

Notes for Machine Learning

Jack Wang

April 11, 2015

Contents

| | | |
|----------|---|----------|
| 1 | Introduction (1) | 1 |
| 1.1 | Definition of Machine Learning (ML) | 1 |
| 1.2 | Types of learning | 1 |
| 1.3 | Todo | 2 |
| 2 | Linear Regression (2,3,4,5) | 2 |
| 2.1 | Regression | 2 |
| 2.2 | Model and Cost function | 3 |
| 2.3 | Parameter learning | 3 |
| 2.4 | Demo of Gradient Descent | 3 |
| 2.4.1 | $J(\text{# of iterations})$ | 3 |
| 2.4.2 | $J(\theta_0, \theta_1)$ in 3D | 3 |
| 2.4.3 | $J(\theta_0, \theta_1)$ in contour | 3 |

1 Introduction (1)

1.1 Definition of Machine Learning (ML)

TEP A computer program is said to learn from experience **E** with respect to some task **T** and some performance measure **P** if its performance on **T**, as measured by **P**, improves with experience E.

Machine learning the field of study that gives computers the ability to learn without being explicitly programmed.

1.2 Types of learning

Supervised learning Given {Features, Label} to estimate the label for new data.

- regression (target variable is continuous)
- classification (target variable is discrete)

Un-supervised learning Given {Features} to do something, e.g, clustering. **No labels are given.**

- clustering (image processing to 3D, separate sounds from two person(SVD) ¹⁾)

¹Is it necessary to have two speakers? What if only one speaker and two person?

Reinforcement learning inspired by behaviorist psychology, concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward

- robot

1.3 Todo

How to seprate sound from two person?

2 Linear Regression (2,3,4,5)

2.1 Regression

| Notation | Meaning |
|--|--|
| $\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{(n+1) \times 1}$ | input variables, features |
| y $\hat{y} = h(\mathbf{x}) = h_{\boldsymbol{\theta}}(\mathbf{x})$ | output/target variable hypothesis function of \mathbf{x} with parameter $\boldsymbol{\theta}$ |
| $\mathbf{x}^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix}$ | i^{th} observation of \mathbf{x} where $1 \leq i \leq m$ |
| $y^{(i)}$ | i^{th} observation of y where $1 \leq i \leq m$ |
| $\mathbf{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times 1}$ | m observations of y |
| $\mathbf{X} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(m)})^T \end{bmatrix} \in \mathbb{R}^{m \times (n+1)}$ | m observations of \mathbf{x} |
| m | # of traing examples |
| n | # of traing features (excluding addtional 1 vector) |
| $\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ | parameter of hypothesis funtion(model) |

Table 1: Notation of Math symbols

TEP defintion of regression:

- E: Given m observations: (\mathbf{X}, \mathbf{y}) ,

- T: Predict y using model $h(\mathbf{x}) = h_{\boldsymbol{\theta}}(\mathbf{x})$,
- P: which minimizes loss(error) function $J(\boldsymbol{\theta})$

After selecting the model (e.g, linear model), T becomes to find $\arg \min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ where:

$$J(\boldsymbol{\theta}) = \sum_{i=1}^m (h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)})^2 \quad (\text{SSE}) \quad (1)$$

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

$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \cdots & x_n^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \cdots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & x_2^{(1)} & \cdots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & x_2^{(2)} & \cdots & x_n^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_0^{(m)} & x_1^{(m)} & x_2^{(m)} & \cdots & x_n^{(m)} \end{bmatrix} = \begin{bmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(m)})^T \end{bmatrix} \in \mathbb{R}^{m \times (n+1)} \quad (3)$$

- 解析解
- Iteration method Gradient decent method
- Regression models:
 - linear model $h_{\boldsymbol{\theta}}(\mathbf{x}) = \langle \boldsymbol{\theta}, \mathbf{x} \rangle = \mathbf{x}^T \boldsymbol{\theta}$
 - poly model

2.2 Model and Cost function

2.3 Parameter learning

Training Set \rightarrow Learning Algorithm (Chosen) $\rightarrow h_{\boldsymbol{\theta}}(x) = h(x)$

2.4 Demo of Gradient Decent

All figures 1, 2, 3, 4, 5 and 6 are generated by `../src/demo_gradient_descent.py`

We can find:

- No scaling, Gradient Decent 在不同维度上的陡峭程度不一样 (二维中轮廓图由圆变为椭圆, 等高线越密集, 越陡峭), 越陡峭, 下降降幅度越厉害, 梯度下降法给予该陡峭维度的权重比较大, 所以会沿着陡峭维度的主方向下降, 之后是第二个陡峭的维度成为主方向, 以此类推。与圆形轮廓图相比, 误入“歧图”, 多乖了个“弯儿”
- Gradient Decent can only find the **local minimum** point
- The step size during update of θ will reduce automatically while approaching non-steep surface, even for the same **learning rate**

