# Model and Cost function

Not"s :

  $m$ $\ne$ no. of training ex·s

  $x$'s : input var·s / features

  $y$'s : output var·s / ~~features~~ target var. that we are
                  trying to predict

  $(x, y)$ : one training ex. / observⁿ

  $(x^{(i)}, y^{(i)})$: $i^{th}$ training ex. / observⁿ

  $(x^{(i)}, y^{(i)})$; $i = 1, 2, \ldots, m$ : training set

  $i$  : index into the training set

  $X$  : space of I/P val·s

  $Y$  : space of O/P val·s

  $n$  : no. of features

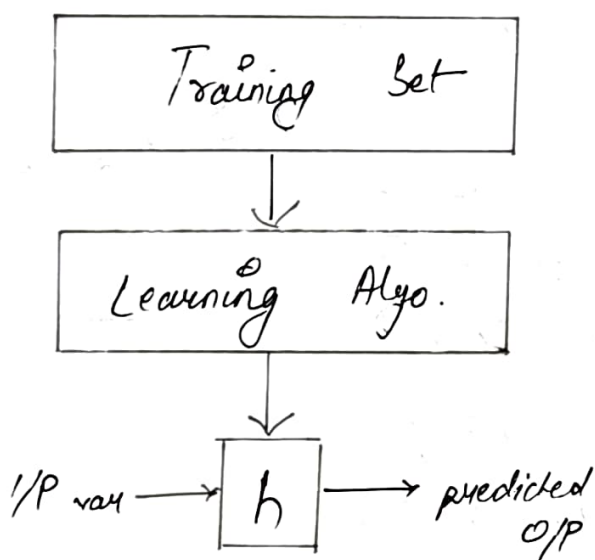$$\boxed{\text{Hypothesis} : h_\theta(x) = \theta_0 + \theta_1 x}$$

$\theta_i$·s  : para. of model

\# $h$ maps from $x$'s to $y$'s

In supervised learning, our goal is: given a training set to learn a fun$^n$ $h: x \rightarrow Y$ so that $h(x)$ is a "good" predictor for the corresponding val. of $y$

How it works: we find a training set (ex: Housing price size in sq. ft on $x$ vs \$ on $y$) into our learning algo. The learning algo. then outputs a fun$^n$ 'h' called hypothesis (terminology)

This fun$^n$ h takes i/p val and predict the desired o/p (here, takes feet$^2$ & predict price)

→ (When the target var. that we are trying to achieve is continuous, such as in our housing problem, we call the learning probl. a **Regression** probl.
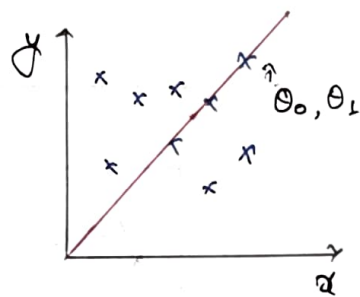
→ (When you can take only a small no. of discrete val.s (such as, if given the living area, we wanted to predict if the dwelling is a house or an apartment), we call it a **Classification** problem.

## Cost Function

let us figure out how to fit the best possible straight line to our data

Idea : choose $\theta_0, \theta_1$ so that $h_\theta(x)$ is close to 'y' for our training ex.s $(x, y)$

\# i.e., $h_\theta(x) - y$ is min for all training set (diff. b/w hypoth. & actual o/p)



\# $\theta_0, \theta_1$ ($\theta_i$'s) are para. of model. we need to find the ~~closest~~ val. of para.s closest to the actual val.s. for this, we will calc. diff b/w our predicted val. & og val for each data set. for that we use **Cost Fun**

We can measure the accuracy of our hypothesis funⁿ by using a cost function.

This takes an avg. difference (actually a fancier ver of an avg) for all results of the hypothesis w/ I/ps from x's and actual o/p, y's.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x_i) - y_i)^2$$

Minimize $J(\theta_0, \theta_1)$
$\theta_0, \theta_1$
$\underbrace{\qquad}_{\text{cost fun}^\text{n}}$

$$h_\theta(x_i) = \theta_0 + \theta_1 x^{(i)}$$

To break it apart, it is $\frac{1}{2} \bar{x}$ where $\bar{x}$ is the mean of squares of $h_\theta(x_i) - y_i$, or the diff. b/w the predicted val. & the actual val.

