# Logistic Regression ⟵ classification
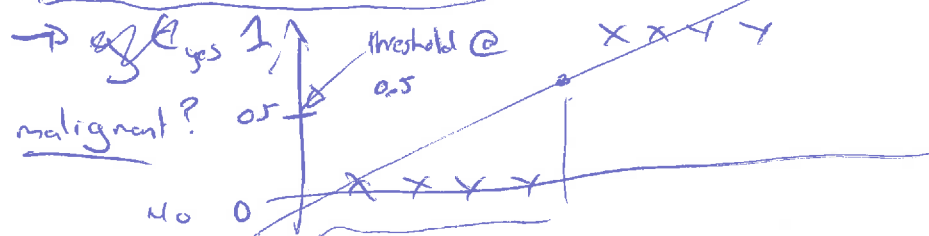
Classification problems - email
- Online Transactions / fradulent etc
- Tumor: Malignant / benign

$y \in (0, 1)$   0: "negative class" (benign tumor)
         1: "positive class" (malignant tumor)

         "absence of something"

"presence of something"

malignant? or
yes 1 ↑ threshold @ 0.5    X X Y Y

No 0

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

⟵ linear regression (hypothesis function)

Threshold classifier output $h_\theta(x)$ at 0.5
  $h_\theta(x) \geq 0.5$, predict "$y=1$"
  $h_\theta(x) \leq 0.5$,  "  "$y=0$"

Linear regression for classification

classification: $y = 0$ or $1$
$h_\theta(x)$ can be $\geq 1$ or $\leq 0$
so lets develop a new algorith such that $0 \leq h_\theta(x) \leq 1$

logistic regression. ⟶ classification algoritm (discrete value 0 or 1)
         ↗ historical (confusing)

## Logistic Regression – hypothesis representation

What is the function that we are going to use to represent our hypothesis when we have a classification problem?

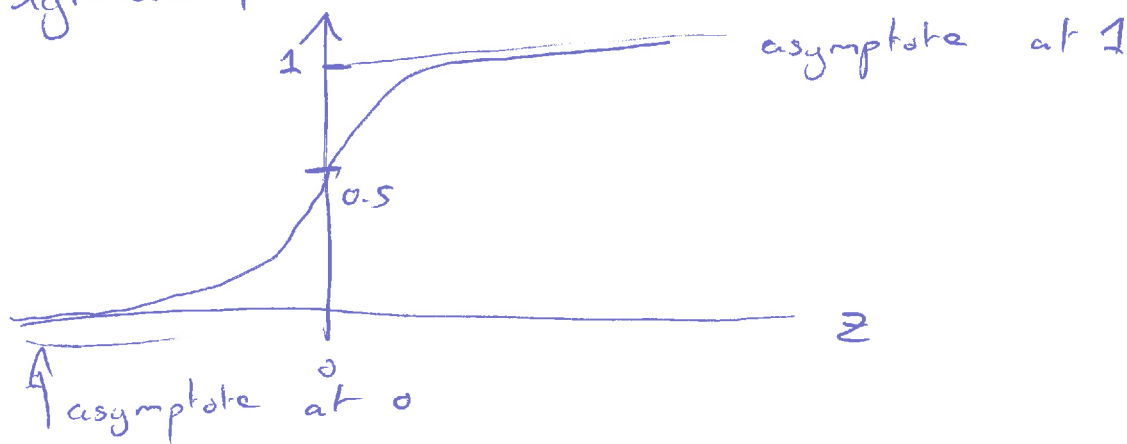### Logistic Regression model

Want $0 \leq h_\theta(x) \leq 1$

$h_\theta(x) = \theta^T x$ ← linear regression

logistic regression

$h_\theta(x) = g(\theta^T x) \to g(z) = \dfrac{1}{1+e^{-z}}$ ← sigmoid function (logistic function)

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

sigmoid function.



asymptote at 1

0.5

$z$

asymptote at 0

Interpretation of hypothesis output $h_\theta(x)$

how we are going to treat the output of our hypothesis function using the sigmoid function.

