Hello everyone, so this will be the beginning of the fourth week of this course on optimization algorithms. In the first three weeks, we learned about an introduction to optimization, then in the second week, an introduction to Python coding, and in the last week, we learned why iterative algorithms are needed to solve optimization problems. We also learned about line search algorithms where you find the step size $\alpha$, and we started with the gradient descent algorithm. We will continue with the gradient descent algorithm.

If you recall where we stopped last week, we choose the negative of the gradient of f at $x_k$ as the descent direction.

We started with two examples: one was the very simple example

$x_1{}^2 + x_2{}^2$,

and for any starting point, we arrived at the answer in just one step. Just as we were beginning to feel happy that this would happen for many other functions, we gave another example where the domain was just $\mathbb{R}^2$ and the function was quadratic as well, but the number of steps taken to converge was about 30 or 31. The end of the last lecture was regarding the condition number. The condition number of any quadratic function depends on the matrix H.

If $\lambda_1$ is the smallest eigenvalue of H and $\lambda_n$ is the largest eigenvalue of H, then $\lambda_1/\lambda_n$ is defined as the condition number.

I gave this result: when $\lambda_1/\lambda_n$ is approximately 1, then the gradient descent algorithm converges to x* very fast.

For example, in the case of $x_1{}^2 + x_2{}^2$, H is [[2, 0], [0, 2]],

which means $\lambda_1 = \lambda_2 = 2$.

So, $\lambda_1/\lambda_2 = 1$, and the gradient descent algorithm converged in just one iteration.

But when $\lambda_1/\lambda_n$ is much less than 1, like in the other example where the eigenvalues were roughly 8.6 and 1.4, then $\lambda_1/\lambda_2$ was more than 6, and gradient descent converged much more slowly.

You might still have the question about where this $\lambda_1/\lambda_n$ comes from.

We saw two examples: in one, $\lambda_1/\lambda_2$ was 1, and in the other, it was less than 1/6, and this result held. But why is this result true?

I am going to mathematically derive why $\lambda_1/\lambda_n$ has such a big impact on the number of iterations, the speed of convergence of the algorithm.

I will first define an error function.

The error function is given by $(1/2)(x_k - x^*)^T H(x_k - x^*)$, where $x^*$ is the solution.

Here, $x^* = -H^{-1}b$ is the solution to the problem min over x of $(1/2)x^T Hx + b^T x + c$.

You take a general quadratic function and minimize over x.

Differentiating and equating to zero, you get $Hx + b = 0$, so $x^* = -H^{-1}b$.

Let us define an error function for any x, say $x_k$, the value of x at the k-th iteration of the gradient descent algorithm:

$E(x_k) = (1/2)(x_k - x^*)^T H(x_k - x^*)$.

Expanding this, we have $(1/2)x_k^T Hx_k - x_k^T Hx^* + (1/2)x^{*T}Hx^*$.

Since $x^* = -H^{-1}b$, we can write this as $(1/2)x_k^T Hx_k + x_k^T b + (1/2)x^{*T}Hx^*$.

Adding and subtracting c, this is just $f(x_k)$ plus some constant.

Note that this constant does not depend on $x_k$.

So the error function is nothing but $f(x_k)$ plus some constant value.

We will understand why this function is needed in some time, but note that $E(x^*) = 0$.

$f(x^*)$ need not be 0 if b and c are non-zero, but $E(x^*)$ is certainly 0 because it is $(1/2)*0^T H*0$.

Another observation is that $H(x_k - x^*) = Hx_k - Hx^*$.

But $Hx^* = -b$, so $Hx_k - (-b) = Hx_k + b = \nabla f(x_k) = g_k$.

So, $H(x_k - x^*) = g_k$. This is another property we will use in our proof.

(Refer Slide Time 9:25)

We will now compute the quantity $[E(x_k) - E(x_{k+1})] / E(x_k)$.

So, $E(x_k) - E(x_{k+1}) = (1/2)(x_k - x^*)^T H(x_k - x^*) - (1/2)(x_{k+1} - x^*)^T H(x_{k+1} - x^*)$.

In the gradient descent algorithm, $x_{k+1} = x_k - \alpha_k g_k$.

