

AI505  
Optimization

# Optimization in Machine Learning

Marco Chiarandini

Department of Mathematics & Computer Science  
University of Southern Denmark

# Simplified Notation

Let  $\xi$  be a random seed or the realization of a single (or a set of) sample  $(\mathbf{x}, y)$ .

For a given  $(\mathbf{w}, \xi)$  let  $f(\mathbf{w}; \xi)$  be the composition of the loss function  $L$  and the prediction function  $h$

Then:

$$R(\mathbf{w}) = \mathbb{E}_{\xi}[f(\mathbf{w}; \xi)] \quad \text{Expected Risk}$$

Let  $\{\xi_{[i]}\}_{i=1}^n$  be realizations of  $\xi$  corresponding to  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and  $f_i(\mathbf{w}) \stackrel{\text{def}}{=} f(\mathbf{w}; \xi_{[i]})$

Then:

$$R_n(\mathbf{w}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{w}) \quad \text{Empirical Risk}$$

# Stochastic vs Batch Optimization Methods

Reduction to minimizing  $R_n$ , with  $\mathbf{w}_0 \in \mathbb{R}^d$  given (deterministic problem)

**Stochastic Approach:** Stochastic Gradient (Robbins and Monro, 1951)

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k \nabla f_{i_k}(\mathbf{w}_k)$$

$i_k$  is chosen randomly from  $\{1, \dots, n\}$ ,  $\alpha_k > 0$ .

- very cheap iteration only on one sample.
- $\{\mathbf{w}_k\}$  is a stochastic process determined by the random sequence  $\{i_k\}$ .
- the direction might not always be a descent but if it is a descent direction in **expectation**, then the sequence  $\{\mathbf{w}_k\}$  can be guided toward a minimizer of  $R_n$ .

**Batch Approach:** batch gradient, steepest descent, full gradient method:

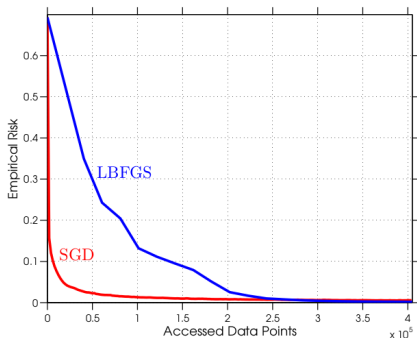
$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \alpha_k \nabla R_n(\mathbf{w}_k) = \mathbf{w}_k - \frac{\alpha_k}{n} \sum_{i=1}^n \nabla f_i(\mathbf{w}_k)$$

- more expensive
- can use all deterministic gradient-based optimization methods
- the sum structure opens up to parallelization

Analogues in simulation: stochastic approximation (SA) and sample average approximation (SAA)

# Stochastic Gradient

- In case of redundancy using all of the sample data in every iteration is inefficient
- Comparison of the performance of a batch L-BFGS method on number of evaluations of a sample gradient  $\nabla f_{i_k}(\mathbf{w}_k)$ . Each set of  $n$  consecutive accesses is called an **epoch**.
- The batch method performs only one step per epoch while SG performs  $n$  steps per epoch.



the fast initial improvement achieved by SG, followed by a drastic slowdown after 1 or 2 epochs, is common in practice

SG more sensitive to  $\alpha_k$  and starting point

if more epochs, batch may become better

# Theoretical Motivations

- a batch approach can minimize  $R_n$  at a fast rate; e.g., if  $R_n$  is strongly convex. A batch gradient method, then there exists a constant  $\rho \in (0, 1)$  such that, for all  $k \in \mathbb{N}$ , the training error follows **linear convergence**

$$R_n(\mathbf{w}_k) - R_n^* \leq \mathcal{O}(\rho^k),$$

- rate of convergence of a basic stochastic method is slower than for a batch gradient; e.g., if  $R_n$  is strictly convex and each  $i_k$  is drawn uniformly from  $\{1, \dots, n\}$ , then for all  $k \in \mathbb{N}$ , SG satisfies the **sublinear convergence property**

$$\mathbb{E}[R_n(\mathbf{w}_k) - R_n^*] = \mathcal{O}(1/k).$$

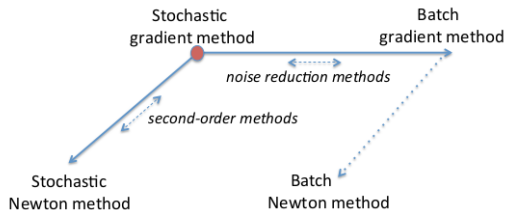
neither the per-iteration cost nor the right-hand side depends on the sample set size  $n$

- in a stochastic optimization setting, SG yields for the expected risk the same convergence rate once substituted  $\nabla f_{i_k}(\mathbf{w}_k)$  replaced by  $\nabla f(\mathbf{w}_k; \xi_k)$  with each  $\xi_k$  drawn independently according to the distribution  $P$

$$\mathbb{E}[R(\mathbf{w}_k) - R^*] = \mathcal{O}(1/k).$$

If  $n \gg k$  up to iteration  $k$  minimizing  $R_n$  same as minimizing  $R$

# Beyond SG: Noise Reduction and Second-Order Methods



- on horizontal axis methods that try to improve rate of convergence
- on vertical axis, methods that try to overcome non-linearity and ill-conditioning

**Minibatch Approach** small subset of samples, call it  $\mathcal{S}_k \subseteq \{1, \dots, n\}$ , chosen randomly in each iteration:

$$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k - \frac{\alpha_k}{|\mathcal{S}_k|} \sum_{i \in \mathcal{S}_k} \nabla f_i(\mathbf{w}_k)$$

due to the reduced variance of the stochastic gradient estimates, the method is easier to tune in terms of choosing the stepsizes  $\{\alpha_k\}$ .

dynamic sample size and gradient aggregation methods, both of which aim to improve the rate of convergence from sublinear to linear