

Stochastic Gradient Descent

Shao Tang

February 24, 2017

Kaggle: Avito Context Ad Clicks

- Objective: predict if context ads will earn a user's click.
- Response (binary): if there was a click on this ad.
- Features (numerical and categorical).

Main Challenges:

- The data set ($>50\text{GB}$) contains over 200 millions of observations.
- The number of categories for each categorical variable is not known in advance.
- large n and large p .

Problem Formulation

- Cost function over the data-generating distribution p_{data} :

$$J^*(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim p_{data}} L(f(\mathbf{x}; \boldsymbol{\theta}), y).$$

- Cost function over the training set:

$$J(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{p}_{data}} L(f(\mathbf{x}; \boldsymbol{\theta}), y) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)}),$$

where L is the per-example loss function, $f(\mathbf{x}; \boldsymbol{\theta})$ and y are the predicted and observed output given input \mathbf{x} , respectively, and \hat{p}_{data} is the *empirical distribution*.

Stochastic Gradient Descent (SGD)

- Exact gradient:

$$\begin{aligned} \mathbf{g} &= \nabla_{\theta} J^*(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim p_{data}} \nabla_{\theta} L(f(\mathbf{x}; \theta), y) \\ &= \sum_{\mathbf{x}} \sum_y p_{data}(\mathbf{x}, y) \nabla_{\theta} L(f(\mathbf{x}; \theta), y). \end{aligned}$$

- An *unbiased* estimator of the exact gradient of the generalization error: (iid \mathbf{x})

$$\hat{\mathbf{g}} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}),$$

where $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$ is a minibatch of examples sampled from the data-generating distribution p_{data} with m as the sample size.

Examples

- Linear regression:

$$\hat{\mathbf{g}}_i = [\mathbf{x}^{(i)T} \boldsymbol{\theta} - y^{(i)}] \mathbf{x}^{(i)},$$

we can verify that the gradient obtained from all samples can be summed up to the full gradient:

$$\sum_{i=1}^n \hat{\mathbf{g}}_i = \sum_{i=1}^n (\mathbf{x}^{(i)} \mathbf{x}^{(i)T} \boldsymbol{\theta} - y^{(i)} \mathbf{x}^{(i)}) = \mathbf{X}^T (\mathbf{X} \boldsymbol{\theta} - \mathbf{y}).$$

- This holds for any **additive** loss function, i.e.,

$$L(f(\mathbf{x}; \boldsymbol{\theta}), y) = \sum_{i=1}^n L(f(\mathbf{x}^{(i)}; \boldsymbol{\theta}), y^{(i)}).$$

Algorithm 1 Basic SGD. (random shuffle all samples before the start of each epoch)

Require: Learning rate ε_k

Require: Initial parameter θ

- 1: **while** not converged **do**
- 2: Sample a minibatch of m examples from the training set as $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$.
- 3: Compute gradient estimate:

$$\hat{\mathbf{g}} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}).$$

- 4: Update: $\theta \leftarrow \theta - \varepsilon_k \hat{\mathbf{g}}$.
 - 5: **end while**
-

Why SGD?

- If the dataset is highly redundant ($N \gg p$):
 - Rather than computing the full gradient, it is better to update the coefficients based on a sub-sample of original observations.
 - The extreme version of this approach updates coefficients after each observation (so called “online”).
- Mini-batches are usually better than online.
 - Computing the gradient for many cases uses matrix-matrix multipliers are very efficient, especially on GPUs.