$h_\theta(x) =$ estimated probability that $y=1$ on input $x$

example: if $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumor size} \end{bmatrix}$

$h_\theta(x) = 0.7$.

Tell patient 70% chance of tumor being malignant.

$h_\theta(x) = P(y=1 | x; \theta)$ → "probability that $y=1$ given $x$, parameterized by $\theta$"

Logist Regression - hypothesis function contd.

$$P(y=0 \mid x_i; \theta) + P(y=1 \mid x_i; \theta) = 1$$

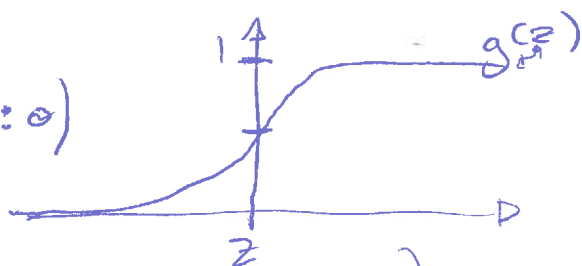$$\hookrightarrow P(y=0 \mid x; \theta) = 1 - P(y=1 \mid x; \theta)$$

# Notes - Week 3

## Logistic Regression - Decision boundary

Logistic regress.

$$h_\theta(x) = g(\theta^T x) = P(y=1 | x : \theta)$$

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

(graph of sigmoid function $g(z)$ with axis $z$, asymptote at 1)

---

Suppose predict "$y=1$" if $h_\theta(x) \geq 0.5$

predict "$y=0$" if $h_\theta(x) < 0.5$

$g(z) \geq 0.5$
whenever
$z \geq 0$
$\parallel$
$\theta^T x$

Similarly

If: $g(z) < 0.5$ then.

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

$\theta^T x < 0$

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

whenever $\theta^T x \geq 0$

---

## Decision Boundary

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2) \quad \therefore \theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

$$\parallel \qquad \parallel \qquad \parallel$$
$$-3 \qquad 1 \qquad 1$$

$y=1$ region

$y=0$ region.    decision boundary

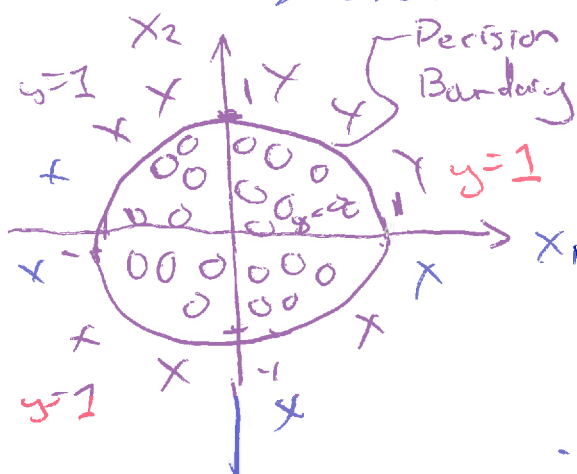ok so....

Predict "$y=1$" if $\overbrace{-3 + x_1 + x_2}^{\theta^T x} \geq 0$

$\longrightarrow x_1 + x_2 \geq 3$.

decision boundary is a property of the hypothesis (and the parameters there of) and not of the data set.

## Logistic Regression - Decision Boundary Contd...



Non linear decision boundaries

consider our hypothesis looks like this

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

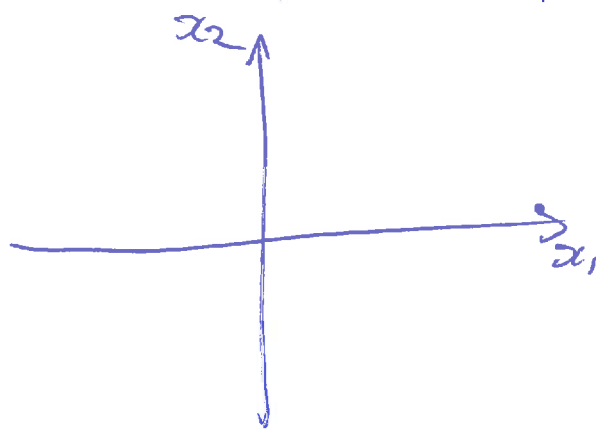$$\theta_0 = -1 \qquad \theta_2 = 0 \qquad \theta_3 = 0 \qquad \theta_4 = 1 \qquad = 1$$

