In [1]:	<pre># import libraries import numpy as np import pandas as pd  import matplotlib.pyplot as plt import seaborn as sns  import warnings warnings.filterwarnings('ignore')</pre>
In [2]•	Gradient Discent (in general)  It is an optimization algorithm to find the minimum of a function. We start with a random point on the function and move in the negative direction of the gradient of the function to reach the local/global minima.  Lets say the function is $y=(x-5)^2$
In [2]:	<pre>X = np.arange(-5,20) y = (X-5)**2 print(X) print(y)  [-5 -4 -3 -2 -1  0  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19] [100  81  64  49  36  25  16  9  4  1  0  1  4  9  16  25  36  49 64  81 100 121 144 169 196]</pre>
In [3]:	<pre># lets visualize the function plt.plot(X,y) plt.ylabel('f(x) = (x+5) 2') plt.xlabel('X') plt.show()</pre>
	125 - (S) 100 - (S) 75 - (S) 7
	Gradient discent formula $\theta^t = \theta^{t-1} - \eta \frac{\partial J}{\partial \theta^{t-1}}$
	where $\theta^t                                     $
In [4]:	J -> function which we have to minimise $\frac{\partial J}{\partial \theta^{t-1}} \text{ -> partial derivative of that function or called Gradient}$ # lets find the minimum of that function $y = (X-5)**2$ # lets start from $x = 15$
	<pre>X = 15 # starting point lr = 0.1 # learning rate  for i in range(100): # lets take 100 iteration     grad = 2*(X-5) # darivative of that function or gradient     X=X-lr*grad # applying formula print(X)  5.000000002037036</pre>
In [5]:	<ul> <li>conclusion</li> <li>So approximately 5 is our global minima for the function y = (X-5)**2.</li> <li>That means at the value of X = 5 our function in minimum.</li> </ul>
	<pre>for i in range(100):     grad = 2*(X-5)     X=X-lr*grad     y = (X-5)**2     plt.scatter(X,y) print(X)</pre> 5.000000002037036
	60 - 50 - 40 - 30 -
	Gradient Discent in Linear Regression
<pre>In [6]: Out[6]:</pre>	<pre># lets take a dataset called housing.csv housing = pd.read_csv(r'E:\linkdin post project\Gradient Discent\housing.csv') housing.head()</pre> <pre>price area bedrooms bathrooms stories mainroad guestroom basement hotwaterheating airconditioning parking prefarea full 1 12250000 8960 4 4 4 1 0 0 0 0 1 3 0</pre> <pre> 1 12250000 8960 4 4 4 1 0 0 0 0 1 3 0</pre>
In [7]:	<pre>2 12250000 9960</pre>
<pre>In [8]: Out[8]:</pre>	housing.drop(['furnishingstatus'],axis=1,inplace=True)
	0       4.562174       1.045766       1.402131       1.420507       1.376952       0.405251       -0.464888       -0.733865       -0.219063       1.471267       1.516299       1.8032         1       4.000809       1.755397       1.402131       5.400847       2.529700       0.405251       -0.464888       -0.733865       -0.219063       1.471267       2.676950       -0.5535         2       4.000809       2.216196       0.047235       1.420507       0.224204       0.405251       -0.464888       1.360148       -0.219063       -0.678439       1.516299       1.8032         3       3.982096       1.082630       1.402131       1.420507       0.224204       0.405251       -0.464888       1.360148       -0.219063       1.471267       2.676950       1.8032         4       3.551716       1.045766       1.402131       -0.569663       0.224204       0.405251       2.147110       1.360148       -0.219063       1.471267       1.516299       -0.5538
In [9]:	# Assign feature variable or independent variable X  X = housing['area']  # Assign response variable or dependent variable to y
In [10]:	<pre>y = housing['price']  # lets visualize relationship between the features and the response using scatterplots sns.scatterplot(X,y) plt.show()</pre>
	For linear regression we use a cost function known as the mean squared error or MSE.
