

# **Alternate Directive method of Multipliers and its Applications**

*Arrun Sivasubramanian(CB.EN.U4AIE19103)*

*B. Tech CSE AI*

*Amrita School of Engineering, Ettimadai, Coimbatore*

## **Abstract:**

ADMM is one of the best methods known at present for solving optimisation problems involving equality constraints. It can handle general constraints and is an effective method to solve problems arising in machine learning including lasso minimization, SVM and some others. It's very effective in statistics and machine learning problems dealing with very large data because it is well suited to distributed(convex) optimisation.

In this project, I have described about ADMM in detail on how it's formulated and its assorted applications in practical areas depending on the type of problem. In the latter parts, I have elaborated some basic concepts needed for speeding up the algorithm by approximately 80 times of the original approach's speed.

## **Introduction to ADMM:**

ADMM combines the decomposability of dual ascent with the superior convergence properties of the method of multipliers. The basic idea is to split the function  $F$  into two independent parts as  $F(x, z) = f(x) + g(z)$  and consider the new minimisation problem subject to  $Ax + Bz = c$ .

Instead of jointly performing minimization, ADMM first performs an  $x$ -minimization step and then a  $z$ -minimization step which accounts for its name alternative direction. Because the Lagrangian is augmented, some mild conditions on  $A$  and  $B$  guarantee termination. Convergence of ADMM is under the following assumptions:

1. The functions  $f(x)$  and  $g(z)$  are proper and closed functions and domain of  $f(x)$  and  $g(x)$  is not equal to null.

2. The square matrices  $A^T A$  and  $B^T B$  of size  $m$  and  $n$  respectively are invertible. Equivalently, the matrices  $A$  and  $B$  of dimension  $(p \times m)$  and  $(p \times n)$  have ranks  $m$  and  $n$  respectively.
3. The un-augmented Lagrangian term:

$$L_0 = f(x) + g(z) + \lambda^T (Ax + Bz - c) \text{ has a saddle point}$$

$$\Rightarrow \text{There exists } x^*, z^*, \lambda^* (\text{not necessarily unique}) \text{ such that:}$$

$$L_0(x^*, z^*, \lambda) \leq L_0(x^*, z^*, \lambda^*) \leq L_0(x, z, \lambda^*) \quad \forall x, z, \lambda$$

In this project, I experimented with variants of ADMM applied to :

1. Basic Pursuit
2. Least Absolute Deviations (LAD)
3. Linear and Quadratic Programming
4. Sparse Covariance Matrices
5. Hubber Loss
6. Regressor Selector
7. Total Variable Minimisation
8. LASSO

## 1. Basic Pursuit:

Basis pursuit is the equality-constrained L1 minimization problem. It is often used as a heuristic for finding a sparse solution to an underdetermined system of linear equations. It plays a central role in modern statistical signal processing, particularly the theory of compressed sensing.

Problem formulation is as follows:

$$\text{minimize } \|x\|_1 \quad \text{subject to } Ax = b$$

*ADMM* formulation:

$$\text{minimize } f(x) + g(z) \quad \text{subject to } x - z = \mathbf{0} \text{ vector}$$

where,

$$f(x) = \|x\|_1$$

$$g(z) = \begin{cases} 0, & Az = b \\ \infty & \text{otherwise} \end{cases}$$

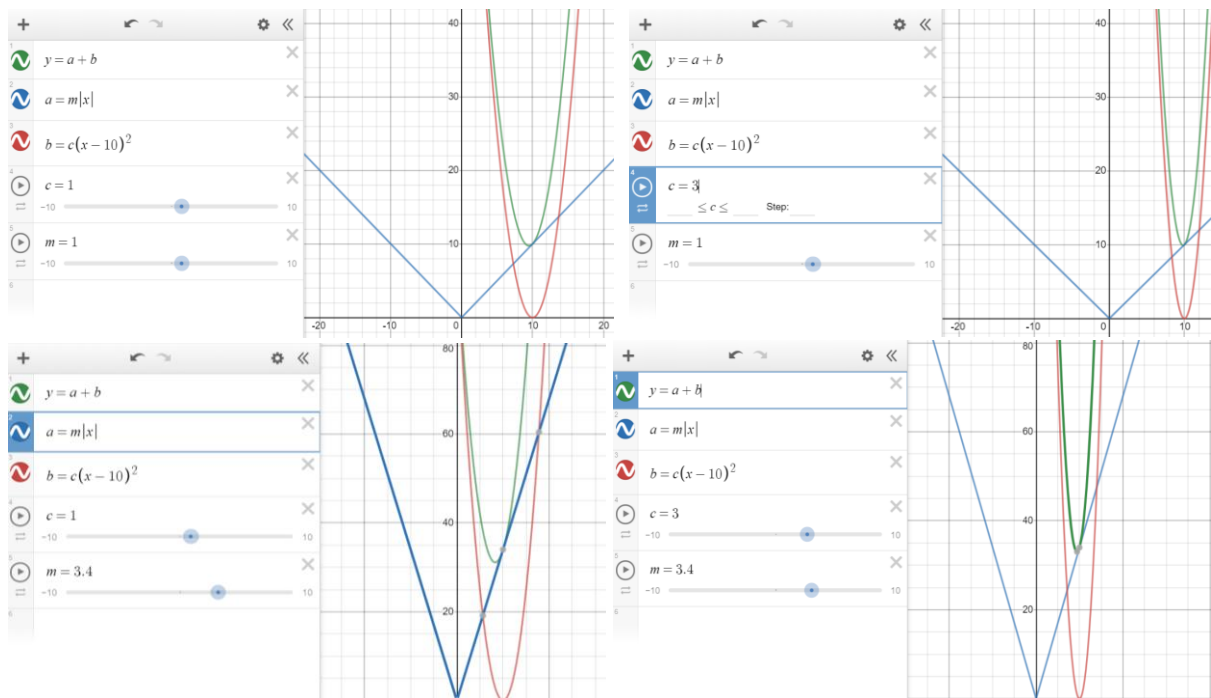
Formation of Lagrangian and solving the problem is as follows:

$$\begin{aligned}
 L(x,z,u) &= \|x\|_1 + g(z) + \frac{1}{2\lambda} \|x - z + u\|_2^2 \\
 x^{k+1} &= \arg \min_x \|x\|_1 + \frac{1}{2\lambda} \|x - z^k + u^k\|_2^2 \\
 &= S_\lambda(z^k - u^k) \\
 z^{k+1} &= \arg \min_z g(z) + \frac{1}{2\lambda} \|x^{k+1} - z + u^k\|_2^2 \\
 z^{k+1} &= P_{Az=b}(x^{k+1} + u^k) = (x^{k+1} + u^k) + \text{pinv}(A) \times (b - A(x^{k+1} + u^k)) \\
 u^{k+1} &= u^k + (x^{k+1} - z^{k+1})
 \end{aligned}$$

**Function to be aware:** The latter part of the added term in x-update is called the penalty or the regularisation function. The function  $S(y)$  is called as the shrinkage function as it shrinks  $y$  towards zero. The exact form of shrinkage function depends on the penalty function.

