# Aantekeningen nn + colleges

Neurons and the brain

Neuron network mimic the brain

Auditory cortex → ...

Neuron in the brain → neuron network

Sigmoid (logistic) activation function

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

$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$   $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$   weights = parameters

$\to h_\theta(x)$

Neural network

$a = x$

layer 1 input layer   layer 2 hidden layer   layer 3 output layer

$a_1^2 = g(\theta_{10}^1 x_0 + \theta_{11}^1 x_1 + \theta_{12}^1 x_2 + \theta_{13}^1 x_3)$

$\theta_j$ : matrix of weights controlling function of mapping from layer $j$ to layer $j+1$

Forward propagation. Vectorial implementation

$a_1^{(2)} = g(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$

$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$   $z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$

$a^{(2)} = g(z^{(2)})$

$z^{(2)} \in \mathbb{R}^3$   $a^{(2)} \in \mathbb{R}^3$

$a_0^{(2)} = 1 \to a^{(2)} \in \mathbb{R}^4$

$z^{(3)} = \theta^{(2)} a^{(2)}$   $g(z^{(3)})$

$h_\theta(x) = a^{(3)} = g(z^{(3)})$
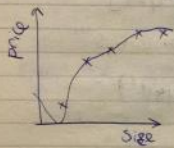
$a = x$

video 1 : introduction , kies juiste algorithme

$$J_{cv} = J_{test}$$

Stel je hebt Linieor regresion Algo met:
costfunction : $J(\theta)$

hypothesis test op nieuwe dataset → error
op voorspelling te hoog wat kun je doen:
- meer trainings sampels
- smaller features $X_1, X_2, X_3 \dots X_N$
- try to get additional features.
- adding polynomial $(X_1^2, X_2^2, X_1 X_2, etc.)$
- decrease lamda
- undecrease lamda

video 2 : Evaluating a hypothsis



Fails to generaliz to new exaples not in trainings set.

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

plotten niet mogelyk bij veel features

data spliten in 2 delen
- 70% traing set $(x^i, y^i)$
- 30% test set $(x_{test}^i, y_{test}^i)$
  $m_{test}$

heel erg overfit → lower training error
→ higher test error

beter om te shuffelen

train /test procedure   $J(\theta)$ 70%
- learn $\theta$ from training data (minimalise traing error)
- compute test set error

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{c=1}^{m_{test}} (h_\theta(x_{test}^{(c)}) - y_{test}^{(c)})^2$$

## traing /test procedure logistic regression

→ learn parameter $\theta$ from traing data
→ compute test set error:

$$J_{test}(\theta) = -\frac{1}{n_{test}} \sum_{i=1}^{n_{test}} y_{test}^{(i)} \log(h_\theta(x_{test}^{(i)})) + (1-y_{test}^{(i)}) \log(h_\theta(x_{test}^{(i)}))$$

→ misclassification error (0/1 classification):

$$err(h_\theta(x), y) = \begin{cases} 1 & \text{if } h_\theta(x) \geq 0.5, \ y=0 \\ \text{or if } h_\theta(x) < 0.5 & y=1 \end{cases} \text{error}$$
$$0 \text{ otherwise}$$

test error: $\frac{1}{n_{test}} \sum_{i=1}^{n_{test}} err(h_\theta(x_{test}^{(i)}), y_{test}^{(i)})$

---

→ for applying machine learning

↳ model selection problems (buide wits kiezen)

overfitting
↳ traing set error gaan goeie
  overspelling hiervoor hypothesis

Stel je met kiezen welke graads
  polynoom je moet weten

$d=1 \to 1: h_\theta(x) = \theta_0 + \theta_1 x \to \theta^{(1)}$
$d=10 \to 10: h_\theta(x) = \theta_0 + \theta_1 x + \ldots + \theta_{10} x^{10} \to \theta^{(10)}$

$d = \text{degree of polynomial} \to \text{niet zorijk als je op test set}$
  does

$\theta^{(1)} \to J_{test}(\theta^{(1)})$
$\theta^{(2)} \to J_{test}(\theta^{(2)})$
$\theta^{(10)} \to J_{test}(\theta^{(10)})$

kijken hoe goed polynoom past

$(x^{(1)}, y^{(1)})$
$(x^{(m)}, y^{(m)})$

indeler in 3
  stukken
  data set

→ training set
  60%

cross
validation
(cv)
20%

test set
20%

$(x_{cv}^{(1)}, y_{cv}^{(1)})$
$(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

$(x_{test}^{(1)}, y_{test}^{(1)})$
$(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

