# LINEAR ALGEBRA AND PROBABILITY TAKE HOME EXAM.

Student Name - Chetana Sharma
Program / Year - MSDS / Spring 2020
NUID        -   001055979

1.

$$f(\lambda x + (1-\lambda)y) < \lambda f(x) + (1-\lambda)f(y)$$

Assuming $f(x) = x^T P x$ is a convex function

$$\Rightarrow f(\lambda x + (1-\lambda)y) < \lambda(x^T P x) + (1-\lambda)(y^T P y)$$

$$\Rightarrow (\lambda x + (1-\lambda)y)^T P(\lambda x + (1-\lambda)y) < \lambda x^T P x + (1-\lambda)y^T P y$$

$$\Rightarrow \lambda x^T P(\lambda x + (1-\lambda)y) + (1-\lambda)y^T P(\lambda x + (1-\lambda)y) <$$
$$\lambda x^T P x + (1-\lambda)y^T P y$$

$$\Rightarrow \lambda^2 x^T P x + \lambda(1-\lambda)x^T P y + \lambda(1-\lambda)y^T P x + (1-\lambda)^2 y^T P y <$$
$$\lambda x^T P x + (1-\lambda)y^T P y$$

\# Taking $x^T P x$ and $y^T P y$ terms from L.H.S to R.H.S

$$\Rightarrow \lambda(1-\lambda)x^T P y + \lambda(1-\lambda)y^T P x < \lambda(1-\lambda)x^T P x + \lambda(1-\lambda)y^T P y$$

Removing $\lambda(1-\lambda)$ from both L.H.S and R.H.S

$$\Rightarrow \quad x^T P y + y^T P x < x^T P x + y^T P y$$

$$\Rightarrow \quad x^T P(y-x) + y^T P(x-y) < 0$$

$$\Rightarrow \quad x^T P(y-x) - y^T P(y-x) < 0$$

$$\Rightarrow \quad (x^T P - y^T P)(y-x) < 0$$

$\Rightarrow$ $(x^T P - y^T P)(x - y) > 0$

$\Rightarrow$ $(x^T - y^T) P (x - y) > 0$

$=)$ $(x - y)^T P (x - y) > 0$

Let $x - y = z$

$=)$ $z^T P z > 0$

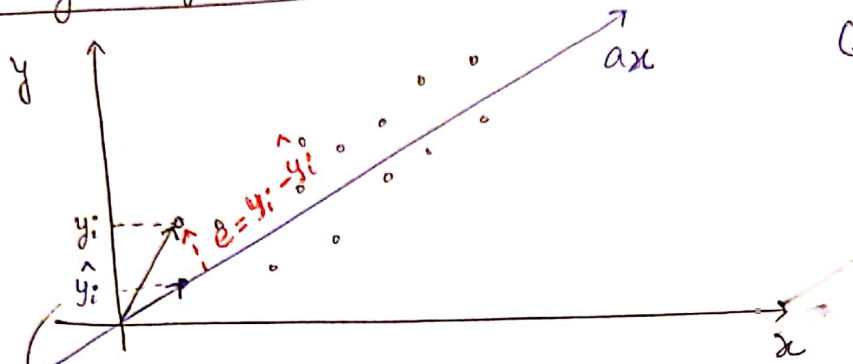$\Rightarrow$ $f(x)$ is convex function when

$z^T P z > 0$

And $z^T P z > 0$ is a property when
P is a positive - definite matrix.

# DETERMINISTIC APPROACH

Given : sets of data $y_i \in R^{n_y}$, $x_i \in R^{n_x}$ where

$$i = 1, 2, \cdots n$$

## Using Projection



Matrix General Form :

$$\hat{Y} = XA$$

$\rightarrow$ Let $\hat{y}_i$ denote the predicted value

We need to minimise the distance between $y_i$ and $\hat{y}_i$. That is possible only when $e'$ is perpendicular to the line $ax$.

