

* NumPy, Matplotlib, Pandas & other necessary libraries.

NumPy

import numpy as np

Creates an array of one with defined shape

Ones = np.ones(shape, dtype=int);

** np.hstack((array_weWant_to_stack, target_array));

↳ hstack takes a tuple of arrays as arguments)

alternatives:

① np.c_[ones, original_array]; → Does column wise concatenation

② np.insert(original_array, index, value, axis=1); → More control.

Pandas

import pandas as pd;

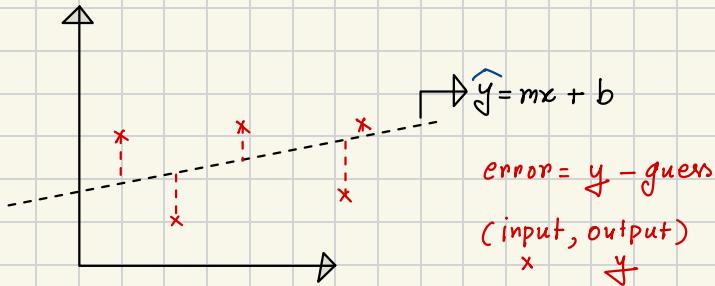
df = pd.read_csv("file-name", header=None, sep=",");

↳ Loads without headers

columns = ['a', 'b', 'c', 'd', 'e', 'f']

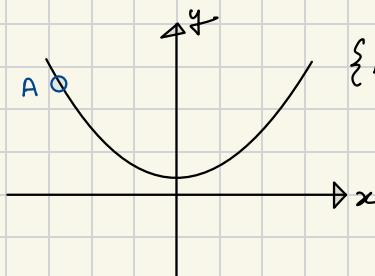
df.columns = columns → Inserts columns

Linear Regression with GD



Cost: $\sum_{i=0}^n (\hat{y}_i - y_i)^2 \rightarrow$ we want to minimise this function.

\therefore Cost = $\sum (\hat{y}_i - y_i)^2$ this is kind of like $y = x^2$



{ Now, we need to find the value of x for that y is lowest.

Suppose, we're at point A & if we want to reduce the function, we need to find its gradient. But how do we go down?

→ We calculate the slope.

$$m = m + \Delta m$$

$$b = b + \Delta b$$

$$\therefore J(m, b) = \frac{1}{n} \sum_{i=0}^n (mx + b - y)^2$$

$$\frac{\delta}{\delta m} J = 2(mx + b - y)m \quad \& \quad \frac{\delta}{\delta b} J = 2(mx + b - y)(1)$$

* $\frac{1}{n}$ in $J(m, b) = \frac{1}{n} \sum (\hat{y}_i - y_i)^2$?

→ We want to take the average error & minimise it.

→ Makes it scale irrelevant so that same learning rate works in different datasets.

We also sometimes take $\frac{2}{n}$, this is purely because of algebraic convenience.

$$\text{Update rule: } m^{(t+1)} = m - \alpha \frac{\partial}{\partial m} J$$

$$b^{(t+1)} = b - \alpha \frac{\partial}{\partial b} J$$

* Example:

$(x, y) = \{(1, 2), (2, 3), (3, 5)\}$ predict y when $x=7$

Here, $n=3$, model, $\hat{y} = mx + b$ let, $\hat{y}_1 = \hat{y}_2 = \hat{y}_3 = 0$ {because, we assumed $m=0$ & $b=0$ }

$$\begin{aligned} J &= \frac{1}{n} \sum_0^3 (\hat{y}_i - y_i)^2 = \{(0-2)^2 + (0-3)^2 + (0-5)^2\} / 3 \\ &= \frac{4+9+25}{3} = \frac{38}{3} \approx 12.67 \end{aligned}$$

let, $\alpha = 0.01$

Now, we'll minimise loss by updating m & b

$$\text{Given, } J = \frac{1}{n} \sum_0^n (mx_i + b - y_i)^2$$

$$\therefore \frac{\delta J}{\delta m} = \frac{1}{n} \sum (mx_i + b - y_i) \times 2x_i$$

