

What is machine learning?

A program learns from experience E with respect to tasks T and some performance measure P if its performance at tasks in T, as measured by P, improves with experience E.

} supervised learning : teach machine what it is and how to do.

} unsupervised learning : learn by itself.

supervised learning :

Given a data set and already know what our correct output should look like, having the idea that there is a relationship between input and output.

} regression : predict results within a continuous output, meaning that we are trying to map input variables to some continuous function.

} predict house price.

} predict human's age given an image.

} classification : predict results in a discrete output (discrete categories).

} given a patient with tumor, predict it is malignant or benign.

unsupervised learning :

allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables.

} clustering : k-means.

} non-clustering : "Cocktail Party Algorithm".

Training dataset :

m : Number of training examples

x 's : input variable / features. X : input space.

y 's : output variable / features. Y : output space.

(x, y) : single training example.

$(x^{(i)}, y^{(i)})$: refer to the i th training example.

We are trying to predict.

Training set.



Learning algorithm



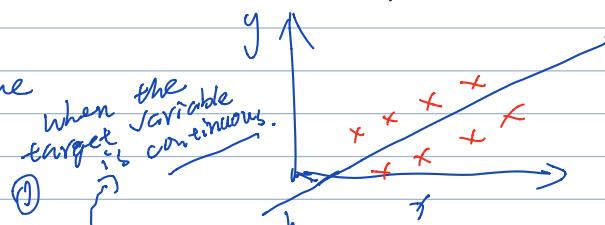
size of house $\rightarrow h \rightarrow$ Estimated price.
 x_1 hypothesis \rightarrow estimated value of y .

h maps from x 's to y 's.
 $h = x \rightarrow y$

How do we represent h ?

$$h(x) = \theta_0 + \theta_1 x. \text{ shorthand } h(x).$$

A simple example : h is a linear function.



linear regression with one variable.

or univariable linear regression.

a fancy way of saying one variable.

Q When the target variable is only

a small number of discrete values.

called classification problem.

Loss functions : $h_0(x) = \theta_0 + \theta_1 x$.

how to choose θ_0 & θ_1 .

Idea : choose θ_0, θ_1 so that $h_0(x)$ is close to y for our training examples (x, y) .

$$\text{minimize} \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

of training examples.

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

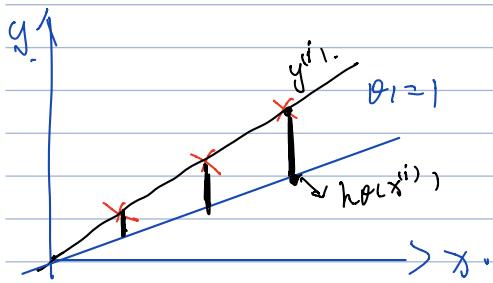
$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad J(\theta_0, \theta_1)$$

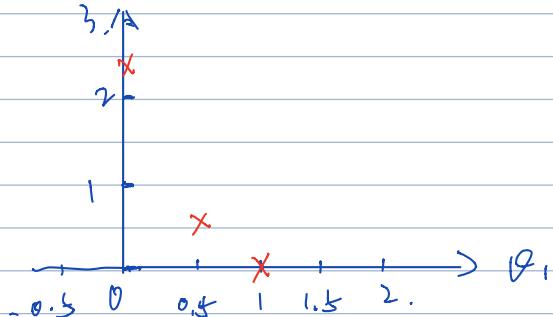
Cost Squared error function \hookrightarrow Loss Function

A simple example: $h_{\theta}(x) = \theta_1 x$.

for fixed θ_1 , h is the function of x .



$J(\theta_1)$:
a function of parameter θ_1 .



for $\theta_1 = 1$

$$\begin{aligned} J(\theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta_1}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 \\ &= \frac{1}{m} (1^2 + 2^2 + 3^2) = 0. \end{aligned}$$

for $\theta_1 = 0$. $h_{\theta_1}(x) = 0$.

$$J(\theta_1 = 0) = \frac{1}{m} [1^2 + 2^2 + 3^2] = \frac{14}{3}.$$

So, for $\theta_1 = 1$, $J(\theta_1) = 0$.

Objective: choose θ_1 to minimize the

for $\theta_1 = 0.5$:

$J(\theta_1)$.

$$\begin{aligned} J(\theta_1 = 0.5) &= \frac{1}{m} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \\ &= \frac{1}{m} \left(\frac{1}{4} + 1 + \frac{9}{4} \right) = \frac{1}{2}. \end{aligned}$$

find the parameters to fit the data well.

Contour plot / figure.: contains many contour lines of a two variable function has a constant value at all points of the same line.

When the value in the contour plot get closer to the centre thus reducing the cost function error.

Gradient descent.

Have a function $J(\theta_0, \theta_1)$, minimize $J(\theta_0, \theta_1)$

outline:

① start with some θ_0, θ_1 , initial.

② keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$.

algorithm:

repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (j = 0 \text{ and } j = 1). \quad \{$$

assignment learning rate. how big step (fast) updating parameters.

$a := b$ set b equal to a
no matter what is a .

