1. LSE:

For convenience, consider two dimenional data $(x, y)$.

Ideally, we would expect that $A\vec{w} = \vec{b}$:

$$\underbrace{\begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^m \\ 1 & x_2 & x_2^2 & \cdots & x_2^m \\ & & & \vdots & \\ 1 & x_n & x_n^2 & \cdots & x_n^m \end{bmatrix}}_{n \times (m+1)} \underbrace{\begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}}_{(m+1) \times 1} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{n \times 1}$$

$$\underset{A}{\phantom{x}} \qquad \underset{\vec{w}}{\phantom{x}} \qquad \underset{\vec{b}}{\phantom{x}}$$

Then $\vec{w}$ can be calculated by $\vec{w} = A^{-1}\vec{b}$.

However, $A$ is not guaranteed to be invertible.

In fact, $A$ doesn't even need to be a square matrix.

Therefore, we change our goal from finding $\vec{w}$ such

that $A\vec{w} = \vec{b}$ to minimizing $\|A\vec{w} - \vec{b}\|^2$.

Then $\|A\vec{w} - \vec{b}\|^2 = (A\vec{w} - \vec{b})^T (A\vec{w} - \vec{b})$

$$= (\vec{w}^T A^T - \vec{b}^T)(A\vec{w} - \vec{b})$$

$$= \vec{w}^T A^T A \vec{w} - \underline{\vec{w}^T A^T \vec{b}} - \underline{\vec{b}^T A \vec{w}} + \vec{b}^T \vec{b}$$

$\because \vec{w}^T A^T \vec{b}$ is a scalar

$\therefore (\vec{w}^T A^T \vec{b})^T = \vec{b}^T A \vec{w} = \vec{w}^T A^T \vec{b}$

$$= \vec{w}^T A^T A \vec{w} - 2 \vec{w}^T A^T \vec{b} + \vec{b}^T \vec{b}.$$

$$\frac{\partial(\vec{w}^T A^T A \vec{w})}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial}{\partial w_0} \\ \frac{\partial}{\partial w_1} \\ \vdots \\ \frac{\partial}{\partial w_m} \end{bmatrix} \begin{bmatrix} W_0 & W_1 & \dots & W_m \end{bmatrix} \begin{bmatrix} a_{0,0} & a_{0,1} & \cdots & a_{0,m} \\ a_{1,0} & a_{1,1} & \cdots & a_{1,m} \\ & & \vdots & \\ a_{m,0} & a_{m,1} & \cdots & a_{m,m} \end{bmatrix} \begin{bmatrix} W_0 \\ W_1 \\ \vdots \\ W_m \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial}{\partial w_0} \\ \frac{\partial}{\partial w_1} \\ \vdots \\ \frac{\partial}{\partial w_m} \end{bmatrix} \begin{bmatrix} W_0(a_{0,0}W_0 + a_{0,1}W_1 + \dots + a_{0,m}W_m) \\ + W_1(a_{1,0}W_0 + a_{1,1}W_1 + \dots + a_{1,m}W_m) \\ + \dots \\ + W_m(a_{m,0}W_0 + a_{m,1}W_1 + \dots + a_{m,m}W_m) \end{bmatrix}$$

<span style="color:red">1 × 1</span>

$$= \begin{bmatrix} (a_{0,0}W_0 + a_{0,1}W_1 + \dots + a_{0,m}W_m) + (W_0 a_{0,0} + W_1 a_{1,0} + \dots + W_m a_{m,0}) \\ (a_{1,0}W_0 + a_{1,1}W_1 + \dots + a_{1,m}W_m) + (W_0 a_{0,1} + W_1 a_{1,1} + \dots + W_m a_{m,1}) \\ \vdots \\ (a_{m,0}W_0 + a_{m,1}W_1 + \dots + a_{m,m}W_m) + (W_0 a_{0,m} + W_1 a_{1,m} + \dots + W_m a_{m,m}) \end{bmatrix}$$

<span style="color:red">(m+1) × 1</span>

$$= \underline{(A^T A)\vec{w}} + \underline{(A^T A)^T \vec{w}}$$

<span style="color:red">$(A^T A)^T = A^T (A^T)^T = A^T A$</span>

$$= 2(A^T A)\vec{w}$$

Similarly, $\dfrac{\partial(\vec{w}^T A^T \vec{b})}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial}{\partial w_0} \\ \frac{\partial}{\partial w_1} \\ \vdots \\ \frac{\partial}{\partial w_m} \end{bmatrix} \begin{bmatrix} W_0 & W_1 & \dots & W_m \end{bmatrix} (A^T \vec{b})$

$$= A^T \vec{b}$$

Hence, define $f(\vec{w}) = \vec{w}^T A^T A \vec{w} - 2 \vec{w}^T A^T \vec{b} + \vec{b}^T \vec{b}$

$$\frac{\partial f}{\partial \vec{w}} = 2 A^T A \vec{w} - 2 A^T \vec{b} .$$

$\frac{\partial f}{\partial \vec{w}} = 0 \implies \vec{w} = (A^T A)^{-1} A^T \vec{b}$ . (*)

Hence $f(\vec{w})$ attains its minimum at $\vec{w} = (A^T A)^{-1} A^T \vec{b}$.

