



→ SUPERVISED Machine : Regression & Classification LEARNING

Supervised vs. Unsupervised Learning

Quick Quest! If Arthur Samuel's checkers-playing program had been allowed to play only 10 games against itself, how would this affected its performance compared to when it was allowed to play over 10k games?
→ would have made it worse

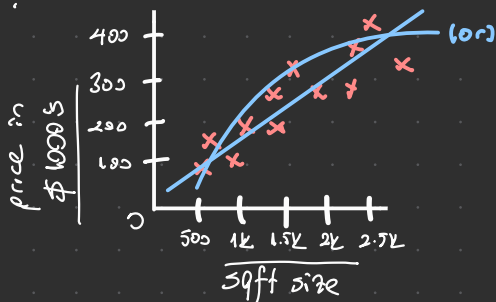
⇒ Supervised Learning: learns from data labeled with "right answers"



e.g.

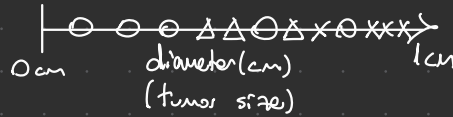
email	→	spam? (0/1)	⇒ spam filtering
audio	→	text transcripts	⇒ speech recognition
ad, user info	→	click? (0/1)	⇒ online advertising
image, radar info	→	pos. of other cars	⇒ self-driving car

→ Regression: Predict a number infinitely many possible outputs

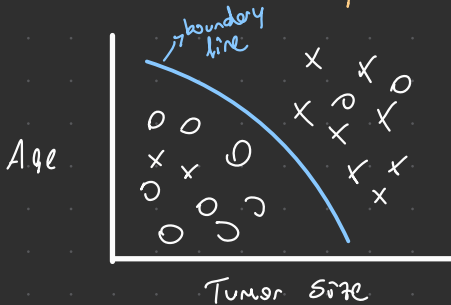


→ Classification: Predict categories, small number of possible outputs

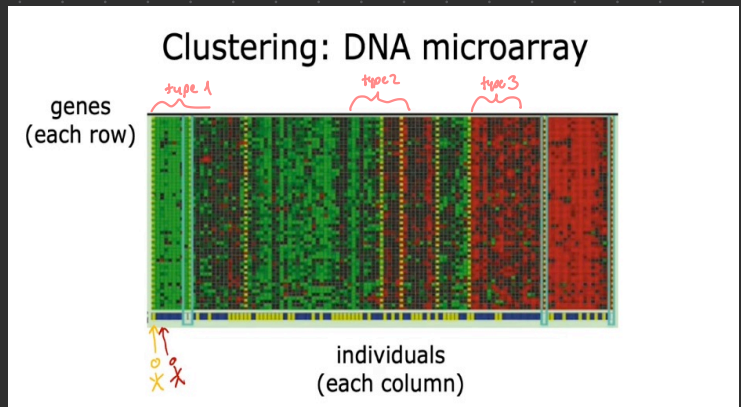
e.g. ○ benign
X malignant 1
△ malignant 2

