

Lecture 1

Cezar Ionescu

11/05/2019

DEPARTMENT FOR
CONTINUING
EDUCATION



- course materials on GitHub
 - lectures notes will generally be available after the lecture, though drafts may appear before
 - exercises will also be posted on GitHub

Administrative matters

- register for CATS points
 - assessment: homework to be handed in at the next meeting (or sent to me by email, or to PPWeekly)
 - every piece of homework is pass or fail

Administrative matters

- course takes place Saturdays 10:00-12:30 at Ewert House
 - exceptions: **no class** on the **4th of May** and **1st of June!**

- main text: *Machine Learning*, Tom Mitchell, 1997
 - there are a couple of copies available in the ContEd library
 - you can buy it used for under 15 GBP on Amazon
 - but you should be able to complete the course using only the lecture notes

Types of learning

- types of learning:
 - by rote
 - conditioning
 - from experience
 - any others?

Why “machine learning”?

- reasons for *machine learning*:
 - programming is hard, it would be better if computers learnt by themselves
 - to study human learning (and intelligence)
 - perhaps we could then improve our own abilities to learn and to teach

Approaches to ML

- two main approaches:
 - modelling how we think and learn, without caring about the underlying physiological mechanisms
 - modelling the underlying physiological mechanism, without caring how they lead to thinking and learning

- **Definition** (Mitchell, p. 2) A computer program is said to **learn** from experience E with respect to some class of tasks T and performance measure P , if its performance at tasks in T , as measured by P , improves with experience E .

Sets and functions

Notation:

- $f : \text{In} \rightarrow \text{Out}$
- $f(x)$
- $x \in \text{In}, f(x) \in \text{Out}$

Black boxes versus functions

- is `Random.choice(seq)` a function?
- we use \mapsto instead of \rightarrow when we need to distinguish black boxes from “real” functions

Mathematical representation

- the set of tasks: Task
- the set of experiences: Experience
- measure of performance: $\text{perf} : \text{Task} \rightarrow \mathbb{R}$
- machine learning system: $\text{learn} : (\text{Task}, \text{Experience}) \rightarrow \text{Task}$
 - learning means improving with experience (according to 'perf):

for all $t \in \text{Task}$, for all $e \in \text{Experience}$:
 $\text{perf}(t) \leq \text{perf}(\text{learn}(t, e))$

Example: checkers learning

- Task = Board \rightarrow Move
- experience : Task \rightarrow List Game
 - we assume play : (Play, Play) \rightsquigarrow Game
 - note the squiggly arrow in the type of play!

experience(t) = [play(t, t), play(t, t), ..., play(t, t)]

- perf : (Task, List Task) $\rightarrow \mathbb{R}$
- we need a function score : Game $\rightarrow \{0, 1\}$

perf (learner, [adv₁, ..., adv_n]) =
(sum [score(play(learner, adv₁)), ...,
score(play(learner, adv_n))]) * 100 / n

Example: self-driving car

- the task is to give steering commands based on sensor input:
 $\text{Task} = \text{Sensor} \rightarrow \text{Command}$
- the set of experiences: $\text{Experience} = \text{List} (\text{Sensor}, \text{Command})$
- performance: $\text{perf} : (\text{Task}, \text{Itinerary}) \rightarrow \text{Time}$
 - $\text{perf} (\text{learner}, \text{itinerary})$ measures how long the learner drives along the given itinerary before making a mistake

Homework

- Give a similar interpretation for the handwriting recognition problem (Mitchell, page 3):
 - Task T : recognizing and classifying handwritten words within images
 - Performance measure P : percent of words correctly classified
 - Training experience E : a database of handwritten words with given classifications
- You have to fill in
 - Task =
 - Experience =
 - perf :
 - something about how perf is computed

Concept learning

- idea: acquiring general concepts from examples
 - e.g., learn to recognise cats from images of animals
- what is a concept?
 - *nominalistic* view: the set of instances of the concept

Mathematical representation of concepts

- mathematically, we can identify a concept with a subset
 - e.g., X is the set of all images of animals, ' $C \subseteq X$ ' is the subset of images of cats
- subsets are in one-to-one correspondence with boolean-valued functions
 - $C \subseteq X$ can be replaced by $c : X \rightarrow \text{Bool}$ such that

$$\forall x \in X \quad c(x) = 1 \quad \text{iff} \quad x \in C$$

- Mitchell uses the functional view and defines:
 - **Concept learning:** inferring a boolean-valued function from training examples of input and output
- **Exercise:** Give an interpretation of concept learning as a learning task (i.e., identify the task, the experience, and the performance measure).