$$\& \frac{\delta J}{\delta b} = \frac{1}{n} \sum (mx_i + b - y_i)$$

$$m^{(t+1)} = m^t + \alpha \frac{\delta}{\delta m} J$$

$$b^{(t+1)} = b^t + \alpha \frac{\delta}{\delta b} J$$

1st step,

$$\begin{aligned} \frac{\delta J}{\delta m} &= \frac{1}{3} \{(0 \times 1 + 0 - 2) \times 2(1) + (0 \times 2 + 0 - 3)(2 \times 2) \\ &\quad + (-5)(2 \times 3)\} = -15.33 \end{aligned}$$

$$\therefore m^1 = m^0 - \alpha \left(\frac{\delta J}{\delta m} \right) = 0 - 0.01(-15.33) = 0.1533$$

let's get b^1

$$\frac{\delta J}{\delta b} = \frac{2}{n} \sum_0^n (mx_i + b - y_i) = \frac{2}{3} (-2 - 3 - 5) = -6.667$$

$$\therefore b^1 = b^0 - \alpha \left(\frac{\delta J}{\delta b} \right) = 0 - 0.0667 = 0.0667$$

So, model updates to,

$$\hat{y}_i = 0.1533x_i + 0.0667$$

$$\therefore \hat{y}_1 = 0.22 \quad \hat{y}_2 = 0.3733 \quad \hat{y}_3 = 0.5266$$

Current error: $J = \frac{1}{n} \sum_0^n (\hat{y}_i - y_i)^2$

$$= \frac{1}{3} ((0.22 - 2)^2 + (0.3733 - 5)^2 + (0.5266 - 5)^2)$$
$$= \{(-1.78)^2 + (-4.6267)^2 + (-4.4734)^2\}/3$$
$$= 10.022067$$

So, we've reduced error in step 1 but we have two more steps to go. It took a lot of manual calculations, let's simplify a bit using matrices.

* Dataset $\{(1,2), (2,3), (3,5)\}$ & prediction $\hat{y}_i = mx_i + b$

$$\theta = \begin{pmatrix} m \\ b \end{pmatrix} \text{ so, } \hat{y} = \bar{x}\bar{\theta}$$

* Design matrix:

$$\bar{x} = \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \quad \bar{\theta} = \begin{pmatrix} m \\ b \end{pmatrix} \quad \hat{y} = \bar{x}\bar{\theta} = \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \begin{pmatrix} m \\ b \end{pmatrix}$$
$$= \begin{pmatrix} m + b \\ 2m + b \\ 3m + b \end{pmatrix}$$

as $m, \theta = 0, 0$

$n = 3$

$$\therefore \bar{y} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \quad \text{Given, } \bar{y} = \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix}$$

$$\bar{e} = \bar{y} - \bar{y} = \begin{pmatrix} -2 \\ -3 \\ -5 \end{pmatrix} \quad \bar{e}^2 = (\bar{y} - \bar{y})^2 = \|\bar{y} - \bar{y}\|_2 \\ = (\bar{y} - \bar{y})^T (\bar{y} - \bar{y})$$

* Cost Function,

$$J = \frac{\bar{e}^2}{n} = 38/3 = 12.667 \quad = (2 \ 3 \ 5) \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix} \\ = (4 + 9 + 25) = 38$$

* Gradient,

$$J(\theta) = \frac{1}{n} (\bar{y} - \bar{y})^T (\bar{y} - \bar{y}) = \frac{1}{n} (\bar{x}\theta - \bar{y})^T (\bar{x}\theta - \bar{y}) \\ = \frac{1}{n} (\bar{x}^T \theta^T - \bar{y}^T) (\bar{x}\theta - \bar{y}) \\ = \frac{1}{n} (\theta^T \bar{x}^T \bar{x}\theta - \bar{x}^T \bar{\theta}^T \bar{y} \\ - \bar{y}^T \bar{x}\theta - \bar{y}^T \bar{y}) \\ = \frac{1}{n} (\theta^T \bar{x}^T \bar{x}\theta - 2\bar{x}^T \bar{\theta}^T \bar{y} - \bar{y}^T \bar{y})$$