∴ our parameters matrix $\theta$ looks like

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

**problem:** predict "$y=1$" if $-1 + x_1^2 + x_2^2 \geq 0$.

$$\boxed{x_1^2 + x_2^2 = 1}$$

**more complex example:**



suppose

$$h_\theta(x) = g(\theta_0 + \theta_1 x + \theta_2 x_2 + \theta_3 x_1^2$$
$$+ \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 +$$
$$\theta_6 x_1^3 x_2 + \dots)$$

# Week 3- Notes

## Logistic Regression Model — Cost function

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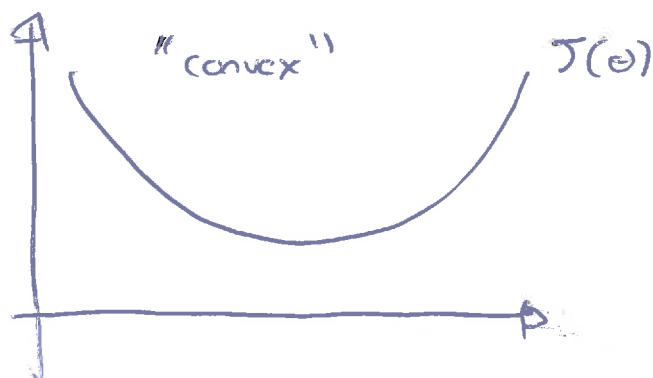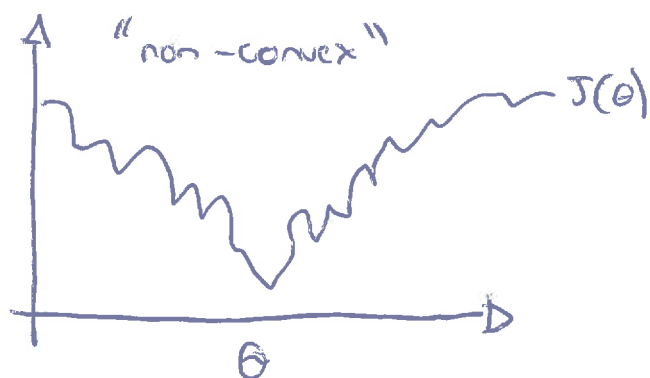
m examples

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \cdots \\ x_n \end{bmatrix} \qquad x_0 = 1, \; y \in \{0, 1\} \qquad h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameters $\theta$?

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

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



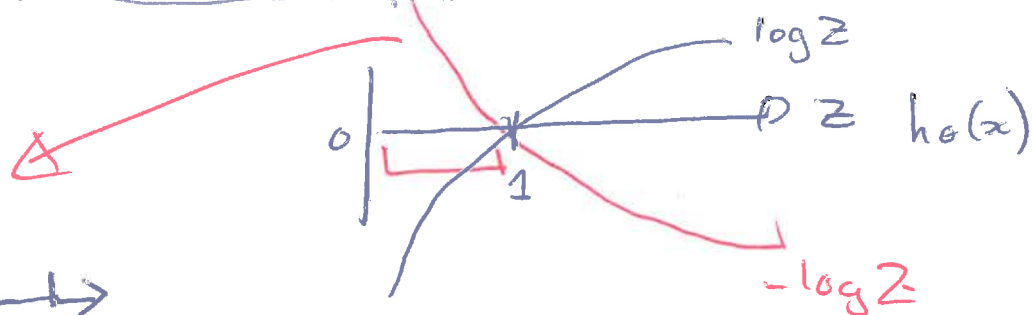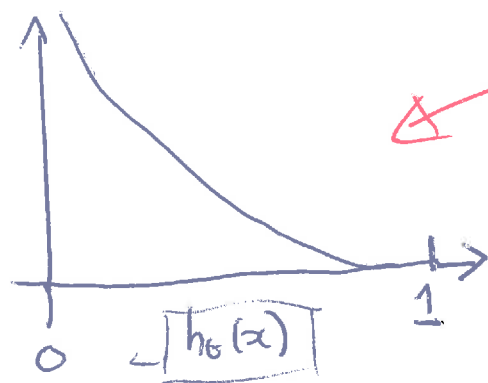"non-convex" $J(\theta)$



