

Machine Learning - Week 1

1. ML: A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P . If its performance at tasks in T as measured by P , improves with experience E .

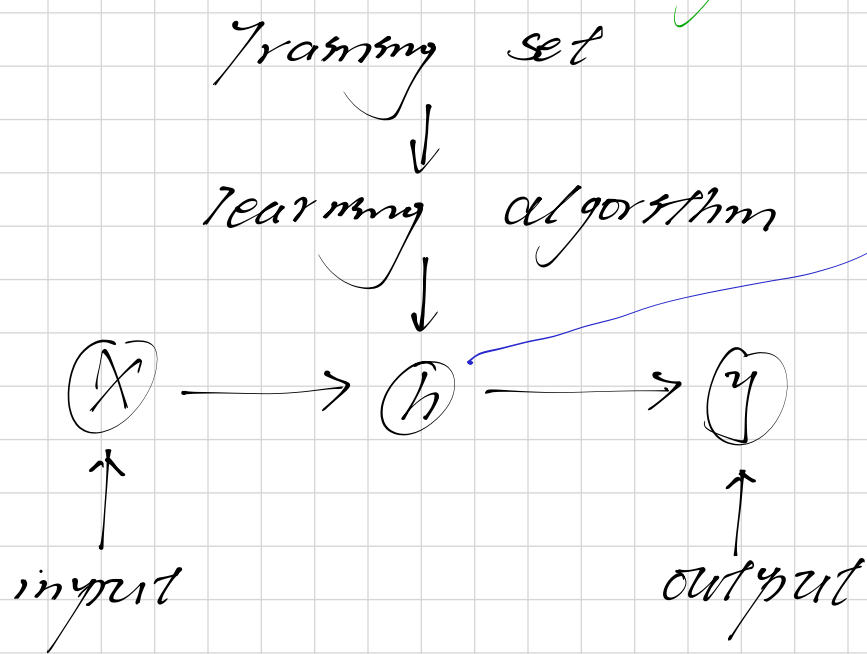
ML \rightarrow Supervised Learning

- ① regression: results within a continuous output
- ② classification: results within a discrete output

\rightarrow Unsupervised Learning

- ① no label
- ② structure of data

2.



hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$



$$h_{\theta}(x) = \sum_{i=0}^n \theta_i x_i = \theta^T x$$

Cost function: $J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$ $\xrightarrow{\text{training examples}}$