# Photo 1

training error
$\quad \hookrightarrow J_{train}(\theta) = \frac{1}{2m}\sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]^2$

cross validation error
$\quad \hookrightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}}\sum_{i=1}^{m_{cv}} (h_\theta(x^{(i)}) - y^{(i)})^2$

test error
$\quad \hookrightarrow J_{test}(\theta) = \frac{1}{2m_{test}}\sum_{i=1}^{m_{test}} [h_\theta(x^{(i)}_{test}) - y^{(i)}_{test}]^2$

model selection

1. $h_\theta(x) = \theta_0 + \theta_1 x \rightarrow \min_\theta J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
   $h_\theta(x) = \theta_0 + \theta_1 x, \theta_2 x^2 \rightarrow \min_\theta J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
   $\vdots$
   $\theta^{10}: J_{cv}(\theta^{(10)})$

Estimate generalization error for test $J_{test}(\theta^{(4)})$

cross validation nodig je een extra parameter toe
$J_{cv}(\theta)$ test is hoger dan $J_{test}$

# Photo 2

regressing bias vs. variance problems

training error
$\frac{1}{2}[h_\theta(x^{(i)}) - y^{(i)}]$

high bias → "just right" high variance
(underfit) → (overfit)



cross validation error $J_{cv}(\theta) = \frac{1}{2m_{cv}}\sum_{i=1}^{m_{cv}} [h_\theta(x^{(i)}) - y^{(i)}]^2$
$J_{train}(\theta) = \frac{1}{2m}\sum_{i=1}^m [h_\theta(x^{(i)}) - y^{(i)}]^2$

if $J_{cv}(\theta)$ or $J_{test}(\theta)$ too high

bias (underfit problem):
$J_{train}(\theta)$ = high
$J_{cv}(\theta)$ = high
$J_{cv}(\theta) \approx J_{train}(\theta)$

variance (overfit problem):
$J_{train}(\theta)$ = low
$J_{cv}(\theta) \gg J_{train}(\theta)$



# Photo 3

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

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

large $\lambda$          unavailable $\lambda$        small $\lambda$
$\lambda = 10000$        "just right"                $\lambda = 0$ overfit
$\theta_1 \approx 0, \theta_2 \approx 0$
underfit

how to choose
e.g. $\lambda = 0, 0.01, 0.02, 0.04 \ldots$
$\lambda = 0.001, 0.002, 0.004 \ldots$
$\lambda = 10$



$J_{train}(\theta)$
$J_{cv}(\theta)$
$J_{test}(\theta)$

## learning curves

→ $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

→ $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x^{(i)}) - y^{(i)})^2$

$J_{cv}(\theta)$

$J_{train}(\theta)$
m (training set size)

$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2$

### high bias

as m increases fit not gonna get better

$J_{cv}(\theta)$

$J_{train}(\theta)$
m (training set size)

high bias → $J_{train}(\theta)$ en $J_{cv}(\theta)$ dicht bij elkaar
↳ more words once eraan
als m groter wordt tight eraan feeling steeds, getting more training data will not help

---

## high variance

$J_{cv}(\theta)$

gap

$J_{train}(\theta)$
m (training set)

$h_\theta(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{...}$

if a learning algorithm is suffering from high variance
getting more training data is likely to help.

Summary → high bias:
  ↳ low training set: $J_{train}(\theta)$ low, $J_{cv}(\theta)$ high
  ↳ large training set: $J_{train}(\theta)$ en $J_{cv}(\theta)$ high
  → high variance:
  ↳ low training set: $J_{train}(\theta)$ low, $J_{cv}(\theta)$ high
  ↳ high training set: $J_{train}(\theta) < J_{cv}(\theta)$ work.
    $J_{cv}$ decrease $J_{train}(\theta) < J_{cv}(\theta)$

---

## debugging learning algorithm

- get more training examples → fixes high variance
- smaller set features → fixes high variance
- getting additional features → fixes high bias
- adding polynomial features → fixes high bias
- try decreasing $\lambda$ → fixes high bias
- try increasing $\lambda$ → fixes high variance

### Neural network and overfitting

→ small (fewer parameters)
  more prone to underfitting

→ Large (more parameters)
  more prone to overfitting
  use regularization $(\lambda)$ to address overfitting

Machine learning system design
- strategie (error analysis)
  ↳ planeren werk
    voorbeeld: spam classifier (log. regression)
    spam (1)  non. spam (0)
    features x: 100 words, indicate specific spam
    ↳ deal, buy, discount
    $x = \begin{cases} 1 & \text{if word } j \text{ appears} \\ 0 & \text{otherwise} \end{cases}$