The projected value is nothing but a projection of the actual point on the linear line / plane or hyperplane.

So if $\hat{Y} = XA$

error $E = Y - \hat{Y}$

And $\langle X, E \rangle = 0$

$\Rightarrow \quad X^T(Y - \hat{Y}) = 0 \quad \Rightarrow X^T(Y - XA) = 0$

$\Rightarrow \quad X^T Y = X^T X A$

$\Rightarrow \quad \boxed{A = (X^T X)^{-1} X^T Y}$

# 2. PROBABILISTIC APPROACH

**Given:** $y_i \in R^{n_y}$, $x_i \in R^{n_x}$, $i \in 1,2,3,--n$

$$y = Ax$$

### for a scalar case:

$$y_i = x_i a \quad \text{or} \quad Y = XA$$

Assuming that each $y_i$ is Gaussian distributed with mean $x_i a$.

$$y_i = N(x_i a, \sigma^2)$$

$$P(Y|X,A,\sigma) = \prod_{i=1}^{n} p(y_i|x_i, A, \sigma)$$

$$= \prod_{i=1}^{n} (2\pi\sigma^2)^{-1/2} e^{-1/2\sigma^2 (y_i - x_i A)^2}$$

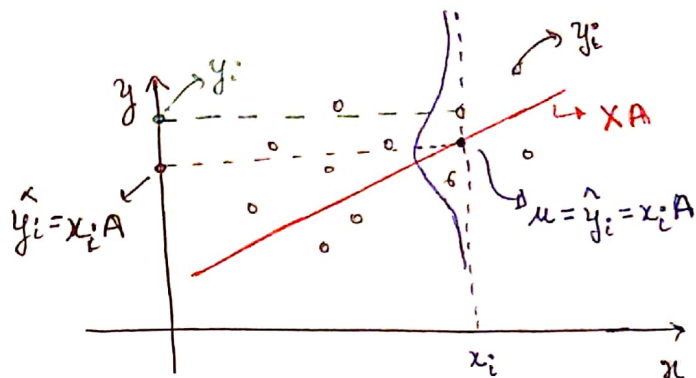$$= (2\pi\sigma^2)^{-n/2} e^{-1/2\sigma^2 \sum_{i=1}^{n} (y_i - x_i a)^2}$$

$$= (2\pi\sigma^2)^{-n/2} e^{-1/2\sigma^2 (Y - XA)^T (Y - XA)}$$

Taking log on both sides

$$\log(P(Y|X,A,\sigma)) = -\frac{n}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(Y-XA)^T(Y-XA)$$

We need to find A that maximises the likelihood of seeing the training data Y given X.

So we differentiate w.r.t. A & equate to zero.

$$\ell(A) = \frac{-n}{2}(2\pi\sigma^2) - \frac{1}{2\sigma^2}(Y - XA)^T(Y - XA)$$

$$\frac{\partial \ell}{\partial A} = \frac{-1}{2\sigma^2}\frac{\partial}{\partial A}\left(Y^TY - Y^TXA - A^TX^TY + A^TX^TXA\right)$$

$$= \frac{-1}{2\sigma^2}\left(-Y^TX - X^TY + 2X^TXA\right)$$

$$\frac{\partial \ell}{\partial A} = 0$$

$$\Rightarrow \frac{-1}{2\sigma^2}\left(-2X^TY + 2X^TXA\right) = 0$$

$$\Rightarrow A = \frac{X^TY}{X^TX}$$

$$\Rightarrow \boxed{A = (X^TX)^{-1}X^TY}$$

3. $x_0 \in R^n$, $\boxed{x_{k+1} = A x_k - ①}$ → Given

Using ①

$x_1 = A x_0$

$x_2 = A x_1$ ⟹ $x_2^{&} = A(A x_0)$ ⟹ $x_2 = A^2 x_0$

⟹ $\boxed{x_k = A^k x_0}$

__Expectation of $x_1$ and $x_k$.__

$x_1 = A x_0$

Applying expectation operator both sides.

