

softmax 回归模型.

模型形式为 $h_{\theta}(x^{(i)}) = \begin{bmatrix} P(y^{(i)}=1 | x^{(i)}; \theta) \\ P(y^{(i)}=2 | x^{(i)}; \theta) \\ \vdots \\ P(y^{(i)}=k | x^{(i)}; \theta) \end{bmatrix} = \frac{1}{\sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \begin{bmatrix} e^{\theta_1^T x^{(i)}} \\ e^{\theta_2^T x^{(i)}} \\ \vdots \\ e^{\theta_k^T x^{(i)}} \end{bmatrix}$

代价函数为

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=1}^k \mathbf{1}\{y^{(i)}=j\} \cdot \log \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}}} \right]$$

直观语义就是使正确输出单元的激活值 $h_{\theta}(x^{(i)})_{y^{(i)}}$ 尽可能大
跟 LR 相似; LR 本质上就是二分类的 softmax. 详见 UFLDL.

求梯度. 注意 θ_j 是向量.

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{\partial -\frac{1}{m} \sum_{i=1}^m \log \frac{e^{\theta_{y^{(i)}}^T x^{(i)}}}{\sum}}{\partial \theta_j}$$

$$= -\frac{1}{m} \cdot \frac{\partial \sum_{i=1}^m \theta_{y^{(i)}}^T x^{(i)} - \log \sum_{l=1}^k e^{\theta_l^T x^{(i)}}}{\partial \theta_j}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[\mathbf{1}\{y^{(i)}=j\} \cdot x^{(i)} - \frac{1}{\sum} \cdot e^{\theta_j^T x^{(i)}} \cdot x^{(i)} \right]$$

$$= -\frac{1}{m} \sum_{i=1}^m x^{(i)} \left(\mathbf{1}\{y^{(i)}=j\} - \frac{e^{\theta_j^T x^{(i)}}}{\sum_{l=1}^k e^{\theta_l^T x^{(i)}}} \right)$$

$$= -\frac{1}{m} \sum_{i=1}^m x^{(i)} \left(\mathbf{1}\{y^{(i)}=j\} - h_{\theta}(x^{(i)}) \right)$$

参考 UFLDL. softmax 是个过度参数化的模型,
最优解不唯一. 可以通过 Weight Decay 解决

$$J(\theta) \triangleq -\frac{1}{m} \left[\sum_{i=1}^m \sum_{j=1}^k 1\{y^{(i)}=j\} \cdot \log \frac{e^{\theta_j^T x^{(i)}}}{\sum_l e^{\theta_l^T x^{(i)}}} \right] + \frac{\lambda}{2} \sum_{i=1}^m \sum_{j=1}^k \theta_{ij}^2$$

引入罚项后 $J(\theta)$ 变成严格凸函数, 有唯一-最优解. 且 Hessian 可逆.

$$\nabla_{\theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[x^{(i)} (1\{y^{(i)}=j\} - h_{\theta}(x^{(i)})) \right] + \lambda \theta_j$$

代码如下:

```
class Softmax:
```

```
    # --init-- (self, k=10) 省略
```

```
    def fit(self, x, y, n-epoch=200,  $\lambda=0.01$ ,  $\eta=1e-4$ ):
```

```
        m = x.shape[0] 样本个数
```

```
        self. $\theta$  = zeros((self.k, x.shape[1]))
```

```
        y = one-hot(y, range(self.k))
```

```
        for _ in range(n-epoch): 注意 n-epoch 不能用科学计数法, float
```

```
            h = self.predict_proba(x)
```

```
            grad = (h - y.T) @ x / m +  $\lambda$  * self. $\theta$ 
```

```
            self. $\theta$  -=  $\eta$  * grad
```

→ 建议琢磨一下向量化的运算细节.
注意公式中的 $-\frac{1}{m}$. 负号和除 m 缺失

```
    def predict_proba(self, x): 不可. 否则梯度爆炸. 出现 inf.
```

```
        return l1-norm(exp(self. $\theta$  @ x.T), axis=0)
```

第 i 行 j 列表示第 j 个样本取值 i 的概率.

```
    def predict(self, x):
```

```
        return argmax(self.predict_proba(x), axis=0)
```

以上代码.

koxo922@qq.com

2018.04.20