How do we compute $\nabla f(x_i; w)$ with respect to w ?

Possibilities

→ Write down a code for $w \mapsto \nabla f(x_i; w)$

⊖ Requires time and effort

⊖ Will change every time the model changes
(e.g. add a layer to a NN architecture)

→ Use a symbolic differentiation tool (Wolfram Alpha)

⊖ Not necessarily easy to use

⊖ Need to store the formula, and regenerate it every time the architecture changes

⊖ Provides a formula for the derivative, that still needs to be evaluated

→ Use approximation formulas (finite differences)

$$\nabla f(w) \approx \left[\frac{f(w + t e_i) - f(w)}{t} \right]_{i=1}^d$$

$e_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i$, $t > 0$

⊕ Only requires evaluations of f , accuracy of the approximation depends on t

⊖ Requires $d+1$ evaluations of f for one gradient approximation (often too expensive in a deep learning context)

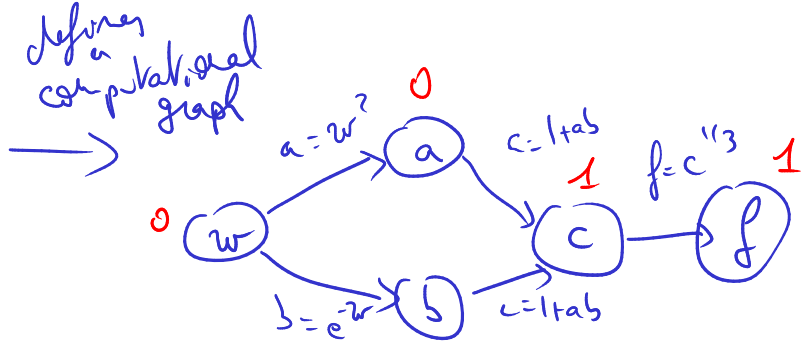
⇒ The most popular approach in Neural Network training is automatic (or algorithmic) differentiation, that uses a graph to represent how $f(x_i; w)$ is computed, in order to evaluate $\nabla f(x_i; w)$ without looking at the data point again.

Toy example $d=1$ $f(w) = (1 + w^2 e^{-w})^{1/3}$

Goal: $\nabla f(0) = \frac{\partial f}{\partial w}(0)$

1) Decompose $f(w)$ in "elementary" operations

$a := w^2$
 $b := e^{-w}$
 $c := 1 + ab$
 $f := c^{1/3}$



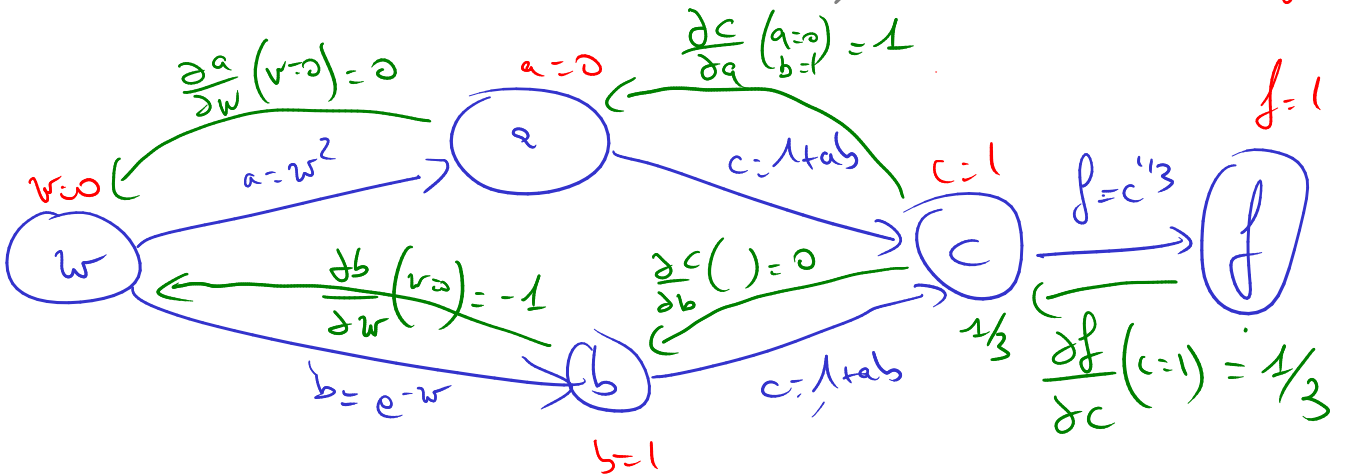
$f(0)$ can be computed by a forward pass on the graph