$E[x_1] = E[A x_0]$

⟹ $\boxed{E[x_1] = A E[x_0]}$

For $x_k = A^k x_0$.

$\boxed{E[x_k] = A^k E[x_0]}$ — ②

__Considering $n = 1$__ and writing in terms of eigen values

$A x_0 = \lambda x_0$ where $\lambda$ = eigen value and $x_0$ = eigen vector of A.

$x_1 = A x_0$ [Using ①]

$E[x_1] = E[\lambda x_0]$

⟹ $E[x_1] = \lambda E[x_0]$

Similarly, we can write

$\boxed{E[x_k] = \lambda^k E[x_0]}$

So, if $|\lambda| < 1$ and $k \to \infty$

$\lambda^k \to 0$

∴ Expectation of $x_k \to 0$ for $k \to \infty$ if Eigen value of A between (-1) and (1).

## For general 'n'

Let $S$ be a matrix where all columns are eigen vectors of $A$.

$$S = \begin{bmatrix} | & | & & | \\ v_1 & v_2 & -- & v_n \\ | & | & & | \end{bmatrix}$$

Let $\Lambda$ be the eigen value matrix of $A$.

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & & 0 \\ 0 & \lambda_2 & --- & \vdots \\ \vdots & \vdots & & \vdots \\ 0 & 0 & & \lambda_n \end{bmatrix}$$

We know that, $Av_1 = \lambda_1 v_1$ , $Av_2 = \lambda_2 v_2$ $--$ $Av_n = \lambda_n v_n$

$\therefore$ $AS = S\Lambda$

Multiplying with $S^{-1}$ both sides

$$\boxed{A = S\Lambda S^{-1}}$$

Also, $A^2 = (S\Lambda S^{-1})(S\Lambda S^{-1}) = S\Lambda^2 S^{-1}$

$\therefore$ $\boxed{A^k = S\Lambda^k S^{-1}}$ $\quad$ — ③

Substituting the above value in ②

$$\boxed{E[x_k] = S\Lambda^k S^{-1} E[x_0]}$$

If $|\lambda_i| < 1$ and $k \to \infty$

then each element of $E[x_k] \to 0$

Else $E[x_k] \to \infty$

## VARIANCE

$$\text{var}(x_o) = E[x_o^2] - (E[x_o])^2$$

For a vector $x_o$

$$\text{var}(x_o) = \Sigma_{x_o} = E[x_o x_o^T] - E[x_o]E[x_o]^T \qquad - \text{④}$$

Variance expression for $x_1$

$$x_1 = A x_o$$

$$\text{var}(x_1) = E[x_1 x_1^T] - E[x_1]E[x_1]^T$$

$$\Rightarrow \text{var}(x_1) = E[(Ax_o)(Ax_o)^T] - E[Ax_o]E[Ax_o]^T$$

$$\Rightarrow \text{var}(x_1) = E[Ax_o x_o^T A^T] - AE[x_o]E[x_o]^T A^T$$

$$\Rightarrow \text{var}(x_1) = AE[x_o x_o^T]A^T - AE[x_o]E[x_o]^T A^T$$

$$\Rightarrow \text{var}(x_1) = A\left(E[x_o x_o^T] - E[x_o]E[x_o]^T\right)A^T$$

Using ④

$$\boxed{\text{var}(x_p) = A\Sigma_{x_o}A^T}$$

Variance expression for $x_k$

$$x_k = A^k x_o$$

$$\text{var}(x_k) = E[x_k x_k^T] - E[x_k]E[x_k]^T$$

$$\Rightarrow \text{var}(x_k) = E[A^k x_o x_o^T (A^k)^T] - E[A^k x_o]E[A^k x_o]^T$$

$$\Rightarrow \text{var}(x_k) = A^k E[x_o x_o^T](A^k)^T - A^k E[x_o]E[x_o]^T(A^k)^T$$

$$\Rightarrow \boxed{\text{var}(x_k) = A^k \Sigma_{x_o}(A^k)^T}$$

## Considering $n=1$.

$$\text{var}(x_k) = A^k(\text{var}(x_0))(A^k)^T$$

for $n=1$  $\quad Ax_0 = \lambda x_0 \quad$ where $\lambda$ = eigen value.

