

Logistic Regression

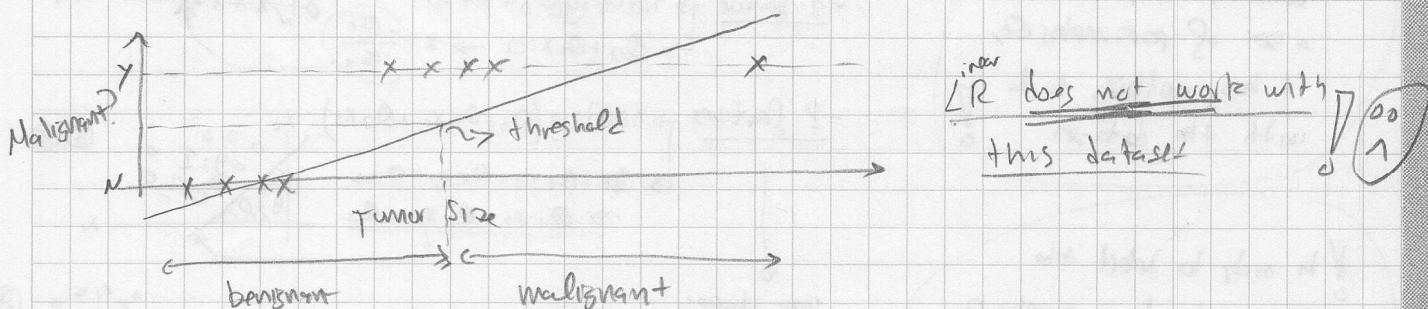
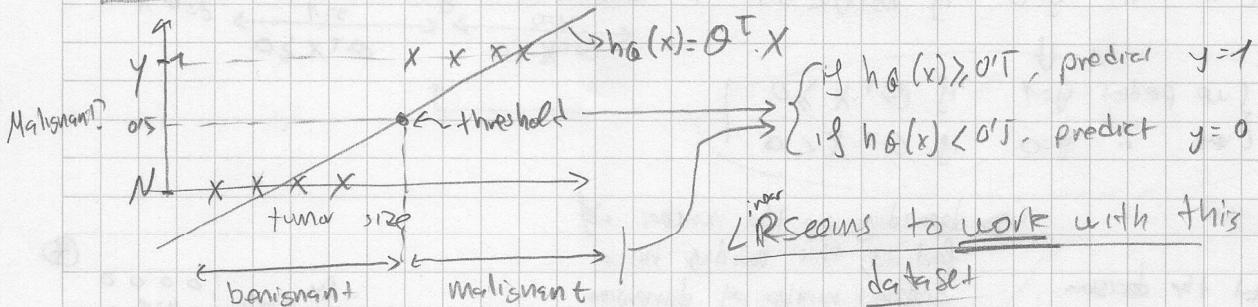
- Algorithm used for classification problems

Classification Problems

- Email: spam / not spam
 - Bank Transactions: fraudulent / not fraudulent
 - Tumors: Malignant / Benignant
 - ⋮
- arbitrarily, the negative class for absence (0)
the positive class for presence (1)
- $\underbrace{1}_{\text{(pos.)}}$ $\underbrace{0}_{\text{(neg.)}}$
- $y \in \{0, 1\}$ 0: Negative class 1: Positive class
- we want to predict the variable y
that takes two values {0, 1} true, false

(there are multi-class classification problems)
where $y \in \{0, 1, 2, 3, \dots\}$

Classifying using Linear Regression?



Another problem with using LR:

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

$h_{\theta}(x)$ can be > 1 or $< 0 \rightarrow$ confusing!

we will use:

Logistic Regression: $0 \leq h_{\theta}(x) \leq 1$

it is confusing to have 'regression' in the name of a classification algorithm

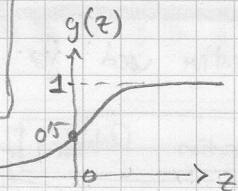
Hypothesis in Logistic Regression

Model

We want
 $0 \leq h_{\theta}(x) \leq 1$

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

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



Sigmoid function

or
logistic function

When looking for best θ , this
 $h_{\theta}(x)$ with the $g(z)$ will be used

Interpretation

If $h_{\theta}(x) = 0.7$, mean that we have a 70% (0.7) probability that

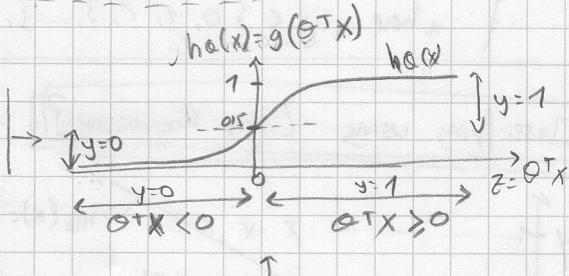
the output is 1 ($y=1$)

