# Linear Neural Networks

## 3.1 Linear Regression - Using Deep learning Book

- Regression - Methods of ==modeling== relationship
  between one or more ==independent==
  ==variable== and a ==dependent== variable.

- The purpose of ==Regression== in ==ML== is for ==prediction==.

eg / Predict a numerical value.
  - Predicting prices (of homes, & stocks)
  - Predicting length of stay (for patients in hospital).
  - Demand forecasting (for retail sales).

## 3.1.1 (Basic Element of Linear Regression).

example :-

- Estimate the ==prices of houses== (in dollars)
  based on ==area== (in square feet) and ==age== (in years).

→ Dataset = Training dataset / Training Set.

→ Each Row = Example / data point, data instance,
  (corresponding to                    Sample.
  one sale).

1⟶ Predict (Price) = Label / target.

2⟶ Independant Variables = Age & Area.

$\underbrace{\qquad\qquad}$

based on the independant variable

called [features (or covariates)]

3⟶ n = Number of examples in our datset.

examples are indexed such as :-

$$x^{(i)} = \begin{bmatrix} x_1^{(i)}, & x_2^{(i)} \end{bmatrix}^T \quad \text{— total} \atop \text{corresponding } y^{(i)}$$

## Linear Model

$$price = \underbrace{W_{area}} \cdot area + \underbrace{W_{age}} \cdot age + \underbrace{b}.$$

weights

bias.

(offset / intercept).

- weights determine the influence of each feature on our prediction.

- bias gives what value of the predicted price should be when all of the features take value 0.

- Without bias we will limit expressivity for our model.

$$\hat{y} = w_1 x_1 + \ldots + w_d x_d + b$$

hat
estimates

Can be expressed as .

$$\hat{y} = w^T x + b$$

Recap    Vectors

- Vectors are list of scalar values

- it is denoted by lowercase letters $x, y, z$.

- Refering element :-

    for $x$    we can    use    $x_i$ - Scalar.

    (in/consist of)

$$x \in \mathbb{R}^n$$

Vector $x$

Real values scalar

$$\hat{y} = Xw + b$$

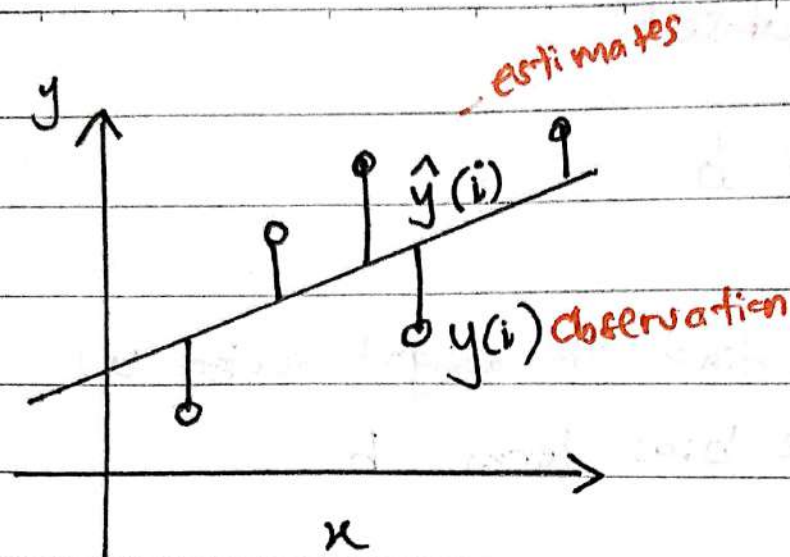- Goal  - To find the weight vector $w$
  - The bias term $b$

## Loss Function

Loss function - measure the distance between the
real & predicted value.

- Loss  → usually non-negative number
  - perfect predictions incur a loss of $0$.

- Squared error

$$\ell^{(i)}(w,b) = \frac{1}{2}\left(\hat{y}^{(i)} - y^{(i)}\right)^2$$

- $\left(\frac{1}{2}\right)$ will improve notationally, canceling out
  when we take derivates loss.

$y$

estimates

$\hat{y}(i)$

$y(i)$ observation

$x$

(Fit data with a linear model)

- difference between $\hat{y}^{(i)}$ & $y^{(i)}$ increase
  $\leadsto$ Larger loss due to quadratic dependece.

- quadratic dependance -

$$L(w,b) = \frac{1}{n} \sum_{i=1}^{n} \ell^{(i)}(w,b) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{2} \left( w^T x^{(i)} + b - y^{(i)} \right)^2$$

with bias.