	Formula $MSE = rac{1}{n} \sum_{i=1}^n \left(y^{(i)} - \hat{y} ight)^2$
	$MSE = rac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - h_{ heta} x^{(i)}  ight)^2 \ MSE = rac{1}{n} \sum_{i=1}^{n} \left( y^{(i)} - \left( m x^{(i)} + c  ight)  ight)^2$
	where ${\bf n}={\rm number\ of\ items}$ ${\bf y}^{(i)}={\rm actual\ point}$ $\hat{y}=h_{\theta}x^{(i)}=mx^{(i)}+c={\rm predicted\ point}$ ${\bf m}={\rm slope}$
	So our goal is to find the best m and c value where our MSE should be minimum.  As per Gradient Discent Formula we have to differentiate our cost function
	first and then put that in to the formula . Partial Darivation of MSE w.r.t slope(m) and intercept(c) $\frac{\partial MSE}{\partial m} = -\frac{2}{n} \sum_{i=1}^n \left(y^{(i)} - mx^{(i)} - c\right) \left(x^{(i)}\right)$
	$rac{\partial m}{\partial c} = -rac{2}{n} \sum_{i=1}^n \left( y^{(i)} - m x^{(i)} - c  ight)$
In [11]: In [12]:	<pre># now for applying gradient discent we need out x,y variables as numby array X = np.array(X) y = np.array(y)  # lets implement the gradient descent function # lets initialised current m and c to 0</pre>
	<pre>def gradient_discent_simple(X,y):     m = 0     c = 0     n = float(len(y))     iters = 1000 # take 1000 iteration     learning_rate = 0.01     df = pd.DataFrame(columns = ['m', 'c', 'cost']) # make a dataframe to keep track how the costs are minimisin     for i in range(iters):         y_pred = m*X+c # prediction values         cost = sum([i**2 for i in (y-y_pred)])/n         d m = ( 2.75) taum(X**(y, y, pred)) # darminusting to y, pred = m*X + c # prediction values</pre>
	<pre>d_m = (-2/n)*sum(X*(y-y_pred)) #darivative w.r.t m d_c = (-2/n)*sum(y-y_pred) #darivative w.r.t c     # we are doing the derivation for minimising the cost     # after darivation we have to update the m and c     m = m - (learning_rate*d_m)     c = c - (learning_rate*d_c)     df.loc[i] = [m,c,cost] # keep track of each cost in each iteration  df.reset_index().plot.line(x='index', y=['cost']) return f'final slope and intercept after 1000 iteration is {round(m,3),round(c,3)}'</pre>
In [13]: Out[13]:	gradient_discent_simple(X,y)  'final slope and intercept after 1000 iteration is (0.536, 0.0)'  100 - cost 0.95 - cost
	0.80 - 0.80 - 0.75 - 0.70 - 0 200 400 600 800 1000
In [14]:	• We can clearly see that after 200 iterations the cost was not decreasing much which means we get our global minima.  # lets visualize how the line was fitted y_pred = 0.536*X+0.0 plt.scatter(X, y) plt.plot(X, y_pred, color = 'red')
	<pre>plt.xlabel('area') plt.ylabel('price') plt.show()</pre>
	Note  • learning rate is the speed at which we want to move towards negetive of the gradient .  • it's always a good practice to choose a small value of learning rate and slowly move towards the nagative of the gradient .  Multiple Linear Regression
In [15]:	Applying Gradient Descent for Multiple (>1) Features  # Assigning feature variable X # lets now take 2 features area and bedrooms X = housing[['area', 'bedrooms']]  # Assigning response variable y y = housing['price']
	Note • when we have more then one features the model now fit a hyperplane instade of line .
	• The cost function here is slightly different from MSE we just dividing the MSE by half (1/2) to get a nice interpretation . ### $J(\theta_0,\theta_1)=\frac{1}{2n}\sum_{i=1}^n\left(y^{(i)}-h_\theta\big(x^{(i)}\big)\right)^2$ ###
	$J(\theta_0,\theta_1,\dots,\theta_i) = \frac{1}{2n} \sum_{i=1}^n \left(y^{(i)} - \left(\theta_0 + \theta_1 x_1 + \theta_2 x_2 \dots + \theta_i x_i\right)\right)^2$ ### we have to first differentiate w.r.t. all thetas and then update each thetas . ### $\frac{\partial J(\theta_0,\theta_1,\dots,\theta_i)}{\partial \theta_0} = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - h_\theta(x^{(i)})\right) \left(x_0^{(i)}\right)$