This funⁿ is otherwise called the "squarred error funⁿ" or "mean squarred error".

The funⁿ is halved ($\frac{1}{2}$) as a convenience for the computⁿ of gradient descent, as the derivative term of the sq. funⁿ will cancel out the $\frac{1}{2}$ term.

# Cost function - Intuition I

If we try to think of it in visual terms, our training data set is scattered on the x-y plane. We are trying to make a straight line (defined by $h_\theta(x)$) wh. passes through these scattered pts.

Our <u>objective</u> is to get the best <u>possible line</u>.

The best possible line would be such so that the avg. sq.ed vertical distances of the scattered pt.s from the line will be the least.

ideally, the line should pass through all the pt.s in our data set. In such a case, $J(\theta_0, \theta_1) = 0$.

Now mathematically speaking, each val. of $\theta_1$ corresponds to a diff. straight line fit or a diff. ex. hypoth ($h_\theta(x)$) and for each val. of $\theta_0, \theta_1$ you can derive some some val. of $J(\theta_0, \theta_1)$.

So we for visualizing what cost fun" does & why we use it, we calc. val.s of $J(\theta_0, \theta_1)$ for diff. val. of hypoth (i.e; $h_\theta(x) = \theta_0 + \theta_1 x$) and plot them in a graph. We then look at the graph and figure out for what val. of $\theta_0, \theta_1$ is $J(\theta_0, \theta_1)$ lowest i.e; we minimize $J(\theta_0, \theta_1)$, wh. is our goal.
$\theta_0, \theta_1$

ex:    let,   $\theta_0 = 0$

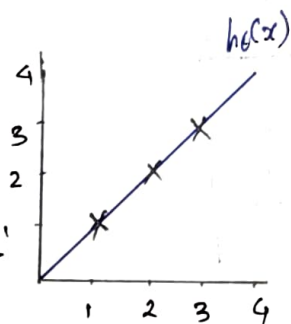so,   $h_\theta(x) = \theta_0 + \theta_1 x$   $\Rightarrow$   $h_\theta(x) = \theta_1 x$

Now,

| $h_\theta(x) = \theta_1 x$ | $J(\theta_1)$ |
|---|---|
| (for fixed $\theta_1$, it is a fun$^n$ of $x$) | (fun$^n$ of para $\theta_1$) |

# so graph of this fun$^n$ will be the slope $h_\theta(x)$ plotted in $x/y$ axes

# so the graph of this fun$^n$ will be plotted in $\theta_1/J(\theta_1)$ axes (only when $\theta_0 = 0$)

for $\theta_1 = 1$

$J(\theta_1) = \frac{1}{2m}$

$\sum_{i=1}^{M} (h_\theta(x^{(i)}) - y^{(i)})^2$

$= \frac{1}{2m} \sum_{i=1}^{M} (\theta_1 x^{(i)} - y^{(i)})^2$
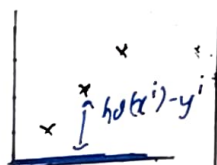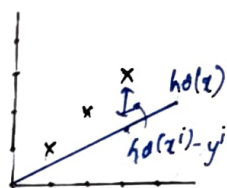
$= \frac{1}{2 \times 3} ((0-0)^2 + (0-0)^2 + (0-0)^2$

$J(\theta_1) = 0 \Rightarrow J(1) = 0$



for $\theta_1 = 0.5$

$J(\theta_1) = \frac{1}{2\times3} \begin{bmatrix} (0.5-1)^2 + \\ (1-2)^2 + (1.5-3)^2 \end{bmatrix}$