- formally: $h_{\theta}(x) = P(y=1 | x; \theta)$ → probability that $y=1$, given x , and parametrized by θ

Decision Boundary

{ we predict $y=1$ if $h_{\theta}(x) \geq 0.5$
 .. " " $y=0$ if $h_{\theta}(x) < 0.5$

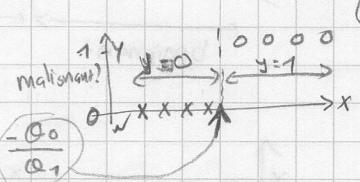
{ we predict $y=1$ if $\theta^T x \geq 0$
 .. " " $y=0$ if $\theta^T x < 0$



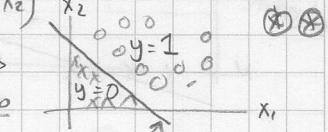
Note that the decision boundary depends on a set of parameters θ , it has nothing to do with the data-set!

Depending on the number of features, this boundary has a different number of dimensions

- 1 feature $\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1)$
 $\rightarrow \theta_0 + \theta_1 x_1 = 0 \rightarrow x_1 = -\frac{\theta_0}{\theta_1}$



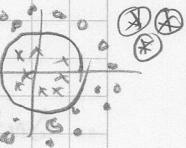
- 2 features $\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$
 $\rightarrow \theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0 \rightarrow \theta_1 x_1 + \theta_2 x_2 = -\theta_0$



any shape:

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

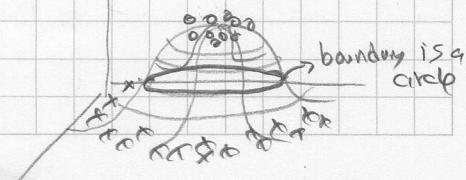
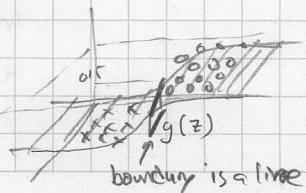
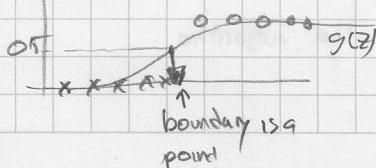
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \dots)$$



In order to select the best set of parameters θ , we will use the different threshold

threshold

I think working with $g(z)$, when looking for the best θ , we are looking at something like:



Cost function - How to choose θ

Training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$ m examples

$$X = \begin{bmatrix} 1 \\ x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{feature vector}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \text{parameter vector}$$

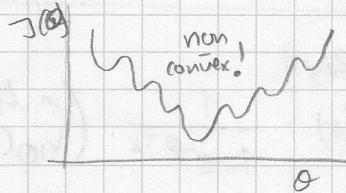
$$y \in \{0, 1\} \quad \text{output}$$

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

How to choose θ ?

• We still want to minimize a $J(\theta)$ function $J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_{\theta}(x_i), y_i)$

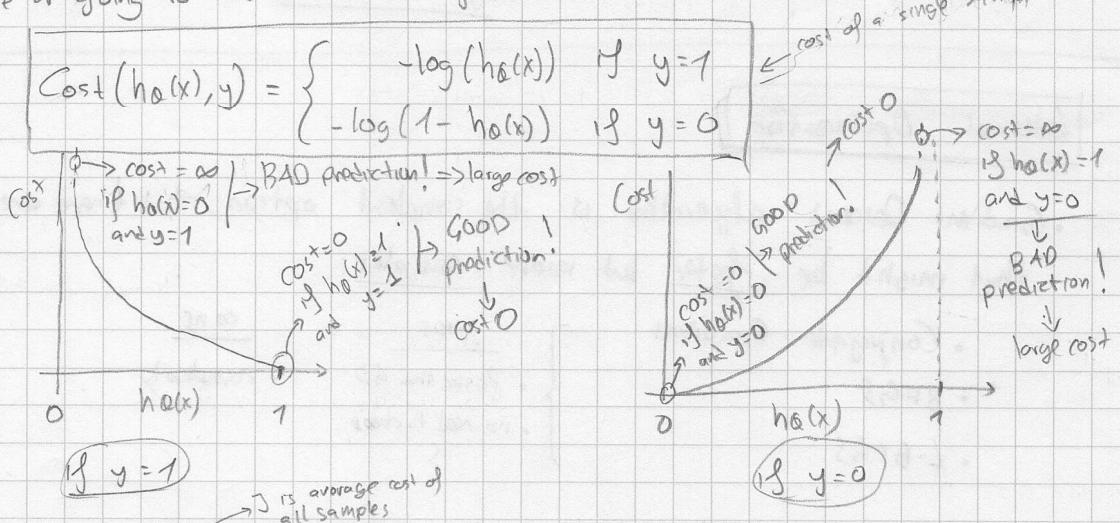
• If we use the same cost function as in LinReg ($\text{cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$), it won't work properly because $J(\theta)$ becomes non-convex:



