

VII. LOGISTIC REGRESSION

Lecture 33: Classification

We will start with a binary classification problem:

→ Email: Spam/Not Spam?

→ Online transactions: Fraudulent (Yes/ No)?

→ Tumor: Malignant / Benign?

We can represent the output as

$$y \in \{0, 1\}$$

"Negative class"

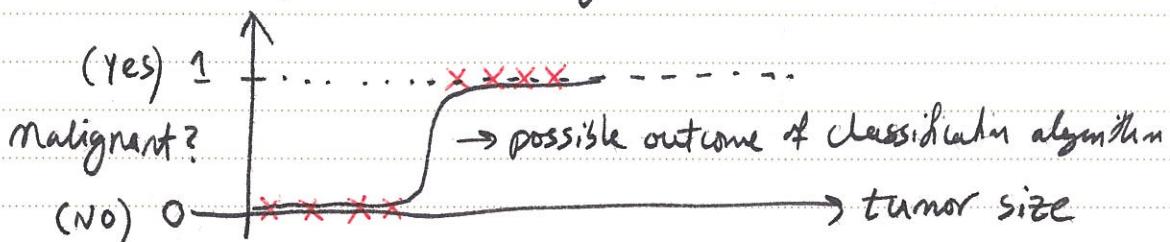
e.g. benign tumor

"Positive class"

e.g. malignant tumor.

→ Later on we will generalize to multi-classification problems where e.g. $y \in \{0, 1, 2, 3\}$

Here's an example of a training set:



- * Any curve that would be a good fit to these points ought to be non-linear. Of course, a non-linearity of this sort is very hard to capture with a polynomial, as it would require many terms. It is fairly clear, therefore, that linear regression is not going to do.
- * It would make sense that in a classification problem we predict probabilities e.g. what is the probability that a 2mm tumor is malignant?

Then in this light, the fitting curve drawn in the previous page can be interpreted as the probability, as a function of tumor size, that the tumor is malignant. Choosing $y = \{0, 1\}$ to correspond to $\{\text{benign, malignant}\}$ also makes more sense now.

- * If we have a probabilistic interpretation, we would want our hypothesis function to be always between 0 & 1.

Lecture 34: Hypothesis Representation

- Logistic Regression is a type of classification problem.
- The hypothesis function used in logistic regression is:

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



* Note that $0 \leq h_\theta(x) \leq 1$, as desired.

* As we mentioned before: $h_\theta(x)$ = estimated probability that $y=1$ on input x .

$$\boxed{h_\theta(x) = P(y=1|x; \theta)}$$

Probability of $y=1$ given x & parameterized by θ .

$$\begin{aligned} * \text{ Obviously } P(y=0|x; \theta) &= 1 - P(y=1|x; \theta) \\ &= 1 - h_\theta(x) \end{aligned}$$

Lecture 35: Decision Boundary

2D Although we're interested in a probabilistic prediction in classification problems, often we'd want to make concrete decisions based on the predicted probability once we've encountered with some new input. The simplest thing to do is:

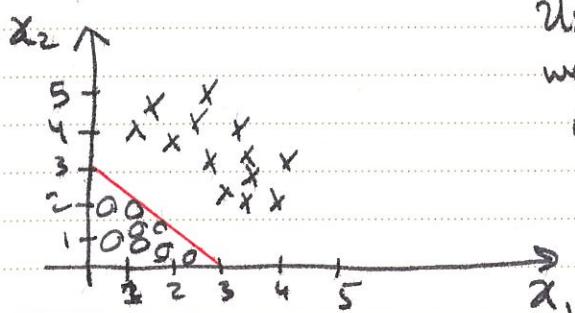
$$\begin{cases} y=1 & \text{if } P(y=1|x; \theta) \geq 0.5 \\ y=0 & \text{if } P(y=1|x; \theta) < 0.5 \end{cases}$$

$$\text{Equivalently: } \begin{cases} y=1 & \text{if } h_\theta(x) \geq 0.5 \\ y=0 & \text{if } h_\theta(x) < 0.5 \end{cases}$$

$$\text{Now since } h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} : \begin{cases} h_\theta(x) \geq 0.5 \Leftrightarrow \theta^T x \geq 0 \\ h_\theta(x) < 0.5 \Leftrightarrow \theta^T x < 0. \end{cases}$$

* The decision boundary is then defined via $\theta^T x = 0$ *

Let's say we have two feature x_1 & x_2 & our training set is as follows:



Using $h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$, we (somehow) fit to the data & find $\theta_0 = -3$, $\theta_1 = \theta_2 = 1$.