* Matrix Calculus Identities

- ① $\nabla_{\theta} (\theta^T A \theta) = 2A\theta$ { When A is symmetric & if $A = x^T x$ then A is always symmetric }
- ② $\nabla_{\theta} (C^T \theta) = C$
- ③ $\nabla_{\theta} (\text{constant}) = 0$

→ Appendix 1

* Applying these in the cost function,

$$\textcircled{1} \quad \nabla_{\theta} (\bar{x}^T \bar{x}\theta) = 2\bar{x}^T \bar{x}\theta$$

$$\textcircled{2} \quad \nabla_{\theta} (2\bar{y}^T \bar{x}\theta) = 2\bar{x}^T \bar{y}$$

$$\textcircled{3} \quad \nabla_{\theta} (\bar{y}^T \bar{y}) = 0$$

$$\therefore \nabla_{\theta} J = (2\bar{x}^T \bar{x}\theta - 2\bar{x}^T \bar{y})/n$$

\downarrow
Gradient
(w.r.t
m & b)

$$= \frac{2\bar{x}^T}{n} \underbrace{(\bar{x}\theta - \bar{y})}_{\text{Error Vector}}$$

Update rule : $\theta := \theta - \alpha \nabla_{\theta} J$

Let's verify !

$$\{(1, 2), (2, 3), (3, 5)\}$$

$$\begin{aligned}\overline{x} &= \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} & \overline{\theta} &= \begin{pmatrix} m \\ b \end{pmatrix} & \overline{y} &= \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix} & \widehat{y} &= \overline{x} \overline{\theta} & \overline{e} &= \widehat{y} - \overline{y} \\ A & & B & & C & & & & \\ & & & = \begin{pmatrix} 0 \\ 0 \end{pmatrix} & & & & & \\ \overline{e} &= \begin{pmatrix} -2 \\ -3 \\ -5 \end{pmatrix} & & & & & & & \\ & & & & & & & & \\ & & & & & & & & \end{aligned}$$

* Cost calc

$$\begin{aligned}J(\theta) &= \frac{1}{n} (\overline{x}\theta - y)^T (\overline{x}\theta - y) \\ &= \frac{1}{3} \begin{pmatrix} -2 & -3 & -5 \end{pmatrix} \begin{pmatrix} -2 \\ -3 \\ -5 \end{pmatrix} = \frac{1}{3} (2^2 + 3^2 + 5^2) = 38/3 = 12.67\end{aligned}$$

$$* \nabla_{\theta} J(\theta) = \frac{2}{n} \overline{x}^T (\overline{x}\theta - y)$$

$$\begin{aligned}&= \frac{2}{3} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix} \left(\begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} - \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix} \right) \\ &= \frac{2}{3} \begin{pmatrix} -2 \\ -3 \\ -10 \end{pmatrix} = \begin{pmatrix} -46/3 \\ -20/3 \end{pmatrix} = \begin{pmatrix} -15.33 \\ -6.667 \end{pmatrix}\end{aligned}$$

$$\therefore \nabla_{\theta} J(\theta) = \begin{pmatrix} -15.33 \\ -6.667 \end{pmatrix}$$

$$\begin{aligned}\therefore \theta^1 &= \theta - \alpha \begin{pmatrix} -15.33 \\ -6.667 \end{pmatrix} \quad \{ \alpha = 0.01 \} \\ &= \begin{pmatrix} 0.1533 \\ 0.0667 \end{pmatrix}\end{aligned}$$

$$\hat{y}_1 = \bar{x} \bar{\theta}$$

$$= \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \begin{pmatrix} 0.1533 \\ 0.0667 \end{pmatrix} = \begin{pmatrix} 0.22 \\ 0.3733 \\ 0.5266 \end{pmatrix}$$

$$e^1 = \begin{pmatrix} -1.78 \\ -2.626 \\ -4.9734 \end{pmatrix}$$

$$\theta^{(2)} = \theta^{(1)} - \alpha \left(\nabla_{\theta} J(\theta) \right) = \theta^{(1)} - \alpha \left(\frac{2}{n} x^T (x \theta^{(1)} - y) \right)$$

