

V26 Recitation

- derivatives
- Newton's method
- constrained optimization

$$f : \mathbb{R}^N \rightarrow \mathbb{R}^M$$

\uparrow \uparrow \uparrow
func. name n inputs m outputs

$$y = g(a, b) \quad g : (\mathbb{R}^{n_a}, \mathbb{R}^{n_b}) \rightarrow \mathbb{R}^{n_y}$$

This means f takes in a vector of length N , and outputs a vector of length M

$$y = f(x), \quad y \in \mathbb{R}^M$$
$$x \in \mathbb{R}^N$$

partial derivatives: $\frac{\partial}{\partial}$ vs $\frac{d}{dx}$

$$\dot{y} = \frac{dy}{dt}, \quad \frac{\partial y}{\partial x}$$

Here is the Jacobian, where all of the partial derivatives are collected and numbers of rows = number of outputs of f , and numbers of columns = number of inputs. There is only one interpretation of this.

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_N} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_N} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \dots & \frac{\partial y_m}{\partial x_N} \end{bmatrix}$$

\leftarrow matrix of scalars

rows = # outputs = M
cols = # inputs = N

$$y = g(x) \quad x \in \mathbb{R}^N, \quad y \in \mathbb{R}$$

As a convenience, we define the gradient to be the partial derivatives of a scalar-output function to be a column vector, making it the transpose of the Jacobian. This is a very important distinction, and only applies to scalar-output functions.

$$\frac{\partial g}{\partial x} = \left[\frac{\partial g}{\partial x_1} \quad \frac{\partial g}{\partial x_2} \quad \dots \quad \frac{\partial g}{\partial x_N} \right] \in (1 \times N) \text{ row vector}$$

gradient $\nabla_x g = \left(\frac{\partial g}{\partial x} \right)^T \in \mathbb{R}^N \text{ vector}$

Hessian

$$\nabla_x^2 g = \begin{bmatrix} \frac{\partial^2 g}{\partial x_1^2} & \frac{\partial^2 g}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 g}{\partial x_1 \partial x_n} \\ \vdots & & & \\ \frac{\partial^2 g}{\partial x_n \partial x_1} & & & \frac{\partial^2 g}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Symmetric matrix

Hessian is all of the second partial derivatives, or the jacobian of the gradient of a function. It is symmetric for all the sorts of functions that we care about in optimization, and is positive semi definite for convex functions.

Taylor series

$$y = f(x), \quad x \in \mathbb{R}^n \\ y \in \mathbb{R}^m$$

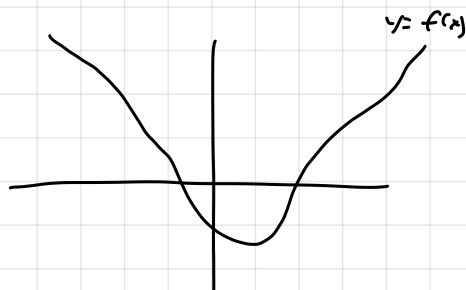
$$f(x) \approx f(\bar{x}) + \left[\frac{\partial f}{\partial x} \right]_{\bar{x}} (x - \bar{x})$$

$(m, 1) \quad (m, 1) \quad (m, n) \quad (n, 1)$

Here we see that the jacobian convention makes sense (size wise) in a Taylor series.

$$g(x) \approx g(\bar{x}) + (\nabla_x g(\bar{x}))^T (x - \bar{x})$$

Newton's Method



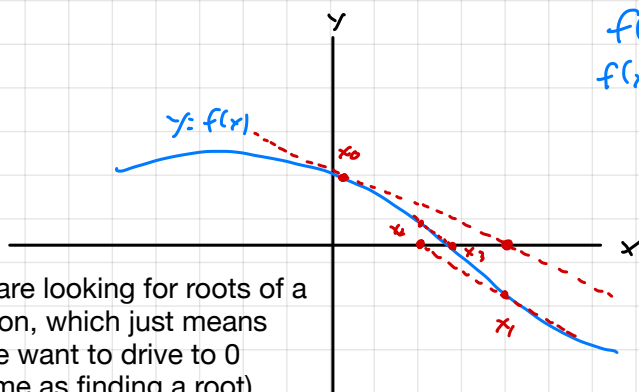
When it comes to solving equations, if there is a solution and the equations are linear, we can solve them. This is the one tool we have that will always work. If the equations are nonlinear, often we have to resort to numerical methods.

We formulate these problems as rootfinding problems, which just means if we want to find $f(x) = y$, we just solve for the roots of $f(x) - y = 0$.