and $A^k x_0 = \lambda^k x_0$

$$\therefore \quad \boxed{\text{var}(x_k) = \lambda^{2k}\,\text{var}(x_0)}$$

For $|\lambda| < 1$ and $k \to \infty$

$\text{var}(x_k)$ tends to zero faster than the expectation.

Otherwise, $\text{var}(x_k) \to \infty$ with $k \to \infty$

## For general '$n$'

$$\text{var}(x_k) = A^k \Sigma_{x_0}(A^k)^T$$

Using ③ i.e. $A^k = S\Lambda^k S^{-1}$

$$\text{var}(x_k) = (S\Lambda^k S^{-1})\Sigma_{x_0}(S\Lambda^k S^{-1})^T$$

If each $|\lambda_i| < 1$ and $k \to \infty$

then $\text{var}(x_k) \to 0$

Else $\text{var}(x_k) \to \infty$

# LINEAR ALGEBRA AND PROBABILITY

## PROBLEM- 4

Let $X$ & $Y$ be two random variables with uniform distribution.

$$Z = \frac{X + Y}{2}$$

If $X = k$, then $Z = z$ only if $Y = 2z - k$

$$\boxed{P(Z=z) = \sum_{-\infty}^{\infty} P(X=k) \cdot P(Y=2z-k)} \quad \leftarrow \begin{array}{l} \text{FOR} \\ \text{DISCRETE} \\ \text{DISTRIBUTION OF } X \text{ \& } Y. \end{array}$$

Consider, throw of a fair six-sided dice.

$$P(X=x) = 1/6 \text{ for } x \in [1, 6]$$

Then $P(Z=z) = \sum_{k=1}^{6} P(X=k) P(Y=2z-k)$

$$\Rightarrow P(Z=z) = \frac{1}{6} \sum_{k=1}^{6} P(Y=2z-k)$$

$$P(Z=1) = \frac{1}{6} \sum_{k=1}^{6} P(Y=2-k) = \frac{1}{6} \left( P(Y=2-1) \right) = \frac{1}{6} \cdot \frac{1}{6} = \frac{1}{36}$$

$$P(Z=1.5) = \frac{1}{6} \sum_{k=1}^{6} P(Y=3-k) = \frac{1}{6} \left( P(Y=2) + P(Y=1) \right) = \frac{2}{36}$$

Similarly, we can find for other values.

$$P(Z=2) = \frac{3}{36} \quad , \quad P(Z=2.5) = \frac{4}{36} \quad , \quad P(Z=3) = \frac{5}{36}$$

$$P(Z=3.5) = \frac{6}{36} \quad , \quad P(Z=4) = \frac{5}{36} \quad , \quad P(Z=4.5) = \frac{4}{36}$$

$$P(Z=5) = \frac{3}{36} \quad , \quad P(Z=5.5) = \frac{2}{36} \quad , \quad P(Z=6) = \frac{1}{36}$$

## FOR A CONTINUOUS UNIFORM DISTRIBUTION.

$$Z = \frac{X+Y}{2}. \quad \text{Let } f(x) \text{ be the PDF for } X \text{ and } Y.$$

