

WEEK 3

Logistic Regression

Regularization

WEEK 3

Classification

$$y = 0 \quad \text{or} \quad 1$$

Logistic regression: $0 \leq h_\theta(x) \leq 1$

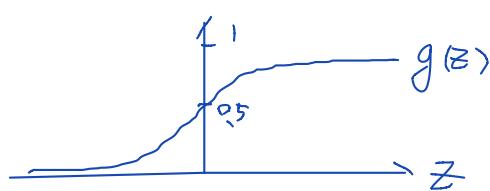
Logistic regression — hypothesis representation

model: want $0 \leq h_\theta(x) \leq 1$

$$h_\theta(x) = g(\theta^T x)$$

$$= \frac{1}{1 + e^{-z}} \leftarrow \text{Si } i \quad \circ \text{ is } -1 \quad - \text{not } 0$$

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



for binary classification: $h_{\theta}(x) = P(y=1|x; \theta)$
 $P(y=0|x; \theta) + P(y=1|x; \theta) = 1$

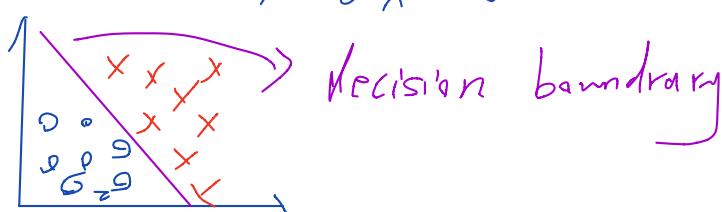
Logistic regression — decision boundary

$$h_{\theta}(x) = P(y=1|\theta; x) \begin{cases} \text{predict } y=1 \text{ if } h_{\theta}(x) \geq 0.5 \\ \dots \quad y=0 \text{ if } h_{\theta}(x) < 0.5 \end{cases}$$

$$g(z) = \frac{1}{1+e^{-z}}$$

$$g(z) \geq 0.5 \text{ when } z \geq 0 \quad \hookrightarrow \theta^T x \geq 0$$

$$g(z) < 0.5 \text{ when } z < 0 \quad \hookrightarrow \theta^T x < 0$$



Logistic regression — cost function

Train set: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad x_0 = 1, \quad y \in \{0, 1\}$$

m train sets, n features

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \quad \text{how to choose parameters } \theta ?$$

Cost function

Linear regression: $J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

* if we use the above same cost function for Logistic regression,
the cost function will be "non-convex" (has many local minimums)



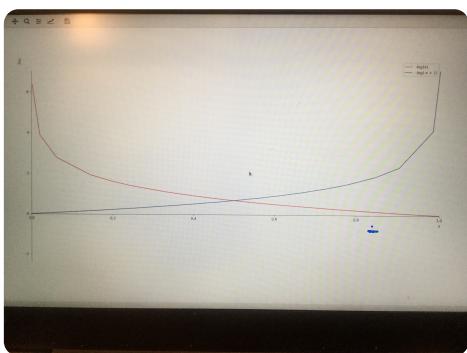
which is not guaranteed to converge to global minimum

Thus we have to modify the cost function for logistic regression

definition:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x_i), y_i)$$

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



Plot of cost

Simplified Cost function and gradient descent

Logistic regression cost function: $J_\theta = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$

Simplification:

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

then:

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

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

To fit parameters θ : $\min_{\theta} J(\theta)$

To make a prediction giving new x : $h(\theta) = \frac{1}{1+e^{-\theta^T x}}$

Gradient descent

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

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

then:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad [\text{Sim update all } \theta_j]$$

Notice Ng did put $\frac{1}{m}$ here, which I think is a mistake

Looks exactly the same as linear regression!

Except the definition of hypothesis $h_\theta(x)$ is diff.

Feature scaling also works for logistic regression!