The shrinkage function sets all values of  $y$  less than linear multiplier of  $x$  ( $\lambda$  in our case) to zero. Thus geometrically, we can define these functions as thresholding functions. But the penalty term can have an additional parameter ( $\alpha$  in our case) which geometrically decides the shape (concavity of the function).

Graphical representation of change in shape and intercept of new curve is as follows:



Figures that help in visualising change in function with change in initial parameters. Here  $c$  denotes  $\alpha$  i.e. the coefficient of the penalty function. With variation in  $c$ , concavity of function changes and with variation in  $M$  (the pulling parameter), the intercepts of the function change.

### Coding basic pursuit is as follows:

```
clc;clear all;close all;
n = 30;
m = 10;
A = randn(m,n);

x = sprandn(n, 1, 0.1*n);
b = A*x;

xtrue = x;
[X,Z,U] = basis_pursuit(A, b, 1.0, 1.0);

function [x,z,u] = basis_pursuit(A, b, rho, alpha)
% basis_pursuit Solve basis pursuit via ADMM
%
% Solves the following problem via ADMM:
% minimize ||x||_1
% subject to Ax = b
%
% The solution is returned in the vector x.
% rho is the augmented Lagrangian parameter and alpha is the
over-relaxation parameter
    MAX_ITER = 1000;
    [m n] = size(A);
    x = zeros(n,1);
    z = zeros(n,1);
    u = zeros(n,1);

    % precompute static variables for x-update (projection on
to Ax=b)
    AAt = A*A';
    P = eye(n) - A' * (AAt \ A);
    q = A' * (AAt \ b);

    for k = 1:MAX_ITER
        % x-update
        x = P*(z - u) + q;

        % z-update with relaxation
        zold = z;
        x_hat = alpha*x + (1 - alpha)*zold;
        z = shrinkage(x_hat + u, 1/rho);
```

```

        u = u + (x_hat - z);
    end
end

function y = shrinkage(a, kappa)
    y = max(0, a-kappa) - max(0, -a-kappa);
end

```

## 2. The Least Absolute Deviations Problem:

Least Absolute Deviations (LAD) is an alternative to traditional Least Squares by using cost function. Least absolute deviations provides a more robust fit than least squares when the data contains large outliers, and has been used extensively in statistics and econometrics. It is similar to an OLS regression model, but it minimizes the absolute deviation instead of the sum of squares. LAD is equivalent to the median regression, a special case of the quantile regression models. LAD is a robust regression technique in the sense that the estimated coefficients are insensitive to outliers.

This problem differs slightly from Basis Pursuit in that we do not have the constraint  $Ax = b$  anymore. We are given observation in the rows of  $A$  and corresponding noisy measurements  $b$ . Our goal here is to actually recover the original, minimal signal  $x$  despite the noise in  $b$ . For this, we don't need to denoise it.

The problem statement is:

$$x^* = \underset{x}{\operatorname{argmin}} \|Ax - b\|_1, \quad \left( \|Ax - b\|_1 \equiv \sum_{i=1}^n |a_i^T x - b_i|, \quad a_i^T \text{ is } i\text{th row of } A \right)$$

Formulation and solving:

$$\begin{aligned} \min \quad & \|z\|_1 \\ \text{subj } & Ax - b - z = 0 \end{aligned}$$

$$\begin{aligned} L(x, y, z) &= \|z\|_1 + \frac{\beta}{2} \|Ax - b - z + u\|^2 \\ \Delta L_x(x, y, z) &= -A^T (Ax - b - z + u) = 0 \\ \Rightarrow x_{k+1} &= (A^T A)^{-1} A^T (b + z^k - u^k) \end{aligned}$$

Instead of naively solving this system, we can use a more efficient approach by finding the Cholesky Decomposition using an upper and lower triangular matrix.

$$\Delta L_z(x, y, z) = 1 - \beta(Ax - b - z + u) = 0$$

$$\hat{Ax} = \alpha Ax_{k+1} + (1 - \alpha)(z_k + \beta)$$

$$\Rightarrow z_{k+1} = S_\lambda \left( \hat{Ax} - b + u^k \right)$$

$$\text{where } S_\lambda(c) = \begin{cases} c - \lambda(\text{sign}(c)) & \lambda \leq c \\ 0 & \lambda > c \end{cases}$$

Updating  $z$  with relaxation, here  $S$  is elementwise shrinkage of  $\lambda$  towards zero.

$$\Delta L_u(x, y, z) = Ax^{k+1} - b - z^{k+1} = 0$$

$$u_{k+1} = u_k + Ax^{k+1} - b - z^{k+1}$$

**MATLAB code to execute this is as follows:**

```
function [x,z,u] = lad(A, b, rho, alpha)

% lad Least absolute deviations fitting via ADMM
% Solves the following problem via ADMM:
%
% minimize ||Ax - b||_1
%
% rho is the augmented Lagrangian parameter and alpha is the over-
relaxation parameter

MAX_ITER = 1000;
[m n] = size(A);
x = zeros(n,1);
z = zeros(m,1);
u = zeros(m,1);

for k = 1:MAX_ITER
    if k > 1
        x = R \ (R' \ (A'*(b + z - u)));
    else
        R = chol(A'*A);
        x = R \ (R' \ (A'*(b + z - u)));
    end

    zold = z;
    Ax_hat = alpha*A*x + (1-alpha)*(zold + b);
    z = shrinkage(Ax_hat - b + u, 1/rho);
```

```

        u = u + (Ax_hat - z - b);
    end
end

function y = shrinkage(a, kappa)
    y = max(0, a-kappa) - max(0, -a-kappa);
end