"convex" $J(\theta)$

# Week 3 - Notes

## Logistic Regression Model — Cost function

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

if $y=1$

cost function $y=1$   $\boxed{-\log(h_\theta(x))}$   $\log z$



$\text{Cost}=0$ if $y=1$, $h_\theta(x)=1$
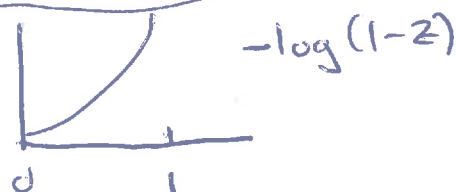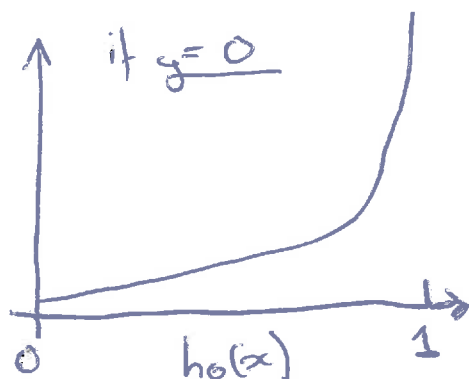
but as $h_\theta(x) \rightarrow 0$

    cost $\rightarrow \infty$

Captures intuition that if $h_\theta(x)=0$

(predict $P(y=1 | x; \theta) = 0$ but $y=1$
   & probability

we'll penalize learning algorithm by a very large cost.

If $y=0 \rightarrow \boxed{-\log(1-h_\theta(x))}$ — cost

$-\log(1-z)$

if $y=0$

# Week 3 – Notes

## Logistic Regression – Simplified cost function and gradient descent

logistic regression cost function

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

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

this is an easier way of writing the following

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

$$Cost(h_\theta(x), y) = -y\log(h_\theta(x)) - (1-y)\log(1-h_\theta(x))]$$

switches terms on and off when needed

derived from stratistics using the principal of maximum likelyhood estimation.

To fit parameter $\theta$:

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

To make a prediction given new $x$:

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

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

( probability that $y=1$, given input $x$ and parameterized by $\theta$)

# Week 3 - Notes

## Logistic Regression - Simplified cost function and gradient descent

so using gradient descent on our cost function:

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

want min $J(\theta)$:

repeats

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

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

(simultaneously ydate all $\theta_j$)

feature scaling works for logistic regression too

# Week 3 - Notes

## Logistic Regression - Advanced optimization

ok so we have the situation

cost function $J(\theta)$ Want $\min_\theta J(\theta)$

Given $\theta$, we have code that can compute

$\longrightarrow J(\theta)$ ✓

$\longrightarrow \left|\dfrac{\partial}{\partial \theta_j} J(\theta)\right|$     (for $j = 0, 1, \ldots, n$)

Gradient descent. (again)

Repeat {

$\longrightarrow \theta_j := \theta_j - \alpha \left|\dfrac{\partial}{\partial \theta_j} J(\theta)\right|$

}

---

Optimization algorithms

- Gradient descent
- Conjugate gradient
- BFGs
- L-BFGs

   ↳ do not need to understand to apply....