$$= \theta^{(1)} - 0.01 \left(\frac{2}{3} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix} \left(\begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \begin{pmatrix} 0.1533 \\ 0.0667 \end{pmatrix} - \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix} \right) \right)$$

$$= \theta^{(1)} - 0.01 \left(\frac{2}{3} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} -1.78 \\ -2.626 \\ -4.9734 \end{pmatrix} \right) e^1$$

$$\therefore \theta^{(2)} = \begin{pmatrix} 0.289698 \\ 0.125896 \end{pmatrix}$$

$$e = x\theta - y$$

$$\theta^3 = \theta^{(2)} - 0.01 \times \frac{2}{3} \begin{pmatrix} 1 & 2 & 3 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} -1.584 \\ -2.274 \\ -9.005 \end{pmatrix}$$

$$= \begin{pmatrix} 0.289698 \\ 0.125896 \end{pmatrix} - \begin{pmatrix} -0.121263 \\ -0.052563 \end{pmatrix} = \begin{pmatrix} 0.410912 \\ 0.1734 \end{pmatrix}$$

$$J(\theta)^{(3)} = |(\bar{x}\theta^{(3)} - y)|_2 \times \frac{1}{n} \quad (n=3)$$

$$= 6.29$$

*How long will we go?

→ Our error term is reducing, so we can go as low as we want till error=0

But, we can assign few other stopping conditions

- (1) if $J(\theta) \approx 0$
- (2) iteration count limit
- (3) $\|e\|_2 \approx 0$
- (4) $\|J(\theta)^{(t+1)} - J(\theta)^{(t)}\| < \epsilon$ or ≈ 0

But no matter what we do, we need computational power because datasets can become very large!

→ We'll use python to create our linear regress with GD algorithm, calculate accuracy.

* Appendix 1

(1) $\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix} \in \mathbb{R}^{d \times 1}$ $\nabla_{\theta} f(\theta) = \begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \\ \vdots \\ \frac{\partial f}{\partial \theta_d} \end{pmatrix}$

* $\frac{d}{d(\text{vector})} (\text{vector}) = (\text{vector})$

(2) $\nabla_{\theta} (c^T \theta) = c$

$$c^T \theta = \sum_0^d c_i \theta_i \quad \frac{\partial}{\partial \theta_j} (c^T \theta) = c_j$$

Stacking all in C

(3) $\nabla_{\theta} (k) = 0$

Constant rule

(4) $\nabla_{\theta} \theta^T A \theta = 2A\theta$ if A is symmetric.

Linear Regression = Quadratic Optimization

Convex = Single Global Minima

Normal Equation (for any number of features)

$\theta = (x^T x)^{-1} x^T y$

$$m = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Stochastic Gradient Descent

- Regular GD uses all datapoints to calculate the gradient but this becomes costly over time as n increases.
- SGD uses one datapoint at a time.

*Steps:

① Define Per-Sample Loss:

$$l_i(\theta) = (x_i^T \theta - y_i)^2$$

$$J\theta = \frac{1}{n} \sum_{i=1}^n l_i(\theta)$$

② Gradient of One Sample (x_i, y_i)

$$\nabla_{\theta} l_i(\theta) = 2(x_i^T \theta - y_i)(x_i)$$

③ Update rule

$$\theta \leftarrow \theta - \underbrace{\alpha \cdot 2x_i (x_i^T \theta - y_i)}_{\nabla_{\theta} l_i(\theta)}$$

Example:

$$\hat{y} = mx + b \quad x_i = [x_i, 1] \quad \theta = \begin{bmatrix} m \\ b \end{bmatrix}$$

i	x	y
1	2	2
2	2	3
3	3	5

$$\theta^{(0)} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\alpha = 0.1$$

