

IAM Data Science Pre-Workshop

Linear Algebra lecture #2

Many optimization algorithms have linear system solves as a key ingredient.

① Optimization Basics.

Many continuous optimization problems can be put into the following form:

(A) optimize $F(\underline{x})$ $\underline{x} \in \mathbb{R}^n$
(max or min) \uparrow
given scalar function.

That is, find \underline{x}^* so that
 $F(\underline{x}^*) \leq F(\underline{x})$ for all $\underline{x} \neq \underline{x}^*$
(at least locally).

(B) possibly subject to $\underline{g}(\underline{x}) = \underline{0}$
 $\underline{g} \in \mathbb{R}^m$ ($m < n$). (equality constraints).

(C) and also possibly subject to
 $\underline{h}(\underline{x}) \geq \underline{0}$, $\underline{h} \in \mathbb{R}^l$. (inequality constraints)
 \uparrow
component-wise

As a reminder, you have seen all of this in Math 200.

(★) 3

for (A) look for critical points \underline{x}^* , $\nabla f(\underline{x}^*) = \underline{0}$.
If the eigenvalues of the Hessian of f at \underline{x}^* are all positive, then \underline{x}^* is a local minimum.

for (B) Lagrange multipliers.

for (C) the inequality constraints specify a domain, and you have to check all the boundary points - with Lagrange multipliers if necessary.

We'll forget (B) & (C) for the rest of this lecture - no constraints!

Notes: + In applications, you may not need to find the minimum - if you can decrease f by just a little bit, maybe that is enough to make millions of dollars.

+ If f , g and h are affine, it is a linear programming problem.

+ If f is quadratic and g & h are affine, it is a quadratic programming problem.

+ If f is a specific type of quadratic with no constraints, $\swarrow \searrow$ A, b given

$$f(\underline{x}) = \|\underline{Ax} - \underline{b}\|^2$$

A $m \times n$, $m > n$, A full rank (n), then it is a linear, least squares problem.

② Linear Least Squares.

Find \underline{x} that minimizes $\|\underline{Ax} - \underline{b}\|^2$

Notation $(\underline{Ax})_i = \sum_{j=1}^n a_{ij} x_j$.

$$F(\underline{x}) = \|\underline{Ax} - \underline{b}\|^2 = \sum_{i=1}^m \left[\left(\sum_{j=1}^n a_{ij} x_j \right) - b_i \right]^2$$

To find the minimum,

$$\frac{\partial F}{\partial x_l} = 0$$

$$l = 1, \dots, n.$$

only nonzero when $j=l$

$$2 \sum_{i=1}^m \left[\sum_{j=1}^n a_{ij} x_j - b_i \right] \underbrace{\sum_{j=1}^n a_{ij} \frac{\partial x_j}{\partial x_l}}_{a_{il}}$$

$$\sum_{i=1}^m a_{il} \underbrace{\sum_{j=1}^n a_{ij} x_j}_{(\underline{Ax})_i} - \sum_{i=1}^m a_{il} b_i = 0.$$

$$\underbrace{\sum_{i=1}^m a_{il} (\underline{Ax})_i}_{(\underline{A}^T \underline{Ax})_l} - \underbrace{\sum_{i=1}^m a_{il} b_i}_{(\underline{A}^T \underline{b})_l} = 0$$

$\begin{matrix} \nearrow & \nearrow & \nearrow \\ n \times m & m \times n & n \times 1 \end{matrix}$

when considering matrix equations, it is good to check sizes.

$$\text{so } \underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b}. \quad (1)$$

Note: If $m > n$, \underline{A} is $m \times n$, \underline{A} has rank n , the $\underline{A}^T \underline{A}$ is $n \times n$ and has rank n , so (1) is

solvable.

+ often the least squares problem is not solved in the form (1), but rather using the Singular Value decomposition (which we will discuss next lecture). This process is better conditioned for iterative approximation.

+ $A^T A$ is symmetric and positive definite (all positive eigenvalues - next lecture).

There are good solver options for this class of matrix (like CG).

③ Multiple linear regression.

Note: some confusing change of notation coming up.

Suppose that a process can be described by

not the form before. \uparrow

$$F(\underline{x}) = a_1 x_1 + \dots + a_n x_n = \underline{a} \cdot \underline{x}$$

$\in \mathbb{R}^n$

Here, the a 's are the unknowns.

m measurements are taken at conditions \underline{x}_i .

$$f(\underline{x}_i) = \underline{a} \cdot \underline{x}_i = \sum_{j=1}^n x_{ij} a_j = (\underline{A} \underline{x})_i \approx y_i$$

row vector. \uparrow

Pick the coefficients \underline{a} so that the error to the measurements is minimized

that is, minimize.

$$\|Xa - y\|^2$$

can be written
 $\|A^T X^T - y^T\|^2$

5.

It is a least squares problem!

$$X^T X a = X^T y.$$

Ex linear regression.

$$f(x) = ax + b.$$

$\uparrow \uparrow$
 a & b to be determined

m data

$$y_i \approx ax_i + b.$$

$\uparrow \uparrow$
given.

$$X = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_m & 1 \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

"train" the model on the data. Useful if the form of the model is appropriate.

$$X^T X = \begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & m \end{bmatrix}$$

$$X^T y = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

→ Computational examples.