$J(0.5) = 0.68$



for $\theta_1 = 0$

$J(0) = \frac{1}{2\times3} (1^2 + 2^2 + 3^2)$

$J(0) = \frac{14}{6} = 2.3$



plotting $J(1), J(0.5), J(0)$ in $J(\theta_1)/\theta_1$ axes
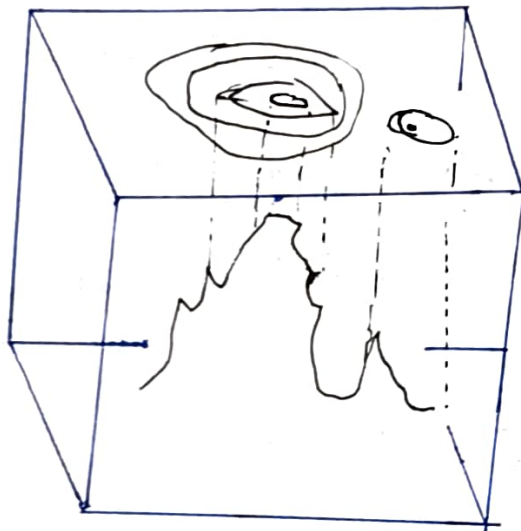↑ (only bc $\theta_0 = 0$)

← on further plotting



Clearly,   $J(\theta_1)$ is min at when   $\theta_1 = 1$.
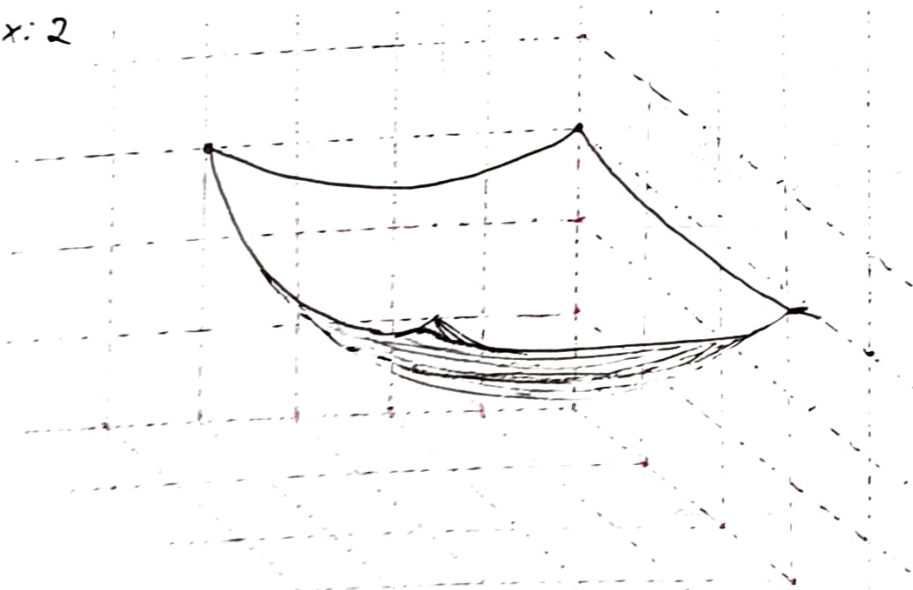
Objective completed

# Cost function - Intuition II

#  a Contour plot is a graph that contains many contour lines. A contour line of a 2 var. funⁿ has a const. val. at all pt.s of the same line. a Contour plot is a graphical technique for plotting a 3-D surface by plotting const. z slices called 'contours' on a 2-D format.
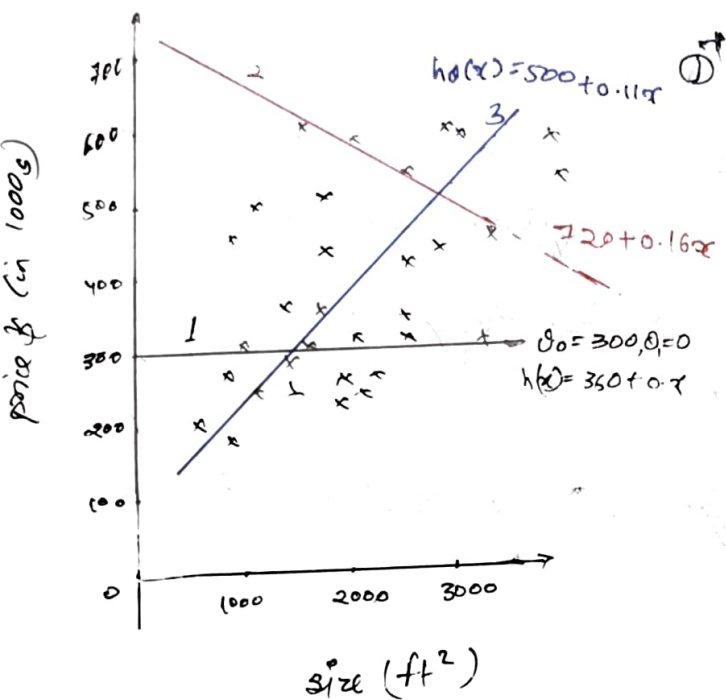
ex: (1)