One random point $(x_2, y_2) = (2, 3)$

Design Matrix $x_2 = \begin{pmatrix} 2 & 1 \end{pmatrix}$

$$\hat{y}_2 = x_2 \theta = \begin{pmatrix} 2 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \end{pmatrix} = 0 + 0 = 0$$

$$e_2 = 0 - 3 = -3 = x_2 \theta - y_2$$

$$G = \text{Gradient} = \nabla_{\theta} l_2(\theta) = 2x_2^T e_2$$

$$= 2 \begin{pmatrix} 2 \\ 1 \end{pmatrix} (-3) = \begin{pmatrix} -12 \\ -6 \end{pmatrix}$$

$$\text{Update: } \theta^{(1)} = \theta - \alpha G = \begin{pmatrix} 0 \\ 0 \end{pmatrix} - (0.1) \begin{pmatrix} -12 \\ -6 \end{pmatrix} = \begin{pmatrix} 1.2 \\ 0.6 \end{pmatrix}$$

*2nd Generation $\{\text{Random point } (x_1, y_1) = (1, 1) \text{ Design Matrix } \begin{pmatrix} 1 & 1 \end{pmatrix}\}$

$$x_1 = (1 \ 1) \quad y_1 = 2$$

$$y_1 = (1 \ 1) \begin{pmatrix} 1.2 \\ 0.6 \end{pmatrix} = 1 \cdot 2 + 0 \cdot 6 = 1.8$$

$$e = 1.8 - 2 = -0.2$$

$$G = 2x_1^T e = 2 \begin{pmatrix} 1 \\ 1 \end{pmatrix} (-0.2) = \begin{pmatrix} -0.4 \\ -0.4 \end{pmatrix}$$

$$\theta^{(2)} = \theta^{(1)} - \alpha G$$

$$= \begin{pmatrix} 1.2 \\ 0.6 \end{pmatrix} - 0.1 \begin{pmatrix} -0.4 \\ -0.4 \end{pmatrix} = \begin{pmatrix} 1.24 \\ 0.69 \end{pmatrix}$$

* Although SGD never truly settles

$\alpha_t = \frac{\alpha_0}{1+t}$ so we need to sometimes implement learning rate decay.

Takeaways

- ① SGD uses one sample.
- ② Each update is cheap but noisy.
- ③ Over many updates noise averages out.
- ④ Learning rate must be smaller than VGD.

Mini Batch Gradient Descent

→ Similar concept as VGD but this time we work with a sub-sample.

Algorithm:

- ① Select sample size.
- ② Select sub-dataset randomly.
- ③ Find gradient of that sub-dataset.
- ④ Shuffle dataset.
- ⑤ Found desired expected loss? (NO). Yes! Stop.
No! → Goto step 2.

Adaptive Gradient Descent AdaGrad

→ Look we're basically calculating how fast objects move in which direction. We are talking about shape of the loss function. Hessian simply describes that shape.

→ Gradient tells about the direction of the steepest descent, but tells nothing about 1) how curved the surface is.
2) whether the direction is safe to move.

Two very different surfaces can have the same gradient at a point.

Appendix 2

Suppose a 1D function $f(x) = \alpha x^2$
 $\alpha = 100 \rightarrow$ very steep.
 $\alpha = 0.01 \rightarrow$ very flat.

$f'(x) = 2\alpha x \rightarrow$ This is the gradient

But what determines how fast the gradient changes?

$f''(x) = 2\alpha \rightarrow$ Thin!

Hessians are basically second derivatives arranged into a matrix.

In gradient descent

$$\theta^{t+1} = \theta^t - \alpha \nabla J(\theta^t)$$

AdaGrad basically gives each parameter its own learning rate based on historical gradients.

① $g_t = \nabla_{\theta} J(\theta^t) \in \mathbb{R}^d \rightarrow$ Gradient of loss function

② Accumulating squared gradients

$$G_t = \sum_{i=1}^t g_i \odot g_i = \bar{g}_t^T \cdot \bar{g}_t \quad \odot \rightarrow \text{Piecewise multiplication}$$
$$G_t \in \mathbb{R}^d$$



↓
Old

③ Update rule:

$$\theta_{t+1} = \theta_t - \eta \frac{g_t}{\sqrt{G_t} + \epsilon} ; \epsilon \text{ is a small constant (so we don't get } \infty \text{ error)}$$

Division is element-wise.

Why squared gradients?

- i) Makes all contributions positive.
- ii) Penalizes large gradients disproportionately.