④ Vector Newton's Method.

Let's go back to Ⓐ $\nabla f(x^*) = \underline{0}$.

In general, this is n nonlinear equations in n unknowns.

Let us consider the more general case $g(x^*) = \underline{0}$

6
Aside: Scalar Newton's method on the board.

The Jacobian matrix of derivatives is

$$Df, \quad d_{ij} = \frac{\partial g_i}{\partial x_j}.$$

Theory: If J has full rank at zero of g , the zero is locally unique.

at an arbitrary point \underline{x} , we can do a linear approximation.

$$g_i(\underline{x}^*) \approx g_i(\underline{x}) + \sum_{j=1}^n \frac{\partial g_i}{\partial x_j} (\underline{x}_j^* - \underline{x}_j).$$

in vector notation

$$g(\underline{x}^*) \approx g(\underline{x}) + J(\underline{x})(\underline{x}^* - \underline{x}).$$

so if \underline{x} were given and $g(\underline{x})$ were not zero, we could use

$$g(\underline{x}) + J(\underline{x}^* - \underline{x}) \approx \underline{0}$$

and solve for \underline{x}^* to approximate the root.

$$\underline{x}^* = \underline{x} - \underbrace{J^{-1}(\underline{x})}_{\underline{\delta}} g(\underline{x}).$$

Remember we never find J^{-1} and multiply by it. This is notation for $\underline{\delta}$ that solves

$$J \underline{\delta} = -g(\underline{x}).$$

Ex Nonlinear boundary value problem. 7.

Consider a FD approx of

$$-u'' + u - u^3 = f(x).$$

$$-ID_2 \underline{U} + \underline{U} - \underline{U}^{\wedge 3} = \underline{F}.$$

↑
pointwise cube.

$$\text{So } g(\underline{U}) = -ID_2 \underline{U} + \underline{U} + \underline{U}^{\wedge 3} - \underline{F}.$$

Jacobian matrix

$$J = -ID_2 + I + 3 \text{diag}(U_j^2)$$

→ computational demo.
aside: continuation on the board.

If $g(x) = \nabla f$, then the Jacobian is
the Hessian of f

$$J_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}.$$

In this case, J is symmetric, and will
be positive definite near a local minimum.

Ex Rosenbrock function

$$f(\underline{x}) = 100 \sum_{j=1}^{n-1} (x_{j+1} - x_j^2)^2 + \sum_{j=1}^{n-1} (x_j - 1)^2$$

↑
 $\in \mathbb{R}^n$

→ computational examples.

⑤ Steepest Descent.

8.

What if the problem is too big to solve

$$\text{H } \underline{d} = \underline{\nabla f}?$$

Well, ∇f is the direction of maximum increase of f , so going in the direction $-\nabla f$ will reduce f .

Iterate $\{\underline{x}^n\}$.

$$\underline{x}^n = \underline{x}^{n-1} - \underset{\uparrow}{\rho} \underline{\nabla f}$$

constant (small) or do a line search, reject if f increases.

Final note: Finding a global minimum may be impossible.

Suggested problems:

9.

#1 Solve the linear programming problem below (find appropriate software).

$$\text{maximize } x_1 + 2x_2 + 3x_3 + 4x_4 + 5$$

subject to the constraints

$$4x_1 + 3x_2 + 2x_3 + x_4 \leq 10$$

$$x_1 - x_3 + 2x_4 = 2$$

$$x_1 + x_2 + x_3 + x_4 \geq 1.$$

and $x_1 \geq 0, x_3 \geq 0, x_4 \geq 0$.

#2. Solve the following linear programming problem with N^2 variables g_{ij} .

$$\text{minimize } \sum_{i=1}^N \sum_{j=1}^N g_{ij} (i-j)^2$$

subject to $\sum_{i=1}^N g_{ij} = 1$ for all j

$$\sum_{j=1}^N g_{ij} = 1 \text{ for all } i.$$

$$0 \leq g_{ij} \leq \delta/N. \text{ for all } i, j. (\delta > 1 \text{ is a parameter}).$$

Notes: + as $N \rightarrow \infty$ this approximates a continuum optimal transport problem.

+ you should see interesting behaviour around $\delta = 2$.

+ as $\delta \rightarrow \infty$ nonzero g values should cluster around the diagonal.

#3. Export the python exponential data points to R and do the regression in that platform.

#4. Investigate ways to limit the influence of the outlier in the linear regression. There are statistical methods to identify outliers. You can use l_1 error estimation. You could also try the following:

$$\text{minimize } \sum_{i=1}^m r(mx_i + b - y_i)$$

$$\text{where } r(z) = \frac{z^2}{1 + z^2/\sigma^2}.$$

as $\sigma \rightarrow \infty$ this tends to the least squares problem. You could take $\sigma \approx$ noise level to do the "right" thing. If you do the nonlinear solve yourself, you may need to do continuation in σ from σ large (use least squares to start).

#5. Apply gradient descent to the Rosenbrock function using different constant ϵ or "strategies".