```

### 3. Basic ADMM for Linear and Quadratic Programming:

Quadratic Programming is the process of solving a specific type of mathematical problem (specifically linear constrained) i.e. optimizing a quadratic function of several variables subject to linear constraints on these variables. A quadratic problem having  $n$  variables and  $m$  constraints is formulated as follows:

$$\begin{aligned}
 &n \times 1 \text{ vector } q \\
 &n \times n \text{ symmetric matrix } P \\
 &m \times n \text{ real matrix } A \\
 &m \times 1 \text{ vector } b
 \end{aligned}$$

For linear programming (a special case of quadratic programming with  $P=0$ ), We reformulate the primal LPP problem to solve using ADMM as follows:

$$\begin{aligned}
 &\min c^T x_1 \\
 &\text{subj to: } Ax_1 = b; \\
 &x_1 - x_2 = 0; x_2 \geq 0
 \end{aligned}$$

The augmented Lagrangian Function is as follows

$$L_p(x_1, x_2, y, s) = c^T x_1 - y^T (Ax_1 - b) - s^T (x_1 - x_2) + \frac{\beta}{2} (\|Ax_1 - b\|^2 + \|x_1 - x_2\|^2)$$

Using the first order gradient condition, we get:

$$\nabla L_{p, x_1}(x_1, x_2, y, s) = c^T - A^T y - s + \beta(A^T (Ax_1 - b) + (x_1 - x_2)) = 0$$

This can also be written as:

$$x_1^{k+1} = \arg \min_{x: Ax=b} f(x_1) + g(x_2) + \frac{1}{2\lambda} \|x_1^k - x_2^k + u^k\|^2;$$

$$i.e. \text{prox}_f(x_2^k - u^k)$$

$$\left( \text{here } u^k = \frac{1}{\rho} y^k; \rho = \frac{1}{\lambda} \right)$$

We obtain first update step for  $x_1$  as:

$$x_1^{k+1} = (A^T A + I)^{-1} + \frac{1}{\beta} (A^T y^k + s^k - c + A^T b^k + x_2^k)$$

Similarly, we update  $x_2$  as:

$$x_2^{k+1} = \arg \min_{x: Ax=b} \left( f(x_1) + g(x_2) + \frac{1}{2\lambda} \|x_1^{k+1} - x_2^k + u^k\|^2 \right);$$

$$i.e. \text{prox}_g(x_1^k - u^k)$$

This expression is same as using shrinkage function used in SVM Updating  $y$  and  $s$  using steepest ascent, we get:

$$y^{k+1} = y^k + \beta (Ax_1 - b)$$

$$s^{k+1} = s^k - \beta (x_1^{k+1} - x_2^{k+1})$$

We terminate the iterative when solution meets primal feasible i.e.

$$Ax_1^{k+1} - b \approx 0 \text{ and } x_1^{k+1} - x_2^{k+1} \approx 0$$

For Quadratic Programming, we formulate:

$$\min \left( \frac{1}{2} \right) x^T P x + q^T x$$

$$subj: Ax = b; x \geq 0$$

$$\text{here } f(x) = \min \left( \frac{1}{2} \right) x^T P x + q^T x;$$

$g(z)$  is Indicator Function for  $z \geq 0$

$$L_\rho(x, z, u) = f(x) + g(z) + \left( \frac{1}{2\lambda} \right) \|x - z + u\|_2^2 + v^T (Ax - b)$$

$$\frac{\partial L}{\partial x} = Px^{k+1} + q + \frac{1}{\lambda} (\|x^{k+1} - z^k + u^k\|) + A^T v$$



$$\left(P + \frac{1}{\lambda}I\right)x^{k+1} + q + A^T v - \frac{1}{\lambda}(z^k - u^k) = 0$$

$$\frac{\partial L}{\partial v} = 0 \Rightarrow Ax^{k+1} - b = 0$$

In the form of a matrix, it can be written and solved as follows:

$$\begin{bmatrix} P + \frac{1}{\lambda}I & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} x^{k+1} \\ v \end{bmatrix} + \begin{bmatrix} q - \frac{1}{\lambda}(z^k - u^k) \\ -b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x^{k+1} \\ v \end{bmatrix} = \begin{bmatrix} P + \frac{1}{\lambda}I & A^T \\ A & 0 \end{bmatrix}^{-1} \begin{bmatrix} \frac{1}{\lambda}(z^k - u^k) - q \\ b \end{bmatrix}$$

From:

$$x_2^{k+1} = \arg \min_{x: Ax=b} \left( f(x_1) + g(x_2) + \frac{1}{2\lambda} \|x^{k+1} - z^k + u^k\|^2 \right)$$

$$u^{k+1} = \arg \min_{x: Ax=b} f(x_1) + g(x_2) + \frac{1}{2\lambda} \|x^{k+1} - z^{k+1} + u^k\|^2$$

We finally arrive at:

$$z^{k+1} = (x^{k+1} + u^k)_+$$

$$u^{k+1} = u^k + x^{k+1} - z^{k+1}$$

Another key point to note is that both Linear and Quadratic programming can be solved using the same Quadratic programming problem where  $P=0$  as mentioned earlier.

**Code to perform both LP and QP is as follows:**

```
clc;clear all;close all;

n = 500; % dimension of x
m = 400; % number of equality constraints

c = rand(n,1) + 0.5; % create nonnegative price vector with mean 1
x0 = abs(randn(n,1)); % create random solution vector

A = abs(randn(m,n)); % create random, nonnegative matrix A
b = A*x0;
[X,Z,U] = linprog(c, A, b, 1.0, 1.0); %function call
```

```

function [x,z,u] = linprog(c, A, b, rho, alpha)
% Linear Programming Solves the following problem via ADMM:
%
%   minimize      c'*x
%   subject to    Ax = b, x >= 0
% rho is the augmented Lagrangian parameter and alpha is the over-
relaxation parameter

    MAX_ITER = 1000;

    [m n] = size(A);
    x = zeros(n,1);
    z = zeros(n,1);
    u = zeros(n,1);

    for k = 1:MAX_ITER

        % x-update
        tmp = [ rho*eye(n), A'; A, zeros(m) ] \ [ rho*(z - u) - c; b ];
        x = tmp(1:n); %extract 1st 500 values of temp

        % z-update with relaxation
        zold = z;
        x_hat = alpha*x + (1 - alpha)*zold;
        z = pos(x_hat + u);

        u = u + (x_hat - z);
    end
end

```

## 4. Interior point ADMM for sparse covariance:

Given a dataset consisting of samples from a zero mean Gaussian distribution in  $\mathbb{R}^n$ , consider the task of estimating the covariance matrix  $\Sigma$  under the prior assumption that the inverse of it is sparse. Since  $\Sigma^{-1}_{ij}$  is zero if and only if the  $i^{\text{th}}$  and  $j^{\text{th}}$  components of the random variable are conditionally independent, this problem is equivalent to the *structure learning* problem of estimating the topology of the undirected graphical model representation of the Gaussian. Determining the sparsity pattern of the inverse covariance matrix  $\Sigma^{-1}$  is also called the *covariance selection problem*.