Because

Large gradients \rightarrow Dominate the sum \rightarrow shrink step size fast
vice-versa.

$H \approx \text{diag}(\sum g^2)$ \rightarrow AdaGrad estimates Hessians basically because
calculating Hessians are 1) Expensive
2) Hard 3) Unstable

Weakness

if $G_t \uparrow$ monotonically then $\frac{1}{\sqrt{G_t}} \downarrow$ & eventually learning rate becomes 0.

Example: (1,2), (2,3), (3,5). Simulate LR with AdaGrad.

$$\rightarrow \text{model: } y = \theta_1 x + \theta_0 \quad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \end{pmatrix}$$

$$\begin{aligned} J(\theta) &= \frac{1}{2m} \sum_{i=0}^m (\hat{y}^{(i)} - y^{(i)})^2 \\ &= \frac{1}{2m} (\hat{y} - y)^T (\hat{y} - y) \\ &= \frac{1}{2m} (x\theta - y)^T (x\theta - y) \\ &= \frac{1}{2m} (x^T \theta^T - y^T)(x\theta - y) \\ &= \frac{1}{2m} (\theta^T x^T x \theta - 2y^T x \theta + y^T y) \end{aligned}$$

$$\text{Design matrix: } \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \quad x$$

$$\begin{aligned} \theta &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \hat{y} &= \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \\ e &= \hat{y} - y \\ &= x\theta \end{aligned}$$

$$\nabla_{\theta} J(\theta) = \frac{1}{2m} (2x^T x \theta - 2x^T y) = \frac{1}{m} (x^T x \theta - x^T y)$$

$$G_t = g_u^T \cdot g_u$$

$$\theta_{t+1} = \theta_t - \eta \frac{g_t}{\sqrt{G_t} + \epsilon}$$

RMSProp

Why Adagrad fails?

$$\textcircled{1} \quad g_t = \sum_{k=1}^t g_k \odot g_k \quad \& \quad \theta_t = \theta_0 - \eta \frac{g_t}{\sqrt{g_t} + \epsilon}$$

g_t increases monotonically

We don't need information about all past gradients rather we need information only about recent curvature.

$s_t \rightarrow$ Running estimate of curvature magnitude per parameter

$s_{t,j} = \mathbb{E}[g_{t,j}^2]$ "How large are the gradients for parameter θ_j , on average, recently?"

$$\text{RMS}(g_j) = \sqrt{s_{t,j}}$$

$$g_t = \nabla_{\theta} J(\theta_t) = \frac{1}{m} x^T (x \theta_t - y)$$

$$s_t = \rho s_{t-1} + (1-\rho) g_t \odot g_t \quad \rho \in (0, 1)$$

$$\theta_{t+1} = \theta_t - \eta \frac{g_t}{\sqrt{s_t} + \epsilon} \quad \theta_t \in \mathbb{R}$$

$$J(\theta) \quad \longrightarrow \quad \theta_{t+1} = \theta_t - \eta (\text{diag}(\sqrt{s_t})^{-1}) g_t$$

* RMSProp performs adaptive diagonal second order scaling.

* Numerical Example

$$(1, 2), (2, 3), (3, 5)$$

* initialisation

$$x = \begin{pmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{pmatrix} \quad y = \begin{pmatrix} 2 \\ 3 \\ 5 \end{pmatrix}$$

$$\begin{aligned} \eta &= 0.1 & \theta_0 &= 0 \\ \rho &= 0.9 & \epsilon &\approx 0 \end{aligned}$$

$$J(\theta) = \frac{1}{n} (x\theta - y)^T (x\theta - y)$$

$$= \frac{1}{n} (\theta^T x^T x \theta - x^T \theta^T y - x \theta y^T + y^T y)$$

$$g_t = \nabla_{\theta} J(\theta) = (2x^T x \theta - 2x^T y - 0)/n$$

$$= \frac{2}{n} \underbrace{x^T (x\theta - y)}_{\text{Error}}$$

$s_t = \rho \cdot s_{t-1} + (1-\rho) g_t \odot g_t$

$$\theta_t = \theta_{t-1} - \eta \frac{g_t}{\sqrt{s_t} + \epsilon}$$

$$g_t = \begin{pmatrix} -15.33 \\ -6.667 \end{pmatrix}$$

$$s_t = \rho \cdot s_{t-1} + (1-\rho) g_t \odot g_t$$

$$= 0.9 \times \begin{pmatrix} 0 \\ 0 \end{pmatrix} + (0.1) \begin{pmatrix} (-15.33)^2 \\ (-6.667)^2 \end{pmatrix}$$