Our decision boundary is then

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

$$\Leftrightarrow x_1 + x_2 = 3$$

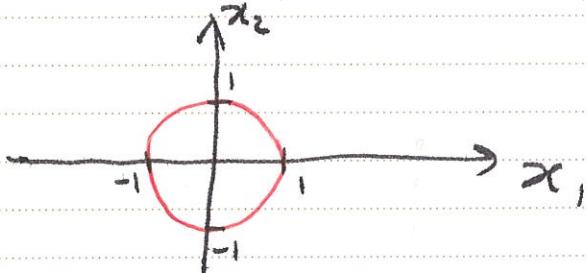
$$\left\{ \begin{array}{l} y=1 : -3 + x_1 + x_2 \geq 0 \\ y=0 : -3 + x_1 + x_2 \leq 0 \end{array} \right.$$

→ Note that the decision boundary is a property of the hypothesis function & the values of θ , which are of course determined by fitting the model to the training set.

→ With more complicated models, one can have more complicated decision boundaries:

$$h_0(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

$$\& \theta_0 = -1 ; \theta_1 = \theta_2 = 0 ; \theta_3 = \theta_4 = 1$$



Lecture 36: Cost Function \oplus Lecture 37: Simplified Cost function and Gradient Descent

→ How do we fit the parameters θ in logistic regression?

→ Remember that $h_\theta(x)$ is interpreted as the probability that $y=1$ at x . What is the probability of obtaining the training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$?

$$x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \quad \leftarrow y^{(i)} \in \{0, 1\}$$

The probability of $(x^{(i)}, y^{(i)})$ happening is:

$$P^{(i)} = \begin{cases} h_\theta(x^{(i)}) & \text{if } y^{(i)} = 1 \\ 1 - h_\theta(x^{(i)}) & \text{if } y^{(i)} = 0 \end{cases}$$

or more condensely $P^{(i)} = h_\theta(x^{(i)})^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}}$

→ The probability of obtaining the whole training set is then:

$$P(\theta) = \prod_{i=1}^m P^{(i)} = \prod_{i=1}^m [h_\theta(x^{(i)})^{y^{(i)}} (1 - h_\theta(x^{(i)}))^{1-y^{(i)}}]$$

* we can now use the Principle of maximum likelihood to find the value of θ which maximizes P . *

Note that if we let $\tilde{P}(\theta) = -\ln P(\theta)$

$$\rightarrow \frac{\partial \tilde{P}}{\partial \theta_k} = -\frac{1}{P(\theta)} \frac{\partial P}{\partial \theta_k} \Rightarrow \frac{\partial \tilde{P}}{\partial \theta_k} = 0 \Leftrightarrow \frac{\partial P}{\partial \theta_k} = 0 \quad \text{since } P(\theta) > 0$$

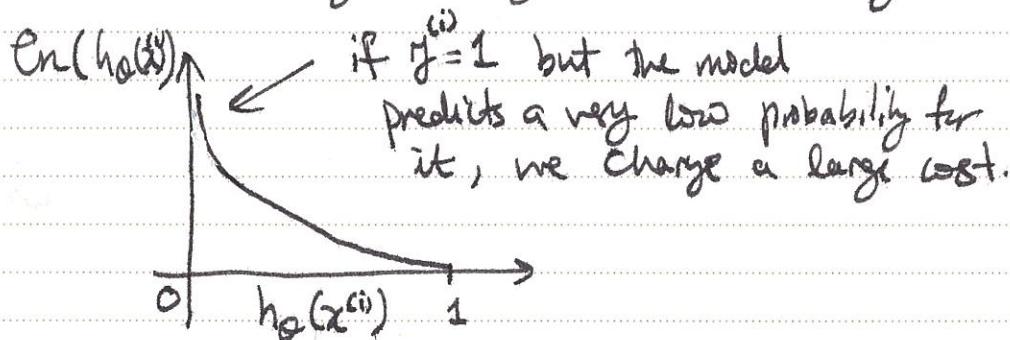
Therefore, all extrema of P & \tilde{P} are the same. Moreover, if θ_{\max} maximizes P , it minimizes \tilde{P} . It is more convenient to minimize \tilde{P} , & this is what we'll do:

→ The cost function for logistic regression is defined as:

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

$$\begin{aligned} \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) &\equiv -y^{(i)} \ln(h_{\theta}(x^{(i)})) - (1-y^{(i)}) \ln(1-h_{\theta}(x^{(i)})) \\ &= \begin{cases} -\ln(h_{\theta}(x^{(i)})) & y^{(i)}=1 \\ -\ln(1-h_{\theta}(x^{(i)})) & y^{(i)}=0 \end{cases} \end{aligned}$$

is the cost we charge the algorithm for training set i .