← note that in Log Reg $h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$
this fact makes it non-convex

← Gradient-descent might end up in a local minimum

• We are going to see this cost function:



• $J(\theta) = \frac{1}{2} \sum_{i=1}^m (\text{cost}(h_{\theta}(x_i), y_i))$, with this cost function, becomes convex! !!

(not demonstrated, beyond objective of the course)

Reformulation:

→ these terms are either 0 or 1.

$$\text{(single sample)} \rightarrow \text{cost}(h_{\theta}(x), y) = -y \cdot \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

$$\text{(all samples)} \rightarrow J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(h_{\theta}(x)) + (1-y^{(i)}) \log(1-h_{\theta}(x))]$$

• To fit parameters θ

$$\min_{\theta} J(\theta) \rightarrow \text{get } \theta$$

• To predict given new x

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

Logistic Regression & Gradient Descent

We will use Gradient Descent to minimize $J(\theta)$ and obtain the best θ parameters.

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

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

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

(simultaneously update all θ_j)

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

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

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

$$\theta_j := \theta_j - \alpha \cdot \sum_{i=1}^m (h_\theta(x) - y) \cdot x_j$$

(simultaneously update all θ_j)

! same algorithm as in LinReg!

The difference is that now $h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$ (in LinReg it was $h_\theta(x) = \theta^T x$)

- Since we apply Gradient Descent, feature scaling is also useful here.

Advanced Optimization

- Gradient Descent algorithm is the simplest option, but there are others that might be faster but more complex:

- Conjugate Gradient
- BFGS
- L-BFGS

	pros	cons
• faster than GD		• complexity
• no need to choose α		

- Can be used as a 'black box' without understanding what they really do.
- you just need to pass a cost function to these algorithms...

Multiclass classification

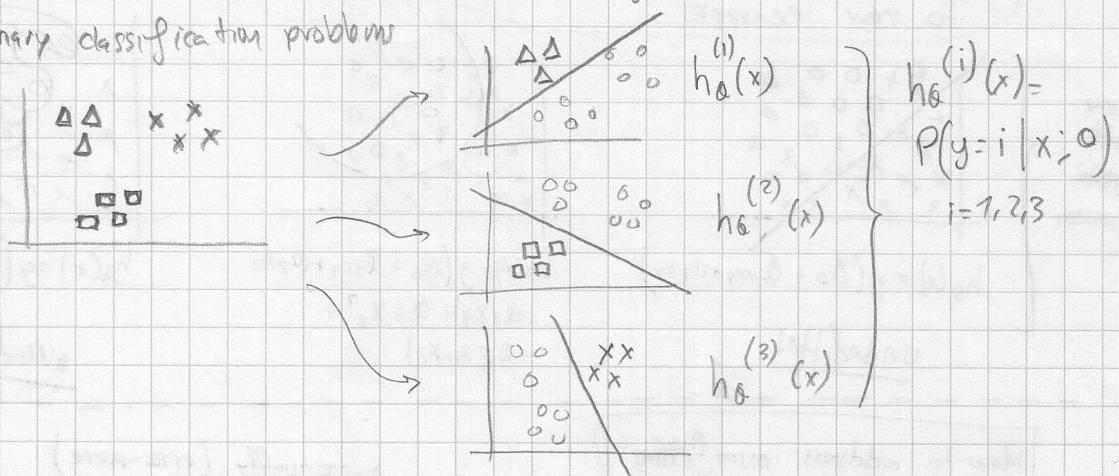
- Email filtering: work, friends, family, hobby,
 $y=1 \quad y=2 \quad y=3 \quad y=4$

- Medical diagnosis: not ill, cold, flu
 $y=0 \quad y=1 \quad y=2$

:

One-vs-all (or One-vs-rest)

- The idea is to convert the multiclass classification problem to several binary classification problems



- So we separate and try to predict the probability that an example belongs to a ~~given~~ particular class; and we compare this probability for each class.

- When we predict, we predict the class that gives the maximum probability.

Newton's Method for Log Reg

- An alternative algorithm to minimize $J(\theta)$ that works well with Logistic Regression
- Quite fast and technical description in the video...