ex: 2



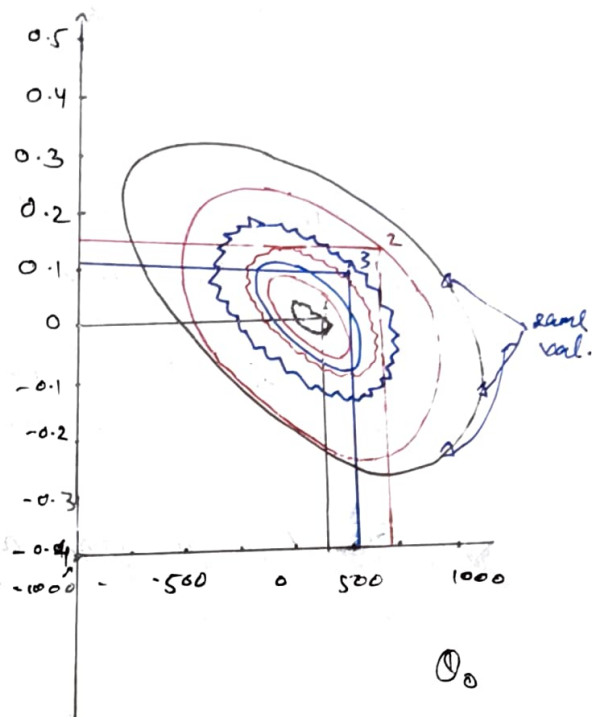3-D ui. is drawn in the following table

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a
fun$^n$ of $x$)

$J(\theta_0, \theta_1)$

(fun$^n$ of para.s $\theta_0, \theta_1$)



Left plot: price $£$ (in 1000s) vs size (ft$^2$). Lines labeled 1, 2, 3.
$h_\theta(x) = 500 + 0.11x$ ①
$720 + 0.16x$
$\theta_0 = 300, \theta_1 = 0$
$h(x) = 350 + 0.x$

Right plot: contour plot of $J(\theta_0, \theta_1)$, axis $\theta_0$.
same val.

$\underset{1}{.} \ (\theta_0 = 300, \theta_1 = 0)$

$2 \ \ \hat\theta_0 = 720, \theta_1 = 0.16$

$3 \ (\theta_0 = 500, \theta_1 = 0.11)$

# all pts. on the same contour line
have same val. so, all 3 $\Delta$s
have equal val (some val. of $J(\theta_0, \theta_1)$)

# ⌇⌇⌇⌇ .... — are linear lines like —
/ or \ only. They are drawn such only
bc I have 3 colour pens.

# Gradient Descent for Linear Regression

here we will put together grad. desc. w our cost fun^n and that will give us an algo for Linear Regression, or putting a straight line to our data.

grad. desc. algo :

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j = 0$ & $j = 1$)

}

Linear Regression Model:

$$h_\theta(x) = \theta_0 + \theta_1 x$$

Linear hypoth.