Randomized coordinate descent

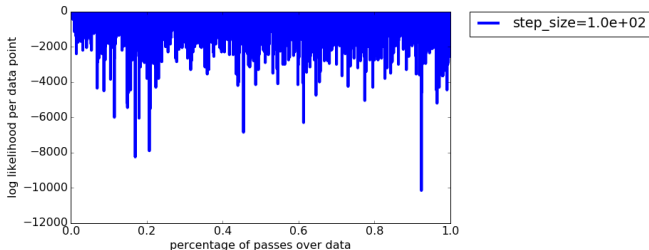
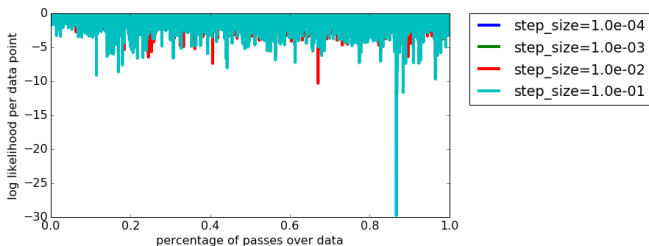
- Instead of sampling from the n observations, we can sample the coordinates of the whole gradient vector $\nabla_{\theta} L(f(\mathbf{x}; \theta), y) \in \mathbb{R}^p$ when p is extremely large.
- In the situation of large n and large p , we can do a two-way sampling, to reduce the computational burden.

Effects of the Step Size

- SGD gradient estimator introduces a source of **noise** that does *NOT* vanish even at a minimum.
- “The learning rate(step size) may be chosen by trial and error, but it is usually best to choose it by monitoring learning curves that plot the objective function as a function of time(number of iterations). This is more of an **art** than a science, and most guidance on this subject should be regarded with some skepticism.”

Effects of the Step Size

Logistic regression on a real world data: amazon review prediction. ($n \sim 50000$, $p \sim 200$)



Effects of the Step Size

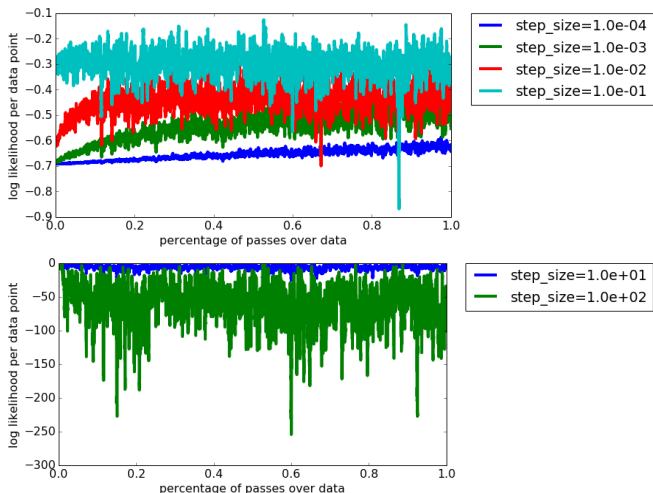


Figure: Batch size $m = 1$, smoothing width 100. (Smoothing is for data visualization only.)

Effects of the Step Size

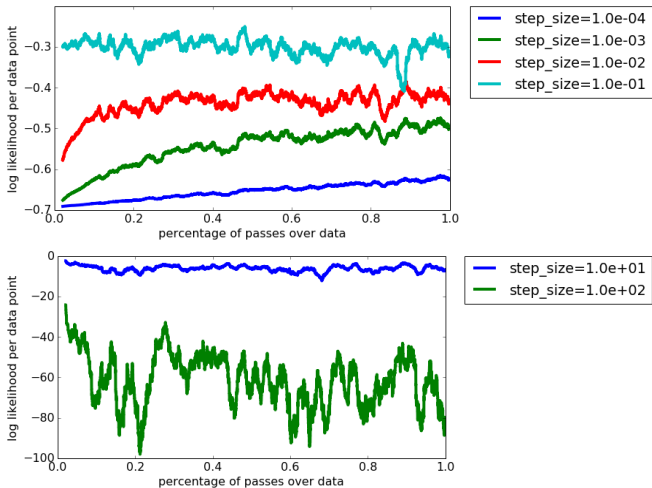


Figure: Batch size $m = 1$, smoothing width 1000.

