

Machine Learning 101

➤ Supervised Learning

"right answers" given

task: estimate more "right answers"

[Regression Problem: continuous value answer

[Classification Problem: discrete valued output

➤ Unsupervised Learning

no "right answer"

can find structure in data

→ clusters

➤ Linear Regression

Notation

- Training set \leadsto dataset for given problem
- $m \leadsto$ number of examples in training set
- x 's \leadsto "input" variables
- y 's \leadsto "output" / "target" variable
- $(x^{(i)}, y^{(i)}) \leadsto$ i -th training example
- hypothesis \leadsto function given by learning algorithm
 - \hookrightarrow maps from x 's to y 's

Cost Function

given: $h_{\theta}(x) = \theta_0 + \theta_1 x$

\hookrightarrow hypothesis $\quad \quad \quad \hookrightarrow$ input $\quad \quad \quad \hookrightarrow$ parameters

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

minimize this:
get θ values that "cost less"
∴ fit data best

➤ Gradient Descent

- * Can reach different local min depending on starting parameters
 - \hookrightarrow works best in cases in which there is only one min point (elliptic paraboloid)

Algorithm:

while !convergence \hookrightarrow learning rate

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

\hookrightarrow univariate case

} * Simultaneous ~~update~~ update ∴ in literal algorithm, next state must be calculated and stored into temp values all at once, then updated all at once or well

→ Learning Rate (α)

too small: gradient descent is too slow

too big: "big steps", may lead to overshooting and failure to converge

→ Gradient Descent for Linear Regression

model $\begin{cases} h_\theta(x) = \theta_0 + \theta_1 x \\ J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \end{cases}$ → apply GD to get lowest cost param

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \rightarrow j=0 \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$
$$\rightarrow j=1 \rightarrow \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x^{(i)} \end{cases}$$

* "Batch" Gradient Descent:
each iteration goes through
the entire training set

Week 2

→ Multivariate Linear Regression

X_n feature contribute to output value Y

Notation: $x^{(i)}$ is vector containing all features in i -th example

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

$$\hookrightarrow \theta_0 = 1 \rightarrow h_\theta(x) = \sum_{i=0}^n x_i \theta_i = \theta^T x$$

vector with all θ_i vector with all x_i

Gradient Descent:

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

→ Feature Scaling

Features being on a similar scale can make gradient descent work faster

get every feature in range $[-1, 1]$ → divide by maximum value

Mean Normalization: $x_i = \frac{x_i - \text{average value}}{\text{maximum range}}$

→ Convergence method

establish a threshold ϵ and if $J(\theta)$ decreases by less than ϵ in one iteration it has converged * hard to determine ϵ

plot no. iterations $\times J(\theta)$ → helpful when choosing α

* It's possible to create new features that enable the fitting of polynomial functions to the data