

Note: This problem set is about linear regression. In particular, we hope to show why traditional methods for solving linear regression problems don't work well as the number of data points n in a data set becomes very large. Two methods (traditional gradient descent and stochastic gradient descent) are proposed to solve these issues and you will analyze their issues. We will talk about more recent developments in the literature on this extremely common problem (matrix sketching and leverage score sampling) in class.

1 Given a data matrix $\mathbf{A} \in \mathbb{R}^{n \times p}$ and an output vector \mathbf{b} , the least squares problem asks us to find the vector \mathbf{x} which minimizes $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$. Here we will ask you to analyze the performance of gradient descent applied to solving this problem. For reference, note that exact, stable solutions to the least squares problem take $2np^2 - \frac{2}{3}p^3$ floating point operations [GolubVanLoan, Algorithm 5.3.2], though these solutions in general don't take advantage of sparsity in the data matrix \mathbf{A} .

- (a) Show that the gradient $\nabla f(\mathbf{x}) = \mathbf{A}^*(\mathbf{A}\mathbf{x} - \mathbf{b})$ and the Hessian $\nabla^2 f(\mathbf{x}) = \mathbf{A}^* \mathbf{A}$. Conclude that f is $\sigma_{\min}(\mathbf{A})$ -strongly convex.
- (b) How many floating point operations does it take to compute $\nabla f(\mathbf{x})$ for an arbitrary input $\mathbf{x} \in \mathbb{R}^p$ if we multiply matrices in the traditional way? Write a bound both in terms of the dimensions n, p , as well as a different bound in terms of the number of non-zero entries in \mathbf{A} (which we denote $\text{nnz}(\mathbf{A})$)
- (c) Give an upper bound on the number of floating point operations it takes to compute an approximate least-squares solution $\tilde{\mathbf{x}}$ with $\|\tilde{\mathbf{x}} - \mathbf{x}\| \leq \epsilon$ via gradient descent which depends on the actual data \mathbf{A} and \mathbf{b} only through the condition number $\kappa = \sigma_{\max}(\mathbf{A})/\sigma_{\min}(\mathbf{A})$, $\sigma_{\min}(\mathbf{A})$, and $\|\mathbf{A}^* \mathbf{b}\|_2$. Give another bound in terms of $\text{nnz}(\mathbf{A})$ which is tighter for sparse data matrices. Assume that we start the algorithm at the zero vector $\mathbf{x}_0 = 0$, and utilize the 'optimal' learning rate to minimize your bounds.
- (d) Detail when you would, and would not, want to use this gradient-descent based approach to solving least squares problems over the traditional exact approach.

- (c) Use Problem 3 from the previous problem set. Write $k = \frac{\kappa-1}{\kappa+1}$ and $\ell = \sigma_{\min}(\mathbf{A})$. You should get something like

$$\text{flops} \sim 4np \log_{1/q} \left(\frac{\|\mathbf{A}^* \mathbf{b}\|_2}{\ell \epsilon} \right).$$

■

2 Observe that we can find the least squares solution \mathbf{x} from Problem 1 by minimizing the rescaled function

$$f(\mathbf{x}) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{a}_i^* \mathbf{x} - b_i)^2, \quad \nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \mathbf{a}_i (\mathbf{a}_i^* \mathbf{x} - b_i)$$

where \mathbf{a}_i are the rows of \mathbf{A} . To mitigate the linear dependence on n in the previous example we might note that randomly sampling i from $\{1, 2, \dots, n\}$ and giving an approximate gradient descent update

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \alpha_{t-1} \underbrace{\mathbf{a}_i (\mathbf{a}_i^* \mathbf{x}_{t-1} - b_i)}_{g_k}$$

gives the original gradient descent update in expectation while ignoring completely the number of samples n . Thus it is natural to think that the algorithm given by replacing the original gradient descent update with this *stochastic gradient* update should perform well, and you will prove that this is true. To do so, we will assume that the ℓ^2 -norm $\sqrt{\mathbb{E}\|g_k\|_2^2} \leq G$ and f is ℓ -strongly convex (i.e. $\sigma_{\min}(\mathbf{A}) \geq \ell$, or $f(\mathbf{x}) - f(\mathbf{y}) \geq \nabla f(\mathbf{y})^* (\mathbf{x} - \mathbf{y}) + \frac{\ell}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$).

- (a) Prove via strong convexity that $\langle \nabla f(\mathbf{x}_t), \mathbf{x}_t - \mathbf{x} \rangle \geq \ell \|\mathbf{x}_t - \mathbf{x}\|_2^2$.
- (b) Prove $\mathbb{E}\|\mathbf{x}_{t+1} - \mathbf{x}\|_2^2 \leq (1 - 2\ell\alpha_t)\mathbb{E}\|\mathbf{x}_t - \mathbf{x}\|_2^2 + \alpha_t^2 G^2$.
- (c) Prove that $\mathbb{E}\|\mathbf{x}_t - \mathbf{x}\|_2^2 \leq \frac{\max\{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2, G^2/\ell^2\}}{t}$ when we set $\alpha_t = \frac{1}{\ell(t+1)}$.
- (d) Give an upper bound on the number of floating point operations it takes to compute an approximate least-squares solution $\tilde{\mathbf{x}}$ with $\mathbb{E}\|\tilde{\mathbf{x}} - \mathbf{x}\|_2 \leq \epsilon$ via stochastic gradient descent like in Problem 1. Can sparsity of \mathbf{A} help?
- (e) In which situations might you prefer to use this method over, say, gradient descent or a traditional exact solver?

- (a) Play with the definition of strong convexity for the points \mathbf{x} and \mathbf{x}_t .
- (b) Expand $\|\mathbf{x}_{t+1} - \mathbf{x}\|_2^2$ deterministically by conditioning on the previous \mathbf{x}_t before applying the expectation to reduce further.
- (c) Induction.
- (d) Reduce to part (c) using Jensen's inequality. You should get something like

$$\text{flops} \sim 4p \frac{\max\{\|\mathbf{x}\|_2^2, G^2/\ell^2\}}{\epsilon^2}.$$

■