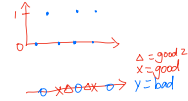Supervised learning
· given input and its corresponding output, after learning, they can take brand new Input to predict output

Input (x)       output (Y)
English    → spanish

regression - predict a number from infinitely many possible outputs


regression line

two or more Input

X X X ← good
age
Tumor size
bad

classification - predict on small outputs / it predicts categories


△ = good 2
X = good
Y = bad

---

unsupervised learning
· Find something interesting from unlabeled data
· Data comes with input x, but not output y & algorithm have to find structure in the data
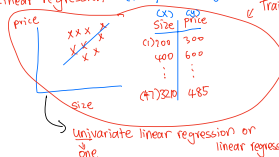

clustering - takes data without labels and automatically group them into clusters

Anomaly detection - Find unusual data points

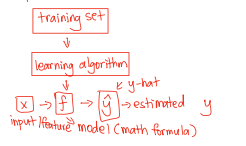Dimensionality reduction - compress data using fewer numbers

---

Linear regression - example of regression model
Training set: Data used to train the model



| (x) Size | (y) Price |
|----------|-----------|
| (1) 200  | 300       |
| 400      | 600       |
| ⋮        | ⋮         |
| (47) 3210| 485       |

univariate linear regression or one
linear regression with one variable
single input

Notation
X = input variable
y = output/target variable
m = number of training examples
(x,y) = single training examples
$(x^{(i)}, y^{(i)})$ = ith training example

supervised model

┌─────────────┐
│ training set │
└─────────────┘
      ↓
┌──────────────────┐
│ learning algorithm │
└──────────────────┘
      ↓
$x$ → $f$ → $\hat{y}$ → estimated  $y$
input /feature  model (math formula)   y-hat

---

cost function

Squared error cost function / difference

$J(w,b) = \frac{1}{2m} \sum_{i=1}^{m} (\hat{y}^{(i)} - y^{(i)})^2$
m = number of training set
goal: minimize J by changing w, b, the smaller J, the better

Imagine a model
$f_{w,b}(x) = wx + b$ or $f(x) = wx + b$

$w, b$ = parameters (variables you can change to improve algorithms)
coefficients
weights

Example:
$f_w(x) = wx$
$J(w) = \frac{1}{2m} \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2$
$= \frac{1}{2m} \sum_{i=1}^{m} (wx^{(i)} - y^{(i)})^2$

---

Gradient descent algorithm
· find the value of w, b that best fit the question or called to find the smallest value for J / $\min_{w,b} J(w,b)$

· $w = w - \alpha \frac{d}{dw}(J(w,b))$ → $\frac{1}{m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$
   learning rate

· $b = b - \alpha \frac{d}{db} J(w,b)$ → $\frac{1}{m} \sum_{i=1}^{m} (f_{w,b}(x^{(i)}) - y^{(i)})$

Just derivative it
$f_{w,b}(x^{(i)}) = wx^{(i)} + b$

· If $\alpha$ is too small, Gradient descent may be slow
· If $\alpha$ is too large, Gradient descent may · overshoot, never reach minimum
                                                · fail to converge, diverge

· You have to update w, b simultaneously
· repeat the w, b function until convergence

correct

$temp\_w = w - \alpha \frac{d}{dw} J(w,b)$
$temp\_b = b - \alpha \frac{d}{db} J(w,b)$

$w = tmp\_w$
$b = tmp\_b$

incorrect

$temp\_w = w - \alpha \frac{d}{dw} J(w,b)$
$w = temp\_w$

$temp\_b = b - \alpha \frac{d}{db} J(w,b)$
$b = temp\_b$

batch gradient descent
· Each step of gradient descent uses all the training examples