    $x: \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \end{bmatrix} \begin{matrix} \text{deal} \\ \text{discount} \\ \text{android} \\ \text{buy} \end{matrix} \quad x \in \mathbb{R}$

    ↳ affect in the error. (feature vector)
    - collect lots of data         (→ honeypot projects
                                      ↳ not interesting)
    - develop sophisticated features based on email
      routing information
    - develop sophisticated features message · routing
    - spelling

- Error analysis
  ↳ Start simple algorithm → test op CV
    ↳ plot learning curves
    ↳ Error analyse: waar is de error
  Nu: 500 examples
    ↳ 100 misclassifieds → manually examine cases
      ↳ 1) what type is it ... → welke algo slaat
      ↳ 1) categorieën: pharma, 2, replica 4, steal/pass (53)
      ↳ 2) spelling 5, email routing 16, spam punctuation, 32
  - numerical evaluation
    ↳ gives back how goed je [stemming]
                                  → algo is
                                    zeker weten
                                    als resultaat
      without stemming  5% error
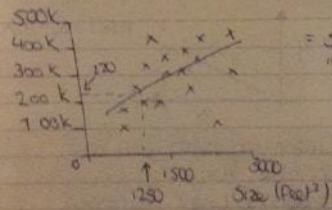      with stemming     3% error  ← wel goed van

# What is machine Learning.

↳ A Computer program is said to learn from ex[perience]
E with respect to some class of tasks T and perf[ormance]
measure P. if its performance at task in T, as
measured by P. improves with experience E.

- ▷ Supervised → given dataset en already know [the]
  output
- ▷ unsupervised → approach problems with Little or [no]
  idea what results look like

Regression → continuous valued output (example = p[rice])
Classification → Discrete valued output (example = 0[...])



price ($)
500k
400k
300k
200k
100k
↑ 1500
1250     3000
Size (feet²)

= Supervised learning, given th[e]
"right answer" for each e[xample]
= regression (real-val[ued])

data set = trainingset
m = number of training examples. (aantal rijen)
x's = input / features
y's = output / target
(x,y) = single training example
(x^i, y^i) = i^th training example (i = rij van de trainings[et])

---

Trainingset

Learning Algorithm
→ basically x's to y's
→ x's to y's

function made by Algorithm
→ Estimated price
[hypothesis]

how do we represent h?

x → h → y

h_θ(x) = θ_0 + θ_1 x
↳ shorthand = h(x)

y
x x x x x
x x x
h_θ(x)

- ▷ Linear regression with one
  variable = h(x)
  ↳ univariate Linear regression
  → one variable

how do we represent h?

$h_\theta(x) = \theta_0 + \theta_1 x$
↳ shorthand = $h(x)$

cost function.
hypothesis: $h(x) = \theta_0 + \theta_1 x$
$\theta_i$'s : parameters
how to choose $\theta_i$'s?

$h(x) = 1.5 - 0 \cdot x$

$h(x) = 0.5x$

$h_\theta(x)$

$\theta_0 = 1.5$      $\theta_0 = 0$      $\theta_0 = 1$
$\theta_1 = 0$        $\theta_1 = 1.5$    $\theta_1 = 0.5$

← how to make this line

→ how to make this line
[choosing $\theta_0, \theta_1$]

we choose $\theta_0, \theta_1$ so that $h_\theta$
is close to $y$ for our training s
$(x, y)$ → training examples

minimize $\theta_0, \theta_1$: $\frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$
$(x^{(i)}, y^{(i)})$

cost function $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

$\underset{\theta_0, \theta_1}{\text{minimize}}\ J(\theta_0, \theta_1)$

$\frac{1}{2m}$ worst hat niet heel groot

cost function

$J(\theta_0, \theta_1) = \frac{1}{2m} \sum (\hat{y}_i - y_i)^2 = \frac{1}{2m} \sum_{x, i} (h_\theta(x) - y_i)^2$

cost function (or) squared e

for fixed $\theta_1$, this is a
function of $x$

$h_\theta(x)$

function of the parameter

$J(\theta_1)$

$h_\theta(x) = \theta_1 x$     $\theta_1 = 1$

$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

$= \frac{1}{2m} \sum_{i=1}^{m} (\theta_1 x^{(i)} - y^{(i)})^2$

$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$

$J(1) = 0$

---

$h_\theta(x)$     $\theta_1 = 0.5$

$J(0.5) = \frac{1}{2m} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2]$
$= \frac{1}{2 \cdot 3} (3.5) = \frac{3.5}{6} \approx 0$ ca 58

$\theta_1 = 0$

$J(0) = ?$

$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2)$
$= \frac{1}{6} \cdot 14 = 2.3$

goal = minimize $J(\theta_1)$     $\theta_1 = 1$

→ gaat vraag blockgle

J functie lijkt welke het best past.