For a very small value of  $n$ , it is feasible to search over all sparsity patterns in  $\Sigma^{-1}$  since for a fixed sparsity pattern, determining the maximum likelihood estimate of  $\Sigma$  is a tractable (convex optimization) problem. A good experiment that scales to much larger values of  $n$  is to minimize the negative log-likelihood (with respect to  $\Sigma^{-1}$ ) with an  $l_1$  regularization term to promote sparsity of the estimated inverse covariance matrix.

If  $c$  is the empirical covariance matrix  $\frac{1}{N} \sum_{i=1}^N a_i a_i^T$ ,

The estimation problem can be written as:

$$\min(c^T x) - \log(\det(x)) + \lambda \|x\|_1.$$

This is a special case of the  $l_1$  regularised problem with the initial part  $(c^T x) - \log(\det(x))$  as the convex loss function.

Hence we reformulate primal LPP using log barrier function as:

$$\begin{aligned} \min_{x_1, x_2} \quad & c^T x_1 - \mu \sum_j \ln(x_{2,j}) \\ \text{subj.} \quad & Ax_1 = b; \quad x_1 - x_2 = 0 \{\text{Implicit}\} \end{aligned}$$

Thus, Lagrangian function with augmented term and equations on partial derivative can be written as:

$$L_{P,\mu}(x_1, x_2, y, s) = c^T x_1 - \mu \sum_j \ln(x_{2,j}) - y^T (Ax_1 - b) - s^T (x_1 - x_2) + \frac{\beta}{2} (\|Ax_1 - b\|^2 + \|x_1 - x_2\|^2)$$

$$\Delta L_{P,\mu,x_1}(x_1, x_2, y, s) = c^T - A^T y - s + \beta(A^T (Ax_1 - b) + (x_1 - x_2)) = 0$$

$$\Delta L_{P,\mu,x_2}(x_1, x_2, y, s) = s - \frac{\mu}{x_2} + \beta(x_2 - x_1)$$

Update becomes:

$$x_1^{k+1} = (A^T A + I)^{-1} + \frac{1}{\beta} (A^T y^k + s^k - c + A^T b^k + x_2^k)$$

$$x_2 = \frac{1}{2\beta} \left( \beta x_1 - s_j + \sqrt{\beta^2 x_{2,i,j}^2 - 2\beta s_j x_{1,j} + s_j^2 + 4\beta u} \right) \text{ for } j=1, \dots, n$$

The final variable update is same as the primal LPP:

$$y^{k+1} = y^k + \beta (Ax_1 - b)$$

$$s^{k+1} = s^k - \beta (x_1^{k+1} - x_2^{k+1})$$

Here  $\beta$  is the step size

At the end of iteration we check for  $\mu_{k+1} = \omega \mu_k$ ;  $\omega < 1$

## MATLAB Snippet is as follows:

```
clc;clear all;close all;
n = 100; % number of features
N = 10*n; % number of samples

% generate a sparse positive definite inverse covariance matrix
Sinv = diag(abs(ones(n,1)));
idx = randsample(n^2, 0.001*n^2);
Sinv(idx) = ones(numel(idx), 1);
Sinv = Sinv + Sinv'; % make symmetric
if min(eig(Sinv)) < 0 % make positive definite
    Sinv = Sinv + 1.1*abs(min(eig(Sinv)))*eye(n);
end
S = inv(Sinv);

% generate Gaussian samples
D = mvnrnd(zeros(1,n), S, N);
[X] = covsel(D, 0.01, 1, 1);

function [Z] = covsel(D, lambda, rho, alpha)
% covsel Sparse inverse covariance selection via ADMM
% Solves the following problem via ADMM:
%
% minimize trace(S*X) - log det X + lambda*||X||_1
%
% with variable X, where S is the empirical covariance of the data
% matrix D (training observations by features).
% rho is the augmented Lagrangian parameter and alpha is the over-
% relaxation parameter

    MAX_ITER = 1000;

    S = cov(D);
    n = size(S,1);

    X = zeros(n);
    Z = zeros(n);
    U = zeros(n);

    for k = 1:MAX_ITER
        % x-update
        [Q,L] = eig(rho*(Z - U) - S);
        es = diag(L);
        xi = (es + sqrt(es.^2 + 4*rho))./(2*rho);
        X = Q*diag(xi)*Q';

        % z-update with relaxation
        Zold = Z;
        X_hat = alpha*X + (1 - alpha)*Zold;
        Z = shrinkage(X_hat + U, lambda/rho);

        U = U + (X_hat - Z);
    end
end

function y = shrinkage(a, kappa)
    y = max(0, a-kappa) - max(0, -a-kappa);
end
```

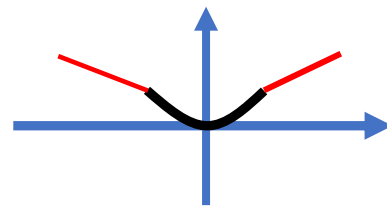
## 5. ADMM for Huber loss:

Huber loss is a loss function that lies between least square error and LAD. In brief, this function is such that it considers absolute error like LAD if absolute value lies after a particular threshold based on the problem(usually 1), and squared error like LSQR for absolute values below the same threshold.