Correct simultaneous update:

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp}_0$$

$$\theta_1 := \text{temp}_1$$

} partial derivative.
derivative

tangent weights

Incorrect:

$$\text{temp}_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

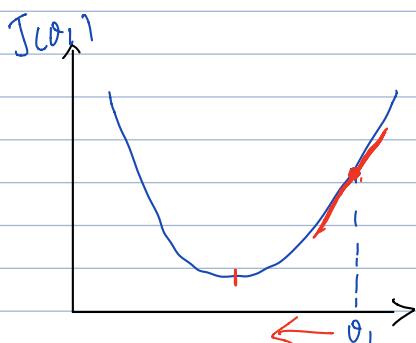
$$\theta_0 := \text{temp}_0$$

$$\text{temp}_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) \leftarrow$$

$$\theta_1 := \text{temp}_1$$

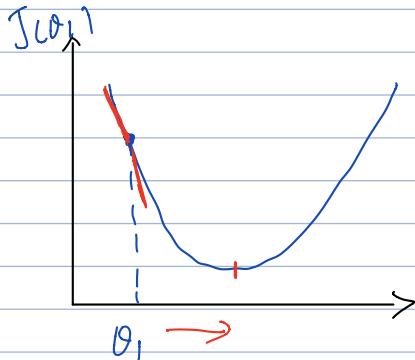
Simple examples.

$$\theta_1 := \theta_1 - \alpha \boxed{\frac{\partial}{\partial \theta_1} J(\theta_1)} \rightarrow \begin{array}{l} \text{the slope of} \\ \text{the tangent.} \end{array}$$



$$\theta_1 := \theta_1 - \alpha \text{ (slope/positive number)}.$$

$\theta_1 \downarrow$. decreasing.

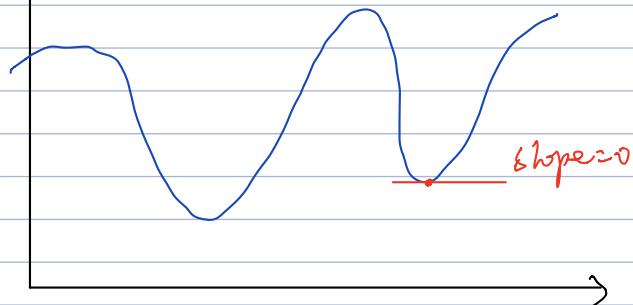
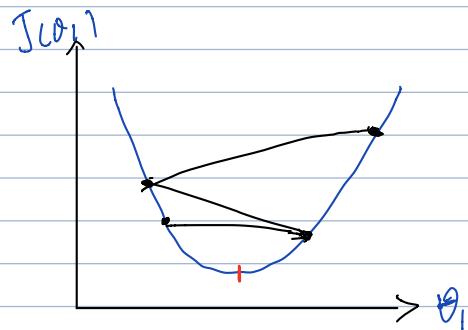
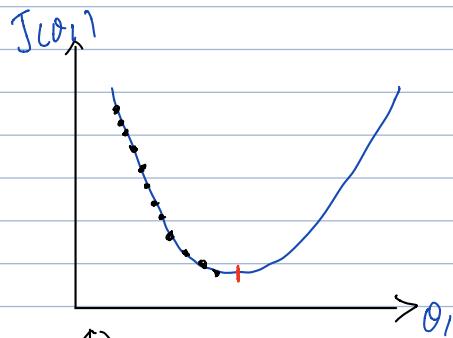


$$\theta_1 := \theta_1 - \alpha \text{ (negative slope)}$$

$\theta_1 \uparrow$ increasing.

learning rate α :

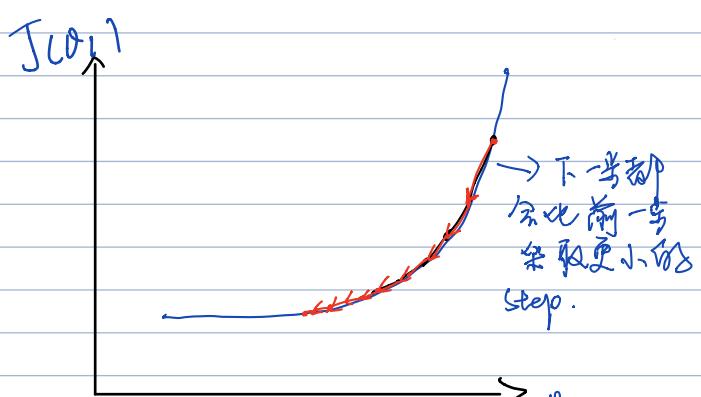
If α is too small, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



If the initial parameter is optimal,

leave it unchanged.

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



gradient descent can converge to a local minimum even with the fixed α .

As we approach a local minimum, gradient descent will automatically take smaller steps. So no need to decrease α over time.

Notice: 虽然是固定的 learning rate α , 但
下一步的 slope 总是比上一步更小, 更
贴近零, 所以每一次更新的距离也会减小。

Gradient descent for linear regression.

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

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

$$\hat{\theta}_0 := \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})$$

$$\hat{\theta}_1 := \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \cdot x^{(i)}$$

repeat until convergence {

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

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

"Batch" Gradient Descent

Each step of gradient descent uses all training samples.

Notice:

for the specific choice of

cost function $J(\theta_0, \theta_1)$

used in linear regression,

there are no local optima.

only one global optimum.

so the cost function always converges to the global optimum.