Let $f(x) = \begin{cases} 1 & 0 \le x \le 1 \\ 0 & \text{otherwise} \end{cases}$

then $f(z) = \int_{-\infty}^{\infty} f(x) \, f(2z-x) \, dx$

$$= \int_0^1 f(2z-x) \, dx.$$

To get the limits:
$$\Rightarrow \quad 0 \le 2z - x \le 1 \quad \Rightarrow \quad -2z \le -x \le 1 - 2z$$

$$\Rightarrow \quad \boxed{2z - 1 \le x \le 2z.}$$

$$f(z) = \begin{cases} \int_0^{2z} dx & 0 \le z \le 1/2 \\[2mm] \int_{2z-1}^1 dx & 1/2 < z \le 1 \\[2mm] 0 & \text{otherwise} \end{cases} \quad \Rightarrow \quad f(z) = \begin{cases} 2z, & 0 \le z \le 1/2 \\ 2(1-z), & 1/2 < z \le 1 \\ 0 & \text{otherwise} \end{cases}$$

Average of



PDF for the average of 2 independent random variables with distribution $f(x)$

# Problem-4 Part-B:

```
In [74]: from matplotlib import pyplot as plt
         from functools import reduce
         import seaborn as sns

         import numpy as np
         import pandas as pd
         import scipy

         from scipy import signal
```
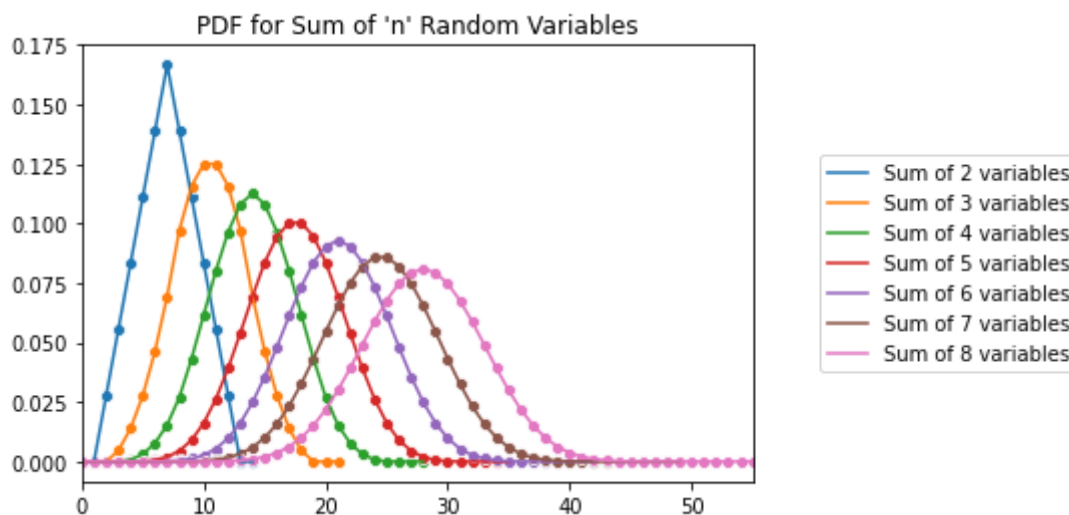
```
In [75]: data_sample = list([0,1/6,1/6,1/6,1/6,1/6,1/6,0])
```

```
In [76]: def n_convolve_sum(n, distribution):
             return np.array(reduce((lambda x, y: signal.convolve(x, y)), [data_s
         ample]*n))

         n =9
         for i in range(2,n):
             conv = n_convolve_sum(i, data_sample)
             sns.lineplot(data = conv ,label = "Sum of " + str(i) +" variables")
             sns.scatterplot(data = conv )
             plt.xlim(0, 6*(n)+1)
             plt.legend(loc='center right', bbox_to_anchor=(1.5, 0.5), ncol=1)
             plt.title("PDF for Sum of 'n' Random Variables")
```
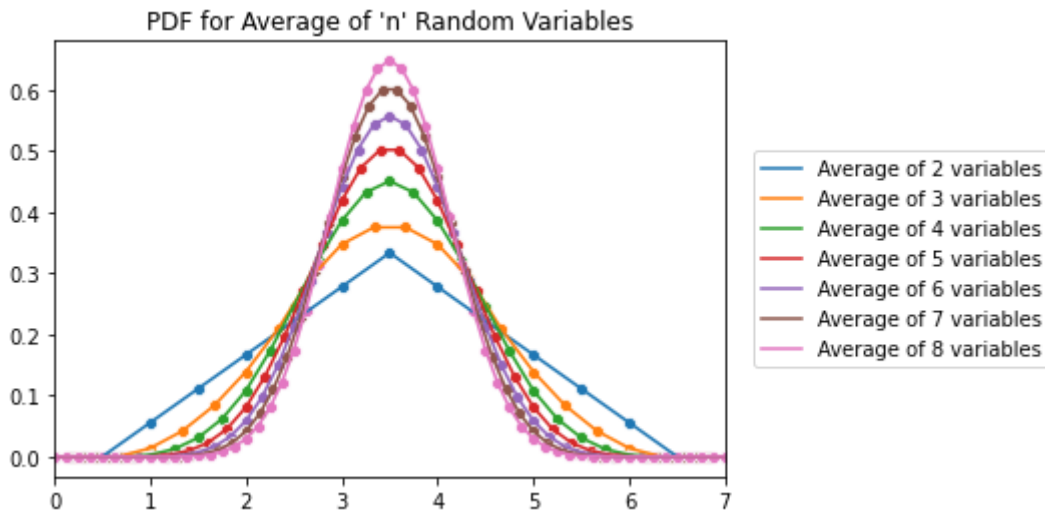
```
In [77]: def n_convolve_average(n, distribution):
             y_convolve = np.array(reduce((lambda x, y: signal.convolve(x, y)), [
         data_sample]*n))
             average_values = {key*(1/n):value/(1/n) for key,value in enumerate(y
         _convolve)}
             return list(average_values.keys()),list(average_values.values())


         n =9
         for i in range(2,n):
             average_values ,average_convolve = n_convolve_average(i, data_sample
         )
             sns.scatterplot(x = average_values, y = average_convolve)
             sns.lineplot(x = average_values, y = average_convolve,label = "Avera
         ge of " + str(i) +" variables")

             plt.xlim(0, 7)
             plt.legend(loc='center right', bbox_to_anchor=(1.5, 0.5), ncol=1)
             plt.title("PDF for Average of 'n' Random Variables")
```