Note that we only have $\det(A^T A) \geq 0$, it's still possible that $(A^T A)$ is not invertible.

Hence we can add the regularized term $\lambda I$ ($\lambda > 0$) to $(A^T A)$, and then $\det(A^T A + \lambda I) > 0$

$$\implies (A^T A + \lambda I) \text{ is invertible.}$$

(*) becomes $\vec{w} = (A^T A + \lambda I)^{-1} A^T \vec{b}$.

This is also called rLSE.


Next we'll illustrate on LU decomposition and how to apply it to finding the inverse of a given matrix.

$\underset{n \times n}{A} = \underset{n \times n}{L} \underset{n \times n}{U} = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ l_{2,1} & 1 & 0 & \cdots & 0 \\ l_{3,1} & l_{3,2} & 1 & & \vdots \\ & & & \ddots & 0 \\ l_{n,1} & l_{n,2} & \cdots & l_{n,n-1} & 1 \end{bmatrix} \begin{bmatrix} u_{1,1} & u_{1,2} & \cdots & u_{1,n} \\ 0 & u_{2,2} & \cdots & u_{2,n} \\ 0 & 0 & & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & u_{n,n} \end{bmatrix}$

$A = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,n} \\ a_{2,1} & & \ddots & \\ \vdots & & & \\ a_{n,1} & & & a_{n,n} \end{bmatrix}$

$$\implies a_{1,1} = u_{1,1} \quad , \quad a_{1,2} = u_{1,2} \quad \cdots \quad a_{1,n} = u_{1,n}$$

$$a_{2,1} = l_{21} \cdot u_{1,1} \quad , \quad a_{3,1} = l_{3,1} \cdot u_{1,1} \quad \cdots \quad a_{n,1} = l_{n,1} \cdot u_{1,1}$$

$$\implies l_{2,1} = \frac{a_{2,1}}{u_{1,1}} \quad , \quad l_{3,1} = \frac{a_{3,1}}{u_{1,1}} \quad , \quad \cdots \quad l_{n,1} = \frac{a_{n,1}}{u_{1,1}}$$

$$a_{2,2} = l_{2,1} \cdot u_{1,2} + u_{2,2} \quad , \quad a_{2,3} = l_{2,1} \cdot u_{1,3} + u_{2,3} \quad \cdots$$

$$\implies u_{2,2} = a_{2,2} - l_{2,1} \cdot u_{1,2} \quad , \quad u_{2,3} = a_{2,3} - l_{2,1} \cdot u_{1,3} \quad \cdots$$

$$a_{3,2} = l_{3,1} \cdot u_{1,2} + l_{3,2} \cdot u_{2,2} \implies l_{3,2} = \frac{a_{3,2} - l_{3,1} \cdot u_{1,2}}{u_{2,2}}$$

$$\vdots$$

Suppose that we already have $A = LU$.
$\underset{n \times n}{}$

Define $A^{-1} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix}$ where $x_1, x_2, \ldots$ are column vectors.

Then $A(A^{-1}) = \begin{bmatrix} Ax_1 & Ax_2 & \cdots & Ax_n \end{bmatrix} = I = \begin{bmatrix} e_1 & e_2 & \cdots & e_n \end{bmatrix}$

$$\implies LUx_1 = e_1 \, , \, LUx_2 = e_2 \, \ldots \, LUx_n = e_n.$$

We solve for $\begin{cases} Ux_1 = y_1 \, , \, Ly_1 = e_1 \\ Ux_2 = y_2 \, , \, Ly_2 = e_2 \\ \quad \vdots \\ Ux_n = y_n \, , \, Ly_n = e_n \end{cases}$

Then $A^{-1}$ can be calculated ✷

## 2. Steepest descent method:

The formula of steepest descent, aka gradient descent method, can be written as:

$$X_{t+1} = X_t - \eta \nabla f(X_t)$$

where $f$ is the loss function.

↳ e.g. MSE

Assume that $f$ is Lipschitz continuous with constant $L > 0$.

Then $\|\nabla f(X) - \nabla f(y)\| \leq L \|X - y\|$ for any $X, y$.