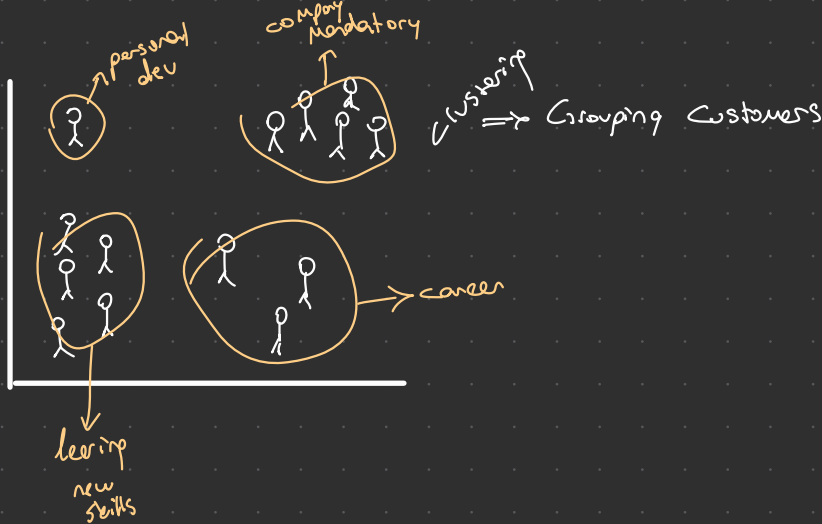


Two or more inputs



⇒ Unsupervised Learning: find something (pattern etc.) interesting in unlabeled data





→ Data only comes with input X , but not output labels Y
 ↳ Algorithm has to find structure in the data

Clustering
 Group similar data points together

Anomaly Detection
 Find unusual data points

Dimensionality Reduction
 Compress data using fewer numbers

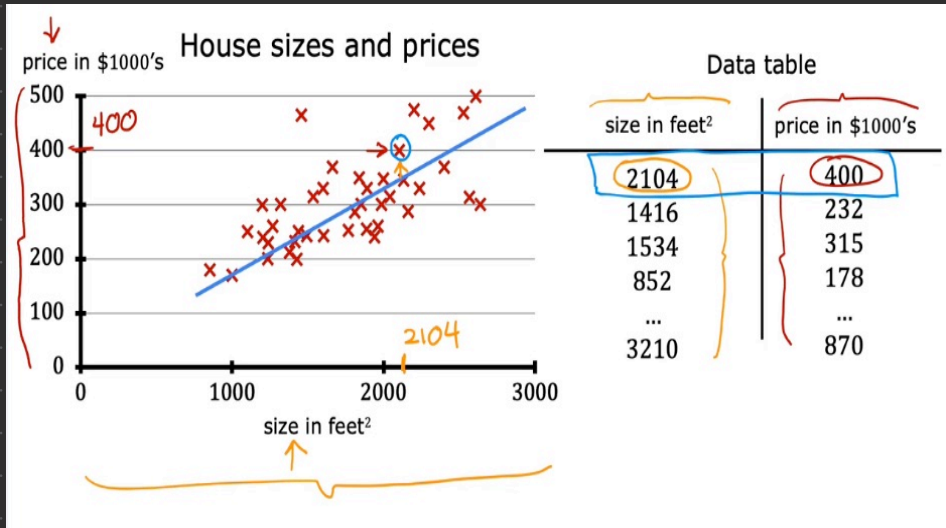
Question

Of the following examples, which would you address using an **unsupervised** learning algorithm?

- ☐ Given email labeled as spam/not spam, learn a spam filter.
- ☒ Given a set of news articles found on the web, group them into sets of articles about the same story.
- ☒ Given a database of customer data, automatically discover market segments and group customers into different market segments.
- ☐ Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not

Regression Model → predicts numbers

⇒ Linear Regression!



→ Terminology!

Training Set: Data used to train model

x	y
sqft size	price
(1) 2104	400
(2) 1416	232
(3) 1534	315
⋮	
(47) 3210	870

$M = 47$

NOTATION:

x = "input" variable
"feature"