PDF for Average of 'n' Random Variables

Legend:
- Average of 2 variables
- Average of 3 variables
- Average of 4 variables
- Average of 5 variables
- Average of 6 variables
- Average of 7 variables
- Average of 8 variables

## Inference:

It can be observed that as the number of random variables increases the PDF of average of random variables follow a Gaussian distribution irrespective of the initial distribution. Also the variance of the distribution decreases as the number of random variables increases. This inference aligns with the Central Limit Theorem as well.

PROBLEM-5

$$y = a_1 f_1(x) + a_2 f_2(x) + a_3 f_3(x)$$

$x, y \in$ scalars , Given: $(x_i, y_i)$ , $i \in \{1, 2, - - n\}$

(a) find parameter vector $(\hat{a}_1, \hat{a}_2, \hat{a}_3)$ that best fits.

In vector notation,

Let $f_1(x) = x_1$ for $x_i, i \in \{1 - n\}$

$$\hat{y} = \hat{a}_1 x_1 + \hat{a}_2 x_2 + \hat{a}_3 x_3$$

$f_2(x) = x_2$
$f_3(x) = x_3$

$$\Rightarrow \hat{y} = [x_1 \ x_2 \ x_3] \begin{bmatrix} \hat{a}_1 \\ \hat{a}_2 \\ \hat{a}_3 \end{bmatrix}$$

$$\Rightarrow \boxed{\hat{y} = XA}$$

The best $A$ would be the one that reduces the error i.e. difference between actual '$Y$' and predicted '$\hat{y}$'.