Contour plots

hypothesis    $h_\theta(x) = \theta_0 + \theta_1 x$

parameters    $\theta_0, \theta_1$

               → cross function

cost functions   $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

goal:   $\underset{\theta_0, \theta_1}{minimize}\ J(\theta_0, \theta_1)$

wanneer er een $\theta_0$ en $\theta_1$ is krijg je een 3D-parabo

Contour plot / contour figures



the sum of square minimum is the lowest

Gradient descent    $J(\theta_0, \dots \theta_n)$

                       $\underset{\theta_0 \dots \theta_n}{min}\ J(\theta_0, \theta_n)$

outline
    Start with some $\theta_0, \theta_1$   (say $\theta_0 = 0, \theta_1 = 0$)
    keep changing $\theta_0, \theta_1$ to reduce $J(\theta_0, \theta_1)$

Gradient descent algorithm
    repeat until convergence {
        $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$   (for $j = 0$ and $j = 1$)

temp 0 : $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$
temp 1 : $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   ← update correctly
    $\theta_0$ = temp 0
    $\theta_1$ = temp 1

α · waarde bekijken (optional.)

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

- if α is to small

$J(\theta_1)$

$\theta_1$

- if α is to large
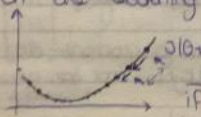
$\theta_1 \rightarrow$

slope = 0

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1) = 0$$
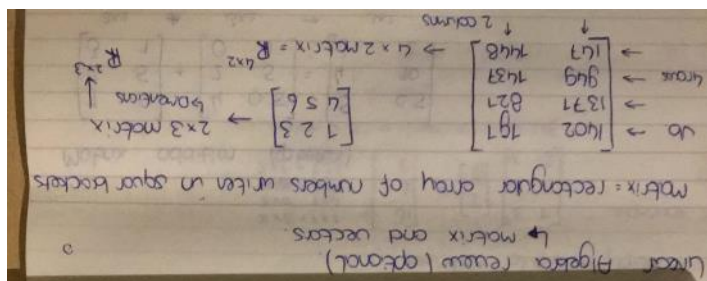
$\theta_1$    $\theta_1$ at local optima

$$\theta_1 = \theta_1 - \alpha \cdot 0 \rightarrow \theta_1 := \theta_1$$

Gradient descent can converge to a local minimum, even
with the learning rate α fixed

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

$J(\theta_1)$

derivatief will get smaller by every step
if we approach local minimum gradient
descent take smaller steps automaticall

Linear Algebra (review / optional)
↳ matrix and vectors

matrix = rectangular array of numbers writen in squar brackets

$\begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$ ← yo
← house
→ 4×2 matrix = $R^{4\times2}$
↑
2 columns

$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$ ← 2×3 matrix $R^{2\times3}$
↑
proventies

herhaling

Gradient descent algorithm ⌐   Linear Regression model
$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$    $h_\theta(x) = \theta_0 + \theta_1 x$
(for $j=1$ and $j=0$)    $J(\theta) = \frac{1}{2m} \sum_{i=1}^{M} (h_\theta(x^{(i)}) - y^{(i)})^2$

$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^{M} (h_\theta(x^{(i)}) - y^{(i)})^2$

$= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum_{i=1}^{M} (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$

$\theta_0$  $J=0: \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{M} \sum_{i=1}^{M} (h_\theta(x^i) - y^i)$
$\theta_1$  $J=1: \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{M} \sum_{i=1}^{M} (h_\theta(x^i) - y^i) \cdot x^{(i)}$

Gradient descent Algorithm
repeat until convergence {
$\theta_0 := \theta_0 - \alpha \left[ \frac{1}{M} \sum_{i=1}^{M} (h_\theta(x^i) - y^i) \right]$ → $\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$
$\theta_1 := \theta_1 - \alpha \frac{1}{M} \left[ \sum_{i=1}^{M} (h_\theta(x^i) - y^i) \cdot x^i \right] \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

Convex functie = de 3D parabol bol vormige.

'Batch' Gradient Descent = Each step of gradient descent
uses all the trainings examples (kijken naar een gehele
trainings set)
↳ $\sum_{i=1}^{M} (h_\theta(x^i) - y^i)$

Aantekeningen Pagina 11