We can perform a quadratic expansion of $f$ around $f(X_t)$ and obtain the following inequality:

$$f(X_{t+1}) \leq f(X_t) + \nabla f(X_t)^T (X_{t+1} - X_t) + \frac{1}{2} \nabla^2 f(X_t) \|X_{t+1} - X_t\|^2$$

$$\leq f(X_t) + \nabla f(X_t)^T (X_{t+1} - X_t) + \frac{1}{2} L \|X_{t+1} - X_t\|^2$$

$$= f(X_t) - \eta \|\nabla f(X_t)\|^2 + \frac{1}{2} L \eta^2 \|\nabla f(X_t)\|^2$$

$$= f(X_t) - (1 - \frac{1}{2} L \eta) \eta \|\nabla f(X_t)\|^2$$

Note that $\|X_{t+1} - X_t\|$ has to be small enough, which implies that $\eta$ also has to be small enough.

Choose $\eta \leq \frac{1}{L}$, then $-\left(1 - \frac{1}{2}L\eta\right) = \frac{1}{2}L\eta - 1$

$$\leq \frac{1}{2}L\left(\frac{1}{L}\right) - 1$$

$$= -\frac{1}{2}.$$

Then $f(x_{t+1}) \leq f(x_t) \underline{- \frac{1}{2}\eta \|\nabla f(x_t)\|^2}$ $\textcolor{red}{(**)}$

$\textcolor{red}{\rightarrow}$ $\textcolor{red}{\text{positive unless } \nabla f(x) = 0}$

Thus the sequence $\{f(x_0), f(x_1), \dots\}$ is indeed decreasing.

Assume that $f$ is convex and $f(x)$ attains its minimum

at $x = x^*$, then we have

$$\textcolor{blue}{f(x^*) \geq f(x) + \nabla f(x)^T(x^* - x) \quad \text{for any } x}$$

$$\textcolor{blue}{\Leftrightarrow f(x) \leq f(x^*) + \nabla f(x)^T(x - x^*).}$$

Then $(**)$ becomes $f(x_{t+1}) \leq f(x^*) + \nabla f(x_t)^T(x_t - x^*)$

$$-\frac{1}{2}\eta \|\nabla f(x_t)\|^2$$

$$\Rightarrow f(x_{t+1}) - f(x^*) \leq \frac{1}{2\eta}\left(2\eta \nabla f(x_t)^T(x_t - x^*) - \eta^2 \|\nabla f(x_t)\|^2\right)$$

$$= \frac{1}{2\eta}\left(\|x_t - x^*\|^2 - \|x_t \eta \nabla f(x_t) - x^*\|^2\right)$$

$$= \frac{1}{2\eta}\left(\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2\right)$$

This implies that the sequence $\{f(x_0), f(x_1)\dots\}$ is bounded.

$\llcorner$(I skip the complicated proof of convergence :c )

In the end, consider $g(\vec{w}) = \|\vec{b} - A\vec{w}\|^2$, aka

the LSE loss, in (1).

Hence $\frac{\partial g}{\partial \vec{w}} = 2A^T A\vec{w} - 2A^T \vec{b}$  $\Rightarrow$ gradient

$\qquad = 2A^T(A\vec{w} - \vec{b})$

On the other hand, for the regularized term in

$L_1$-norm, the gradient of it can be written as

the sign function.   $\text{sign}(w_i) = \begin{cases} 1, & \text{if } w_i > 0 \\ -1, & \text{if } w_i < 0 \\ 0, & \text{if } w_i = 0. \end{cases}$

Thus the gradient in total is $2A^T(A\vec{w} - \vec{b}) + \lambda \cdot \text{sign}(\vec{w})$.

3. Newton's method:

We have the following equation (from Taylor expansion)
$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + O((x - x_0)^2)$$

If we want to find $x$ such that $f(x) = 0$, and $|x - x_0|$ is small enough, then
$$0 = f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$
$$\Rightarrow x \approx x_0 - \frac{f(x_0)}{f'(x_0)}.$$

Then we derive the formula of Newton's method:
$$x_{t+1} = x_t - \frac{f(x_t)}{f'(x_t)}$$

it'll converge to the root of the original equation, and the error is up to your tolerance.

If we want to apply Newton's method to an optimization problem, we may need to solve $f'(x) = 0$:
$$x_{t+1} = x_t - \frac{f'(x_t)}{f''(x_t)}$$
$$\Rightarrow x_{t+1} = x_t - (\nabla^2 f(x_t))^{-1} \nabla f(x_t) \quad \text{in multi-dimensional case.}$$

From (1) LSE and (2) Steepest descent method, we know

that $\nabla f(\vec{w}) = 2A^T(A\vec{w} - \vec{b})$ for $f$: square error.

$\nabla^2 f(\vec{w}) = 2A^T A.$ ※