2) Apply the chain rule $\rightarrow c(w)$ evaluated at $w=0$

$$\frac{\partial f}{\partial w}(0) = \frac{\partial f}{\partial c}(c(0)) \left(\frac{\partial c}{\partial a}(a(0)) \frac{\partial a}{\partial w}(0) + \frac{\partial c}{\partial b}(b(0)) \frac{\partial b}{\partial w}(0) \right)$$

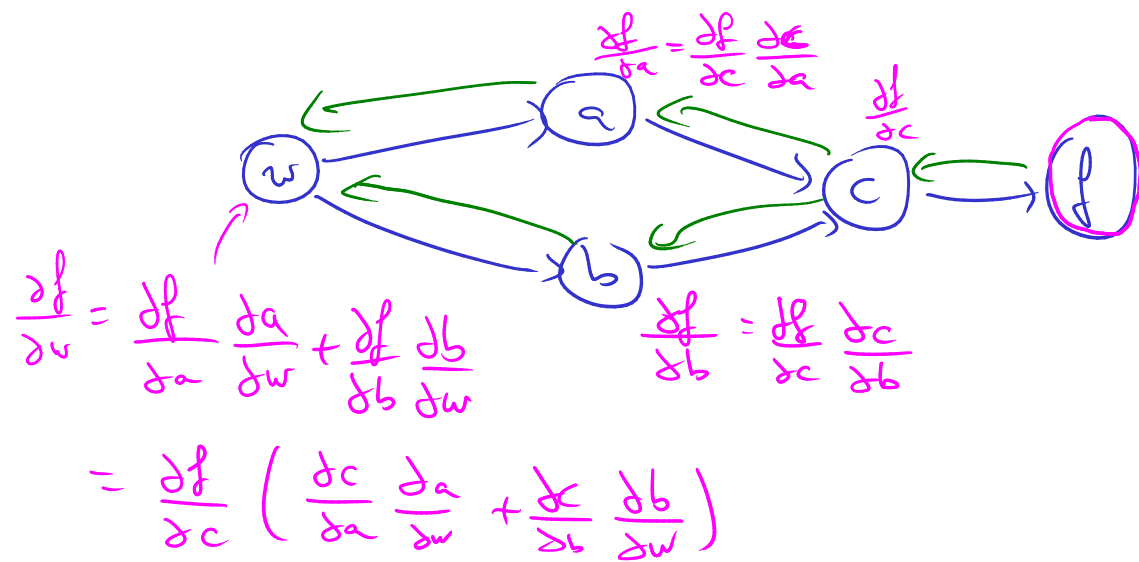
f as a function of c

to go backwards in the graph $\rightarrow \frac{\partial(1+a)}{\partial a}(a=0) = 1$ $f(0) = 1$



— Forward pass

— Backward pass \Rightarrow For every node "backwards" of f , we store the derivative of f w.r.t. that node



Takeaway: In modern ML software (PyTorch, JAX, ...), gradients/derivatives are computed automatically using computational graphs and the chain rule.
 \Rightarrow Also true for generalized derivatives

② Large-scale optimization

minimize $f(w)$ with $d \gg 1$ $f \in \mathcal{C}^1$
 $w \in \mathbb{R}^d$

Setup: Performing calculations with d -dimensional vectors is expensive

Approach: Perform optimization steps involving a small number of coordinates per step \Rightarrow Coordinate descent methods

Basic coordinate descent method ($w_0 \in \mathbb{R}^d$)

Iteration k : \rightarrow Pick a coordinate $j_k \in \{1, \dots, d\}$.
 \rightarrow Define $w_{k+1} = w_k - \alpha_k \nabla_{j_k} f(w_k) e_{j_k}$ for $\alpha_k > 0$
 j_k \uparrow j_k th coordinate of the gradient
 e_{j_k} \rightarrow j_k th coordinate