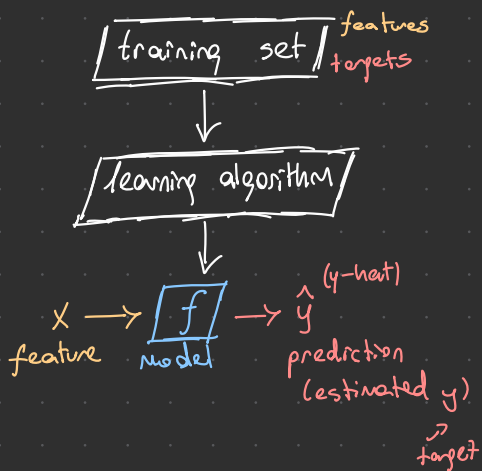
y = "output" variable
"target" variable

n = number of training examples

(x, y) = single training example

$(x^{(i)}, y^{(i)})$ = i th training example

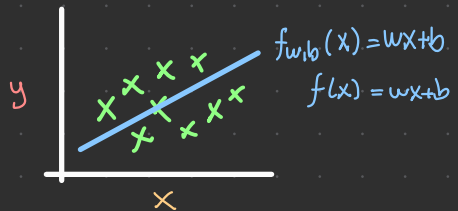
e.g.
 $(x^{(1)}, y^{(1)}) = (2104, 400)$



How to represent f ?

$$f_{w,b}(x) = wx + b$$

$$f(x)$$



linear regression with one variable
univariate linear regression
one variable

\Rightarrow cost function (squared error cost function)

$$J(w, b) = \frac{1}{2M} \sum_{i=1}^M (\hat{y}^{(i)} - y^{(i)})^2$$

error

$$= J(w, b) = \frac{1}{2M} \sum_{i=1}^M (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

→ COST FUNCTION intuition

model: $f_{w,b}(x) = wx + b$

parameters: w, b

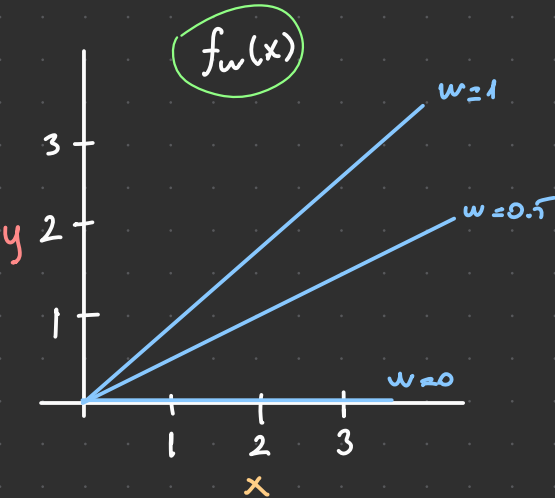
cost function: $\frac{1}{2n} \sum_{i=1}^M (f_{w,b}(x^{(i)}) - y^{(i)})^2$

goal: minimize $J(w,b)$
 w, b

↳ simplified

$f_w(x) = wx$ ($b = 0$)

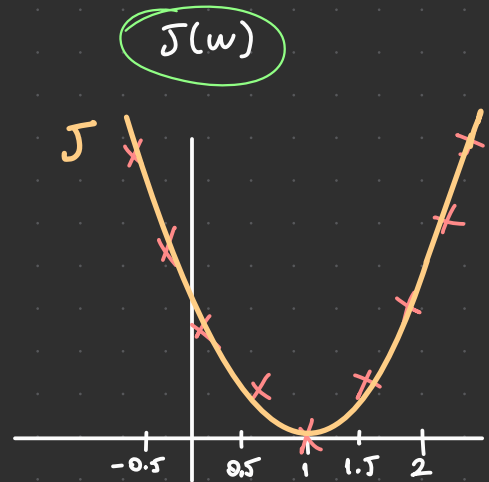
$J(w) = \frac{1}{2n} \sum_{i=1}^M (f_w(x^{(i)}) - y^{(i)})^2$, minimize $J(w)$
 w



$$J(0) = \frac{1}{6} (1^2 + 2^2 + 3^2) = \frac{1}{6} [14]$$

$$J(0.5) = \frac{1}{6} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2) = \frac{1}{6} [3.5]$$