So let's show that for $h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$, $J(\theta)$ is convex.

$$\frac{\partial h}{\partial \theta_k} = \frac{e^{-\theta^T x}}{(1+e^{-\theta^T x})^2} x_k = e^{-\theta^T x} h_{\theta}(x) x_k$$

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

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

$$\begin{aligned} \frac{\partial J}{\partial \theta_k} &= -\frac{1}{m} \sum_{i=1}^m \frac{\partial h_\theta(x^{(i)})}{\partial \theta_k} \left\{ \frac{y^{(i)}}{h_\theta(x^{(i)})} - \frac{(1-y^{(i)})}{h_\theta(x^{(i)})} e^{\theta^T x} \right\} \\ &= -\frac{1}{m} \sum_{i=1}^m e^{-\theta^T x^{(i)}} h_\theta(x^{(i)}) x_k^{(i)} \cdot \frac{y^{(i)} [1+e^{\theta^T x}] - e^{\theta^T x}}{h_\theta(x^{(i)})} \end{aligned}$$

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

$$\boxed{\frac{\partial J}{\partial \theta_k} = \frac{1}{m} \sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}] x_k^{(i)}}$$

$$\frac{\partial^2 J}{\partial \theta_k \partial \theta_l} = \frac{1}{m} \sum_{i=1}^m \frac{\partial h_\theta(x^{(i)})}{\partial \theta_l} x_k^{(i)}$$

$$= \frac{1}{m} \sum_{i=1}^m e^{-\theta^T x^{(i)}} h_\theta(x^{(i)}) x_l^{(i)} x_k^{(i)}$$

$$\rightarrow \nabla^2 J = \frac{1}{m} \sum_{i=1}^m e^{-\theta^T x^{(i)}} h_\theta(x^{(i)})^2 x^{(i)T} x^{(i)} > 0 \quad \checkmark$$

$\Rightarrow J(\theta)$ only has one minima, i.e. the global minima!

Gradient Descent: [repeat until converge {

$$\theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m [\text{h}_\theta(x^{(i)}) - y^{(i)}] x_j^{(i)}$$

$$\text{where } \text{h}_\theta(x) = (1 + e^{-\theta^T x})^{-1}$$

As with linear regression, feature scaling is important!

Lecture 38: Advanced Optimization

Gradient descent is not the only algorithm we can use to minimize our cost function $J(\theta)$. Here are a few more advanced ones:

→ Conjugate Gradient

→ BFGS

→ L-BFGS

Some advantages of these more advanced algorithms:

→ No need to manually pick the learning rate α .

→ Often faster than Gradient Descent.

There are implementations of these in various languages.

They normally take $J(\theta)$ & $\frac{\partial J}{\partial \theta_j}$ as inputs.

In Octave, `fminunc` can be used.

Lecture 39: Multiclass Classification: One-vs-all

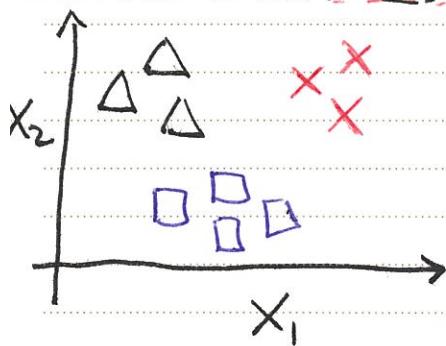
So far we've been discussing binary classification problems

What if we have multiple classes?

