

Presentation in *Convex Optimization*

Mingrui Zhang, Xialiang Dou, Dongming Huang

Dec 22, 2014

Sample size selection in optimization methods for machine learning

Sample size selection in optimization methods for machine learning

Main results: presents a methodology for using varying sample sizes in batch-type optimization methods for large-scale machine learning problems.

Sample size selection in optimization methods for machine learning

Main results: presents a methodology for using varying sample sizes in batch-type optimization methods for large-scale machine learning problems.

- Dynamic sample selection in the evaluation of the function and gradient.
- A practical Newton method that uses smaller sample to compute Hessian vector-products.

The dynamic sample size gradient method

Problem: To determine the values of the parameters $\omega \in \mathbb{R}^m$ of a prediction function $f(\omega; x)$, where we assume:

$$f(\omega, x) = \omega^T x \quad (1)$$

The dynamic sample size gradient method

Problem: To determine the values of the parameters $\omega \in \mathbb{R}^m$ of a prediction function $f(\omega; x)$, where we assume:

$$f(\omega, x) = \omega^T x \quad (1)$$

Common Approach: To minimize the empirical loss function:

$$J(\omega) = \frac{1}{N} \sum_{i=1}^N l(f(\omega; x_i), y_i) \quad (2)$$

where $l(\hat{y}, y)$ is a convex loss function.

The dynamic sample size gradient method

Assumption: The size of the data set N is extremely large, numbered somewhere in the millions or billions, so that the evaluation of $J(\omega)$ is very expensive.

The dynamic sample size gradient method

Assumption: The size of the data set N is extremely large, numbered somewhere in the millions or billions, so that the evaluation of $J(\omega)$ is very expensive.

A gradient-based mini-batch optimization algorithm: At every iteration, chooses a subset $\mathcal{S} \subset \{1, 2, \dots, N\}$ of the training set, and applies one step of an optimization algorithm to the objective function:

$$J_{\mathcal{S}}(\omega) = \frac{1}{|\mathcal{S}|} \sum_{i \in \mathcal{S}} l(f(\omega; x_i), y_i) \quad (3)$$

The dynamic sample size gradient method

Measure of the quality of the sample \mathcal{S} : the variance in the gradient $\nabla J_{\mathcal{S}}$

The dynamic sample size gradient method

Measure of the quality of the sample \mathcal{S} : the variance in the gradient $\nabla J_{\mathcal{S}}$

It's easy to verify that the vector $d = -\nabla J_{\mathcal{S}}(\omega)$ is a descent direction for J at ω if:

$$\delta_{\mathcal{S}}(\omega) \equiv \|\nabla J_{\mathcal{S}}(\omega) - \nabla J(\omega)\|_2 \leq \theta \|\nabla J_{\mathcal{S}}(\omega)\|_2 \quad (4)$$

where $\theta \in [0, 1)$.

Note that

$$\mathbb{E}[\delta_{\mathcal{S}}(\omega)^2] = \mathbb{E}[\|\nabla J_{\mathcal{S}}(\omega) - \nabla J(\omega)\|_2^2] = \|\text{Var}(\nabla J_{\mathcal{S}})\|_1 \quad (5)$$

The dynamic sample size gradient method

By simple calculations, we have:

$$\text{Var}(\nabla J_{\mathcal{S}}(\omega)) = \frac{\text{Var}(\nabla l(\omega; i))}{|\mathcal{S}|} \frac{N - |\mathcal{S}|}{N - 1} \quad (6)$$

where $\frac{N - |\mathcal{S}|}{N - 1} \approx 1$.

The dynamic sample size gradient method

By simple calculations, we have:

$$\text{Var}(\nabla J_{\mathcal{S}}(\omega)) = \frac{\text{Var}(\nabla l(\omega; i))}{|\mathcal{S}|} \frac{N - |\mathcal{S}|}{N - 1} \quad (6)$$

where $\frac{N - |\mathcal{S}|}{N - 1} \approx 1$.

Then we can rewrite the condition in the following format:

$$|\mathcal{S}| = \frac{\|\text{Var}_{i \in \mathcal{S}}(\nabla l(\omega; i))\|_1}{\theta^2 \|\nabla J_{\mathcal{S}}(\omega)\|_2^2} \quad (7)$$

This is also the criterion that we use to determine the dynamic sample size.