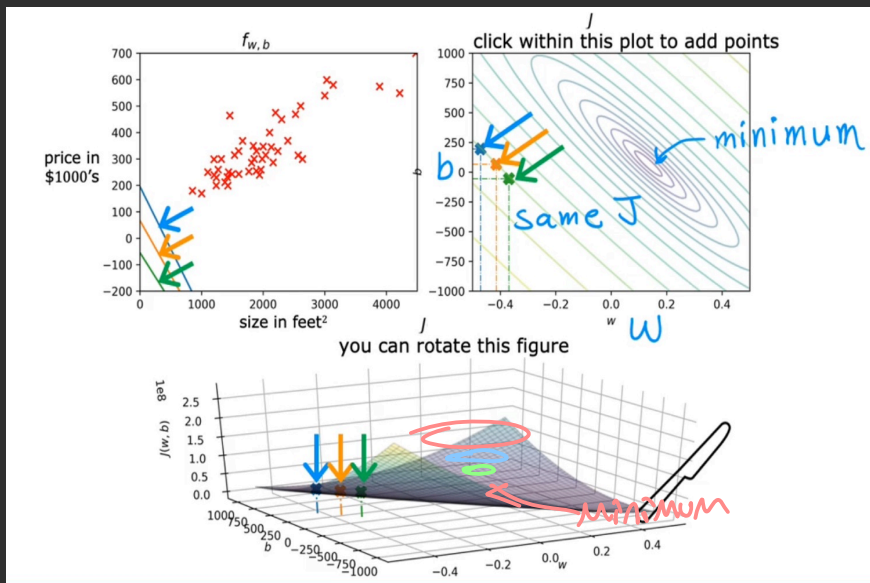
$$J(1) = 0$$



goal of linear regression: minimize $J(w)$
 w

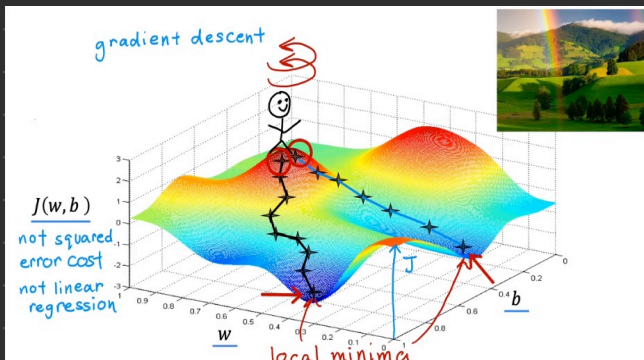
general case: minimize $J(w,b)$
 w, b

⇒ Visualizing
COST - FUNCTION



TRAIN THE MODEL WITH GRADIENT DESCENT

Gradient Descent Algorithm: an Optimization algorithm used to minimize a function by iteratively moving towards the steepest descent



using to fit a
model to holding data

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

assigning \downarrow
 learning rate (taking step) \downarrow
 (partial) derivative \rightarrow

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

Simultaneously update w and b

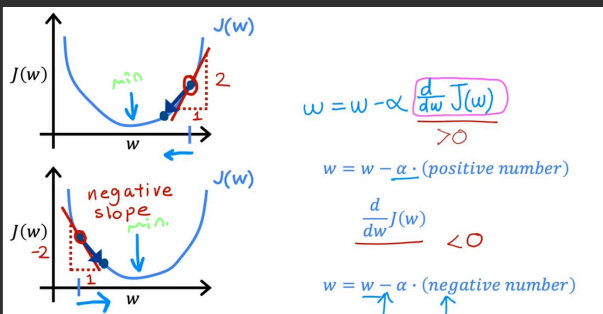
$$tmp_w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$tmp_b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

$$w = tmp_w$$

$$b = tmp_b$$

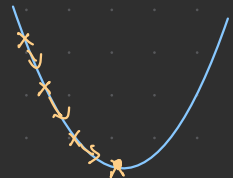
→ Gradient
Descent
Intuition



in both example, gradient descent doing something reasonable as it gets you closer to minimum

→ Learning Rate (α)

→ if learning rate is too small, gradient descent will work, but it will be slower
(think like taking a small step while walking)

