

Classification

$$y = \begin{cases} 0 & : \text{Negative class} \\ 1 & : \text{Positive class} \end{cases}$$

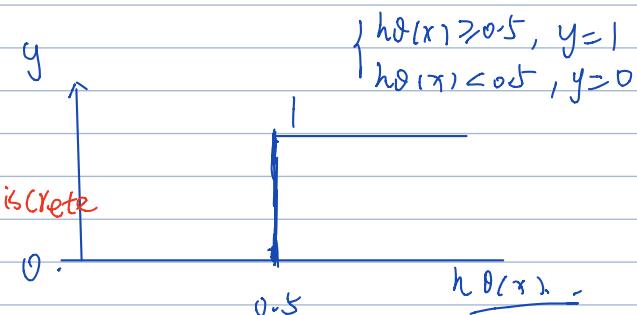
Because sometimes classification is not a linear function, we cannot use linear regression to attempt classification.

multiclass classification problem:

$$y = \{0, 1, 2, 3\}$$

实际是一种分类算法
适用于离散情况 discrete

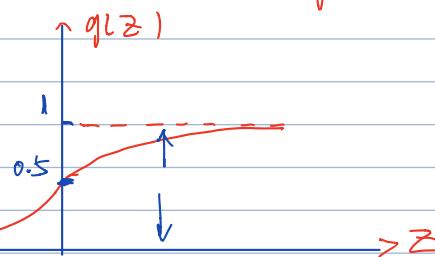
Logistic Regression Model (逻辑回归)



$$\text{Want } 0 \leq h(x) \leq 1$$

$$\left. \begin{array}{l} h(x) = g(z) \\ g(z) = \frac{1}{1+e^{-z}} \\ z = \theta^T x \end{array} \right\} \Rightarrow h(x) = \frac{1}{1+e^{-\theta^T x}}$$

Sigmoid Function (S型)
Logistic Function



$g(z)$ maps any real number into the $(0, 1)$ interval

$h(x) = g(\theta^T x)$
X: 特征向量
g: 逻辑函数

Interpretation of Hypothesis Output

$h(x)$ = estimated probability that $y=1$ on input x

根据选择概率数计算出变量二值的可能性, $h(x) = p(y=1 | x, \theta)$

$$\text{eg. If } x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorsize} \end{bmatrix}$$

that patient with feature x , the probability that $y=1$ is 0.7.

Tell patient that tumor 70% chance of tumor being malignant.

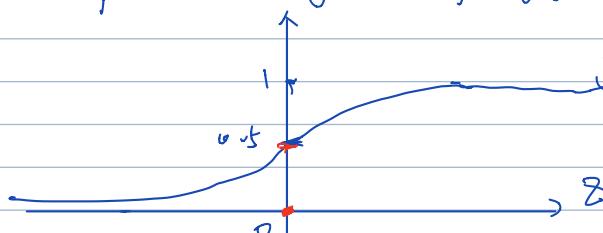
$h(x) = p(y=1 | x; \theta)$ "probability that $y=1$, given x , parameterized by θ .

$$p(y=0 | x; \theta) = 1 - p(y=1 | x; \theta)$$

Decision Boundary: the line separate the area where $y=0$ and where $y=1$.
it is created by hypothesis function.

suppose predict: $y=1$ if $h(x) \geq 0.5$.

predict: $y=0$ if $h(x) < 0.5$. when $g(z) \geq 0.5$, $y=1$

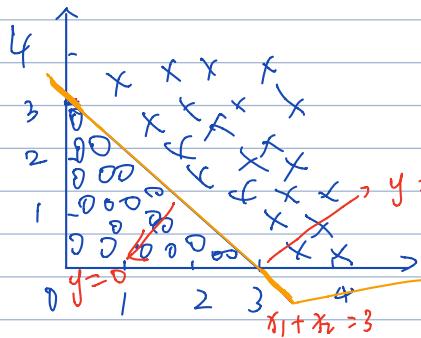


we know that:

$$h(x) = g(\theta^T x) \geq 0.5$$

$$\theta^T x \geq 0$$

$$\therefore \theta^T x \geq 0$$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) = g(-3 + \theta_1 + \theta_2)$$

$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

Remember:

$$\left\{ \begin{array}{l} g(z) = \frac{1}{1+e^{-z}} \\ z = \theta^T x \end{array} \right.$$

