

2.2.1 $h_{\theta}(x_i) = \theta^T x_i$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \theta^T \theta = \frac{1}{m} \sum_{i=1}^m (\theta^T x_i - y_i)^2 + \lambda \theta^T \theta = \frac{1}{m} (X\theta - y)^T (X\theta - y) + \lambda \theta^T \theta = \frac{1}{m} \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2$$

2.2.3 $\nabla_{\theta} J(\theta) = \frac{2}{m} X^T (X\theta - y) + 2\lambda \theta$

2.3.1 - 泰勒展开

$$J(\theta + \eta h) = J(\theta) + (\eta h)^T \nabla_{\theta} J(\theta) + \dots$$

$h^T \nabla_{\theta} J(\theta)$ 只有 h 与 $\nabla_{\theta} J(\theta)$ 方向相同时取得最大值

2.3.2 $\theta^{t+1} = \theta^t - \eta \nabla_{\theta} J(\theta) |_{\theta=\theta^t}$

2.3.4 步长为 0.01 时收敛最快，步长为 0.5, 0.1 时会发散

2.4.1 $X_k = \begin{pmatrix} X_{i1} & X_{i2} & \dots & X_{ind+1} \\ X_{i2} & X_{i2} & \dots & X_{ind+1} \\ \vdots & \vdots & \ddots & \vdots \\ X_{in} & X_{in} & \dots & X_{ind+1} \end{pmatrix} \quad y_k = \begin{pmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{in} \end{pmatrix}$

$$J_{SGD}(\theta) = \frac{1}{n} \sum_{k=1}^n (h_{\theta}(x_{ik}) - y_{ik})^2 + \lambda \theta^T \theta = \frac{1}{n} (X_k \theta - y_k)^T (X_k \theta - y_k) + \lambda \theta^T \theta$$

$$\nabla_{\theta} J_{SGD}(\theta) = \frac{2}{n} X_k^T (X_k \theta - y_k) + 2\lambda \theta$$

2.4.3 批增大，损失函数呈收敛趋势

2.4.5 发现在所取数值中，正则化系数越大越稀疏

3.2 发现 distance 对房价有较大负影响，latitude 对房价有较大正影响，当前将时间运用特征工程加入了模型，且分为年月日，可能会与房屋年龄产生关联，同时也引入过多时间波动，弱化其他指标作用。可以减少时间变量重新估计模型。

把不同变量都进行特征归一化，但变量间方差不同，可能进行标准化会更好。
(最大最小) (正态)

4.1.4 $h_{\theta}(x) = P(y=1|x;\theta) = \frac{1}{1 + e^{-\theta^T x}}$

$$J(\theta) = [y \log(1 + e^{-\theta^T x}) + (1-y) \log(\frac{1 + e^{-\theta^T x}}{e^{-\theta^T x}})]$$

$$\nabla_{\theta} J(\theta) = y \frac{-x e^{-\theta^T x}}{1 + e^{-\theta^T x}} - (1-y) \frac{-x e^{-\theta^T x}}{e^{-\theta^T x} + 1} + (1-y) \frac{-x e^{-\theta^T x}}{1 + e^{-\theta^T x}} = -x e^{-\theta^T x} h_{\theta}(x) + (1-y)x = -x(e^{-\theta^T x} + 1) h_{\theta}(x) + x h_{\theta}(x) + (1-y)x = (h_{\theta}(x) - y)x$$

4.1.5. 发现手动得出的结果与调用算法包得到的结果一样

4.2.2 选取了花瓣长度与花瓣宽度

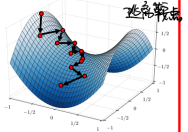
4.2.3 引入基函数 决策边界曲线且拟合效果变好

Stochastic Gradient Descent (SGD)

- In each iteration t ($\leq T$):
 - Randomly sample a minibatch of $m \ll n$ points $\{(x_i, y_i)\}_{i=1}^m$
 - Set $J^t(\mathbf{w}^t) = \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{w}^t; x_i, y_i)$
 - Compute gradient on minibatch:

$$\Delta^t = \nabla_{\mathbf{w}} J^t(\mathbf{w}^t)$$
 - Update parameters with learning rate η :

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \Delta^t$$



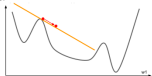
- For linear regression: $\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - 2\eta \mathbf{X}_m^T (\mathbf{X}_m \mathbf{w}^t - \mathbf{y})$

- Gradient shows direction that function varies fastest.

Same dimension with parameter \mathbf{w} 's dimension

$$\mathbf{g} = \nabla_{\mathbf{w}} J(\mathbf{w})$$

$$\mathbf{g}_j = \nabla_{w_j} J(\mathbf{w})$$



- First-order Taylor approximation of the objective function:

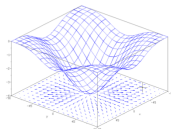
$$J(\mathbf{w}) \approx J(\mathbf{w}_0) + (\mathbf{w} - \mathbf{w}_0)^T \mathbf{g} + \dots$$

- Go along gradient for a step with a small rate η :

$$J(\mathbf{w} - \eta \mathbf{g}) \approx J(\mathbf{w}) - \eta \mathbf{g}^T \mathbf{g}$$

- Reach a point with smaller loss.

$$-\eta \mathbf{g}^T \mathbf{g} \leq 0$$



- Repeat this step, and we get the Gradient Descent (GD) algorithm.

- Compute the gradient and set it to zero

$$\nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w}) = 2\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\lambda \mathbf{w} = 0$$

$$\Rightarrow (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y}$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

有唯一解

λ 小 \Rightarrow 不稳定

- Or using Gradient Descent (GD):

$$\mathbf{w}^{t+1} \leftarrow \mathbf{w}^t - \eta \nabla_{\mathbf{w}} \hat{\epsilon}(\mathbf{w})|_{\mathbf{w}=\mathbf{w}^t}$$

for $\lambda > 0$

$d \times d$ -dim identity matrix