MiniBatch SGD Step Size

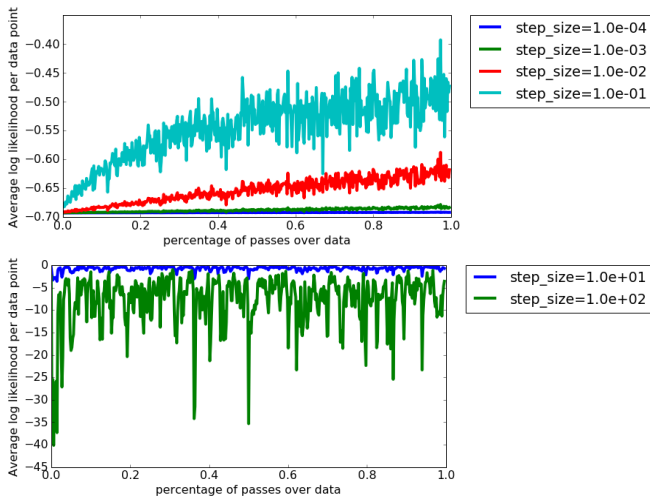


Figure: Batch size $m = 100$, non-smoothing.

MiniBatch SGD and GD

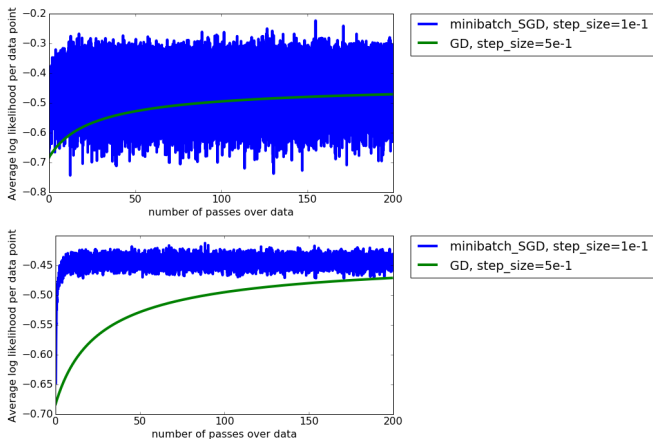


Figure: Batch size $m = 1$ and $m = n$, smoothing width 1 (top) and 50 (bottom).

Empirical Conclusions

- The convergence is strongly affected by the *pre-fixed* step size.
- Increasing the batch size helps function value convergence.
- The minibatch SGD often improves function value faster than the Gradient Descent (GD) method.
- SGD can be directly used on the streamed data (online learning).
- “Don’t trust the last coefficient estimate”.

SGD with momentum

Algorithm 2 SGD with momentum.

Require: Learning rate ε , momentum parameter α

Require: Initial parameter θ , initial velocity v

- 1: **while** not converged **do**
- 2: Sample a minibatch of m examples from the training set as $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$.
- 3: Compute gradient estimate:

$$\hat{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(f(\mathbf{x}^{(i)}; \theta), y^{(i)}).$$

- 4: Velocity update: $v \leftarrow \alpha v - \varepsilon \hat{g}$.
 - 5: Update: $\theta \leftarrow \theta + v$.
 - 6: **end while**
-

Intuition of the momentum

Imagine a ball on the error surface $E(x, y)$. The location of the ball (x, y) represents the coefficient vector.

- Initially, the ball follows the gradient. However, it no longer does steepest descent when it has the velocity.
- Its momentum keeps it going in the previous direction.
- The momentum damps oscillations in directions of high curvature by combining gradients with opposite signs.
- It helps build up speed in directions with a consistent gradient.

SGD with Nesterov momentum

Algorithm 3 SGD with Nesterov momentum.