Examples: (i) Email folder/Hatting: Work, Friends, Family, Hobby

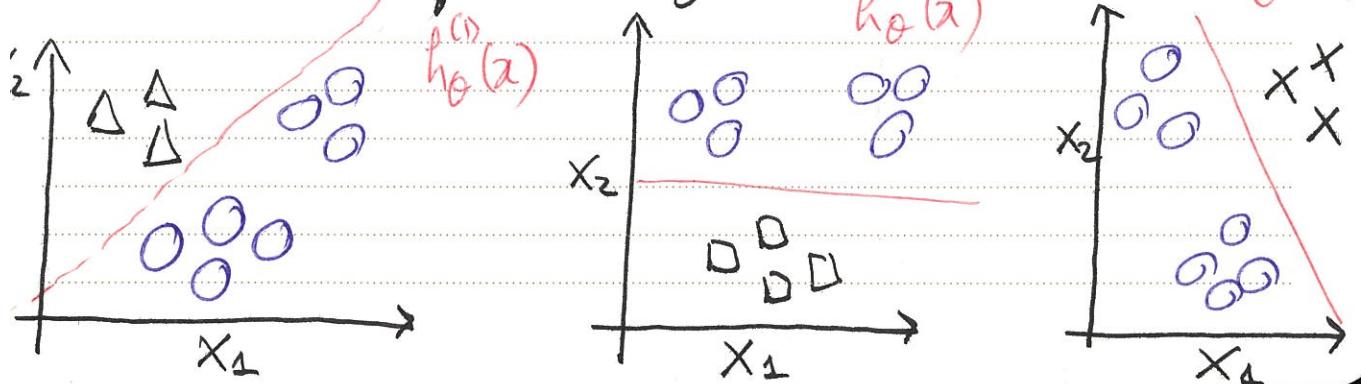
(ii) Weather prediction: Sunny, Cloudy, Rain, Snow

One-vs-all (One-vs-rest):



Class 1 : \triangle
Class 2 : \square
Class 3 : \times

For each class, we solve a binary classification problem where the chosen class is taken as positive, and all other classes are mapped to negative:



On a new input x , to make a prediction, pick the class q . That has the maximum $h_q^{(2)}(x)$:

$$\boxed{\max_k h_k^{(2)}(x)}$$

Note that nothing ensures that $\sum_{q=1}^K h_q^{(2)}(x) = 1$, so in fact we should not really think about $h_q^{(2)}(x)$ as the probability that class = q at x .

I suspect, though, if there the training set is large, the sum-probabilities will be very close to 1. Why?

Assume we're using logistic regression as our binary classifier
Then we would need to solve: (see page 34)

$$\frac{\partial J^{(1)}}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m [h_0^{(1)}(x^{(i)}) - \delta_{y^{(i)}, 1}] x_j^{(i)} = 0$$

$$\vdots \\ \frac{\partial J^{(K)}}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m [h_0^{(K)}(x^{(i)}) - \delta_{y^{(i)}, K}] x_j^{(i)} = 0$$

assuming we have K classes: $y^{(i)} \in \{1, 2, \dots, K\}$ $\forall i$

Summing up these equations & noting that $\sum_{q=1}^K \delta_{q(i), q} = 1$

$$\sum_{i=1}^m \left\{ \sum_{q=1}^K h_\theta^{(q)}(x^{(i)}) - 1 \right\} x_j^{(i)} = 0$$

If there are many training examples, i.e. m is very large, it's likely that solving the equation above would require

$$\sum_{q=1}^K h_\theta^{(q)}(x) = 1.$$

Again, though, there's no guarantee at all.

There's a multiclass classification problem which does predict probabilities that sum to 1, by construction.

This is not covered in the course but I'll look at it briefly below.

Multinomial Logistic Regression

This is also known as the Maximum Entropy Classifier, Softmax regression, and a bunch of other names.

Suppose: $y \in \{1, 2, 3, \dots, K-1, K\}$.

Given a training set $\{(x^{(i)}, y^{(i)})\}$, our goal is to predict the following probabilities for an unseen value of x :

Probability of $y^{(i)}=1$ Probability of $y^{(i)}=2$ Probability of $y^{(i)}=K$

$P_1(x), P_2(x), \dots, P_K(x)$

Of course, we require that

$$\sum_{q=1}^K P_q(x) = 1 \quad \forall x. \quad \text{#}$$

Multinomial logistic regression uses the following parametrization:

$$P_q(x) = \frac{e^{-\theta_q^T x}}{1 + \sum_{q=1}^{K-1} e^{-\theta_q^T x}} \quad q \in [1, K-1]$$

P_K is fixed by # above:

$$P_K(x) = \frac{1}{1 + \sum_{q=1}^{K-1} e^{-\theta_q^T x}}$$

Our job is now to estimate $\{\theta_q\}$ $q \in [1, K-1]$

Note that every θ_k is an $n+1$ dimensional vector.
 \uparrow
 $\#$ of features.