Notation

- the training data $\mathcal{D} = \{((x_1, c(x_1)) \dots, (x_n, c(x_n)))\}$
- the subset of negative training examples
 $\mathcal{D}_0 = \{(x, 0) \mid (x, 0) \in \mathcal{D}\}$
- the subset of positive training examples
 $\mathcal{D}_1 = \{(x, 1) \mid (x, 1) \in \mathcal{D}\}$

Example

- we want to learn the concept enjoyable : $\text{Day} \rightarrow \{0, 1\}$
- days are described via *attributes*:
 - $\text{Day} = (\text{Sky}, \text{Temp}, \text{Humidity}, \text{Wind}, \text{Water}, \text{Forecast})$
 - $\text{Sky} = \{\text{Sunny}, \text{Cloudy}, \text{Rainy}\}$
 - $\text{Temp} = \{\text{Warm}, \text{Cold}\}$
 - $\text{Humidity} = \{\text{Normal}, \text{High}\}$
 - $\text{Wind} = \{\text{Strong}, \text{Weak}\}$
 - $\text{Water} = \{\text{Warm}, \text{Cool}\}$
 - $\text{Forecast} = \{\text{Same}, \text{Change}\}$

Training data

Nr	Sky	Temp	Humidity	Wind	Water	Forecast	Enjoyable
1	Sunny	Warm	Normal	Strong	Warm	Same	Yes
2	Sunny	Warm	High	Strong	Warm	Same	Yes
3	Rainy	Cold	High	Strong	Warm	Change	No
4	Sunny	Warm	High	Strong	Cool	Change	Yes

Two problems

- Day contains 96 elements; there are 2^{96} concepts.
- We can represent a concept by a lookup table, but not if it's too big!
- The training data does not suffice to determine the concept we are looking for.

- The two problems force us to make two decisions:
 - ① How to represent **some** of the concepts (the hypotheses space)
 - ② Which hypothesis to pick.

- The assumptions under which we manage to learn the correct concept form **the inductive bias**.

Hypothesis space for the weather example

- each hypothesis is described by a tuple of
(Sky*, Temp*, Humidity*, Wind*, Water*, Forecast*), where
 $S^* = S \cup \{?, \emptyset\}$
- notation: $s \sim s^*$ iff $s = s^*$ or $s^* = ?$ (s matches s^*)

Let h be described by $(s^*, t^*, u^*, w_i^*, w_a^*, f^*)$. Then

$$h(s, t, u, w_i, w_a, f) = s \sim s^* \text{ and } t \sim t^* \text{ and } u \sim u^* \text{ and } \\ w_i \sim w_i^* \text{ and } w_a \sim w_a^* \text{ and } f \sim f^*$$

- Find-S solves problem 2 by choosing the *most specific* hypothesis that is consistent with the training data.
- Our hypothesis correspond to subsets. We have a natural ordering on subsets: \subseteq .

Find-S algorithm

```
-- input: training data  $\{(x_1, c(x_1)) \dots, (x_n, c(x_n))\}$ 
--          hypothesis set  $H$ 
h = min  $H$  -- "the" (or "a") smallest element of  $H$ 
for i in 1:n
    if  $c(x_i) = 0$ 
        then keep h
    else if  $x_i \in h$  then keep h
        else  $h = \min \{h' \in H \mid h \subseteq h' \text{ and } x_i \in h'\}$ 
-- output: "the" (or "a") most specific hypothesis
--          consistent with  $D$ 
```

Example

Work out how **Find-S** works on the weather example.

- If H contains all possible concepts, then the result of Find-S is D_1 .
- A bad situation for Find-S:
 - $X = \{a, b, c, d\}$, $H = \{\emptyset, \{a, b\}, \{a, c\}\}$, $D_1 = \{a\}$
- Another bad situation:

$X = \{a, b, c, d\}$, $H = \{\emptyset, \{a, b, c\}, \{a, b, d\}\}$,
 $D_1 = \{a\}$, $D_0 = \{c\}$

- The choice of H can avoid these problems, as in the weather example.

Property of Find-S

If Find-S works, then

$s = \text{Find-S}(D_0, D_1, H)$ implies

$D_0 \subseteq \neg s$, $D_1 \subseteq s$, and

for all $h \in H$, $D_0 \subseteq \neg h$ and $D_1 \subseteq h \Rightarrow s \subseteq h$

Exercise

What does Find-S (D_1 , D_0 , H) do?

A better Find-S

```
-- input: training data  $\{(x_1, c(x_1)) \dots, (x_n, c(x_n))\}$ 
--          hypothesis set  $H$ 
S = allMin H -- start with all smallest element of H
repeat until S no longer changes:
  for i in 1:n
    if  $c(x_i) = 0$ 
    then eliminate from S all  $\{s \mid x_i \in s\}$ 
    else for all  $s \in S$ 
      if  $x_i \in s$ 
      then keep s
      else replace s with allMin  $\{h' \in H \mid h \subseteq h' \text{ and } x_i \in h'\}$ 
-- output: all most specific hyp consistent with D
```


- why do we need to repeat the for loop?
- the algorithm terminates (why?)

Avoiding repeat

- we need to keep a record of the negative examples
- idea: do that in the same form as the record we keep for positive examples!
- this leads to the **Candidate-Elimination** algorithm