Substituting that, we have:

$= (1/2)(x_k - x^*)^T H(x_k - x^*) - (1/2)(x_k - \alpha_k g_k - x^*)^T H(x_k - \alpha_k g_k - x^*)$

$= (1/2)[ (x_k - x^*)^T H(x_k - x^*) - (x_k - x^* - \alpha_k g_k)^T H(x_k - x^* - \alpha_k g_k) ]$

Expanding the second term:

$= (x_k - x^*)^T H(x_k - x^*) - 2\alpha_k g_k^T H(x_k - x^*) + \alpha_k^2 g_k^T H g_k$

So, the difference becomes:

$(1/2)[ 2\alpha_k g_k^T H(x_k - x^*) - \alpha_k^2 g_k^T H g_k ] = \alpha_k g_k^T H(x_k - x^*) - (\alpha_k^2/2) g_k^T H g_k$

But $H(x_k - x^*) = g_k$, so this becomes $\alpha_k g_k^T g_k - (\alpha_k^2/2) g_k^T H g_k$.

The denominator is $E(x_k) = (1/2)(x_k - x^*)^T H(x_k - x^*) = (1/2) g_k^T H^{-1} g_k$ (since $x_k - x^* = H^{-1} g_k$).

So, $[E(x_k) - E(x_{k+1})] / E(x_k) = [ \alpha_k g_k^T g_k - (\alpha_k^2/2) g_k^T H g_k ] / [ (1/2) g_k^T H^{-1} g_k ] = [ 2\alpha_k g_k^T g_k - \alpha_k^2 g_k^T H g_k ] / [ g_k^T H^{-1} g_k ]$

Now, recall that in exact line search for quadratic functions, $\alpha_k = (g_k^T g_k) / (g_k^T H g_k)$. Substituting this:

$= [ 2*(g_k^T g_k)^2/(g_k^T H g_k) - (g_k^T g_k)^2/(g_k^T H g_k) ] / [ g_k^T H^{-1} g_k ]$

(Refer Slide Time 16:41)

$$= [ (g_k^Tg_k)^2 / (g_k^THg_k) ] / [ g_k^TH^{-1}g_k ] = (g_k^Tg_k)^2 / [ (g_k^THg_k)(g_k^TH^{-1}g_k) ]$$

(Refer Slide Time 17:26)



Now, we use the Kantorovich inequality.

For any positive definite matrix H and vector $x \neq 0$, we have:

$$(x^Tx)^2 / [ (x^THx)(x^TH^{-1}x) ] \geq 4\lambda_1\lambda_n / (\lambda_1 + \lambda_n)^2$$

where $\lambda_1$ and $\lambda_n$ are the smallest and largest eigenvalues of H.

Applying this to $x = g_k$, we get:

$[E(x_k) - E(x_{k+1})] / E(x_k) \geq 4\lambda_1\lambda_n / (\lambda_1 + \lambda_n)^2$

Therefore, $E(x_{k+1}) / E(x_k) \leq 1 - 4\lambda_1\lambda_n / (\lambda_1 + \lambda_n)^2 = [ (\lambda_1 + \lambda_n)^2 - 4\lambda_1\lambda_n ] / (\lambda_1 + \lambda_n)^2 = (\lambda_n - \lambda_1)^2 / (\lambda_1 + \lambda_n)^2$

So, $E(x_{k+1}) \leq [ (\lambda_n - \lambda_1) / (\lambda_1 + \lambda_n) ]^2 E(x_k)$

This tells us that the error function decreases by a factor of $[ (\lambda_n - \lambda_1) / (\lambda_1 + \lambda_n) ]^2$ each iteration. We want the error function to go to 0. The quicker $x_k$ gets nearer to $x^*$, the better.

If $\lambda_n = \lambda_1$, meaning all eigenvalues are equal, then $E(x_{k+1}) \leq 0$, so $E(x_1) = 0$. We reach the solution in one step.

But if $\lambda_n$ is much larger than $\lambda_1$, for example, $\lambda_n = 100$ and $\lambda_1 = 1$, then $E(x_{k+1}) \leq (99/101)^2 E(x_k)$ $\approx 0.9606\ E(x_k)$.