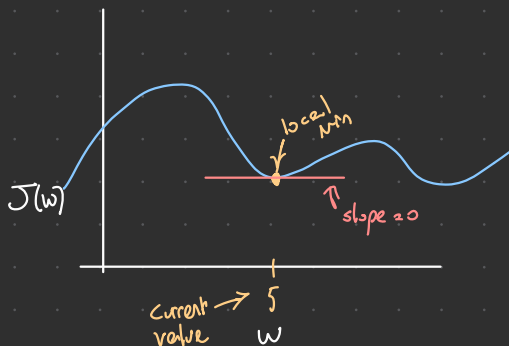


→ if α is too large

- gradient descent may never reach minimum (overshooting)
- fail to converge / diverge



Quick Ques: What happens if you've already reach / at minimum?



slope = 0 \Rightarrow derivative = 0

$$w = w - \alpha \left(\frac{\partial}{\partial w} J(w) \right)$$

$$= w = w$$

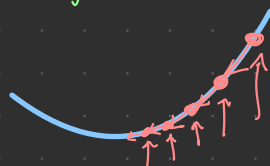
Can reach local minimum with fixed α ?

$$w = w - \alpha \cdot \frac{\partial}{\partial w} J(w)$$

as we get close to local minimum slope gets near zero

so " $\alpha \cdot \frac{\partial}{\partial w} J(w)$ " get smaller in each step

so we can reach min. even with fixed α



\leadsto Gradient Descent
for
Linear Regression

linear regression

$$f_{w,b}(x) = wx + b$$

Cost function

$$J(w,b) = \frac{1}{2N} \sum_{i=1}^N (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Gradient Descent Algo.

repeat until diverge {

$$w = w - \alpha \left(\frac{\partial}{\partial w} J(w,b) \right)$$

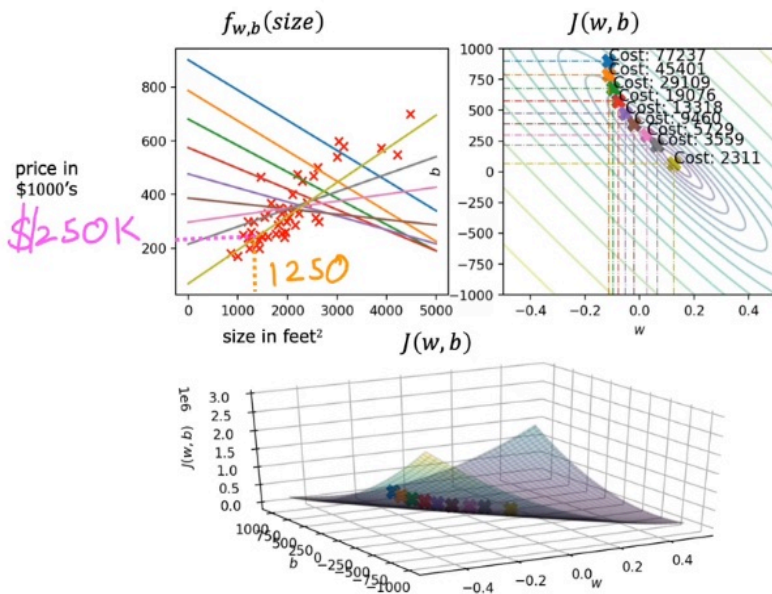
$$b = b - \alpha \left(\frac{\partial}{\partial b} J(w,b) \right)$$

}

$$\rightarrow \frac{1}{N} \sum_{i=1}^N (f_{w,b}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$\rightarrow \frac{1}{N} \sum_{i=1}^N (f_{w,b}(x^{(i)}) - y^{(i)})$$

→ Running Gradient Descent



"Batch" Gradient Descent

↓
Each step of gradient descent
uses all the training examples

$M=47$
prev. example

$$\rightarrow \sum_{i=1}^M (f_{w,b}(x^{(i)}) - y^{(i)})^2$$