2.4.1 $J(\# \text{ of iterations})$

2.4.2 $J(\theta_0, \theta_1)$ in 3D

2.4.3 $J(\theta_0, \theta_1)$ in contour

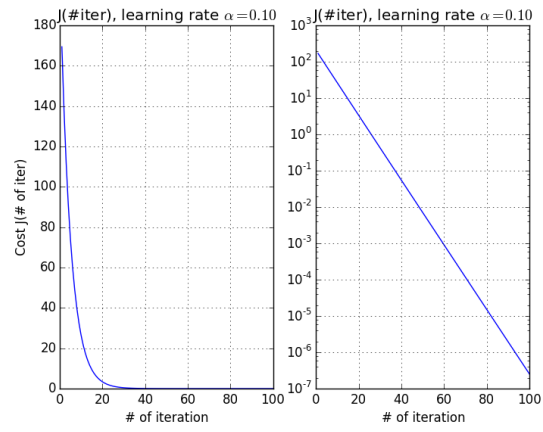


Figure 1: $J(\# \text{ of iterations})$

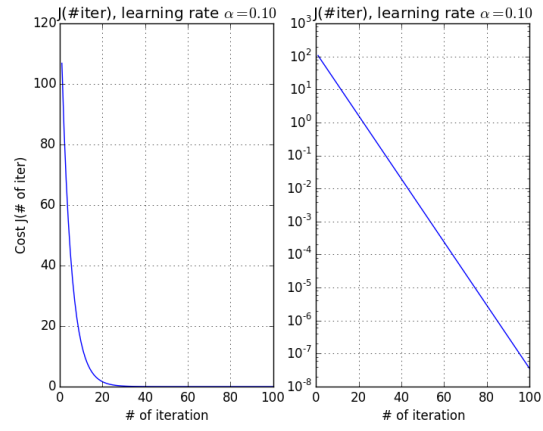


Figure 2: $J(\# \text{ of iterations})$ without scaling

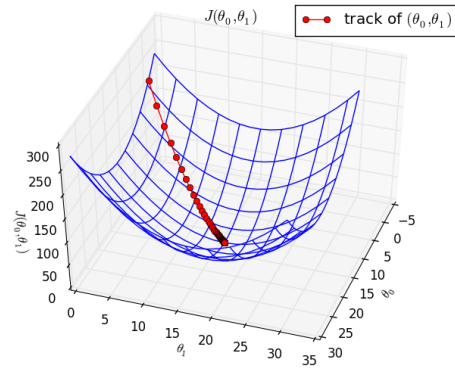


Figure 3: $J(\theta_0, \theta_1)$ in 3D

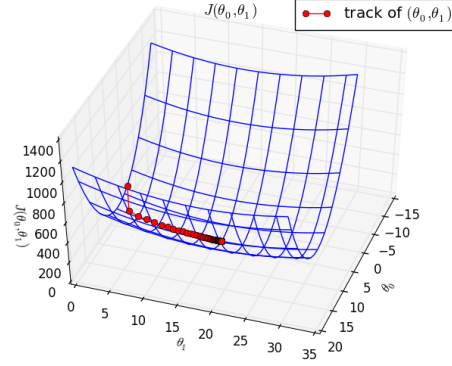


Figure 4: $J(\theta_0, \theta_1)$ in 3D without scaling

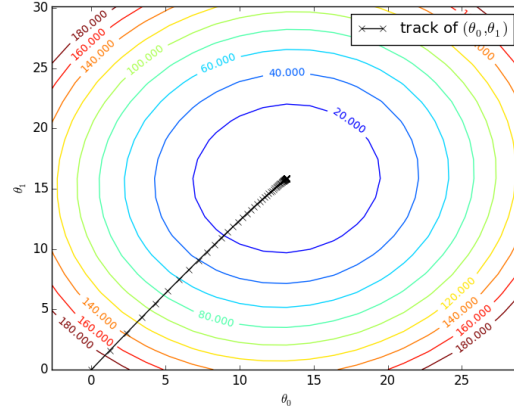


Figure 5: $J(\theta_0, \theta_1)$ in contour

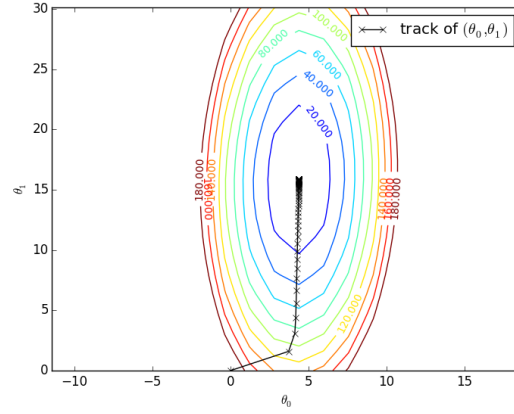


Figure 6: $J(\theta_0, \theta_1)$ in contour without scaling