

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

\hookrightarrow input

\hookrightarrow parameters

$$\hookrightarrow J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad \left. \vphantom{\sum} \right\} \begin{array}{l} \text{minimize this:} \\ \text{get } \theta \text{ values that "cost less"} \\ \therefore \text{fit data best} \end{array}$$

➤ 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

$$\theta_j := \theta_j - \alpha \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \quad \hookrightarrow \text{univariate case}$$

} * Simultaneous ~~update~~ update \circ 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

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

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

* "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

Normal Equations

- Solver for optimum value directly
- Take derivative and equal to zero
- Defining X as a matrix containing all features and y as the solution vector
- $\theta = (X^T X)^{-1} X^T y$
- No need to choose α or to iterate
- Does not scale well to a large number of features

Week 3

Logistic Regression

- Used in classification problems \rightarrow Linear regression is ineffective
- Model (Binary case)
 - $0 \leq h_\theta(x) \leq 1$
 - Logistic function: $g(z) = \frac{1}{1+e^{-z}} \rightarrow h_\theta(x) = g(\theta^T x)$
 - h_θ will output the probability of the correct prediction being 1
- Decision Boundary
 - if $\theta^T x \geq 0 \rightarrow y=1$ else $y=0$
 - $\theta^T x$ can represent all kinds of polynomials by adding extra features
- Cost
 - Original cost function is not convex in logit regression \rightarrow gradient descent wouldn't work
 - Defining:
$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$
 - $\rightarrow \text{case } y=1 = -\log(h_\theta(x))$
 - $\rightarrow \text{case } y=0 = -\log(1-h_\theta(x))$
 - In LR, we have:
 - cost zero when prediction correct
 - cost blows up when prediction incorrect
 - Simplified: $\text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$
- Gradient Descent
 - Same update rule:
$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$
 - Vectorized:
$$\theta := \theta - \frac{\alpha}{m} X^T (y - \vec{y})$$
- Advanced Optimization Algorithms
 - Conjugate Gradient
 - BFGS
 - L-BFGS
 - provided by language $\neq \text{fminunc()}$
- Multiclass Classification
 - One-vs-all: find θ_i for each class and predict whichever maximizes the result