Thus:

$$huber(z_i) = \frac{1}{2}z_i^2, \quad |z_i| \leq 1$$

$$huber(z_i) = |z_i| - \frac{1}{2}, \quad |z_i| > 1$$



Huber Loss function plot

The Optimisation problem now is:

$$\min_{x,z} g(z) = \sum_{i=1}^m huber(z_i)$$

*subject to*  $z = Ax - b$

*ADMM* formulation is

$$\min_{x,z} g(z)$$

*subject to*  $z = Ax - b$

$$\Rightarrow Ax - b - z = \mathbf{0} \text{ vector}$$

*Augmented Lagrangian* is

$$L(x,z,u) = g(z) + \frac{1}{2\lambda} \|Ax - b - z + u\|_2^2$$

Update x:

$$L(x, z, u) = g(z) + \frac{1}{2\lambda} \|Ax - b - z + u\|_2^2$$

$$x^{k+1} = \arg \min_x \|Ax - b - z^k + u^k\|_2^2$$

Solution :

$$\begin{aligned} \text{Let } f(x) &= \|Ax - b - z^k + u^k\|_2^2 \\ &= (Ax - b - z^k + u^k)^T (Ax - b - z^k + u^k) \end{aligned}$$

$$\nabla f(x) = \mathbf{0} \Rightarrow A^T (Ax - b - z^k + u^k) = \mathbf{0}$$

$$A^T Ax = A^T (b + z^k - u^k)$$

$$x^{k+1} = (A^T A)^{-1} A^T (b + z^k - u^k)$$

Update z:

$$L(x, z, u) = g(z) + \frac{1}{2\lambda} \|Ax - b - z + u\|_2^2$$

$$z^{k+1} = \arg \min_z g(z) + \frac{1}{2\lambda} \|Ax^{k+1} - b - z + u^k\|_2^2$$

Solution :  $g(z) = \text{Huber}(z)$  consists of two piecewise continuous functions .

Let us consider first one.

$$\text{Let } g(z) = \frac{z^2}{2}, \text{ then}$$

$$z^{k+1} = \arg \min_z \frac{z^2}{2} + \frac{1}{2\lambda} \|Ax^{k+1} - b - z + u^k\|_2^2$$

$$z - \frac{1}{\lambda} (Ax^{k+1} - b - z + u^k) = 0 \Rightarrow (1 + \lambda)z^{k+1} = (Ax^{k+1} - b + u^k)$$

$$z_i^{k+1} = \frac{1}{1 + \lambda} (Ax^{k+1} - b + u^k)_i \text{ ----- (1)}$$

$$(Ax^{k+1} - b + u^k)_i \text{ is the } i\text{th element of vector } (Ax^{k+1} - b + u^k)$$

*Update z:* Consider second function

Let  $g(z) = |z| - \frac{1}{2}$ , then

$$z^{k+1} = \arg \min_z |z| - \frac{1}{2} + \frac{1}{2\lambda} \|Ax^{k+1} - b - z + u^k\|_2^2$$

$$z^{k+1} = \arg \min_z |z| + \frac{1}{2\lambda} \|Ax^{k+1} - b - z + u^k\|_2^2$$

$$z_i^{k+1} = S_\lambda \left( (Ax^{k+1} - b + u^k)_i \right) = (Ax^{k+1} - b + u^k)_i \pm \lambda \text{ -----(2)}$$

$S_\lambda$  represent shrinkage by  $\pm \lambda$ , sign depending on the sign of  $(Ax^{k+1} - b + u^k)_i$   
 $(Ax^{k+1} - b + u^k)_i$  is the  $i$ th element of vector  $(Ax^{k+1} - b + u^k)$

Then we combine the two formula (1) and (2) for  $z$  update.

*Split  $z_i^{k+1}$  in (2) into two terms, one term containing (1) and express the residual in terms of a Shrinkage function*

$$z_i^{k+1} = (Ax^{k+1} - b + u^k)_i \pm \lambda \text{ -----(2)}$$

$$\begin{aligned} z_i^{k+1} &= (Ax^{k+1} - b + u^k)_i \pm \lambda = \frac{1}{1+\lambda} (Ax^{k+1} - b + u^k)_i + \frac{\lambda}{1+\lambda} \left\{ (Ax^{k+1} - b + u^k)_i \pm (1+\lambda) \right\} \\ &= \frac{1}{1+\lambda} (Ax^{k+1} - b + u^k)_i + \frac{\lambda}{1+\lambda} S_{1+\lambda} (Ax^{k+1} - b + u^k)_i \end{aligned}$$

or

$$z^{k+1} = \frac{1}{1+\lambda} (Ax^{k+1} - b + u^k) + \frac{\lambda}{1+\lambda} S_{1+\lambda} (Ax^{k+1} - b + u^k)$$

## 6. ADMM for non-convex Regression Selector:

As an example, consider the least squares regressor selection or feature selection problem,

$$\begin{aligned} x^* &= \arg \min_x \frac{1}{2} \|Ax - b\|_2^2 \\ &\text{subject to } \text{card}(x) \leq k \end{aligned}$$

which is to find the best fit to  $b$  as a linear combination of no more than  $c$  columns of  $A$ . For this problem, ADMM takes the form above, where the  $x$ -update involves a regularized least squares problem, and the  $z$ -update involves keeping the  $c$  largest magnitude elements of  $x_{k+1} + u_k$ .

This is just like ADMM for the lasso, except that soft thresholding is replaced with 'hard' thresholding. This close connection is hardly surprising, since lasso can be thought of as a heuristic for solving the regressor selection problem. From this viewpoint, the lasso controls the trade-off between least squares error and sparsity through the parameter  $\lambda$ , whereas in ADMM for regressor selection, the same trade-off is controlled by the parameter  $c$ , the exact cardinality desired.

**Problem formulation is as follows:**

$$\min_{x,z} f(x) + g(z) = \frac{1}{2} \|Ax - b\|_2^2 + g(z)$$

$$g(z) = \begin{cases} 0 & \text{if } \text{card}(z) \leq k \\ \infty, & \text{if } \text{card}(z) > k \end{cases}$$

*ADMM* formulation :

*Augmented* Lagrangian is

$$\begin{aligned} L(x, z, u) &= f(x) + g(z) + \frac{1}{2\lambda} \|x - z + u\|_2^2 \\ &= \frac{1}{2} \|Ax - b\|_2^2 + g(z) + \frac{1}{2\lambda} \|x - z + u\|_2^2 \end{aligned}$$

x-update is as follows:

$$x^{k+1} = \arg \min_x \frac{1}{2} \|Ax - b\|_2^2 + \frac{1}{2\lambda} \|x - z^k + u^k\|_2^2$$