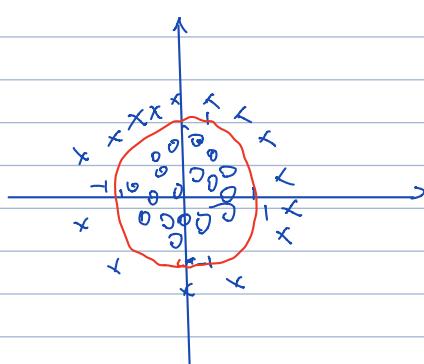
$$\Rightarrow z = 0, g(z) = \frac{1}{2}$$

$$z \rightarrow \infty, e^{-z} \rightarrow 0, g(z) = 1$$

$$z \rightarrow -\infty, e^{-z} \rightarrow \infty, g(z) = 0.$$

Non-linear Decision Boundaries

polynomial regression =



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

assume $\theta = \begin{bmatrix} -1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$

$$y=1, -1 + x_1^2 + x_2^2 \geq 0.$$

$$x_1^2 + x_2^2 \geq 1$$

Notice: higher order of polynomial, the more complex Decision Boundary.

How to fit the parameters = 定义用来拟合参数的优化目标或者代价函数.

supervised model =

m examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$.

$n+1$ features:

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}, x_0 = 1, y \in \{0, 1\}$$

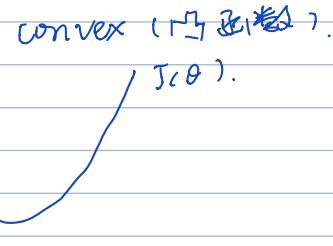
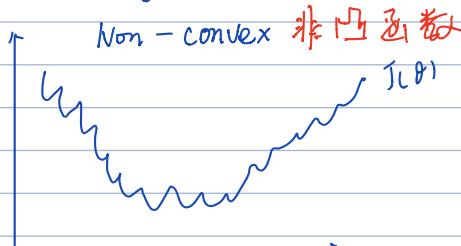
$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$\text{Linear Regression} = J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} \underbrace{(h_{\theta}(x^{(i)}))}_{\text{cost}(h_{\theta}(x^{(i)}), y^{(i)})}^2$$

→ cost function

$$\text{cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)})) - y^{(i)}^2$$

cause
many
local
optimum

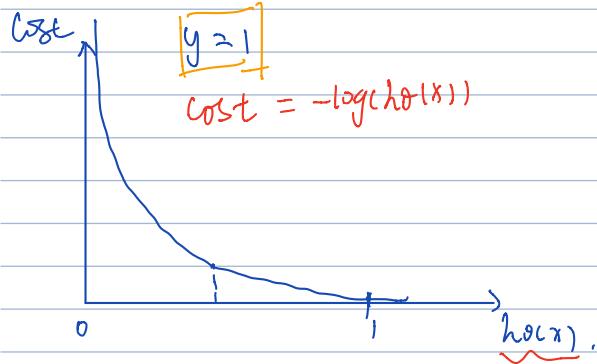


Because $h\theta(x)$ in classification problem is nonlinear ($h\theta(x) = \frac{1}{1+e^{-\theta^T x}}$) ;

$J(\theta)$ is non-convex; if use gradient descent, there is no global optimization.

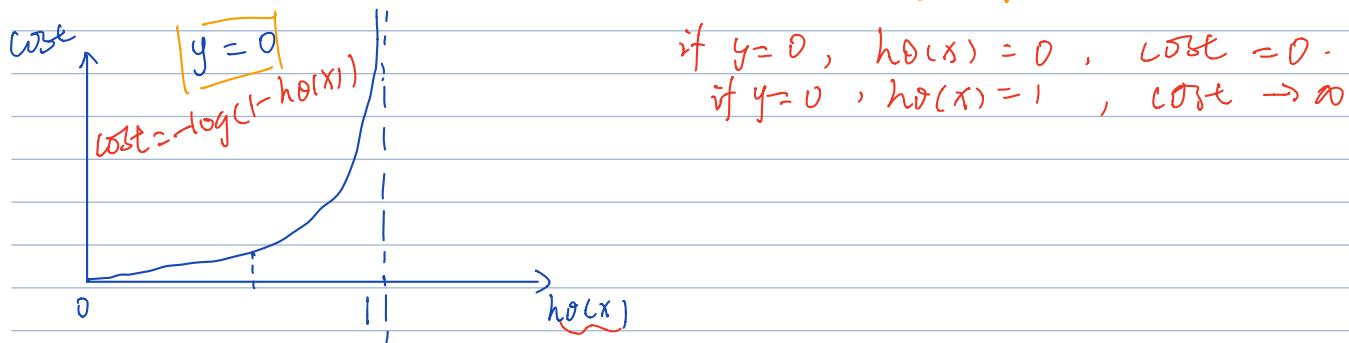
Logistic regression cost function: $J\theta = \frac{1}{m} \sum \text{Cost}(h\theta(x), y)$