Require: Learning rate ε , momentum parameter α

Require: Initial parameter θ , initial velocity v

- 1: **while** not converged **do**
- 2: Sample a minibatch of m examples from the training set as $\{(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}$.
- 3: Update $:\tilde{\theta} \leftarrow \theta + \alpha v$.
- 4: Compute gradient estimate:

$$\hat{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\tilde{\theta}} L(f(\mathbf{x}^{(i)}; \tilde{\theta}), y^{(i)}).$$

- 5: Velocity update: $v \leftarrow \alpha v - \varepsilon \hat{g}$.
- 6: Update: $\theta \leftarrow \theta + v$.
- 7: **end while**

SGD with Adaptive Learning Rates

- AdaGrad: individually adapts the learning rates of all model parameters by scaling them inversely proportional to the square root of the sum of all of their historical squared values.
- Adadelta: is an extension of AdaGrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size w .

SGD with Adaptive Learning Rates

- RMSProp: modifies AdaGrad to perform better in the non-convex setting by changing the gradient accumulation into an exponentially weighted moving average.
- Adam (adaptive moments): In addition to storing an exponentially decaying average of past squared gradients like Adadelta and RMSprop, Adam also keeps an exponentially decaying average of past gradients.

Follow-the-Regularized-Leader

- However, the online learning is not particularly effective at producing sparse models.
- Let $L_i(\theta) \triangleq L(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$, and $\mathbf{g}_i = \nabla_{\theta} L_i(\theta_i)$.
Considering

$$\theta_t = \operatorname{argmin}_{\theta} \sum_{i=1}^{t-1} \langle \mathbf{g}_i, \theta \rangle + \frac{1}{2\varepsilon} \|\theta\|_2^2, \quad (\star)$$

we have

$$\theta_t = -\varepsilon \sum_{i=1}^{t-1} \mathbf{g}_i.$$

It is easy to see that

$$\theta_{t+1} = \theta_t - \varepsilon \mathbf{g}_t,$$

which is the online learning and can be viewed as optimization with the linearization trick.

Follow-the-Regularized-Leader

- With $u_i(\boldsymbol{\theta}) = \langle \mathbf{g}_i, \boldsymbol{\theta} \rangle$ and $R(\boldsymbol{\theta}) = \frac{1}{2\varepsilon} \|\boldsymbol{\theta}\|_2^2$, (\star) belongs to the Follow-the-Regularized-Leader algorithm:

$$\boldsymbol{\theta}_t = \operatorname{argmin}_{\boldsymbol{\theta}} \sum_{i=1}^{t-1} u_i(\boldsymbol{\theta}) + R(\boldsymbol{\theta}).$$

- Adding the ℓ_1 penalty, $R(\boldsymbol{\theta}) = \frac{1}{2\varepsilon} \|\boldsymbol{\theta}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_1$, we have

$$\boldsymbol{\theta}_t = \operatorname{argmin}_{\boldsymbol{\theta}} \left\langle \sum_{i=1}^{t-1} \mathbf{g}_i, \boldsymbol{\theta} \right\rangle + \frac{1}{2\varepsilon} \|\boldsymbol{\theta}\|_2^2 + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_1.$$

$$\boldsymbol{\theta}_t = \operatorname{argmin}_{\boldsymbol{\theta}} \|\boldsymbol{\theta} + \varepsilon \mathbf{g}_i\|_2^2 + \lambda \varepsilon \|\boldsymbol{\theta}\|_1,$$

which can be solved by soft-thresholding.

Follow-the-Regularized-Leader

- We have

$$\theta_{t,j} = \begin{cases} 0 & |\sum_{i=1}^{t-1} g_{i,j}| < \lambda \\ -\varepsilon[\sum_{i=1}^{t-1} g_{i,j} - \lambda \text{sign}(\sum_{i=1}^{t-1} g_{i,j})] & \text{otherwise} \end{cases}$$

- In each step, we only need to compute \mathbf{g}_i and update an auxiliary vector $\mathbf{z}_t = \sum_{i=1}^{t-1} \mathbf{g}_i$ to obtain the coefficients update.
- FTRL can be easily extended to **FTRL-Proximal algorithm with coordinate-wise learning rate**.