# Derivation for the equation.

$$\hat{y} = w_1 x_1 + \ldots + w_d x_d + b.$$

$$\hat{y} = w^T x + b \qquad \qquad —①$$

**Loss function**

$$\mathcal{L}^{(i)}(w, b) = \frac{1}{2}\left(\hat{y}^{(i)} - y^{(i)}\right)^2$$

$$= \frac{1}{2}\left(w^T x + b - y^{(i)}\right)^2 \quad —②$$

**Average**

$$\frac{1}{n} \sum_{i=1}^{n} \qquad —③$$

**Combine**

$$\frac{1}{n} \sum_{i=1}^{n} \frac{1}{2}\left(w^T x + b - y^{(i)}\right)^2 \qquad \begin{array}{l} ① \text{ into } ② \\ ② \text{ into } ③ \end{array}$$

$$\#$$

$$w^* b^* = \arg\min \; \mathcal{L}(w, b)$$

$\downarrow$

To minimize the total loss across all training examples.

# ✳ Gradient Descent
---×---

- An **algorithm** that consist of iteratively **reducing** the **error** by updating the parameter in the direction that incrementally **lowers** the **loss function.**

- Average of the losses computed on every single example in the dataset.

- Problem : Extremely slow, we must pass entire dataset before making a single update.

- Solution :- Use a random minibatch of examples every time we need to compute the update.
  ↝ Called minibatch stochastic gradient descent.

→ In each iteration, we sample a minibatch $\beta$
  consist of fixed number of training examples.

→ Compute the derivative (gradient) of the average loss on the minibatch.

→ We multiply the gradient by a predetermined positive value $\eta$.

$$(w, b) \leftarrow (w, b) - \frac{\eta}{|B|} \sum_{i \in B} \partial_{(w,b)} \ell^{(i)}(w, b)$$

$$w - \frac{\eta}{|B|} \sum_{i \in B} \partial_w \ell^{(i)}(w, b) = w - \frac{\eta}{|B|} \sum_{i \in B} x^{(i)} \left( w^T x^{(i)} + b - y^{(i)} \right)$$

- we initialize the values of the model parameters, at random.

$$b - \frac{\eta}{|B|} \sum_{i \in B} \partial_b \ell^{(i)}(w, b) = b - \frac{\eta}{|B|} \sum_{i \in B} \left( w^T x^{(i)} + b - y^{(i)} \right)$$

- we iteratively sample random minibatches from the data.

- Updating the parameters in the direction of negative gradient.

$|B|$ = Number of examples in each minibatch (the batch size)

$\eta$ = Learning rate.

## Linear Regression

Equation : $Y = a + bx$

intercep

explanatory variable.

dependant Variable

slope (gradient)

old eq:
$$y = mx + c$$

## Least-Squares Regression

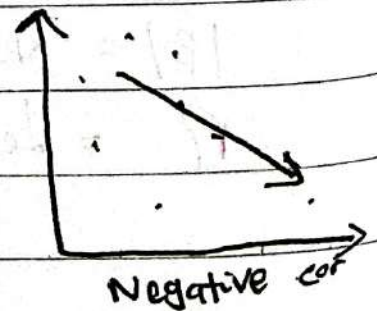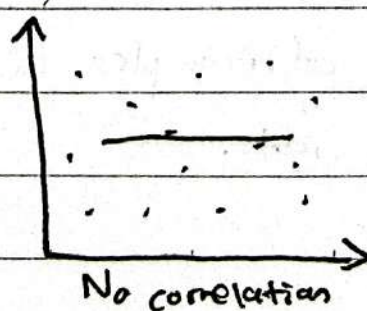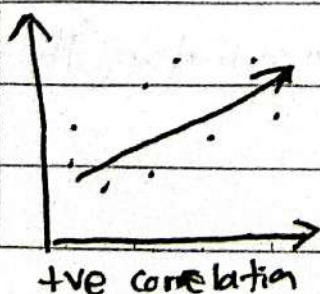-Used for fitting ~~a regression~~ a regression line.

~ calculates the best fitting line

- by minimizing the sum of squares of vertical deviations from each data points to the line.