	#### $rac{\partial J( heta_0, heta_1,\ldots, heta_i)}{\partial  heta_1}=rac{1}{n}\sum_{i=1}^n\left(y^{(i)}-h_ hetaig(x^{(i)}ig)ig)ig(x_1^{(i)}ig)  ight.$ $\dots$
	$\frac{\partial J(\theta_0,\theta_1,\dots,\theta_i)}{\partial \theta_i} = \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - h_\theta\big(x^{(i)}\big)\right) \big(x_i^{(i)}\big)$ • We need to minimise the cost function J( $\theta$ ) One way to do this is to use the batch gradient decent algorithm. In batch gradient decent, the values are updated in each iteration: #### $\theta_0 = \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left(y^{(i)} - h_\theta\big(x^{(i)}\big)\right) \big(x_0^{(i)}\big)$
	#### $ heta_1 =  heta_1 - lpha rac{1}{n} \sum_{i=1}^n ig(y^{(i)} - h_ hetaig(x^{(i)}ig)ig)ig(x_1^{(i)}ig)$ ####
	$ heta_2= heta_2-lpharac{1}{n}\sum_{i=1}^nig(y^{(i)}-h_ hetaig(x^{(i)}ig)ig)ig(x_2^{(i)}ig) \ \ldots \  heta_n= heta_n-lpharac{1}{n}\sum_{i=1}^nig(y^{(i)}-h_ hetaig(x^{(i)}ig)ig)ig(x_n^{(i)}ig)$
In [16]:	• Now here we are taking the help of matrix multiplication because we are deal with more than one features. # Add a columns of 1s as an intercept to X. # The intercept column is needed for convenient matrix representation of cost function X['intercept'] = 1 X = X.reindex(['intercept', 'area', 'bedrooms'], axis=1) X.head()
Out[16]:	intercept         area         bedrooms           0         1         1.045766         1.402131           1         1         1.755397         1.402131           2         1         2.216196         0.047235           3         1         1.082630         1.402131
In [17]:	<pre># Convert X and y to arrays import numpy as np X = np.array(X) y = np.array(y)</pre>
In [18]: In [19]:	theta = np.matrix(np.array([0,0,0])) alpha = 0.01 iterations = 1000
In [20]:	<pre>def compute_cost(X, y, theta):     return np.sum(np.square(np.matmul(X, theta) - y)) / (2 * len(y))  # gradient descent # takes in current X, y, learning rate alpha, num_iters # returns cost (notice it uses the cost function defined above)  def gradient_descent_multi(X, y, theta, alpha, iterations):</pre>
	<pre>theta = np.zeros(X.shape[1]) n = len(X) df = pd.DataFrame( columns = ['coefficients', 'cost'])  for i in range(iterations):     cost = compute_cost(X, y, theta)     derivative = (1/n) * np.matmul(X.T, np.matmul(X, theta) - y) # we are doing derivative for minimizing theta = theta - alpha * derivative # here we have to update our theta means we are going to next step df.loc[i] = [theta,cost]</pre>
In [21]: Out[21]:	gradient_descent_multi(X, y, theta, alpha, iterations)
	997 [3.337635977993293e-16, 0.4916563172711537, 0 0.314176 998 [3.3429324548080186e-16, 0.4916565358164564, 0 0.314176 999 [3.3447658506285006e-16, 0.49165675237880896, 0.314176, 'final coefficients are [3.34476585e-16 4.91656752e-01 2.91844700e-01] ')  0.500
	0.400 - 0.375 - 0.350 - 0.325 - 0.325 - 0.000   0.000
In [22]:	• We can clearly see that after 200 iterations the cost was not decreasing much which means we get our global minima.  EXAMPLE 2  # lets take a dataset called advartising.csv
Out[22]:	<pre>advartising = pd.read_csv(r'E:\\advertising (1).csv') advartising.head()</pre>
In [23]:	TV         Radio         Newspaper         Sales           0         230.1         37.8         69.2         22.1           1         44.5         39.3         45.1         10.4           2         17.2         45.9         69.3         9.3
	<b>0</b> 230.1 37.8 69.2 22.1
Out[23]:	0 230.1 37.8 69.2 22.1 1 44.5 39.3 45.1 10.4 2 17.2 45.9 69.3 9.3 3 151.5 41.3 58.5 18.5 4 180.8 10.8 58.4 12.9  # Normalisisng the data advartising = (advartising - advartising.head()  TV Radio Newspaper Sales 0 0.967425 0.979066 1.774493 1.548168 1 -1.194379 1.080097 0.667903 -0.694304 2 -1.512360 1.524637 1.779084 -0.905135 3 0.051919 1.214806 1.283185 0.858177