Goal: minimize $J(\theta)$

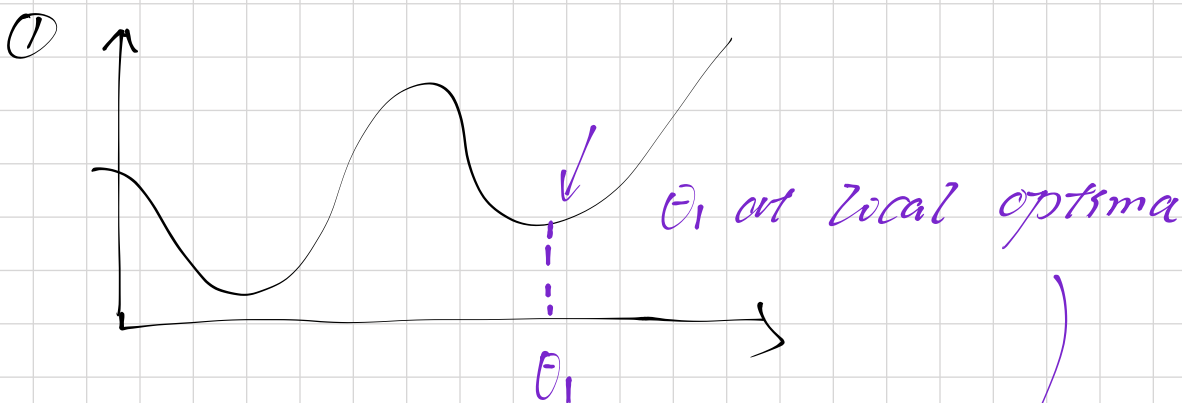
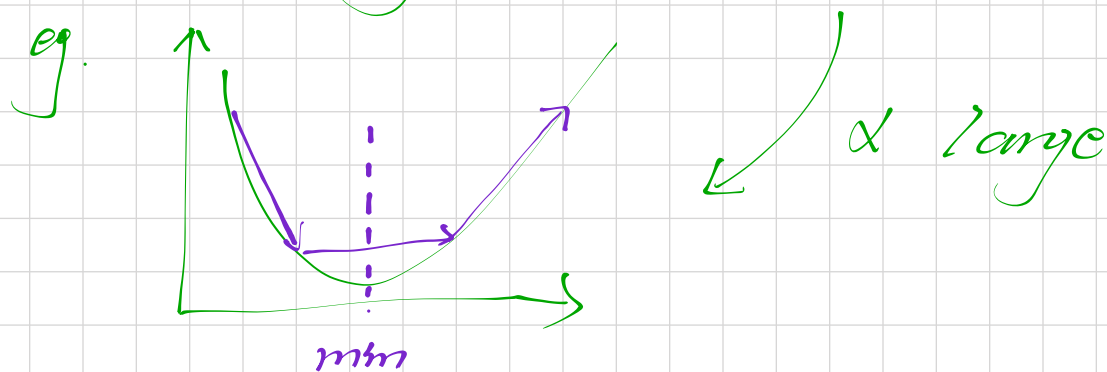
3. $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$

\downarrow

α : Learning rate \nearrow Small: gradient descent \searrow

\searrow Large: gradient descent \nearrow

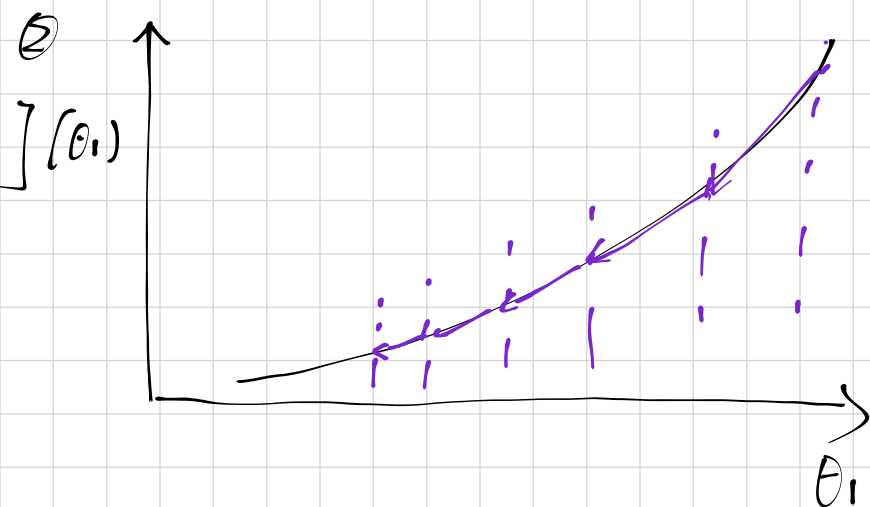
α large: may overshoot the min



$$\therefore \theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$\therefore \theta_1 \rightarrow \text{unchanged}$

gradient descent can converge to a local minimum, even with the learning rate α fixed



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

changed

\therefore approach a local minimum, gradient descent will automatically take smaller steps

\therefore don't need to decrease α

$$4. \begin{cases} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \\ J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \end{cases}$$

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

$$= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(X^{(i)}) - y^{(i)}) X_j$$

$$\Rightarrow \theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(X^{(i)})) X_j^{(i)}$$

① Update $\theta_0, \theta_1, \dots, \theta_j$ at the same time

② "Batch" gradient descent: use all training set