Advantages

- No need to manually pick $\alpha$ ✓
- often faster than gradient descent.

Disadvantages

- More complex.

## Logistic Regression — Advanced Optimization | Cont'd

Advanced optimization example...

$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$  ∴ values that minimise $J(\theta) =$
$\theta_1 = 5, \ \theta_2 = 5$

Matlab code

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

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

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

```
function [jVal, gradient]
   = costFunction (theta)
jVal = (theta(1)-5)^2 + ....
         (theta(2)-5)^2 ;

gradient = zeros (2, 1);
gradient (1) = 2 * (theta(1) - 5);
gradient (2) = 2 * (theta(2) - 5);
```

```
options = optimset ('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros (2,1);
[optTheta, functionVal, exitFlag] ...
     = fminunc (@costfunction, initialTheta, options);
```

octave notation to reference value is!

theta $= \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$ — theta (1)
— theta (2)
— theta (n+1)

```
function [jVal, gradient] = costFunction (theta)
   jVal = [code to compute J(θ)]  ✓
   gradient (1) = [ "   "   "  ∂/∂θ₀ J(θ)]
   gradient (2) = [ "   "   "  ∂/∂θ₁ J(θ)]
   ⋮
   gradient (n+1) = [code to compute ∂/∂θₙ J(θ)];
```

where:

gradient (1) = $\left[\cdots \frac{\partial}{\partial \theta_0} J(\theta)\right]$

gradient (2) = $\left[\cdots \frac{\partial}{\partial \theta_1} J(\theta)\right]$

gradient (n+1) = $\left[\text{code to compute } \frac{\partial}{\partial \theta_n} J(\theta)\right];$

# Week 3 - Notes

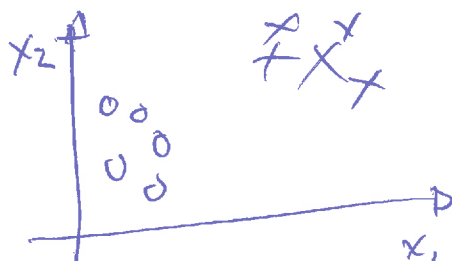## Multiclass classification: One vs all

Example:
let's say you want to automatically tag your email with

Email tagging: Work, Friends, Family, Hobby
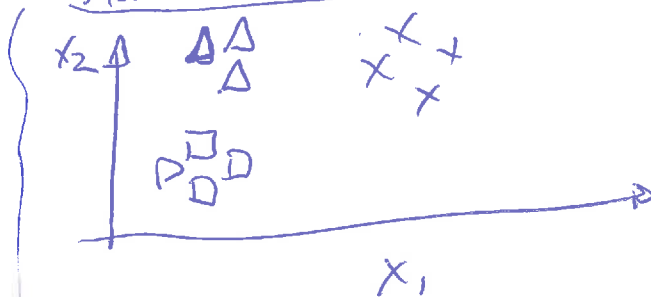$y=1$      $y=2$      $y=3$      $y=4$

Medical diagrams: Not ill, cold, Flu
$y=1$      2      3

Weather: Sunny, Cloudy, Rain, Snow
$y=$   1      2      3      4

---

### Binary Classification



### Multi-class classification



---

### One vs. All



essentially mapping this data to a fake data set like so

assigned value of "1"

essentially $h_\theta^{(1)}(x)$

$\rightarrow P(y=1 \mid x; \theta)$.

assigned value of "0"

class 1: △

class 2: □

class 3: ✕

$h_\theta^{(2)}(x)$

$h_\theta^{(3)}(x)$.