- Changes one coordinate per iteration
 ⇒ Not all coordinates of w_k are necessary to compute w_{k+1} (only $[w_k]_{j_k}$ gets updated)
- Interesting when $\nabla_{j_k} f(w_k)$ can be computed without having to access all coordinates of w_k

vector $\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}_{j_k}$

Ex) minimize $w \in \mathbb{R}^d$ $\frac{1}{2m} \sum_{i=1}^m (x_i^T w - y_i)^2 + \lambda \sum_{j=1}^d [w]_j^2$

Linear least squares with ℓ_2 regularization

separable: every term depends on its own coordinate

If the x_i 's ($\in \mathbb{R}^d$) are sparse, then the gradient of this function with respect to a particular coordinate can be computed without accessing all coordinates ("partial separability")

⇒ The cost of computing $\nabla_{j_k} f(w)$ is proportional to the number of x_i 's that have a non-zero j^{th} coordinate

Variants of coordinate descent depend on how $\{j_k\}$ are chosen:

→ Cyclic: cycling through $\{1, \dots, d\}$

($j_0=1, j_1=2, \dots, j_{d-1}=d, j_d=1, \dots$)

has more guarantees than the previous one → cycling through a random permutation of $\{1, \dots, d\}$ ⇒ randomized cyclic

NR. A cycling approach through $\{1, \dots, d\}$ may not work (counter-example due to Powell)

→ Randomized: draw j_k at random (e.g. w. family) in $\{1, \dots, d\}$ with replacement

\Rightarrow Nice theoretical guarantees

\hookrightarrow Can also use blocks of coordinates instead of a single one (block coordinate descent)

Remark (P) minimize $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$

Stochastic gradient: $w_{k+1} = w_k - \alpha_k \nabla f_{i_k}(w_k)$

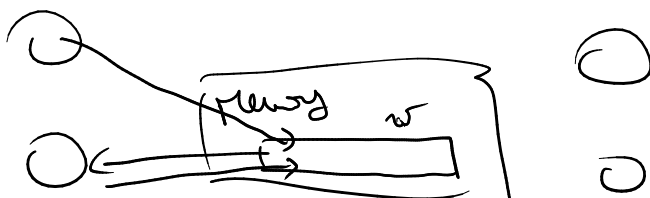
Coordinate descent: $w_{k+1} = w_k - \alpha_k \nabla_{j_k} f(w_k) e_{j_k}$

random data point
random coordinate
 \equiv random parameter among the d parameters represented by w

Coordinate descent actually corresponds to stochastic gradient on a problem related to (P), which is a problem in n variables called the dual problem

(3) Distributed and decentralized optimization

Distributed: \rightarrow Several processors performing calculation
 \rightarrow shared memory \Rightarrow All processors access the same location for the parameter $w \in \mathbb{R}^d$



Coordinate descent methods are designed for distributed setups:

→ At every iteration, every processor performs a coordinate update on its own coordinate

⇒ Efficient, parallel implementation of coordinate descent

→ Can be done in a synchronized way

($w_k \rightarrow$ every processor reads and updates its own coordinate $\rightarrow w_{k+1}$)

or in an asynchronous way

(Every processor runs iterations of coordinate descent and does not wait for updates from the other processors)

⇒ The latter approach is the most efficient and has been used in connection with stochastic gradient (Hogwild!, 2011)

Decentralized : → "Federated Learning" paradigm
→ Also called consensus optimization, network optimization

Idea: → No centralized calculations or memory

→ Instead, a set of agents with their own piece of the data are trying to solve an optimization problem that involves the entire dataset

→ The agents can communicate using a communication graph, and are trying to

reach consensus.

Mathematically:

→ Graph: (V, E)

V set of vertices
 $i=1 \dots n$

E set of edges

V represents the agents with their own portion of data

$\{(i, j) \mid i, j \in \{1, \dots, n\}\}$

E represents the communication between two agents
direct

Optimization problem

minimize $f(w) = \frac{1}{n} \sum_{i=1}^n f_i(w)$, where every f_i is only available to agent i
 $w \in \mathbb{R}^d$

One vector of parameters per agent

⇒ minimize $w^{[1]} \in \mathbb{R}^d$
:
 $w^{[n]} \in \mathbb{R}^d$

$\frac{1}{n} \sum_{i=1}^n f_i(w^{[i]})$
known to agent i

s.t. $w^{[i]} = w^{[j]}$
 $\forall (i, j) \in E$

Constraints: every agent must agree with its neighbors

Consensus optimization problem

⇒ Specific case of an optimization problem with linear constraints, that can be tackled with dual optimization methods like ADMM