To do this iteratively, we resort to our one tool, linear systems, and we linearize this function and solve for the root of the linearized version. It turns out that if our initial guess is within some "basin of attraction" of the solution, linearizing and solving iteratively is extremely fast effective. This is Newton's method.

$$\begin{aligned} 5a + 6b &= 7 && \text{if a solution exists} \\ 6a - 4b &= 11 && \text{I can find it} \end{aligned}$$

$$\begin{aligned} e^a + \cosh(b) &= \sin(a) && \text{not easy...} \\ a^2 + \sin(b) a^e &= -1 \end{aligned}$$



$$\begin{aligned} f(x) &= y \\ f(x) - y &= 0 \leftarrow \text{find "roots"} \end{aligned}$$

Oftentimes, we are looking for roots of a "residual" function, which just means it's a function we want to drive to 0 (which is the same as finding a root).

Root finding

$$f(x) = 0 \quad y \approx f(x_0) + \left(\frac{\partial f}{\partial x} \right)_{x_0} (x - x_0) = 0$$

Solve for root of our linearized version

$$y \approx f(x_0) + \left(\frac{df}{dx} \Big|_{x_0} \right) (x - x_0) = 0$$

Solve for root
of our linearized version

$$f(x_0) = - \left(\frac{df}{dx} \Big|_{x_0} \right) (x - x_0)$$

$$\left(\frac{df}{dx} \Big|_{x_0} \right)^{-1} f(x_0) = x - x_0$$

$$x = x_0 - \left(\frac{df}{dx} \Big|_{x_0} \right)^{-1} f(x_0)$$

$$x_{k+1} = x_k - \left(\frac{df}{dx} \Big|_{x_k} \right)^{-1} f(x_k)$$

We approximate our residual function with a first order Taylor series (making our approximate function linear), then we solve for the root of this linearized function. This is an iteration of Newton's method.

$x = x_k + \Delta x$ I can define delta x to be this

$$x_k + \Delta x - x_k = \Delta x$$

$$f(x_k) + \left(\frac{df}{dx} \Big|_{x_k} \right) (x - x_k)$$

$$f(x_k) + \left(\frac{df}{dx} \Big|_{x_k} \right) \Delta x = 0$$

Which let's me write out the Newton step as the following equivalent expression:

$$\Delta x = - \left(\frac{df}{dx} \Big|_{x_k} \right)^{-1} f(x_k)$$

"Newton step"

$$x_{k+1} = x_k + \alpha \Delta x$$

$$= x_k + \alpha \left(\frac{df}{dx} \Big|_{x_k} \right)^{-1} f(x_k)$$

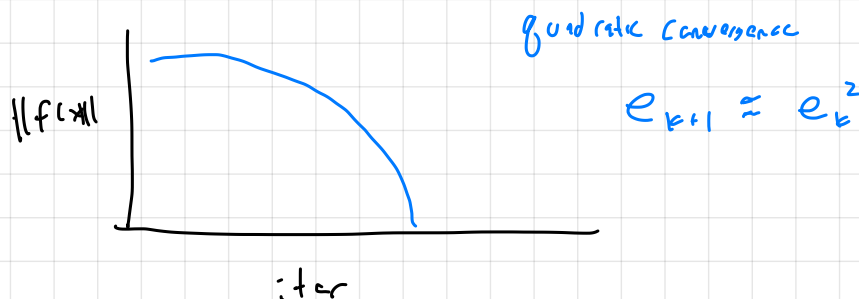
I can apply a step length (alpha) to this step. If my linearization is accurate, alpha=1 is ok, if it is inaccurate, we probably have to use a smaller alpha. We can use a linesearch to find an appropriate alpha to use (more on this later).

input : $f(x)$, x_0 ,

use newton

$$x_{k+1} = x_k + \alpha \left(\frac{\partial f}{\partial x} \Big|_{x_k} \right)^{-1} f(x_k)$$

solved when $\|f(x_k)\| < \text{tolerance}$ ($1e-3, 1e-6, 1e-15$)



$$\min_x f(x)$$

We can easily use Newton's rootfinding method for optimization. First we can take an unconstrained optimization problem, write out the optimality condition (gradient = 0), and then rootfind on this condition. This means our residual function is the gradient of f .

$$\nabla_x f(x) = 0 \quad \leftarrow \text{optimality condition}$$

$$g(x) = \nabla_x f(x) = 0$$

$$\begin{aligned} x_{k+1} &= x_k + \alpha \left(- \frac{\partial g}{\partial x} \Big|_{x_k} \right)^{-1} g(x_k) \\ &= x_k - \alpha \left(\nabla_x^2 f(x_k) \right)^{-1} g(x_k) \end{aligned}$$

If we write out our Newton step, we see that we are taking the jacobian of our gradient, which is the hessian of the original function. This is why you always see people in optimization talk about how Newton's method is really approximating the objection function as a quadratic, and solving for the minimum of this quadratic. This is a really important idea and it's worth spending some time and convincing yourself of this.