So we need to minimise $(Y - \hat{y})^2$.

Squared error $= (Y - XA)^T (Y - XA)$

$$E = Y^T Y - A^T X^T Y - Y^T XA + A^T X^T XA$$

To find $A$ that minimises error, we differentiate.

$$\frac{\partial E}{\partial A} = -X^T Y - Y^T X + 2X^T X A = 0$$

$$\Rightarrow -2X^T Y + 2X^T X A = 0$$

$$\Rightarrow X^T X A = X^T Y$$

$$\Rightarrow \boxed{A = (X^T X)^{-1} X^T Y}$$

# Problem-5

## PART - (b) , (c) , (d)

Picking $f_1$ , $f_2$, $f_3$ as following: $f_1(x) = x$ , $f_2(x) = x^3$ , $f_3(x) = e^{-x}$ .

Using $a_1 = 11.3$, $a_2 = 8.7$, $a_3 = 5.2$ randomly to create the training dataset.

## Data Simulation

Using small deltas from a normal distribution[N(0,1)] with a mean of 0 and std. deviation as 1, we add some random noise in the training data.

```python
In [224]: import matplotlib.pyplot as plt
          import seaborn as sns; sns.set()
          import numpy as np
          from numpy.linalg import inv
```

```python
In [225]: # Simulation to create training dataset of 200 rows
          np.random.seed(12)
          x = 10 * np.random.sample(200)
          deltax = np.random.normal(0,1,200)
          deltay = np.random.normal(0,1,200)
```

```python
In [226]: def functions(x):
              fx1 = x
              fx2 = x*x*x
              fx3 = np.exp(-x)
              return fx1,fx2,fx3
```

```python
In [227]: def get_training_data(x, deltay):
              fx1,fx2,fx3 = functions(x)
              train_y = 11.3*fx1 + 8.7*fx2 + 5.2*fx3 + deltay
              Train_X = np.matrix([[fx1[i],fx2[i],fx3[i]]for i in range(len(x))] )
              Train_Y = np.matrix([[train_y[i]] for i in range(len(train_y))])
              return Train_X, Train_Y
```

```python
In [228]: def calculate_coeff(X, Y):
              first = np.dot(X.T, X)
              second = np.dot(X.T, Y)
              return np.dot(inv(first), second)
```

Below in line 220, added the delta_x to the training data.

```
In [229]: Train_X, Train_Y = get_training_data(x+deltax, deltay)
          plt.scatter(x, np.array(Train_Y), label = "Training data")
          plt.title("Training data")
          plt.xlabel("x + deltax")
          plt.ylabel("y values")
```

Out[229]: Text(0, 0.5, 'y values')



## Predicting using the coefficients calculated from the training data

The coefficients are calculated using A = $(X^{T}X)^{-1}X^{T}Y$

```
In [230]: coeff_A = calculate_coeff(Train_X, Train_Y)
          print("Coefficients calculated are: ", coeff_A)
```

```
Coefficients calculated are:  [[11.28469767]
 [ 8.7003012 ]
 [ 5.16718308]]
```

```
In [231]: fx1,fx2,fx3 = functions(x)
          y_predicted = float(coeff_A[[0]])*fx1 + float(coeff_A[[1]])*fx2 + float(
          coeff_A[[2]])*fx3
          plt.scatter(x, np.array(Train_Y), label = "Training data")
          plt.scatter(x, y_predicted, color = "red", label = "Predicted data")
          plt.legend()
          plt.xlabel("x")
          plt.ylabel("y")
```

Out[231]: Text(0, 0.5, 'y')



**Inference:**

The coefficients calculated from the training data are very close to the actual coefficients chosen. The slight variation is due to the small randomness introduced in each $x_i$ and $y_i$ in the training dataset. The curve for the predicted values follow a generic trend of the training data. Hence we can say that the model proposed for calculation is good.