*Solution* : Let  $\rho = \frac{1}{\lambda}$

$$A^T (Ax - b) + \rho (x - z^k + u^k) = 0 \text{ vector}$$

$$(A^T A + \rho I) x = A^T b + \rho (z^k - u^k)$$

$$x^{k+1} = (A^T A + \rho I)^{-1} (A^T b + \rho (z^k - u^k))$$

z-update is as follows:

$$L(z) = g(z) + \frac{\rho}{2} \|x^{k+1} - z + u^k\|_2^2$$



From the latter term, we get updated  $z$  and project it to get:

$$z^{k+1} = \max(0, x^{k+1} + u^k)$$

{in accordance with  $g(z)$ ->Indicator Func<sup>n</sup>}

The original lagrangian multiplier term is  $u^T(x - z)$

$\therefore$  The gradient w.r.t.  $u$  is  $x - z$

$$\Rightarrow u^{k+1} = u^k + (x^{k+1} - z^{k+1})$$

## MATLAB Code to execute is as follows:

```
clc;clear all;close all;
randn('seed', 0);
rand('seed', 0);

m = 1500;          % number of examples
n = 5000;          % number of features
p = 100/n;         % sparsity density

% generate sparse solution vector
x = sprandn(n,1,p);

% generate random data matrix
A = randn(m,n);

% normalize columns of A
A = A*spdiags(1./sqrt(sum(A.^2))', 0, n, n);

% generate measurement b with noise
b = A*x + sqrt(0.001)*randn(m,1);

xtrue = x; % save solution
x = regressor_sel(A, b, p*n, 1.0);
function z = regressor_sel(A, b, K, rho)
% regressor_sel Solve lasso problem via ADMM
%
% Attempts to solve the following problem via ADMM:
%
% minimize || Ax - b ||_2^2
% subject to card(x) <= K
%
% where card() is the number of nonzero entries.
%
% The solution is returned in the vector x.
%
% rho is the augmented Lagrangian parameter and alpha is the over-
relaxation parameter
MAX_ITER = 1000;
[m, n] = size(A);

% save a matrix-vector multiply
Atb = A'*b;
x = zeros(n,1);
z = zeros(n,1);
```

```

u = zeros(n,1);

% cache the factorization
[L U] = factor(A, rho);
for k = 1:MAX_ITER

    % x-update
    q = Atb + rho*(z - u); % temporary value
    if( m >= n ) % if skinny
        x = U \ (L \ q);
    else % if fat
        x = q/rho - (A'*(U \ (L \ (A*q))))/rho^2;
    end

    % z-update with relaxation
    zold = z;
    z = keep_largest(x + u, K);

    % u-update
    u = u + (x - z);
end

end

function z = keep_largest(z, K)
    [val,pos] = sort(abs(z), 'descend');
    z(pos(K+1:end)) = 0;
end

function [L U] = factor(A, rho)
    [m, n] = size(A);
    if ( m >= n ) % if skinny
        L = chol( A'*A + rho*speye(n), 'lower' );
    else % if fat
        L = chol( speye(m) + 1/rho*(A*A'), 'lower' );
    end

    % force matlab to recognize the upper / lower triangular structure
    L = sparse(L);
    U = sparse(L');
end

```

## 7. ADMM for total variable minimisation: (One of the powerful tools used in Signal and Image processing)

Problem formulation is as follows:

$$\text{minimize } \frac{1}{2} \|x - b\|_2^2 + \lambda \sum_{i=1}^{n-1} |x_{i+1} - x_i|, \quad x \in \mathbb{R}^n$$

$$\Rightarrow \text{minimize } \frac{1}{2} \|x - b\|_2^2 + \lambda \|z\|_1$$

*Subject to*  $z = Dx$

*Augmented Lagrangian* is

$$L(x, z, y) = \frac{1}{2} \|x - b\|_2^2 + \lambda \|z\|_1 + y^T (Dx - z) + \frac{\rho}{2} \|Dx - z\|_2^2$$

$$\underbrace{\begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ z_4 \\ \vdots \\ z_{n-2} \\ z_{n-1} \end{pmatrix}}_z = \underbrace{\begin{pmatrix} -1 & 1 & . & . & . & . & . & . \\ . & -1 & 1 & . & . & . & . & . \\ . & . & -1 & 1 & . & . & . & . \\ . & . & . & -1 & 1 & . & . & . \\ . & . & . & . & -1 & 1 & . & . \\ . & . & . & . & . & . & -1 & 1 \end{pmatrix}}_D \underbrace{\begin{pmatrix} x_1 \\ x_2 \\ x_2 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{pmatrix}}_x$$

Representation of constraint in matrix format.

x-update is as follows:

$$L(x, z^k, y^k) = \frac{1}{2} \|x - b\|_2^2 + (y^k)^T Dx + \frac{\rho}{2} \|Dx - z^k\|_2^2$$

$$x - b + D^T y^k + \rho D^T (Dx - z^k) = 0$$

$$(I + \rho D^T D)x = b - D^T y^k + \rho D^T z^k = b + \rho D^T \left( z^k - \frac{1}{\rho} y^k \right)$$

$$x^{k+1} = (I + \rho D^T D)^{-1} \left( b + \rho D^T \left( z^k - \frac{1}{\rho} y^k \right) \right) \dots (1)$$

z-update is as follows:

*Augmented Lagrangian* is

$$L(x, z, y) = \frac{1}{2} \|x - b\|_2^2 + \lambda \|z\|_1 + y^T (Dx - z) + \frac{\rho}{2} \|Dx - z\|_2^2$$

$$L(x^{k+1}, z, y^k) = \lambda \|z\|_1 - (y^k)^T z + \frac{\rho}{2} \|Dx^{k+1} - z\|_2^2$$

$$L(x^{k+1}, z, y^k) = \|z\|_1 + \frac{\rho}{2\lambda} \left\| Dx^{k+1} - z + \frac{y^k}{\rho} \right\|_2^2$$

$$z^{k+1} = S_{\lambda/\rho} \left( Dx^{k+1} + \frac{y^k}{\rho} \right)$$

Similarly, y(scaled u update) is:

$$y^{k+1} = y^k + Dx^{k+1} - z^k$$

## Code to execute is as follows:

```
clc;clear all;close all;
rand('seed', 0);
randn('seed', 0);
n = 100;
x0 = ones(n,1);

for j = 1:3
    idx = randsample(n,1);
    k = randsample(1:10,1);
    x0(ceil(idx/2):idx) = k*x0(ceil(idx/2):idx);
end

b = x0 + randn(n,1);
lambda = 5;
e = ones(n,1);
D = spdiags([e -e], 0:1, n,n);

x = total_variation(b, lambda, 1.0, 1.0);

function x = total_variation(b, lambda, rho, alpha)
% total_variation Solve total variation minimization via ADMM
% Solves the following problem via ADMM:
%
% minimize (1/2)||x - b||_2^2 + lambda * sum_i |x_{i+1} - x_i|
% where b in R^n.
%
% The solution is returned in the vector x.
%
% rho is the augmented Lagrangian parameter and alpha is the over-
relaxation parameter
    MAX_ITER = 1000;
    n = length(b);

    % difference matrix
    e = ones(n,1);
    D = spdiags([e -e], 0:1, n,n);
    x = zeros(n,1);
    z = zeros(n,1);
    u = zeros(n,1);
    I = speye(n);
    DtD = D'*D;

    for k = 1:MAX_ITER
        % x-update
        x = (I + rho*DtD) \ (b + rho*D'*(z-u));

        % z-update with relaxation
        zold = z;
        Ax_hat = alpha*D*x + (1-alpha)*zold;
        z = shrinkage(Ax_hat + u, lambda/rho);

        % y-update
        u = u + Ax_hat - z;
    end
end

function y = shrinkage(a, kappa)
y = max(0, a-kappa) - max(0, -a-kappa);
end
```

## 8. Lasso using ADMM:

An important special case of regularized linear regression, also called the *lasso*. Lambda is a scalar regularization parameter that is usually chosen by cross-validation. In typical applications, there are many more features than training examples, and the goal is to find a parsimonious model for the data. The lasso has been widely applied, particularly in the analysis of biological data, where only a small fraction of a huge number of possible factors are actually predictive of some outcome of interest. It is prevalent across machine learning, model selection in statistics.

Problem formulation and solution is as follows:

$$\min f(x) + g(x)$$

$$\text{where } f(x) = \frac{1}{2} \|Ax - b\|^2; g(x) = \tau \|z\|_1$$

$$\text{subj } x - z = 0$$

$$L(x, y, z) = (Ax - b)^T (Ax - b) + \tau \|z\|_1 + y^T (x - z) + \frac{\beta}{2} (x - z)^T (x - z)$$

$$L(x, y, z) = (Ax - b)^T (Ax - b) + \tau \|z\|_1 + \frac{\beta}{2} \|x - z + u\|^2$$

$$\Delta L_x(x, y, z) = A^T (Ax - b) + \beta (x - z + u) = 0$$

$$\Rightarrow x_{k+1} = (A^T A + \beta I)^{-1} (A^T b + \beta (z^k - u^k))$$

$$\Delta L_z(x, y, z) = \tau + \beta (x - z + u) = 0$$

$$\Rightarrow z_{k+1} = S_\lambda (x^{k+1} + u^k)$$

$$\text{where } S_\lambda(c) = \begin{cases} c - \lambda(\text{sign}(c)) & \lambda \leq c \\ 0 & \lambda > c \end{cases}$$

$$\Delta L_u(x, y, z) = x^{k+1} - z^{k+1} = 0$$

$$u_{k+1} = u_k + x^{k+1} - z^{k+1}$$

## Code to execute is as follows:

```
function [x,z,u] = lasso(A, b, lambda, rho, alpha)
% lasso Solve lasso problem via ADMM
% Solves the following problem via ADMM:
%
%   minimize 1/2*|| Ax - b ||_2^2 + \lambda || x ||_1
%
% rho is the augmented Lagrangian parameter and alpha is the over-
relaxation parameter

MAX_ITER = 1000;
[m, n] = size(A);

x = zeros(n,1);
z = zeros(n,1);
u = zeros(n,1);

% cache the factorization
[L U] = factor(A, rho);

for k = 1:MAX_ITER

    % x-update
    q = Atb + rho*(z - u);      % temporary value
    if( m >= n )      % if skinny
        x = U \ (L \ q);
    else      % if fat
        x = q/rho - (A'*(U \ ( L \ (A*q) )))/rho^2;
    end

    % z-update with relaxation
    zold = z;
    x_hat = alpha*x + (1 - alpha)*zold;
    z = shrinkage(x_hat + u, lambda/rho);

    % u-update
    u = u + (x_hat - z);
end
end

function z = shrinkage(x, kappa)
z = max( 0, x - kappa ) - max( 0, -x - kappa );
end

function [L U] = factor(A, rho)
[m, n] = size(A);
if ( m >= n )      % if skinny
    L = chol( A'*A + rho*speye(n), 'lower' );
else      % if fat
    L = chol( speye(m) + 1/rho*(A*A'), 'lower' );
end

% force matlab to recognize the upper / lower triangular structure
L = sparse(L);
U = sparse(L');
end
```

## Essential Methods that help in speeding up the process of solving problems using ADMM:

### a) Block Splitting:

The concepts of Computer software were written conventionally for serial computing. This meant that to solve a problem, an algorithm divides the problem into smaller instructions. These discrete instructions are then executed on Central Processing Unit of a computer one by one. Only after one instruction is finished, next one starts.

Parallel Computing is the use of multiple processing elements simultaneously for solving any problem. Problems are broken down into instructions and are solved concurrently as each resource which has been applied to work is working at the same time.

### Primal LPP Block Splitting:

Since the size of data needed to be optimised in ADMM is very huge, we need to divide into B blocks of equal sizes as follows:

$$x_1 = \begin{bmatrix} x_{1,1} \\ \cdot \\ \cdot \\ \cdot \\ x_{1,B} \end{bmatrix}; A = [A_1 \quad \cdot \quad \cdot \quad \cdot \quad A_B]; c = \begin{bmatrix} c_1 \\ \cdot \\ \cdot \\ \cdot \\ c_B \end{bmatrix}$$