## Tips

Correlation . ~ measure how strong a relationship between two variables.

eg



tve correlation          No correlation          Negative cor
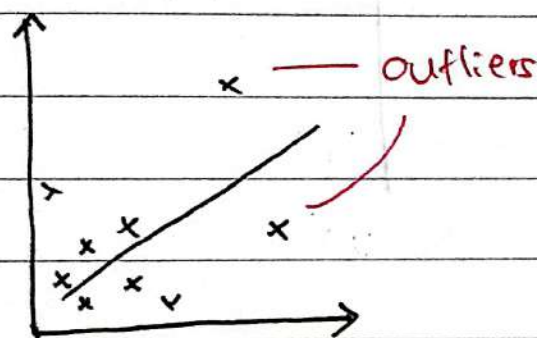
example : dataset 'Television, Physician & Life Expectancy'.
- Contains 40 ~~countries~~ countries.
- 8 removed 32 remaining
- Correlation coefficient, $r = 0.852$
- $r^2 = 0.726$
= 72.6% of variation in one variable.



- Most of the data points are at the **left** and the ones on the right called **outliers**.

- ~~The~~ Point lies for from other data in the horizontal direction :- **Influential observation.**

- After the influential observation removed :-
Correlation $r$ dropped to $0.427$.
$r^2 = 0.182$
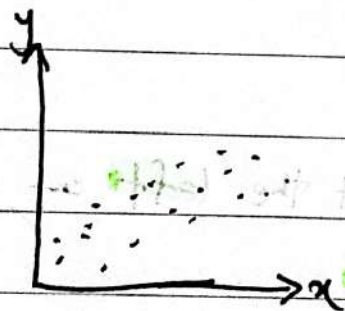= < 20% of variation in number of people per physician

Recall Supervised Learning
— x —

example

Living areas and prices of 47 houses.

| $x$ (Living area) | $y$ (1000 $s) |
|---|---|
| data | data. |

Produce
scatter
plot.
(+ve
correlation)



$x^{(i)}$ = input (living area) called input features

$y^{(i)}$ = output (price) target variable

A pair of $\left( x^{(i)}, y^{(i)} \right)$ – called traing training example.

Used to learn
a list of $n$ training
examples. $\}$.

Training Set.

$$(x^{(i)}, y^{(i)}) \qquad \qquad -①$$

$$\{(x^{(i)}, y^{(i)}) ; i = 1 \dots n\} \qquad -②$$

$$X = \text{Input} \quad ; \quad Y = \text{Output}$$

$$X = Y = \mathbb{R} \quad -\text{real numbers.}$$

Given Training Set $\quad h : X \longrightarrow Y$

$h(x)$ predictor of value $y$.

Process



- Predict Continuous problem $\sim$ Regression Problem.

- $y$ can take only small discrete values

   - we call it as classification problem.

Based on Lecture Notes Andrew NG

## Linear Regression

Using a data set

| Living area (feet)$^2$ | # bedrooms | Price (1000$s) |
|---|---|---|
| | | |
| | | |
| | | |

- $x$'s , two dimensional vector $\mathbb{R}^2$.
- $x_1^{(i)}$ = living area
- $x_2^{(i)}$ = Number of bedrooms.
- Approximate Y as a linear function of $x$.

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 \dots \theta_n x_n^{(i)}$$

assume $x_0 = 1$    $\therefore$
(intercept term)

$\theta_i$  = parameters / weights.

~~We take~~

We take the eq as :

$$h(x) = \sum_{i=0}^{d} \theta_i x_i = \theta^T x.$$

-mate $h(x)$ close to $y$.

**linear eq**

$$y = mx + c$$
$$h(x) = \theta x$$

~~the~~ -We use loss function (least square function)

- square error:

$$\frac{1}{2}\left(\hat{y}^{(i)} - y^{(i)}\right)^2 \qquad -①$$

estimates

In this case :

$$\hat{y} = h_\theta(x) \qquad\qquad ② \quad sub$$

$$J(\theta) = \frac{1}{2}\left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$$

$\Leftarrow$ We find sum $h_\theta(x) = \sum \theta_i x_i$

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

# Least Mean Square LMS Algorithm.

~ To make $J(\theta)$ minimize we need to choose $\theta$.

~ Algorithm that repeatedly changes $\theta$ to make $J(\theta)$ smaller.

~ We use Gradient Descent Algorithm.

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta).$$ ①

$\alpha =$ Learning rate.      $\alpha \in \mathbb{R}^+$ $\left(\begin{array}{c} \text{range between} \\ 0 \text{ } \& \text{ } 1 \end{array}\right)$.