We will use the principle of maximum likelihood:

What is the probability of $(x^{(i)}, y^{(i)})$ occurring?

$$P^{(i)} = P_1(x^{(i)})^{\delta_{y^{(i)}, 1}} P_2(x^{(i)})^{\delta_{y^{(i)}, 2}} \dots P_K(x^{(i)})^{\delta_{y^{(i)}, K}}$$

$$J(\theta) = -\frac{1}{m} \ln \prod_{i=1}^m P^{(i)} = -\frac{1}{m} \sum_{i=1}^m \ln P^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m [\delta_{y^{(i)}, 1} \ln P_1(x^{(i)}) + \dots + \delta_{y^{(i)}, K} \ln P_K(x^{(i)})]$$

$$-\theta_1^T x^{(i)}$$

$$= -\frac{1}{m} \sum_{i=1}^m \left[\delta_{y^{(i)}, 1} \ln \frac{e^{-\theta_1^T x^{(i)}}}{1 + \sum_{k=1}^{K-1} e^{-\theta_k^T x^{(i)}}} + \delta_{y^{(i)}, 2} \ln \frac{e^{-\theta_2^T x^{(i)}}}{1 + \sum_{k=1}^{K-1} e^{-\theta_k^T x^{(i)}}} \right. \\ \left. + \dots + \delta_{y^{(i)}, K-1} \ln \frac{e^{-\theta_{K-1}^T x^{(i)}}}{1 + \sum_{k=1}^{K-1} e^{-\theta_k^T x^{(i)}}} \right]$$

$$-\theta_{K-1}^T x^{(i)}$$

$$+ \delta_{y^{(i)}, K-1} \ln \frac{e^{-\theta_K^T x^{(i)}}}{1 + \sum_{k=1}^{K-1} e^{-\theta_k^T x^{(i)}}}$$

$$+ \delta_{y^{(i)}, K} \ln \frac{1}{1 + \sum_{k=1}^{K-1} e^{-\theta_k^T x^{(i)}}}]$$

$$\begin{aligned} J(\theta) = & -\frac{1}{m} \sum_{i=1}^m \left\{ \delta_{y^{(i)}, 1} \left[-\theta_0^T \boldsymbol{x}^{(i)} - \ln \left[1 + \sum_{\ell=1}^{K-1} e^{-\theta_\ell^T \boldsymbol{x}^{(i)}} \right] \right] \right. \\ & + \dots \\ & \left. + \delta_{y^{(i)}, K-1} \left[-\theta_{K-1}^T \boldsymbol{x}^{(i)} - \ln \left[1 + \sum_{\ell=1}^{K-1} e^{-\theta_\ell^T \boldsymbol{x}^{(i)}} \right] \right] \right. \\ & \left. + \delta_{y^{(i)}, K} \times -\ln \left[1 + \sum_{\ell=1}^{K-1} e^{-\theta_\ell^T \boldsymbol{x}^{(i)}} \right] \right\} \end{aligned}$$

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \left[\sum_{q=1}^{K-1} \delta_{y^{(i)}, q} \theta_q^T \boldsymbol{x}^{(i)} + \ln \left[1 + \sum_{\ell=1}^{K-1} e^{-\theta_\ell^T \boldsymbol{x}^{(i)}} \right] \right]$$

where we have used the fact that $\sum_{q=1}^{K-1} \delta_{y^{(i)}, q} = 1$

We can think about $\{\theta_\ell\}$ as a $(n+1) \times (K-1)$ matrix:

$$\textcircled{H} \underbrace{\delta_{qj}}_{\substack{\text{rows are features} \\ \text{---}}} = [\theta_1]_j \quad \textcircled{H} = [\theta_1 | \theta_2 | \dots | \theta_{K-1}]$$

rows are features classes

$$\theta_\ell^T \boldsymbol{x}^{(i)} = \sum_{j=1}^{n+1} [\theta_\ell]_j \boldsymbol{x}_j^{(i)} = \sum_{j=1}^{n+1} \textcircled{H}_{qj} \boldsymbol{x}_j^{(i)}$$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m [y^{(i)}_j - \frac{e^{-\theta_j^T x^{(i)}}}{1 + \sum_{q=1}^{k-1} e^{-\theta_q^T x^{(i)}}}] x_j^{(i)}$$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m [y^{(i)}_j - p_j(x^{(i)})] x_j^{(i)}$$