A Newton-CG method with dynamic sampling

At each iteration, the subsampled Newton-CG method chooses samples \mathcal{S}_k and \mathcal{H}_k such that $|\mathcal{H}_k| \ll |\mathcal{S}_k|$, and defines the search direction d_k as an approximate solution of the linear system

$$\nabla^2 J_{\mathcal{H}_k}(\omega_k)d = -\nabla J_{\mathcal{S}_k}(\omega_k) \quad (8)$$

A Newton-CG method with dynamic sampling

At each iteration, the subsampled Newton-CG method chooses samples \mathcal{S}_k and \mathcal{H}_k such that $|\mathcal{H}_k| \ll |\mathcal{S}_k|$, and defines the search direction d_k as an approximate solution of the linear system

$$\nabla^2 J_{\mathcal{H}_k}(\omega_k)d = -\nabla J_{\mathcal{S}_k}(\omega_k) \quad (8)$$

Now we turn to create automatic criterion for deciding the accuracy in the solution of (8)

A Newton-CG method with dynamic sampling

$$r_k \equiv \nabla^2 J_{\mathcal{H}_k}(\omega_k)d + \nabla J_{\mathcal{S}_k}(\omega_k) \quad (9)$$

A Newton-CG method with dynamic sampling

$$r_k \equiv \nabla^2 J_{\mathcal{H}_k}(\omega_k)d + \nabla J_{\mathcal{S}_k}(\omega_k) \quad (9)$$

Then we write the residual of the standard Newton iteration as:

$$\nabla^2 J_{\mathcal{S}_k}(\omega_k)d + \nabla J_{\mathcal{S}_k}(\omega_k) = r_k + [\nabla^2 J_{\mathcal{S}_k}(\omega_k) - \nabla^2 J_{\mathcal{H}_k}(\omega_k)]d \quad (10)$$

A Newton-CG method with dynamic sampling

$$r_k \equiv \nabla^2 J_{\mathcal{H}_k}(\omega_k)d + \nabla J_{\mathcal{S}_k}(\omega_k) \quad (9)$$

Then we write the residual of the standard Newton iteration as:

$$\nabla^2 J_{\mathcal{S}_k}(\omega_k)d + \nabla J_{\mathcal{S}_k}(\omega_k) = r_k + [\nabla^2 J_{\mathcal{S}_k}(\omega_k) - \nabla^2 J_{\mathcal{H}_k}(\omega_k)]d \quad (10)$$

If we define:

$$\mathbb{E}[\Delta_{\mathcal{H}_k}(\omega_k; d)^2] \equiv [\nabla^2 J_{\mathcal{S}_k}(\omega_k) - \nabla^2 J_{\mathcal{H}_k}(\omega_k)]d \quad (11)$$

Then we can make the approximation:

$$\mathbb{E}[\Delta_{\mathcal{H}_k}(\omega_k; d)^2] \approx \frac{\|\text{Var}_{i \in \mathcal{H}_k}(\nabla^2 l(\omega_k; i)d)\|_1}{|\mathcal{H}_k|} \quad (12)$$

A Newton-CG method with dynamic sampling

In order to avoid to recompute the variance at every CG iteration, we initialize the CG iteration at the zero vector. From (9), we have that the initial CG search direction is given by

$$p_0 = -r_0 = -\nabla J_{S_k}(\omega_k).$$

A Newton-CG method with dynamic sampling

In order to avoid to recompute the variance at every CG iteration, we initialize the CG iteration at the zero vector. From (9), we have that the initial CG search direction is given by

$$p_0 = -r_0 = -\nabla J_{S_k}(\omega_k).$$

We compute (12) at the beginning of the CG iteration, for $d = p_0$.

A Newton-CG method with dynamic sampling

In order to avoid to recompute the variance at every CG iteration, we initialize the CG iteration at the zero vector. From (9), we have that the initial CG search direction is given by

$$p_0 = -r_0 = -\nabla J_{S_k}(\omega_k).$$

We compute (12) at the beginning of the CG iteration, for $d = p_0$.

The stop test for the $j + 1$ CG iteration is then set as:

$$\|r_{j+1}\|_2^2 \leq \Psi \equiv \left(\frac{\|\text{Var}_{i \in \mathcal{H}_k}(\nabla^2 l(\omega_k; i)p_0)\|_1}{|\mathcal{H}_k|} \right) \frac{\|d_j\|_2^2}{\|p_0\|_2^2} \quad (13)$$

where d_j is the j th trial candidate for the solution of (8) generated by the CG process and the last ratio accounts for the length of the CG solution.