Out[23]:	0 230.1 37.8 69.2 22.1 1 44.5 39.3 45.1 10.4 2 17.2 45.9 69.3 9.3 3 151.5 41.3 58.5 18.5 4 180.8 10.8 58.4 12.9  # Normalisising the data advartising = (advartising - advartising.mean())/advartising.std() advartising.head()  TV Radio Newspaper Sales 0 0.967425 0.979066 1.774493 1.548168 1 -1.194379 1.080097 0.667903 -0.694304 2 -1.512360 1.524637 1.779084 -0.905135 3 0.051919 1.214806 1.283185 0.858177 4 0.393196 -0.839507 1.278593 -0.215143
	0 230.1 37.8 69.2 22.1  1 44.5 39.3 45.1 10.4  2 17.2 45.9 69.3 9.3  3 151.5 41.3 58.5 18.5  4 180.8 10.8 58.4 12.9  # Normalisisng the data advartising = (advartising - advartising.mean())/advartising.std() advartising.head()  TV Radio Newspaper Sales  0 0.967425 0.979066 1.774493 1.548168  1 -1.194379 1.080097 0.667903 -0.694304  2 -1.512360 1.524637 1.779084 -0.905135  3 0.051919 1.214806 1.283185 0.858177  4 0.393196 -0.839507 1.278593 -0.215143  # Assigning feature variable X # lets see how attribute 'TV' affects 'Sales' X = advartising['TV']  # Assigning response variable y
In [24]:	0 2301 37.8 692 221 1 445 393 451 104 2 172 459 693 93 3 1515 413 585 185 4 1808 10.8 584 129  # Normalisisng the data advartising = (advartising - advartising.mean())/advartising.std() advartising.head()  TV Radio Newspaper Sales  0 0.967425 0.979066 1.774493 1548168 1 -1.194379 1.080097 0.667903 -0.694304 2 -1.512360 1.524637 1.779084 -0.905135 3 0.051919 1.214806 1.283185 0.858177 4 0.393196 -0.839507 1.278593 -0.215143  # Assigning feature variable X # lates see how attribute 'TV' affects 'Sales' X = advartising('TV') # Assigning response variable y y = advartising('TV')  # Assigning response variable y y = advartising('Sales') plt.scatter(X,y)  Cmatplotlib.collections.FathCollection at 0x2715a6628e0>
In [24]:	0 230.1 37.8 69.2 22.1 1 44.5 39.3 45.1 10.4 2 17.2 45.9 69.3 9.3 3 151.5 41.3 58.5 18.5 4 180.8 10.8 58.4 12.9  # Mormalisiang the data advartising = (advartising - advartising .mean())/advartising.std() advartising   (advartising - advartising - advartising - advartising .std() advartising .std() advartising 1.080097 0.667903 0.694304 2 -1.512260 1.524637 1.779084 0.905135 3 0.051919 1.214806 1.283185 0.858177 4 0.393196 -0.839507 1.278593 -0.215143  # Assigning feature variable X # lets see how attribute 'TV' affects 'Sales' X = advartising('Sales') plt.scatter(X, y)  Cmatpiolib.collections.FathCollection at 0x2715a6628e0>
In [24]: Out[24]:	0 2301 378 692 221  1 445 393 451 104  2 172 459 693 93  3 1515 413 585 185  4 1808 108 584 129  **Normalising the data advantising = (advantising - advantising, mean())/advantising.std()  **TV Radio Newspaper Sales  0 0967425 097966 1,774493 1.548168  1 -1194379 1.888097 0.667903 -0.694904  2 -1.512360 1524637 1,779684 -0.905135  3 0051919 1.214806 1283185 0.858177  4 0.393196 -0.839507 1,278593 -0.215143  **Assigning feature variable X** slots see how actribute 'TV' affects 'Sales' X** advantising('TV')  **Assigning response variable Y** a advantising('TV')  **Assigning response variable Y** y = advantising('TV')  **Assigning response variable Y** a advantising
In [24]: Out[24]:	0 2001 378
In [24]: Out[24]:	0 220.1 37.8 692 221 1 445 383 451 104 2 172 459 693 93 93 3 1515 413 585 185 4 1938 108 584 129  ***Moccalitation** Use data
In [24]: Out[24]:	0 2261 378 682 221 1 443 383 451 00 2 172 453 683 93 3 1515 413 585 785 4 1828 188 84 129  * Moraticing conditions controlling mean())/advan(ising.std) advantising conditions * Moraticing conditions
In [24]: Out[24]:	0 2031 578 692 201 1 449 903 63 63 63 63 63 63 63 63 63 63 63 63 63
In [24]: Out[24]:	### 1928 ### 622 ### 1924 ### 1925 ###
In [24]:  Out[24]:  In [26]:	0 2003
In [24]:  Out[24]:  In [26]:	2011   2012   2013   2013   2014
In [24]: Out[24]: Out[25]: Out[28]:	### 2015
In [24]:  Out[24]:  Out[25]:  In [26]:	
In [24]:  In [25]:  In [26]:  In [28]:  In [30]:	Security
In [24]:  In [25]:  In [26]:  In [28]:  In [30]:	0 301