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

that is the probability that $y=i$ given $x$ and parameterized by $\theta$

# One-vrs-all

The basic idea is :

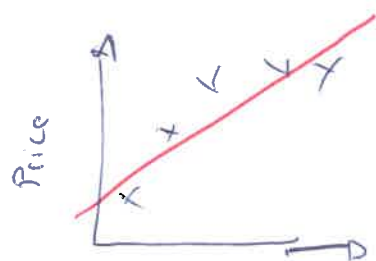Train a logistic regression classifier $h_\theta^{(i)}(x)$ for each class $i$ to predict the probability that $\underline{y=i}$

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

$$\max_i h_\theta^{(i)}(x)$$

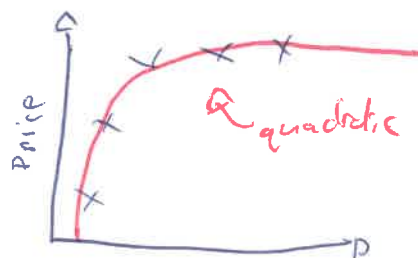here for example run all 3 classifiers on the input $x$, then pick class $i$ that maximizes the 3.

## Regularization — the problem of overfitting

tries 'too hard to fit the training set'
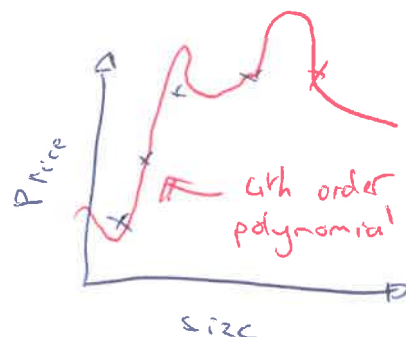
example: linear regression (housing prices)



Price vs Size

$\theta_0 + \theta_1 x$



Price vs Size — quadratic

$\theta_0 + \theta_1 x + \theta_2 x^2$

"Just right"



Price vs size — 4th order polynomial

$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$ ← 4th order polynomial

This example shows what is known as "overfitting" or "high variance"

→ This (the red line) is a bad model of the housing price data, this is known as "underfitting" or "high bias"

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2 \approx 0),$$ but fail to generalize to new examples (predict prices on new examples)

# Week 3 – Notes

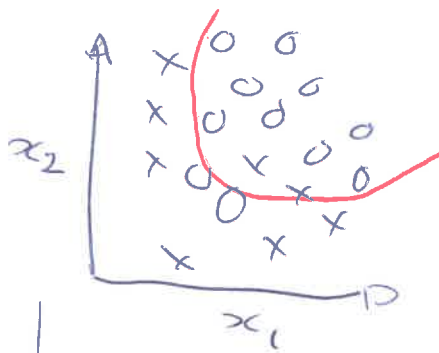## Regularization – the problem of overfitting
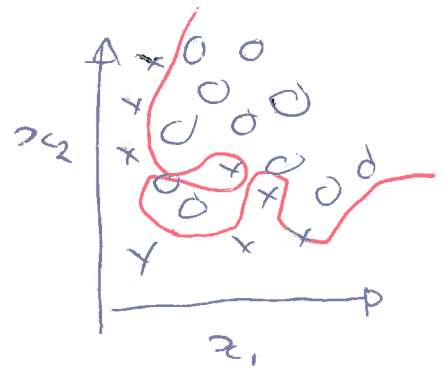
example: logistic regression



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

($g$ = sigmoid function)

↳ "underfit"

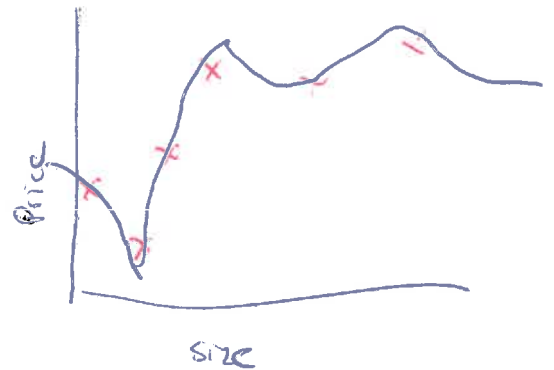$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2)$$

$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \theta_6 x_1^3 x_2 + \dots$$

↳ "overfit"

---

## Addressing overfitting

$x_1$ = size of house
$x_2$ = # of bedrooms
$x_3$ = # of floors
$x_4$ = age of house
$x_5$ = average income of neighborhood
$x_6$ = kitchen size
$\vdots$
$x_{100}$



too many features, not enough training data, overfitting can become a problem.