$$\text{Cost}(h\theta(x), y) = \begin{cases} -\log(h\theta(x)), & \text{if } y=1 \\ -\log(1-h\theta(x)), & \text{if } y=0 \end{cases} \quad h\theta(x) \in [0, 1]$$



Notice: $\text{Cost} = 0$, if $y=1, h\theta(x)=1$
But as $h\theta(x) \rightarrow 0$
 $\text{Cost} \rightarrow \infty$

Captures intuition that if $h\theta(x)=0$,
predict $P(y=1|x, \theta)=0$, but $y=1$.
We will penalize learning algorithm by
a very large cost.



Simplified cost function and Gradient descent.

$$\text{Cost}(h\theta(x), y) = \begin{cases} -\log(h\theta(x)) & , y=1 \\ -\log(1-h\theta(x)) & , y=0 \end{cases}$$

$$\Rightarrow \boxed{\text{Cost}(h\theta(x), y) = -y \log(h\theta(x)) - (1-y) \log(1-h\theta(x))}$$

Verifying: $y=1 \quad \text{Cost}(h\theta(x), y) = -\log(h\theta(x))$

$$y=0 \quad \text{Cost}(h\theta(x), y) = -\log(1-h\theta(x))$$

Logistic Regression Cost Function:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h\theta(x^{(i)}), y^{(i)})$$

$$= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h\theta(x^{(i)}))$$

vectorized implementation

$$h = g(X\theta)$$

$$J(\theta) = \frac{1}{m} \cdot (-y^\top \log(h) - (1-y)^\top \log(1-h))$$

To fit parameters θ :

$$\min J(\theta) \rightarrow \text{Get a set of } \theta$$

To make a prediction given new x :

$$\text{output } h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$J(\theta) = \frac{1}{m} \sum$$

Gradient Descent:

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

Repeat

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta). \Rightarrow \theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\log(h_{\theta}(x^{(i)})) - y^{(i)}) \cdot x_j^{(i)}$$

(Simultaneously update all θ_j)

Notice:

this equation is identical to linear regression. But they have different meaning.

for linear regression:

$$h_{\theta}(x) = \theta^T x.$$

for logistic regression:

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

the vectorized Gradient Descent:

$$\theta := \theta - \alpha \cdot \underbrace{\frac{1}{m} x^T (g(x\theta) - \bar{y})}_{\text{gradient}}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$h_{\theta}(x) = \theta^T x. = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n.$$

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

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} (J(\theta))$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum (h_{\theta}(x^{(j)}) - y^{(j)}) \cdot x_j^{(j)}$$

Advanced Optimization:

optimization algorithms:

} Conjugate gradient
BFGS. 变更度法
L-BFGS. 限制高斯度法.

Pros = - No need to manually pick d.
- often faster than gradient descent.

Cons = - more complex.

Example

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

Logistic Regression:

$$\theta = \begin{bmatrix} \theta_0 & \theta_1 \\ \theta_1 & \theta_2 \\ \vdots & \vdots \\ \theta_n & \end{bmatrix}$$

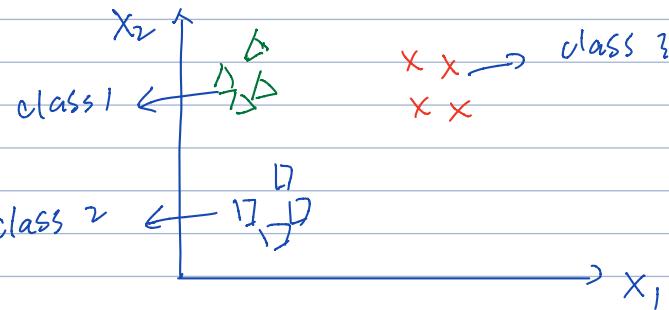
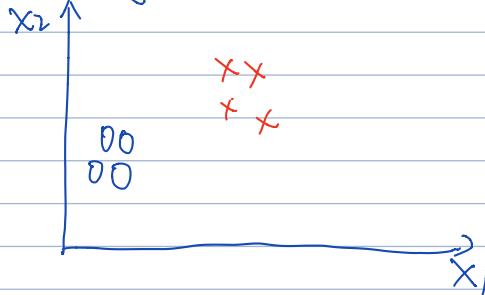
function [JVal, gradient] = costFunction(theta):