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

sq.ed error cost fun

Now, we need to calc. $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^i) - y^i \right)^2$$

$$= \frac{1}{2m} \sum_{i=1}^{m} \left( \theta_0 + \theta_1 x^i - y^i \right)^2$$

so, for

$j = 0$ : $\dfrac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \left( h_\theta(x^i) - y^i \right)$ ⁽ᵖᵃʳᵗⁱᵃˡ ᵈⁱᶠᶠ ʷʳᵗ θ₁⁾

$j = 1$ : $\dfrac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \dfrac{1}{m} \sum\limits_{i=1}^{m} \left( h_\theta(x^i) - y^i \right) \cdot x^i$

putting them in grad. desc. algo

repeat until convergence {

$\theta_0 := \theta_0 - \alpha \dfrac{1}{m} \sum\limits_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)$

$\theta_1 := \theta_1 - \alpha \dfrac{1}{m} \sum\limits_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x^{(i)}$

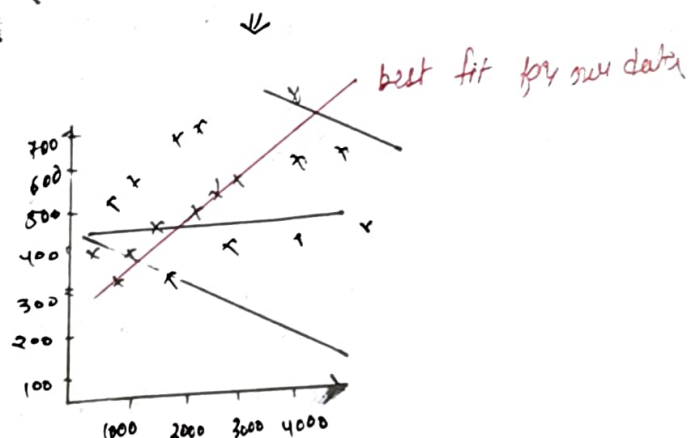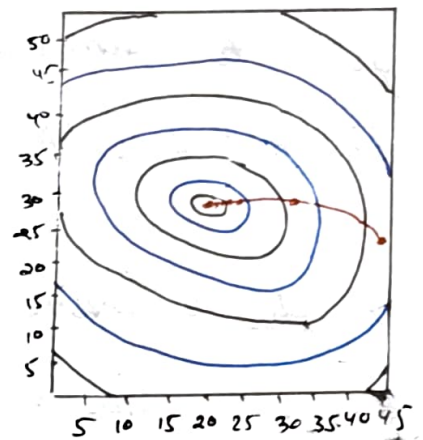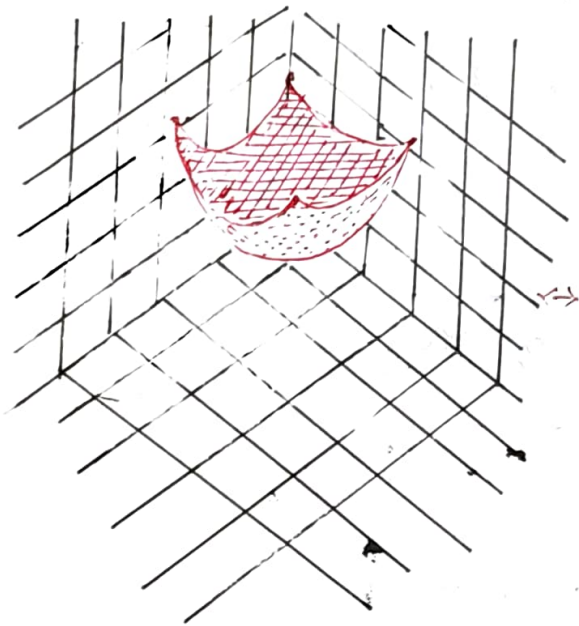# This algo is also k/a Batch grad. desc⁽ʳⁱᵍ⁾ where

Batch means that every step of grad. desc. uses

all training ex.s $\sum\limits_{i=1}^{m} \left( h_\theta(x^i) - y^i \right)$

However, some algo.s also use subsets of training
ex.s and not all of them.

If we start w a guess for our hypoth. &
then repeatedly apply grad. desc. eq's, our hypoth.
will become more accurate.

★

the cost fun for Linear Regression is always
going to be a Convex fun ( bow-shaped fun )
and so, it doesn't have the local optimum
except the one global optimum.
so, the grad. desc. will always converge to the
1 local optimum.