---

## Addressing overfitting:

Options:

① Reduce number of features
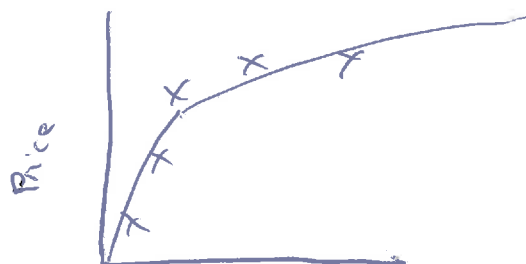→ manually reduce which features to keep
→ model selection algorithm.

② Regularization
- Keep all the features but reduce magnitude / values of parameters $\theta_j$.
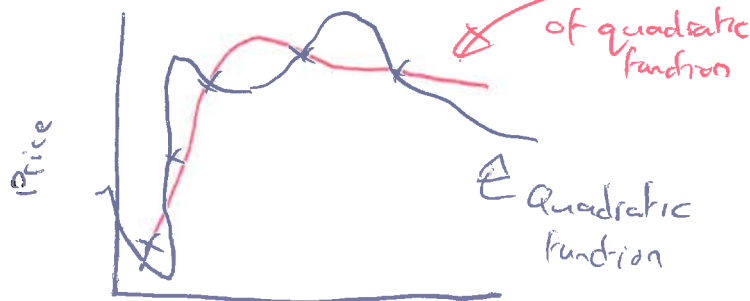- works well when we have a lot of features, each of which contributes a bit to predicting $y$

## Regularization - Cost function

### Intuition



Price (y-axis), Size of house (x-axis)

$$\theta_0 + \theta_1 x + \theta_2 x^2$$

penalized version of quadratic function

Quadratic function

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make $\theta_3, \theta_4$ really small

$$\rightarrow \min_\theta \frac{1}{2m} \sum_{i=1}^m \left(h_\theta(x^{(i)}) - y^{(i)}\right)^2 + 1000\theta_3^2 + 1000\theta_4^2$$

$$\therefore \quad \theta_3 \approx 0 \qquad \theta_4 \approx 0$$

### Regularization

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

$\quad \rightarrow$ "Simpler" hypothesis

$\quad \rightarrow$ Less prone to overfitting.

$$\boxed{\theta_3, \theta_4 \approx 0}$$

Housing:
- Features: $x_1, x_2, \ldots, x_{100}$
- Parameters: $\theta_0, \theta_1, \theta_2, \ldots, \theta_{100}$

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

modify cost function to shrink all parameters

$\theta_1 \ldots \theta_{100}$

with the exception of

$$\boxed{\theta_0}$$

## Regularization - Cost function

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

regularization term
2nd goal

1st goal goal

regularization parameter

$\min_\theta J(\theta)$

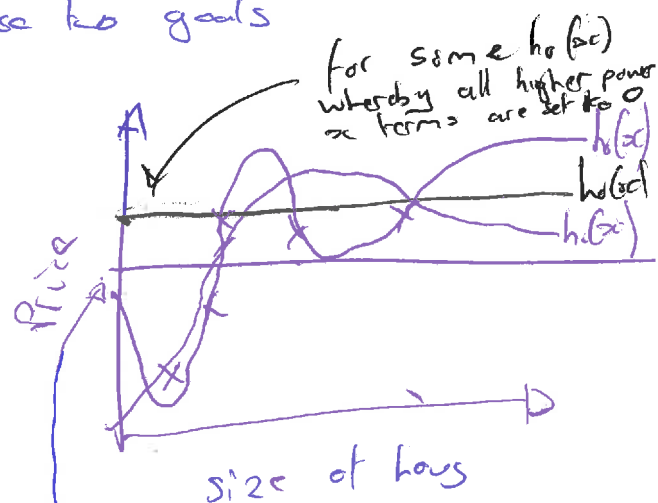$\lambda$ controls a trade off between two different goals.