JVal = [Code to compute $J(\theta)$];
gradient(1) = [Code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];
⋮

gradient(n) = - - - $\frac{\partial}{\partial \theta_n} J(\theta)$;

Multi-class classification: One vs all (one-vs-rest)
multiclass classifications.

Binary

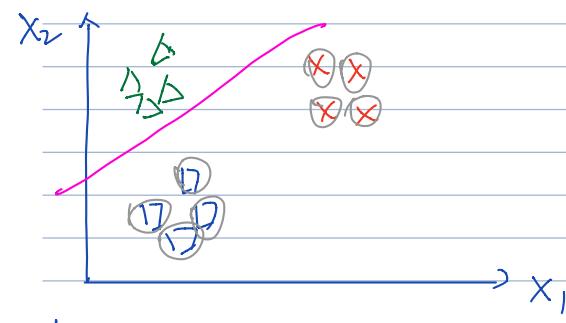


Train a logistic regression classifier $h_i(x^{(i)})$ for each class i to predict the probability that $y=i$. 我们将每个类的每一个标记为正向类，另外的标记为负向类。

On a new input x , to make a prediction, pick the class i that maximizes

$$\max_i h_i(x)$$

输入一个 x 值，在许多类中输入 x 值，然后我们选择一个在 $h_i(x)$ 最大的 i ，即 $\max_i h_i(x)$



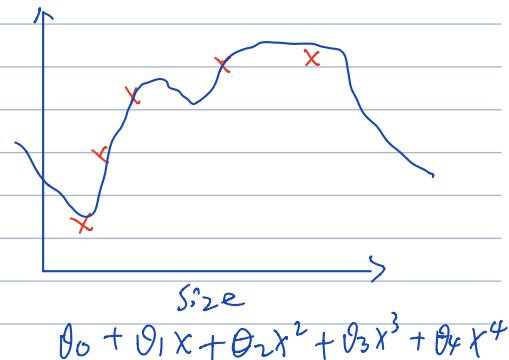
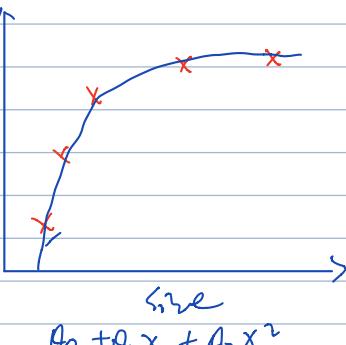
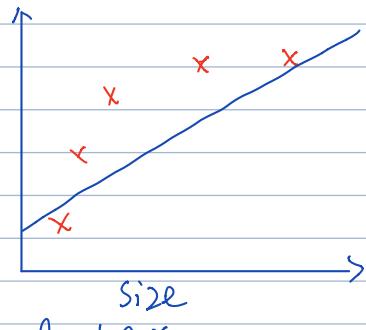
assume, the value of triangles is 1, the circle is 0.

$h(\theta(x))$ to classify the class 1.

Similarly, we can separate the other two classes.

Regularization:

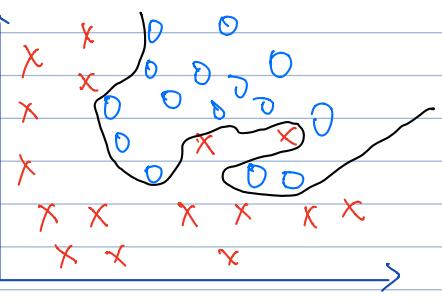
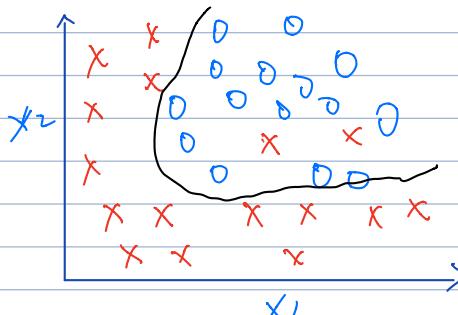
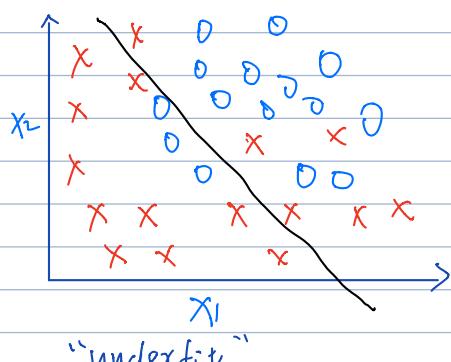
example: linear function (hunting price)



"underfit" "high bias"

"overfit" "High variance"

overfit: If we have too many features, the learned hypothesis may fit the training set very well ($J_\theta = \frac{1}{2m} \sum (h_\theta(x_i) - y_i)^2 \approx 0$), but fail to generalize to new examples.



Underfit: When the form of hypothesis function h maps poorly to trend of the data.

Address overfitting problem.

option 1: reduce the number of features:

- manually select which features to keep.
- use a model selection algorithm (e.g. PCA)

option 2: Regularization

- keep all the features, but reduce the magnitude of parameters θ_j .
- regularization works well when we have a lot of slightly useful features.

regularization:

Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$

- simpler hypothesis \leftarrow penalize.

- less prone to overfitting

Example: Housing:

- Features: x_1, x_2, \dots, x_{100}

- Parameters: $\theta_0, \theta_1, \theta_2, \dots, \theta_{100}$.

We modify the cost function to shrink all parameters.

$$J(\theta) = \frac{1}{2m} \left[\sum (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{100} \theta_j \right]$$

↑ 不知道要对哪些特征进行惩罚

$$\min J(\theta)$$

regularization parameter

control trade off between two different goals:

① fit the training data well;

If λ is very large (if $\lambda = 10^{10}$), -> underfit. ② keep parameters small.

We will penalize $\theta_1, \theta_2, \dots, \theta_{100}$..

then $\theta_1 = 0, \theta_2 = 0, \dots, \theta_{100} = 0$

$$h_\theta(x) = \theta_0 \rightarrow \text{underfit.}$$

正是 $h_\theta(x)$ 中的高次项导致了过拟合的产生，所以我们需要做的就是减小这些 θ 的值。

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

我们在 Cost Function 中设置一点 penalty.

$$J(\theta) = \frac{1}{2m} \sum (h_\theta(x^{(i)}) - y^{(i)})^2 + 1000\theta_3^2 + 1000\theta_4^2$$

regularized linear regression:

We will modify the gradient descent to separate out θ_0 , from the rest of parameters:

$$\theta_0 := \theta_0 - 2 \frac{1}{m} \sum (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

第0特征 θ_0 对应的 $x_0^{(i)}$ 在每个样本中。

$$\theta_j := \theta_j - 2 \left[\frac{1}{m} \sum (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right]$$

$$= \theta_j (1 - 2 \frac{\lambda}{m}) - 2 \frac{1}{m} \sum (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$1 - 2 \frac{\lambda}{m} < 1$: Reducing the value of θ_j by some amount on every update.

$$\text{grad } [\theta] = \frac{1}{m} * X [1, :, 1] * (h_\theta(x^{(i)}) - y^{(i)})$$

Normal equation:

$$\theta = (X^T X + \lambda \cdot I)^{-1} \cdot X^T y.$$

$$I = \begin{bmatrix} 0 & & 0 \\ & 1 & & 0 \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix}_{(n+1) \times (n+1)}$$

If $m < n$, $X^T X$ is non-invertible,
when add the term λI ,
 $X^T X + \lambda I$ becomes invertible.

regularized logistic regression:

Cost Function:

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right] + \frac{\lambda}{m} \sum_{j=1}^n \theta_j^2$$

Gradient descent

Repeat {

$$\theta_0 := \theta_0 - \frac{\lambda}{m} (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_j := \theta_j - \frac{\lambda}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \quad i = 1, 2, \dots, n.$$

Regularized logistic regression:

$$J(\theta) = \frac{1}{m} \sum_i [y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$