For the sake of derivation, we will do it using 2 blocks. So values become:

$$x_1 = \begin{bmatrix} x_{1,1} \\ x_{1,2} \end{bmatrix}; A = [A_1 \quad A_2]; c = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$L_P(x_1, x_2, y) = c_1^T x_{1,1} + c_2^T x_{1,2} - y^T (A_1 x_{1,1} + A_2 x_{1,2} - b) - s_1^T (x_{1,1} - x_{2,1}) \dots$$

$$\dots - s_2^T (x_{1,2} - x_{2,2}) + \frac{\beta}{2} \left( \|A_1 x_{1,1} + A_2 x_{1,2} - b\|^2 + \|x_{1,1} - x_{2,1}\|^2 + \|x_{1,2} - x_{2,2}\|^2 \right)$$

$$\Delta L_{P,x_{1,1}}(x_1, x_2, y) = c_1 - A_1^T y - s_1 + \beta \left( A_1^T (A_1 x_{1,1} + A_2 x_{1,2} - b) + (x_{1,1} - x_{2,1}) \right)$$

$$\Delta L_{P,x_{1,2}}(x_1, x_2, y) = c_2 - A_2^T y - s_2 + \beta \left( A_2^T (A_1 x_{1,1} + A_2 x_{1,2} - b) + (x_{1,2} - x_{2,2}) \right)$$

Our Update statements will be:

$$\Rightarrow x_{1,1}^{k+1} = (A_1^T A_1 + I)^{-1} \left( \frac{1}{\beta} (-c_1 + A_1^T y + s_1) + (A_1^T b + x_{2,1}^k - A_1^T A_2 x_{1,2}^k) \right)$$

$$\Rightarrow x_{1,2}^{k+1} = (A_2^T A_2 + I)^{-1} \left( \frac{1}{\beta} (-c_2 + A_2^T y + s_2) + (A_2^T b + x_{2,1}^k - A_1^T A_2 x_{1,1}^k) \right)$$

The general statement for n blocks will be:

$$x_{1,i}^{k+1} = (A_i^T A_i + I)^{-1} \left( \frac{1}{\beta} (-c_i + A_i^T y + s_i) + \left( A_i^T b + x_{2,i}^k - \sum_{j=1; j \neq i} A_i^T A_j x_{1,j}^k \right) \right)$$

## b) Cholesky Decomposition:

In linear algebra, a matrix decomposition or matrix factorization is a factorization of a matrix into a product of matrices. There are many different matrix decompositions. One of them is Cholesky Decomposition.

The Cholesky decomposition or Cholesky factorization is a decomposition of a Hermitian, positive-definite matrix into the product of a lower triangular matrix and its conjugate transpose. The Cholesky decomposition is roughly twice as efficient as the LU decomposition for solving systems of linear equations.

The Cholesky decomposition of a Hermitian positive-definite matrix  $A$  is a decomposition of the form  $A = LL^T$ , where  $L$  is a lower triangular matrix with real and positive diagonal entries, and  $L^T$  denotes the conjugate transpose of  $L$ . Every Hermitian positive-definite matrix (and thus also every real-valued symmetric positive-definite matrix) has a unique Cholesky decomposition.

$$\begin{bmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_{00} & 0 & 0 \\ L_{10} & L_{11} & 0 \\ L_{20} & L_{21} & L_{22} \end{bmatrix} \begin{bmatrix} L_{00} & L_{10} & L_{20} \\ 0 & L_{11} & L_{21} \\ 0 & 0 & L_{22} \end{bmatrix}$$

Lower Triangular L

Transpose of L



Every symmetric, positive definite matrix A can be decomposed into a product of a unique lower triangular matrix L and its transpose:  $A = LL^T$ . The following formulas are obtained by solving above lower triangular matrix and its transpose.

$$L_{i,j} = \sqrt{A_{j,j} - \sum_{k=0}^{j-1} L_{j,k}^2}$$

$$L_{i,j} = \frac{1}{L_{j,j}} \left( A_{j,j} - \sum_{k=0}^{j-1} L_{i,k} L_{j,k} \right)$$

Example of Cholesky Decomposition is as follows:

$$\begin{bmatrix} 4 & 12 & -16 \\ 12 & 37 & -43 \\ -16 & -43 & 98 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 6 & 1 & 0 \\ -8 & 5 & 3 \end{bmatrix} \begin{bmatrix} 2 & 6 & -8 \\ 0 & 1 & 5 \\ 0 & 0 & 3 \end{bmatrix}$$

### c) LU Decomposition:

L U decomposition of a matrix is the factorization of a given square matrix into two triangular matrices, one upper triangular matrix and one lower triangular matrix, such that the product of these two matrices gives the original matrix. Basically, the L U decomposition method comes handy whenever it is possible to model the problem to be solved into matrix form. Conversion to the matrix form and solving with triangular matrices makes it easy to do calculations in the process of finding the solution.

Representation is as follows:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix}$$

$A \qquad \qquad = \qquad \qquad L \qquad \qquad U$

### Steps for LU Decomposition:

Given a set of linear equations, first convert them into matrix form  $A X = C$  where A is the coefficient matrix, X is the variable matrix and C is the matrix of numbers on the right-hand side of the equations.

Now, reduce the coefficient matrix A, i.e., the matrix obtained from the coefficients of variables in all the given equations such that for 'n' variables we have an  $n \times n$  matrix, to row echelon form using Gauss Elimination Method. The matrix so obtained is U.

To find L, we have two methods. The first one is to assume the remaining elements as some artificial variables, make equations using  $A = L U$  and solve them to find those artificial variables. The other method is that the remaining elements are the multiplier coefficients because of which the respective positions became zero in the U matrix.

Now, we have A (the  $n \times n$  coefficient matrix), L (the  $n \times n$  lower triangular matrix), U (the  $n \times n$  upper triangular matrix), X (the  $n \times 1$  matrix of variables) and C (the  $n \times 1$  matrix of numbers on the right-hand side of the equations).

The given system of equations is  $A X = C$ . We substitute  $A = L U$ . Thus, we have  $LUX = C$ . We put  $Z = U X$ , where Z is a matrix or artificial variables and solve for  $L Z = C$  first and then solve for  $U X = Z$  to find X or the values of the variables, which was required.

### Example of LU Decomposition is as follows:

$$\begin{bmatrix} 1 & 8 & 9 \\ 4 & 6 & 8 \\ -5 & -7 & -2 \end{bmatrix} = \begin{bmatrix} -0.2 & 1 & 0 \\ -0.8 & 0.06 & 1 \\ 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} -5 & -7 & -2 \\ 0 & 6.6 & 8.6 \\ 0 & 0 & 5.9 \end{bmatrix}$$

$A \qquad \qquad \qquad L \qquad \qquad \qquad U$