$J(\theta_0, \theta_1)$ fun of paras $\theta_0, \theta_1$

best fit for our data

# Multivariate Linear Regression.

Linear Regression w multiple var.s is k/a "multivariate linear regression".

→ i.e., we have multiple features/var.s w wh. we try to fit a model to our data.

ex:

| size $(ft^2)$ | no. of bedrooms | no. of floors | age of home (yrs) | price ($1000) |
|---|---|---|---|---|
| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 ⎱ |
| | | | | ⋮ ⎱ ₋47 |

We now introduce nota's for eq's where we can have any no. of i/p var.s

$$m = \text{no. of training ex.s} = 47$$
$$n = \text{no. of features} = 4$$

$x^{(i)} = $ i/p (features) of $i^{th}$ training ex.

$i →$ index to training set

$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$

$x_j^{(i)} = $ val. of feature $j$ is $i^{th}$ training ex.

$x_3^{(2)} = 2$

$x_1^{(4)} = 852$

\# similar to our hypo. for linear regr. $h_\theta(x) = \theta_0 + \theta_1(x)$

The multivariate form of the hypothesis function accomodating these multiple features is as follows:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

★ for convenience, we assume $x_0^{(i)} = 1 \; \forall \, (i \in 1, ..., M)$ in this course.

This allows vec.s '$\theta$' & $x^{(i)}$ match each other element wise $(n+1 \text{ el.s})$

i.e., $x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$ & already $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$

Now,

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n$$

$$= \underbrace{\begin{bmatrix} \theta_0 & \theta_1 & \cdots & \theta_n \end{bmatrix}}_{1 \times n+1} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}}_{n+1 \times 1}$$

$y = \boxed{h_\theta(x) = \theta^T x}$

# Gradient Descent for Multiple Variables

hypoth : $h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n$

Para.s : $\theta_0, \theta_1, \theta_2, \ldots, \theta_n = 0$

$\theta = [\theta_0, \ldots, \theta_1]$ is an $n+1$ dimension vec.

Cost fun$^n$ : $J(\theta_0, \theta_1, \ldots, \theta_n) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$

or, $J(\theta) = \dfrac{1}{2m} \sum\limits_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)^2$

# Gradient Descent

\# for single feature / var,

\# for $n = 1$

Repeat

$\{$

$\qquad \theta_0 := \theta_0 - \alpha \underbrace{\dfrac{1}{m} \sum\limits_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)}_{\frac{\partial}{\partial \theta_0} J(\theta_0)}$

$\qquad \theta_1 := \theta_1 - \alpha \dfrac{1}{m} \sum\limits_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right) x^{(i)}$

$\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \& \; x_j^{(i)}$

$\qquad \qquad$ (simultaneously update $\theta_0, \theta_1$)

$\}$

The gradient descent equation itself is generally the same form, we just have to repeat it for our 'n' features.

$n \geq 1$.

repeat until convergence

{

$$\theta_j := \theta_j^{\circ} - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

$$(\text{simultaneously update } \theta_j \text{ for } j = 0, 1, \ldots, n)$$

}

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_0^{(i)}$$

$$\text{sim. to prev } \theta_0 \text{ (for } n=1) \text{ as}$$

for convenience, we have assumed $x_0 = 1$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_1^{(i)}$$

$$\text{sim to prev } \theta_1 \text{ (for } n=1)$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_2^{(i)}$$

# Features and Polynomial Regression.

We can improve our features and the form of our hypoth. fun $^n$ in a couple diff. ways.

We **combine** multiple features into one.

ox :   Housing Price predic$^n$ :

$$h_\theta(x) = \theta_0 + \theta_1 \times \underset{x_1}{frontage}$$

$$+ \theta_2 \times \underset{x_2}{depth}$$



Instead of using the features we have in hand, we may create another feature wh. helps us create a better model.