The error decreases very slowly each iteration.

(Refer Slide Time 24:03)



This is the mathematical proof of why the condition number $\lambda_1/\lambda_n$ affects the convergence speed. When $\lambda_1/\lambda_n$ is close to 1, gradient descent converges very fast. When $\lambda_1/\lambda_n$ is much less than 1, convergence is very slow.

I will illustrate this with code. For $x_1^2 + x_2^2$, $\lambda_1/\lambda_n = 1$, and it converges in one step.
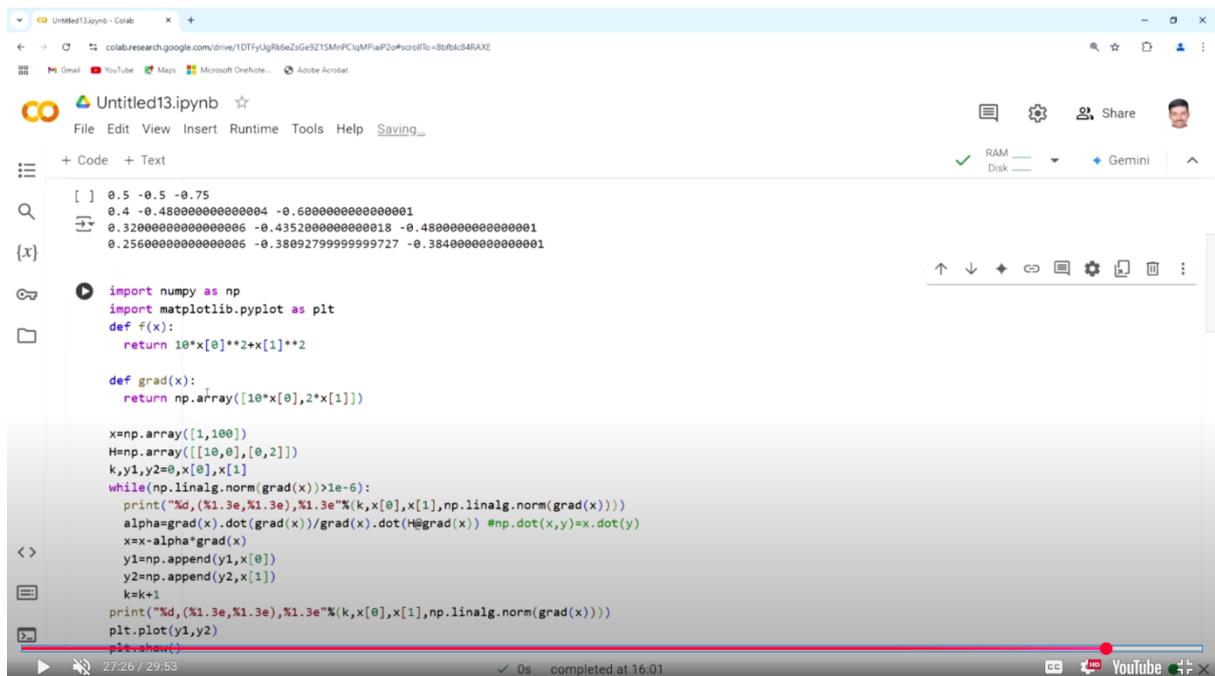
If we change the function to $2x_1^2 + x_2^2$, the condition number changes, and it takes more steps. For example, with initial point (1, 100):

- For $2x_1^2 + x_2^2$, it takes about 5 steps.

- For $5x_1^2 + x_2^2$, it takes about 8 steps.

- For $10x_1^2 + x_2^2$, it takes about 16 steps.

- For $25x_1^2 + x_2^2$, it takes about 68 steps.

- For $50x_1^2 + x_2^2$, it takes about 316 steps.

- For $100x_1^2 + x_2^2$, it takes about 973 steps.

You can see the zigzagging path the algorithm takes from (1, 100) to (0, 0).

The function is just $100x_1^2 + x_2^2$, but the algorithm takes 973 steps because $\lambda_1/\lambda_n = 1/100$ is very small, making convergence very slow.

(Refer Slide Time 27:30-29:00)

We will continue the discussion on this algorithm in the next lecture. Thank you.