General dual approach for linear constraints

↳ Consider minimize $f(w)$ s.t. $Aw = b$
where f is C^1 $w \in \mathbb{R}^d$ $A \in \mathbb{R}^{m \times d}$, $b \in \mathbb{R}^m$

↳ this problem can be reformulated as an unconstrained problem:

$$\text{minimize}_{w \in \mathbb{R}^d} \left[\max_{z \in \mathbb{R}^m} L(w, z) := f(w) + z^T (Aw - b) \right]$$

\uparrow dual variables \downarrow Lagrangian function of the problem

↳ The dual problem is

$$\text{maximize}_{z \in \mathbb{R}^m} \left[\min_{w \in \mathbb{R}^d} L(w, z) \right]$$

⊕ Always equivalent to a convex minimization problem ($z \mapsto -[\min_{w \in \mathbb{R}^d} L(w, z)]$ is convex)

⊕ Can rely on subgradients (implicit in the methods below) to define algorithms

Idea: One can solve the dual problem (in z) to obtain a solution to the original problem by mixing unconstrained optimization steps and dual steps.

Algorithm: Dual ascent method ($w_0 \in \mathbb{R}^d, z_0 \in \mathbb{R}^m$)

$$w_{k+1} \in \underset{w}{\operatorname{argmin}} L(w, z_k) = f(w) + z_k^T (Aw - b)$$

≈ subgradient step $z_{k+1} = z_k + \alpha_k (Aw_{k+1} - b)$ where $\alpha_k > 0$

⚠ The minimization problem may not have a solution
 \Rightarrow One can use augmented Lagrangian (instead of Lagrangian) to guarantee that the subproblem has a (unique) solution.

Augmented Lagrangian ($w_0 \in \mathbb{R}^d$, $z_0 \in \mathbb{R}^m$, $\lambda_0 > 0$)

• $w_{k+1} \in \arg \min_w L(w, z_k, \lambda_k) := f(w) + z_k^T (Aw - b) + \frac{\lambda_k}{2} \|Aw - b\|^2$
 "regularization" of the Lagrangian \nearrow

• $z_{k+1} = z_k + \lambda_k (Aw_{k+1} - b)$

• Set $\lambda_{k+1} \geq \lambda_k$] Typical choice, may also work with $\lambda_k \rightarrow 0$

⊕ Good convergence guarantees to a solution of the original, unconstrained optimization problem

Remaining issue: How can we handle decentralized optimization with that method, given that every agent has its own w ?

\Rightarrow Answer: ADMM (Alternated Direction Method of Multipliers)

ADMM with two "blocks"

Pb: minimize $f_1(w_{[1]}) + f_2(w_{[2]})$ s.t. $A_1 w_{[1]} + A_2 w_{[2]} = b$
 $w_{[1]} \in \mathbb{R}^d$
 $w_{[2]} \in \mathbb{R}^d$
 $A_1 \in \mathbb{R}^{m \times d}$
 $A_2 \in \mathbb{R}^{m \times d}$

ADMM is a rewriting of augmented Lagrangian where the updates for $w^{[1]}$ and $w^{[2]}$ are done in a separate fashion

Algo ($w_{0,[1]} \in \mathbb{R}^d$, $w_{0,[2]} \in \mathbb{R}^d$, $z_0 \in \mathbb{R}^m$, $\lambda_0 > 0$)

- $w_{k+1,[1]} \in \operatorname{argmin}_w L(w, w_{k+1,[2]}, z_k, \lambda_k)$
- $w_{k+1,[2]} \in \operatorname{argmin}_w L(w_{k+1,[1]}, w, z_k, \lambda_k)$
- $z_{k+1} = z_k + \lambda_k (A_1 w_{k+1,[1]} + A_2 w_{k+1,[2]} - b)$
- Choose λ_{k+1}

⇒ Can be extended to more than 2 agents!

$$L(w^{[1]}, w^{[2]}, z, \lambda) = f_1(w^{[1]}) + f_2(w^{[2]}) + z^T (A_1 w^{[1]} + A_2 w^{[2]} - b) + \frac{\lambda}{2} \|A_1 w^{[1]} + A_2 w^{[2]} - b\|^2$$