1st goal: We would like to fit the the training data well, by minimizing the error.

2nd goal: We want to keep the parameters small

$\boxed{\lambda}$ controls the trade off between these two goals

for some $h_\theta(x)$ whereby all higher power $x$ terms are set to 0

$h_\theta(x)$
$h_\theta(x)$
$h_\theta(x)$

Price

size of haus

if $\lambda$ is very very large

$\theta_1, \theta_2, \theta_3, \theta_4 \approx 0$ (approach)

∴ for a hypothesis e.g

$h_\theta(z) = \theta_0 + \cancel{\theta_1 x} \cancel{\theta_2 x^2} \cancel{\theta_3 x^3}$

↳ approximates a straight line ie this is an example of "underfitting"

# Week 3 — Notes

## Regularization — Linear Regression

Gradient descent

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_0^{(i)}$$

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

$$(j = 0, \ 1, 2, 3, \ldots, n)$$

### Combining the two yields:

$$\theta_j := \theta_j \times \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \left[ \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \right] \quad \text{remains relatively the same}$$

$$1 - \alpha \frac{\lambda}{m} < 1$$

▷ So here we are multiplying $\theta_j$ by a number a little bit less than one (for all $\theta_i$)...

## Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}$$

$m \times (n+1)$

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

m — dimensional vector

$$\rightarrow \min_\theta J(\theta) \qquad \frac{\partial}{\partial \theta_j} J(\theta) \stackrel{\text{set}}{=} 0$$

$$\rightarrow \theta = \left( X^T X + \lambda \begin{bmatrix} 0 & 0 & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

$(n+1) \times (n+1)$

E.g  $n = 2$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

minimises $\theta$ when we are not using regularization

# week 3 - Notes

Non -invertibility (optional / advanced)

~~non inver~~

normal equation, suppose $m \leq n$
(#examples)  (#features)

$$\Theta = (X^T X)^{-1} X^T y$$

$\hookrightarrow$ non-invertible / singular

pinv. thv.

if $\lambda > 0$,

$$\Theta = \left( X^T X + \lambda \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix} \right)^{-1} \cdot X^T y$$

(don't get this)

## Regularized Logistic Regression



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

→ regularized boundary.

### Cost function

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

$$+ \frac{\lambda}{2m}\sum_{j=1}^{n}\theta_j^2 \Bigg|$$

$$\underbrace{\qquad}_{\text{regularization term}}$$

$\theta_1, \theta_2, \dots \theta_n$.

### implementing regularized gradient descent (algorithm for regularized logistic regression)

$$\theta_0 = \theta_0 - \alpha \frac{1}{m}\sum_{i=1}^{m}\left(h_\theta(x^{(i)}) - y^{(i)}\right)x_0^{(i)}$$

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

$$j = (0), 1, 2, 3, \dots n. \qquad \to \frac{\partial}{\partial\theta_j}J(\theta)$$

hypothesis is different → using sigmoid function ...

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

## Regularized Logistic Regression

Advanced optimization — Matlab code.

```
function [jVal, gradient] = costFunction (theta)

    jVal = [code to compute J(θ)]
             (cost function with regularized term)

    gradient(1) = [code to compute ∂/∂θ₀ J(θ)];
             └→ algorithm for θ₀

    gradient(2) = [code to compute ∂/∂θ₁ J(θ)];
             ⋮

    gradient(n+1) = [code to compute ∂/∂θₙ J(θ)];
```

need to create a vector:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$ — theta(1)
— theta(2)
— theta(n+1)

cost function needs to return jVal.

fminunc (@costFunction