
CS189 (SPRING 2023) NOTES

A PREPRINT

Chawin Sitawarin *

University of California, Berkeley

1 Discussion 3

1.1 Solving Linear Regression with Gradient Descent

“Why do we care so much about eigenvalues?”

There are two billion reasons, but today we will look at one of them.

Let’s consider a quadratic function of the form

$$f(w) = \frac{1}{2} \|Xw - y\|^2 \quad (1)$$

$$f(w) = \frac{1}{2} w^\top Q w + w^\top a + b \quad (2)$$

Note that in this special setting we have $Q = X^\top X$ and $a = -X^\top y$. Importantly, Q is a positive semidefinite matrix, i.e., $\forall X, Q \succeq 0$. Let’s also further assume that $X^\top X$ is full rank so we have a stronger condition that Q is positive definite, i.e., $Q \succ 0$.

Linear regression wants to solve

$$w^* = \arg \min_w \frac{1}{2} \|Xw - y\|^2 \quad (3)$$

$$= \arg \min_w \frac{1}{2} w^\top Q w + w^\top a + b \quad (4)$$

We already know of this solution (analytically):

$$\nabla_w f(w^*) = Qw^* + a = 0 \quad (5)$$

$$w^* = -Q^{-1}a \quad (6)$$

$$= (X^\top X)^{-1} X^\top y \quad (7)$$

However, inverting an $n \times n$ matrix is $\mathcal{O}(n^3)$, which is rather expensive for a large n , i.e., high-dimensional data. Note that realistically, while Gaussian Elimination involves $\mathcal{O}(n^3)$ operations, it does not sufficiently account for the number of bits needed to represent these intermediate values.

1. Finding matrix inverse with Gaussian Elimination: https://en.wikipedia.org/wiki/Gaussian_elimination#Finding_the_inverse_of_a_matrix
2. Why it is actually $\mathcal{O}(n^5)$ instead of $\mathcal{O}(n^3)$: https://en.wikipedia.org/wiki/Gaussian_elimination#Computational_complexity and <https://rjlipson.wpcomstaging.com/2015/01/14/forgetting-results/>

*Corresponding email: chawins@berkeley.edu

In short, we rarely want to invert a matrix in practice.

Now we will consider the gradient descent algorithm for solving linear regression. We will denote the eigenvalues of \mathbf{Q} as $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$ (since we assume that \mathbf{Q} is full rank).

$$f(w_{t+1}) = f(w_t - \alpha \nabla f(w_t)) \quad (8)$$

$$= f(w_t) - \alpha \nabla f(w_t)^\top \nabla f(w_t) + \frac{\alpha^2}{2} \nabla f(w_t)^\top \nabla^2 f(\zeta_{w_t, \alpha}) \nabla f(w_t) \quad (\text{MVT}) \quad (9)$$

$$= f(w_t) - \alpha \|\nabla f(w_t)\|^2 + \frac{\alpha^2}{2} \nabla f(w_t)^\top \mathbf{Q} \nabla f(w_t) \quad (10)$$

$$\leq f(w_t) - \alpha \|\nabla f(w_t)\|^2 + \frac{\alpha^2 \lambda_1}{2} \|\nabla f(w_t)\|^2 \quad (\text{Max. eigenvalue of } \mathbf{Q}) \quad (11)$$

$$(12)$$

Choosing an optimal step size $\alpha = 1/\lambda$, we get

$$f(w_{t+1}) \leq f(w_t) - \frac{\alpha}{2} \|\nabla f(w_t)\|^2 \quad (13)$$

Note that choosing $\alpha > 1/\lambda$ will mean that the gradient descent algorithm will diverge. So knowing the largest eigenvalue is already important for choosing a good step size and hence ensuring that the gradient descent algorithm will converge.

We can end our proof here and get a convergence rate of $\mathcal{O}(1/T)$, but we can do better under the assumption that \mathbf{Q} is PD.

$$f(w_{t+1}) \leq f(w_t) - \frac{\alpha}{2} \|\nabla f(w_t)\|^2 \quad (14)$$

$$f(w_{t+1}) \leq f(w_t) - \alpha \lambda_n (f(w_t) - f^*) \quad (\text{Polyak-Lojasiewicz condition}) \quad (15)$$

$$f(w_{t+1}) - f^* \leq f(w_t) - f^* - \alpha \lambda_n (f(w_t) - f^*) \quad (16)$$

$$f(w_{t+1}) - f^* \leq (1 - \alpha \lambda_n) (f(w_t) - f^*) \quad (17)$$

$$f(w_T) - f^* \leq (1 - \alpha \lambda_n)^T (f(w_0) - f^*) \quad (18)$$

$$f(w_T) - f^* \leq \exp(-\alpha \lambda_n T) (f(w_0) - f^*) \quad (1 + x \leq e^x) \quad (19)$$

$$= \exp\left(-\frac{\lambda_n}{\lambda_1} T\right) (f(w_0) - f^*) \quad (20)$$

This means that the gradient descent algorithm converges exponentially fast at a rate of $\mathcal{O}(\exp(-T))$. And finally, the closer λ_n is to λ_1 , the faster the gradient descent algorithm will converge.

In general, the *condition number* of the problem is defined as:

$$\kappa = \frac{\lambda_1}{\lambda_n} \quad (21)$$

“Poorly conditioned” problems have a large condition number, and “well-conditioned” problems have a small condition number (close to 1). Even when the rate is exponential but the condition number is very large, the gradient descent algorithm can still converge very slowly (in the worst case).

In general, if we are dealing with other non-quadratic function, we can still replace \mathbf{Q} with the Hessian matrix $\mathbf{H}(w) := \nabla^2 f(w)$, and everything we talk about still holds. λ_1 corresponds to the L -Lipschitz constant of the gradients of f . This generally tells how smooth f is, and larger L means that f is less smooth (gradients change fast in some directions). On the other hand, λ_n corresponds to the “strong” convexity of f , a stronger condition than convexity.

Reference: <https://www.cs.cornell.edu/courses/cs4787/2020sp/lectures/Lecture2.pdf>