$$= \begin{pmatrix} 235.111 \\ 49.999 \end{pmatrix}$$

$$\sqrt{s_t} = \begin{pmatrix} 4.847 \\ 2.108 \end{pmatrix} \quad \frac{g_1}{s_1} = \begin{pmatrix} -15.33/4.847 \\ -6.667/2.108 \end{pmatrix} = \begin{pmatrix} -3.164 \\ -3.163 \end{pmatrix}$$

$$\theta = \begin{pmatrix} 0 \\ 0 \end{pmatrix} - 0.01 \begin{pmatrix} -3.164 \\ -3.163 \end{pmatrix}$$

This part follows lecture on Stanford CS229

No need to go through if you've understood the concepts well

Supervised Learning Setup

Regression (Where does this data fit?)
 Classification (What type of data is this?)

$x \rightarrow y$
 (Input) (Output)

$n \rightarrow$ number of examples in training set.
 (x_i, y_i)

$d \rightarrow x \in \mathbb{R}^d$ dimension of input.

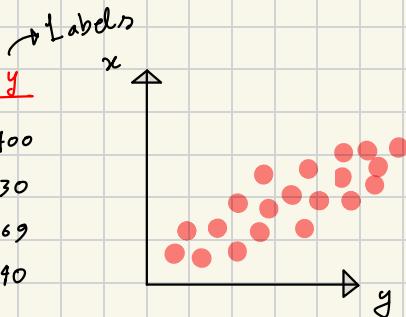
$x^{(i)} \rightarrow$ i th example input.

$y^{(i)} \rightarrow$ i th example output. (label / ground truth)

$(x^{(i)}, y^{(i)}) \rightarrow$ i th example.

* Regression

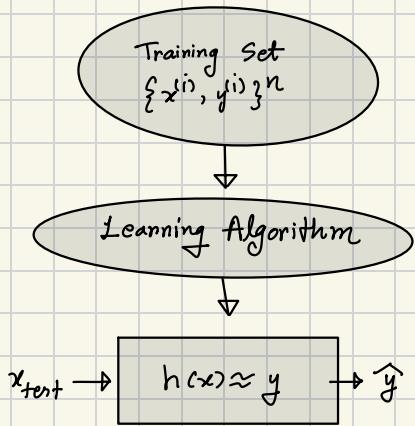
<u>x</u>	<u>y</u>
2104	400
1600	330
2400	369
3000	590



Goal: $f(x) \approx y \rightarrow$ Learnt Model

Or find the best fit line

* Pipeline



* Linear Regression:

$x \in \mathbb{R}^d$; $y \in \mathbb{R}$, n such examples

Hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d$$

$$h_{\theta}(x) = \sum_{i=1}^d \theta_i x_i + \theta_0 ;$$

$x_0 = 1 \rightarrow$ Intercept term

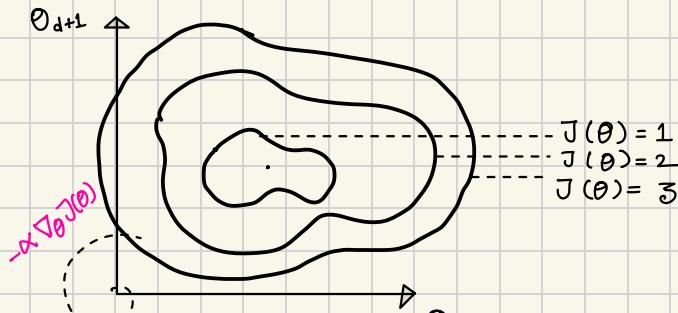
$$h_{\theta}(x) = \theta_0 x_0 + \dots + \theta_d x_d$$

$$= \sum_{i=0}^d x_i \theta_i = \theta^T x$$

$$\text{Cost Function: } J(\theta) = \frac{1}{2} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2 \quad | \quad J(\theta) = f(x)$$

$$\hat{\theta} \rightarrow \arg \min_{\theta} J(\theta) \rightarrow \arg \min_{\theta} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$