Vectorized implementation:

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X^T \theta) - \vec{y})$$

Logistic regression — adv. optimization

given θ , we have code that Compute:

- $J(\theta)$
- $-\frac{\partial}{\partial \theta_j} J(\theta)$ for $j = 0, 1, \dots, n$

optimization algorithms

- Gradient decent

- Conjugate gradient
- BFGS
- L-BFGS

advantages:

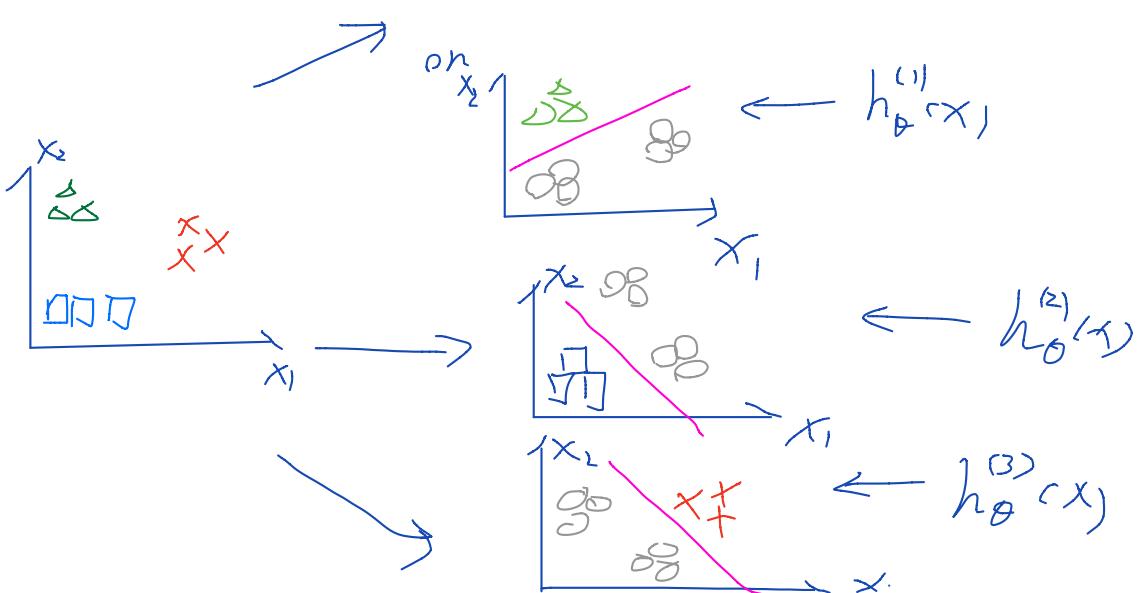
- No need to manually pick α
- often faster than gradient descent

Disadvantages:

- more complex

Multi-class classification

one-vs-all (rest) classification



$$h_{\theta}^{(i)}(x) = P(y=i | x; \theta), i = 1, 2, 3$$

one-vs-all

Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class i to predict the probability that $y = i$

on a new output x , to make a prediction, pick the class i that maximizes: $\max_i h_\theta^{(i)}(x)$

Problem of overfitting

Overfitting:

If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) \approx 0$), but fail to generalize to new example

Addressing overfitting

Options:

1. Reduce number of features

a. Manually select which features to keep

b. Model selection algorithm (later)

2. Regularization

a. keep all features, but reduce magnitude / values of θ_j

↳ works well when having a lot of features, each of which contributes a bit to prediction of

Regularization — cost function

Regularization:

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

↪ "Simpler" hypothesis

↪ less prone to over fitting

regularized Cost function

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

omitted θ_0

convention

and θ_0 doesn't
make any diff.

regularization parameter

λ : too larger \rightarrow underfit

Regularized linear regression:

① gradient decent:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) X_0^{(i)}$$

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

($j = \cancel{0}, 1, 2, 3, \dots, n$)

$$\theta_j := \underbrace{\theta_j}_{\text{---}} \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) X_j^{(i)}$$

usually α, λ small, m large, so $1 - \alpha \frac{\lambda}{m} < 1$

② Normal Equation

$$X = \begin{bmatrix} (x^{(1)})^\top \\ \vdots \\ (x^{(m)})^\top \end{bmatrix}_{m \times n+1}$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$$\min_{\theta} J(\theta)$$

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

Non-invertibility

if $m \leq n$, $[X^T X]$ non-invertible

if $\lambda > 0$, then $[X^T X + \lambda [1, 1]]^{-1}$ is guaranteed to be invertible

Regularized logistic regression

regularize

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

Gradient descent:

Same as linear regression, except that hypothesis func $h_\theta(x)$ is different.