$\theta_j =$ Parameters.

$j = 0, \dots, n.$

$\theta =$

~ This algorithm repeatedly decrease the steepness of $J$.

We take the eq from $J(\theta)$

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

Recap Partial derivative.

Basic differentiation concept

$$\frac{d}{dx}(ax+b)^n = n(ax+b)^{n-1}(a)$$

$$= (n)(a)(ax+b)$$

Continue :- $J(\theta) = \frac{1}{2}(h_\theta(x) - y)^2$

Apply basic concept :

$h_\theta(x) = \theta_1 x_1 + x_0 \dots$

$= \sum \theta_i x_i$

We do this to sub into $\theta_j$

$$\frac{d}{d\theta_j} J(\theta) = \frac{1}{2}(2)(h_\theta(x) - y)$$

$$\frac{d}{d\theta_j} J(\theta) = \frac{1}{2}(2)\left[\frac{d}{d\theta_j}(h_\theta(x) - y) \cdot (h_\theta(x) - y)\right]$$

$$= \left[\frac{d}{d\theta_j}\left(\sum_{i=0}^{d} \theta_j x_j\right) \cdot (h_\theta(x) - y)\right]$$

$$= \left[x_j \cdot (h_\theta(x) - y)\right]$$

$$= (h_\theta(x) - y)x_j \quad \#. \quad ②$$

Sub eq ① - into ②

$$\theta_j = \theta_j - \alpha \frac{d}{d\theta_j} J(\theta) \quad - ①$$

$$\frac{d}{d\theta_j} J(\theta) = \left(h_\theta(x) - y\right) x_j \quad - ②$$

$$\theta_j = \theta_j - \alpha \left(h_\theta(x) - y\right) x_j .$$

For Single training example.

- The rule is called LMS update rule.
  (Least Mean Squares).
- Widrow - Hoff learning rule.
- directly proportional to the error term
  $$\left(y^{(i)} - h_\theta(x^{(i)})\right)$$
- if the training set have more than
  one example :-

- Repeat until convergence :-

$$\theta_j = \theta_j + \alpha \sum_{i=1}^{n} \left(y^{(i)} - h_\theta(x^{(i)})\right) x_j^{(i)}, \text{ for every } j$$

succinct = short / compact.

$$\theta = \theta + \alpha \sum_{i=1}^{n} \left( y^{(i)} - h_\theta(x^{(i)}) \right) x^{(i)}.$$

- method looks at every example in the entire training set on every step.

- Called **batch gradient descent.**

- Gradient descent always converges to the global minimum.

- J is a convex quadratic function

-

# The normal equations

- Another way of optimizing $J$ other than gradient descent.
- This method will minimize $J$ by taking its derivatives with respect to $\theta_j$ & setting them to (zero).

## Least Square ~~Matrix~~ Using Matrix Derivatives

Training Examples :-

$$X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(n)})^T - \end{bmatrix} \qquad \text{design matrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$X\theta = \begin{bmatrix} (x^{(1)})^T \theta \\ \vdots \\ (x^{(n)})^T \theta \end{bmatrix}$$

$$h_\theta\left(x^{(i)}\right) = \left(x^{(i)}\right)^T \theta.$$

$$\cancel{X\theta =}$$

$$X\theta - \vec{y} = \begin{bmatrix} \left(x^{(1)}\right)^T \theta \\ \left(x^{(n)}\right)^T \theta \end{bmatrix} - \begin{bmatrix} y^{(1)} \\ y^{(n)} \end{bmatrix}$$

$$= \begin{bmatrix} h_\theta\left(x^{(1)}\right) - y^{(1)} \\ h_\theta\left(x^{(n)}\right) - y^{(n)} \end{bmatrix}$$

$$\frac{1}{2}\left(X\theta - \vec{y}\right)^T \left(x\theta - \vec{y}\right) = \frac{1}{2}\sum_{i=1}^{n} \left(h_\theta\left(x^{(i)} - y^{(i)}\right)\right)^2$$

$$= J(\theta)$$

$$\cancel{\left(x\theta - y\right)^T \left(x\theta - y\right)}$$

Eq:

$$\left(x\theta - y\right)^T \left(x\theta - y\right) = \sum_{i=1}^{m}\left(h_\theta\left(x^{(i)} - y^{(i)}\right)\right)^2$$

$$J(\theta) = \frac{1}{2}\left(x\theta - y\right)^T \left(x\theta - y\right). \#.$$