*Gradient Descent



Contour plot of the cost function

$\theta^{(0)}$ = Initialisation

$$\theta_j^{(t+1)} = \theta_j^{(t)} - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\therefore \theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta^{(t)})$$

→ In vector form

(we repeat this till convergence)

→ Gradient

$$\text{so, } \theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} J(\theta^{(t)})$$

How to check for convergence

$$1) \|\theta^{(t)} - \theta^{(t+1)}\| \approx 0$$

$$\text{or, } 2) \|\nabla_{\theta} J(\theta^{(t)})\| \approx 0$$

$$\text{or, } 3) |J(\theta^{t+1}) - J(\theta^t)| \approx 0$$

$\theta \rightarrow$ Set of all parameters.

Gradient Descent on Linear Regression

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} [\text{Cost function of LR}]$$

$$= \theta^{(t)} - \alpha \nabla_{\theta} \left(\frac{1}{2} \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right)$$

$$= \theta^{(t)} - \alpha \nabla_{\theta} \left(\frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2 \right)$$

$$= \theta^{(t)} - \alpha \left(\frac{1}{2} \sum_{i=1}^n \cancel{\times} (\theta^T x^{(i)} - y^{(i)}) x^{(i)} \right)$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \left(\sum_{i=1}^n (\underbrace{\theta^T x^{(i)}}_{\text{Scalar}} - \underbrace{y^{(i)}}_{\text{vector}}) \underbrace{x^{(i)}}_{\text{vector}} \right)$$

* Stochastic Gradient Descent - SGD (I find it stupid)

- For each small progress we need to scan through all of the dataset.
- To make it a bit simple we can do,

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} \tilde{J}(\theta)$$

$$\tilde{J}(\theta) = \frac{1}{2} (\theta^T x^{(u)} - y^{(u)})^2$$

* Only linear regression has a closed form solution of GD. In other cases, the solution varies with the problem.

We defined $J(\theta)$ as,

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2$$

Design Matrix

$$\bar{X} = \begin{pmatrix} n \times d+1 \\ x^{(1)} \\ \vdots \\ x^{(i+1)} \\ \vdots \\ x^{(n)} \end{pmatrix} \quad \bar{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_{d+1} \end{pmatrix}$$

$$\begin{aligned} & x\theta - y \\ \mathbb{R}^{n \times d+1} \mathbb{R}^{d+1} - \mathbb{R}^n &= \mathbb{R}^n \\ \therefore x\theta - y &= \begin{pmatrix} x^{(1)\top} \theta & -y^{(1)} \\ x^{(i)\top} \theta & -y^{(2)} \\ x^{(n)\top} \theta & -y^{(n)} \end{pmatrix} \end{aligned}$$

$$\text{so, } J(\theta) = \frac{1}{2} (x\theta - y)^T (x\theta - y)$$

$$\text{Now, } \nabla_{\theta} J(\theta) = 0$$

$$\begin{aligned} & \nabla_{\theta} \frac{1}{2} (x\theta - y)^T (x\theta - y) \\ &= \nabla_{\theta} \frac{1}{2} ((x\theta)^T (x\theta) - (x\theta)^T y - y^T (x\theta) + y^T y) \\ &= \nabla_{\theta} \frac{1}{2} (\theta^T (x^T x) \theta - 2\theta^T (x^T y) + y^T y) \\ &= \frac{1}{2} (2(x^T x)\theta - 2x^T y) = 0 \\ \Rightarrow x^T x \theta &= x^T y \quad \text{or, } \hat{\theta} = (x^T x)^{-1} x^T y \quad \left\{ \begin{array}{l} \text{As long as} \\ x^T \text{ is invertible} \end{array} \right\} \end{aligned}$$

*Probabilistic Interpretation

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$$\epsilon \approx N(0, \sigma^2)$$