newton for unconstrained opti

$$f(x) \approx f(x_k) + (\nabla_x f(x_k))^T \Delta x + \frac{1}{2} \Delta x^T (\nabla_x^2 f(x_k)) \Delta x$$

$$\min_{\Delta x} f(x_k) + (\nabla_x f(x_k))^T \Delta x + \frac{1}{2} \Delta x^T (\nabla_x^2 f(x_k)) \Delta x$$

Notice how if we approximate our objective function with it's second order Taylor expansion, then minimize this approximation, we get the exact same Newton step.

$$\nabla_x f(x_k) + (\nabla_x^2 f(x_k)) \Delta x = 0$$

$$\Delta x = -(\nabla_x^2 f(x_k))^{-1} \nabla_x f(x_k)$$

Root finding = linearize + solve

$$\begin{aligned} \min_x \quad & f(x) \\ \text{s.t.} \quad & c(x) = 0 \end{aligned}$$

$$L(x, \lambda) = f(x) + \lambda^T c(x)$$

Newtons method for rootfinding can also be used for constrained optimization in the exact same way. Write out the optimality conditions as a root finding problem -> use Newton to solve this rootfinding problem.

Our optimality conditions use lagrange multipliers to enforce the constraints, so we are solving for both the variable x (primal variable) and the lagrange multipliers λ (dual variable).

$$\nabla_x L = 0 = \begin{bmatrix} \nabla_x f(x) + \left(\frac{\partial c}{\partial x}\right)^T \lambda \\ c(x) \end{bmatrix} = 0 \quad \text{optimality (KKT)}$$

$$z = \begin{bmatrix} x \\ \lambda \end{bmatrix} \quad r(z) = \begin{bmatrix} \nabla_x f(x) + \left(\frac{\partial c}{\partial x}\right)^T \lambda \\ c(x) \end{bmatrix} = 0$$

Stacking these conditions into a residual function like this lets us use Newton normally.

$$z = \begin{bmatrix} x \\ \lambda \end{bmatrix} \quad r(z) = \begin{bmatrix} \nabla_x f(x) + \left(\frac{\partial C}{\partial x}\right)^T \lambda \\ C(x) \end{bmatrix} = 0$$

When we are doing Newton's steps on a constrained optimization problem, we see that the hessian of the lagrangian shows up. We can form this (especially with modern autodiff tools), but it's expensive and can sometimes ruin the nice positive definiteness of our cost hessian. This leads us to a choice for Newton steps when it comes to constrained optimization:

$$\Delta z = - \left(\frac{\partial r}{\partial z} \Big|_{z_k} \right)^{-1} r(z_k)$$

$$\begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = - \begin{bmatrix} \nabla_x^2 L(x_k, \lambda_k) & \left(\frac{\partial C}{\partial x} \Big|_{x_k}\right)^T \\ \frac{\partial C}{\partial x} \Big|_{x_k} & 0 \end{bmatrix}^{-1} \begin{bmatrix} \nabla_x L(x_k, \lambda_k) \\ C(x_k) \end{bmatrix}$$

Full Newton: Do Newton as normal (form the hessian of the lagrangian)

Gauss Newton: Do Newton mostly as normal, but in the jacobian of our residual (the matrix we form), ignore the constraint curvature and just assume $\text{hess}(\text{lag}) = \text{hess}(\text{cost})$.

$$L(x, \lambda) = f(x) + \lambda^T C(x)$$

$$\nabla_x L = \nabla_x f + \left(\frac{\partial C}{\partial x}\right)^T \lambda$$

$$\nabla_x^2 L = \nabla_x^2 f + \frac{\partial}{\partial x} \left(\left(\frac{\partial C}{\partial x}\right)^T \lambda \right)$$

$$\begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = - \begin{bmatrix} H & \left(\frac{\partial C}{\partial x} \Big|_{x_k}\right)^T \\ \frac{\partial C}{\partial x} \Big|_{x_k} & 0 \end{bmatrix}^{-1} \begin{bmatrix} \nabla_x L(x_k, \lambda_k) \\ C(x_k) \end{bmatrix}$$

for full newton: $H = \nabla_x^2 f(x_k) + \frac{\partial}{\partial x} \left(\left(\frac{\partial C}{\partial x} \Big|_{x_k}\right)^T \lambda_k \right) = \nabla_x^2 L$

Gauss-Newton: $H = \nabla_x^2 f(x_k)$

Constraint Curvature: $\frac{\partial}{\partial x} \left(\underbrace{\left(\frac{\partial C}{\partial x} \Big|_{x_k}\right)^T \lambda_k}_{(N, 1)} \right)$

$\lambda \in \mathbb{R}^m$
 $C(x) \in \mathbb{R}^m$
 $x \in \mathbb{R}^N$

$\underbrace{\hspace{10em}}_{(N, N)}$