$$Area$$
$$x = frontage \times depth$$

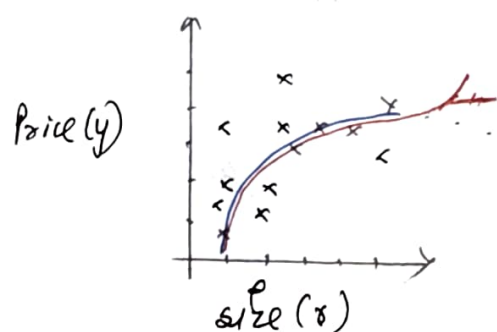then,   $h_\theta(x) = \theta_0 + \theta_1 \underset{\curvearrowleft land\ area}{x}$

Polynomial Regression → how to put polynomial (quadratic, cubic) fun$^n$ into data). Polynomial Regression allows us to use the machinery of Linear Regression to fit very complicated, even non-linear data

Our hypo. fun? need not be a linear (straight line) if it does not fit the data well.

We can change the behaviour or curve of our hypo fun? by making it quadratic, cubic or sq. root fun? (or any other form).

# if you have data set



Price (y)

size (x)

you can fit a model

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

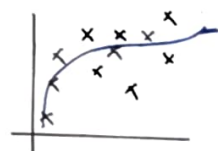← (∵ quad. fun's aze go up then come down)

or, $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$

(3im·season : cubic eqⁿ soes up)

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$= \theta_0 + \theta_1 (size) + \theta_2 (size)^2 + \theta_3 (size)^3$$

or

$$= \theta_0 + \theta_1 (size) + \theta_2 (\sqrt{size})$$



When fitting such model to our data, feature scaling becomes extremely imp. &

ex. $x_1 = size = 1 - 1000 \ ft^2$
$x_2 = size^2 = 1 - 10^6 \ ft^2$
$x_3 = size^3 = 1 - 10^9 \ ft^2$.

# Normal Equation.

Gradient Descent gives 1 way of minimizing J.
Let's discuss a second way of doing so, this
time performing the minimization explicitly & w/o
resorting to an iterative algo.

In the "Normal Equation" method, we will minimize
J by explicitly taking its (partial) derivatives w.r.t.
the $\theta_j$'s, and setting them to 0.
This allows us to find the optimum $\theta$ w/o
iter$^n$.

The normal eq$^n$ formula is given below :

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

→ Intu$^n$ : If 1D ($\theta \in \mathbb{R}$)

to minimize $J(\theta)$

$$\frac{d}{d\theta} J(\theta) \xrightarrow[set]{} 0$$

solve for $\theta$     (gives us optimum val.
of $\theta$)

but in our problem, $\theta \in \mathbb{R}^{n+1}$ ($\theta$ is $n+1$ dimensional vec.)

$$J(\theta_0, \theta_1, \ldots, \theta_m) = \frac{1}{2m} \sum_{i=1}^{m} \left(h_\theta(x^{(i)}) - y^{(i)}\right)$$

so, $\frac{\partial}{\partial \theta_j} J(\theta) \overset{set}{=} 0$  (for every $j$)

solve for $\theta_0, \theta_1, \ldots, \theta_m$  gives us the vals of $\theta$ that minimize the cost fun $J(\theta)$.

→ Examples :

$m = 4$

| $x_0$ | size $(ft^2)$ $x_1$ | no. of bedrooms $x_2$ | no. of floors $x_3$ | age of home (yrs) $x_4$ | Price ($1000) $y$ |
|---|---|---|---|---|---|
| 1 | 2104 | 5 | 1 | 45 | 460 |
| 1 | 1416 | 3 | 2 | 40 | 232 |
| 1 | 1534 | 3 | 2 | 30 | 315 |
| 1 | 852 | 2 | 1 | 36 | 178 |

we added this col. ↓

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}_{m \times (n+1)} \qquad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}_{m \times 1}$$

then,

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

gives us the vals of $\theta$ that minimizes the cost fun.

→ In gen.,

$m$ ex.s : $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \ldots, (x^{(M)}, y^{(M)})$

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in R^{n+1} \quad \& \quad \underset{\substack{\uparrow \\ k/a \\ \text{design matr.}}}{X} = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(m)})^T - \end{bmatrix}$$

ex:

$$\text{if } x^{(i)} = \begin{bmatrix} 1 \\ x_2^{(i)} \end{bmatrix} \quad \Rightarrow \quad X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix}$$

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

then, $\qquad \theta = \cancel{X} \; (X^T X)^{-1} X^T y$.

