# Unit 5: Machine Learning & Al

Dave Abel

March 4th, 2016





#### Demo time!



# Overfitting

- Generalization is the goal of learning.
- Overfitting: when our classifier pays too much attention to noisy details of our data instead of the underlying trends.
- When we do well on classifying training data, but poorly on testing data)



# Overfitting

Michael Littman and Charles Isbell feat. Infinite Harmony Overfitting ML4LIFE Udacity Records















#### observation, reward



Memory

8

#### action



















#### Goal: Maximize reward!





#### Goal: Maximize reward!







Q: How can we provide reward in Mario so that an RL agent would be incentivized to win the game?







Q: How can we provide reward in Mario so that an RL agent would be incentivized to win the game?

Some ideas:

- Give Reward when Mario beats the level
- Give Reward when Mario's score goes up
- Punish when Mario dies
- Punish when Mario loses a mushroom
- Give Reward when Mario kills a Goomba
- etc.





#### Clicker Question!





# Clicker Question!



[A] Every time the agent moves a pawn, get punished

[B] Every time the agent moves a knight, get punished

[C] Every time the agent captures a piece, get reward

[D] If the agent wins, get reward.



## Clicker Answer!



[A] Every time the agent moves a pawn, get punished

[B] Every time the agent moves a knight, get punished

[C] Every time the agent captures a piece, get reward

[D] If the agent wins, get reward.



### Clicker Answer!

Could do more reward for better pieces.



VS.

[D] If the agent wins, get reward.







#### Clicker Answer!

**In general:** a trail of breadcrumbs is much easier for an agent than just getting reward at the end!

If we only give reward at the end, all of its behavior leading up to the end is basically arbitrary!



#### Problem: Maximize Reward







#### Problem: Maximize Reward



- Two things the agent can *learn*:
  - (A) The dynamics of the world!
    - What will my next observation be, if I act in this way?
  - (B) How do observations relate to reward?
    - Where's the positive reward located?









#### Rules:

- Costs one token to pull an arm. Given T tokens.
- Each slot machine will pay some amount of money on average.
- Maximize money earned with fewest arm pulls.
- **NOTE**: Tokens are *not* money. Tokens represent # pulls.





#### Rules:

- Costs one token to pull an arm. Given T tokens.
- Each slot machine pays some money *on average.* Could pay \$5, then \$15, then \$10 if it pays \$10 avg.

- Maximize money earned with fewest arm pulls.

- NOTE: Tokens are *not* money. Tokens represent # pulls.





The Problem: How do I (or the agent) know when I should *explore* (to gain more information), as opposed to *exploiting* the information I have?



3



???



















3



???













































3



\$2?

















|--|



3



\$2?

























???

















\$2?





























**Issue:** How do we trade off *exploring* new strategies with *exploiting* the information we have now?



**Cool note:** general model of decision making, in life!

- What city to live in?
- What apartment to move in to?
- Who to date/marry?
- What field to concentrate in?
- What food to order at a restaurant?
- What instrument to learn?



# Problem: Slot Machines

- INPUT: some number of slot machines, K, and T tokens for playing (so you get to pull T times).
- OUTPUT: Some amount of money made (want to be as high as possible!)





# Algorithm 1: Dora The Explora!



 Pull every arm equally often so you know which arm is probably the best!



**Issue:** How do we trade off *exploring* new strategies with *exploiting* the information we have now?



# Algorithm 1: Dora The Explora!



- Pull every arm equally often so you know which arm is probably the best!
- Explores way too much. (We never exploit!)





# Algorithm 2: Hoyt The Exployter!



- Pull each arm once.
- Then pull the arm that paid the most until you run out of tokens.



# Algorithm 2: Hoyt The Exployter!



- Pull each arm once.
- Then pull the arm that paid the most until you run out of tokens.
- Exploits way too much! (We barely explore..)



### The Goal

Q: How do we *perfectly* balance exploration and exploitation?





#### Biology and Reinforcement Learning





Lee, D., Seo, H., & Jung, M. W. (2012). Neural basis of reinforcement learning and decision making. Annual review of neuroscience, 35, 287.



#### One Scratch Note...



Let's Take a Look!



# Al Reflection

- Learning can be formalized as an algorithm
- Classification is learning rules to determine concepts, based on labeled examples.
- Reinforcement Learning is a (biologically inspired) super general model of learning based on rewards and punishment.
- An exciting road ahead!