→ There is no need to do __feature scaling__.
  w the Normal Eqⁿ.

\# ~~for~~ more complex algo.s like Classificⁿ and
  Logistic Regression cannot be solved using
  Normal Equation. for them we have to use
  grad. desc.

→ Octave: 
$$\text{pinv}(x' * x) * x' * y$$

→ Comparison b/w Gradient Descent & Normal Eq$^n$ for $m$ training ex.s, $n$ features.

| Grad Desc. | Normal Eq$^n$ |
|---|---|
| → need to choose $\alpha$ | → no need to choose $\alpha$ |
| → needs many iter's | → no need to iterate |
| → $O(kn^2)$ | → {$\theta \propto O(n^3)$, need to calc. $(x^Tx)^{-1}$} ← causes |
| → works well when $n$ is large $\rightarrow n > 10^6$ | → slow if $n$ is very large. |

Normal Equation and non-invertibility.

If $x^Tx$ is non-invertible (singular / degenerate), the common causes might be having

• redundant features (linearly dependent) ex: $x_1 = ft^2$  $x_2 = m^2$

• too many features (ex: $m < n$) In this case, del some features or use "regulariz$^n$".

#. In MATLAB/OCTAVE, use pinv() instead of inv() to calc $(x^Tx)^{-1}$

# Regularized Linear Regression.

Optimize objective:

$$J(\theta) = \frac{1}{2m}\left[\sum_{i=1}^{M}\left(h_\theta(x^{(i)}) - y^{(i)}\right)^2 + \lambda \sum_{i=1}^{n} \theta_j^2\right]$$

( regularized lin regr

$$\min_\theta J(\theta)$$

## Grad. Desc.:

Repeat

$$\{ \quad \theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^{M}\left(h_\theta(x^{(i)}) - y^{(i)}\right)x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha\left[\frac{1}{m} \sum_{i=1}^{M}\left(h_\theta(x^{(i)}) - y^{(i)}\right)x_j^{(i)} + \frac{\lambda}{m}\cdot \theta_j\right] \quad j \in \{1, 2, ..., n\}$$

the term $\frac{\lambda}{m}\cdot\theta_j$ performs our regulariz.

Update rule on further manipul:

$$\theta_j := \theta_j^\circ\left(1 - \frac{\alpha\lambda}{m}\right) - \alpha\cdot\frac{1}{m}\sum_{i=1}^{M}\left(h_\theta(x^{(i)}) - y^{(i)}\right)x_j^{(i)}$$

$$\because m > 0, \alpha > 0 \quad \Rightarrow \quad 1 - \frac{\alpha\lambda}{m} < 1$$

(ex 0.99)

What the above expr. means is that $\forall j = 1, ..., n$, for every "iter$^n$" we are updating $\theta_j$ a little by first regularizing (shrinking) it a little - $\left(1 - \frac{\alpha \lambda}{m}\right)$, and then performing a similar

(grad. desc.) update as before $\left(-\alpha \cdot \frac{1}{m} \sum_{i=1}^{\hat{m}} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}\right)$.

#$^\circ$ is just "intu$^n$", mathematically it is just performing grad. desc. on regularized cost fun$^n$.

## Normal Eq$^n$:

To add in regulariz$^n$, the eq$^n$ is same as og, except that we add another term inside parantheses.

$$X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ - (x^{(m)})^T - \end{bmatrix}_{m \times (n+1)} \text{mat}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}_{m \times 1 \text{ vec.}}$$

$$\theta = \left(X^T X + \lambda \cdot L\right)^{-1} X^T y$$

$$\text{where} \quad L = \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}_{(n+1) \times (n+1)}$$

ex: $n = 2$

$$L = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}_{3 \times}$$

# in regularized lin. regr., if $\lambda$ is set to an extremely large val., algo results in Underfitting. bc $h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$.

if $\lambda = 10^4$ (say) $\Rightarrow \theta_1, \theta_2, ..., \theta_n \approx 0 \Rightarrow h_\theta(x) = \theta_0$