# AI in Logistics

**AUTUMN 2021**

**Christina Imdahl**

IE & IS, OPAC

# Schedule of Second Part

| Monday | Tuesday | Wednesday | Thursday | Friday |
|---|---|---|---|---|
| Oct 4 | Oct 5 | Oct 6 | Oct 7 | Oct 8 |
| | Introduction to RL ⭐ | | | Case Study Descriptives and Orientation |
| Oct 11 | Oct 12 | Oct 13 | Oct 14 | Oct 15 |
| | Reinforcement Learning – Key Concept | | | *(Homework: Implement RL)* *Momentum* |
| Oct 18 | Oct 19 | Oct 20 | Oct 21 | Oct 22 |
| | Inventory Management - Heuristics | No Availability | | *Case Study Benchmark* |
| Oct 25 | Oct 26 | Oct 27 | Oct 28 | Oct 29 |
| | Wrap-up / Case Study | | | Case Study |
| Nov 1 | Nov 2 | Nov 3 | Nov 4 | Nov 5 |
| | Submission Case | | | Case Presentation |

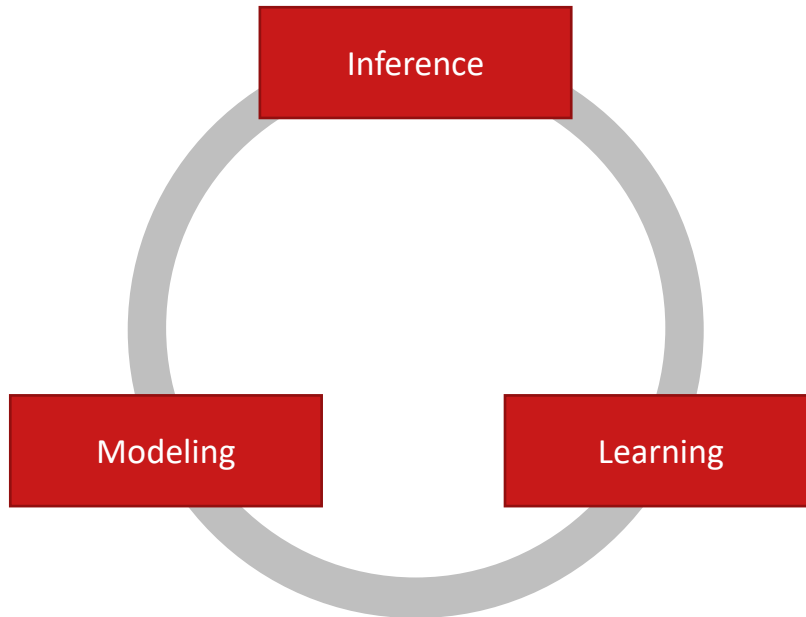AI in Logistics | Christina Imdahl

TU/e

# Objectives of Today

- Differentiate different ML techniques

- Learn about the basics of Reinforcement Learning

- Understand the interaction between model formulation and learning

- INTUITION on reinforcement learning

TU/e

# Agenda
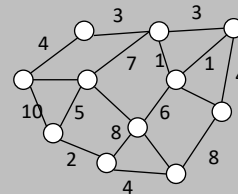
TU/e

# Paradigms

**TU/e**

# Modelling

Real World Problem

Mathematical Model

TU/e

# Inference

Full Model
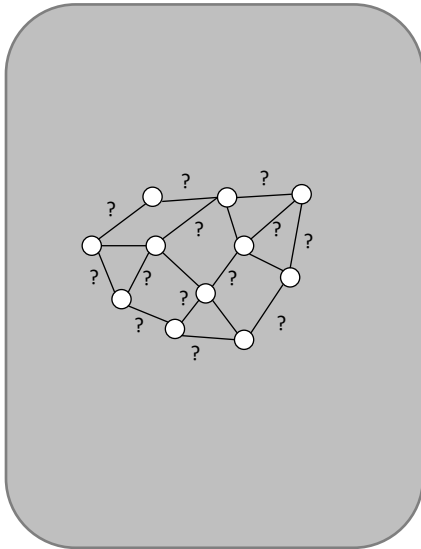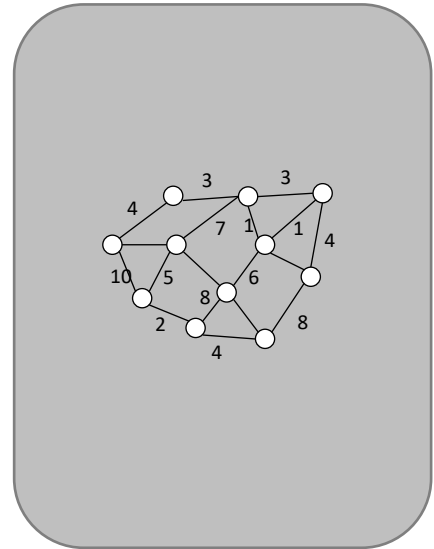
Inference



AI in Logistics | Christina Imdahl

TU/e

# Learning

Model



Learned Model

TU/e

# Overview Machine Learning

**Supervised Learning**

➡️

**Predict a value**
Input: Labeled Training Data
Output: Prediction Model
Type: Regression/Classification

**Unsupervised Learning**

➡️

**Identify patterns**
Input: Unlabeled Training Data
Output: Classes/Associations
Type: Clustering/Associations

**Reinforcement Learning**

➡️

**Find policy to optimize rewards**
Input: No Predefined Data
Output: Policy
Type: Reward-based

AI in Logistics | Christina Imdahl

TU/e

# Supervised Learning

**Labeled Training Set**

0

1

1

0

0

0

1

**Feature Set**

**Training**

Machine Learning Alg.

e.g. Support Vector Machines, Elastic Net, Lasso/Ridge Regression, Neural Networks

**Test Set**

+ Features

Prediction Model

**Predictions**

0.01     0.24     0.98

TU/e

# OM Applications
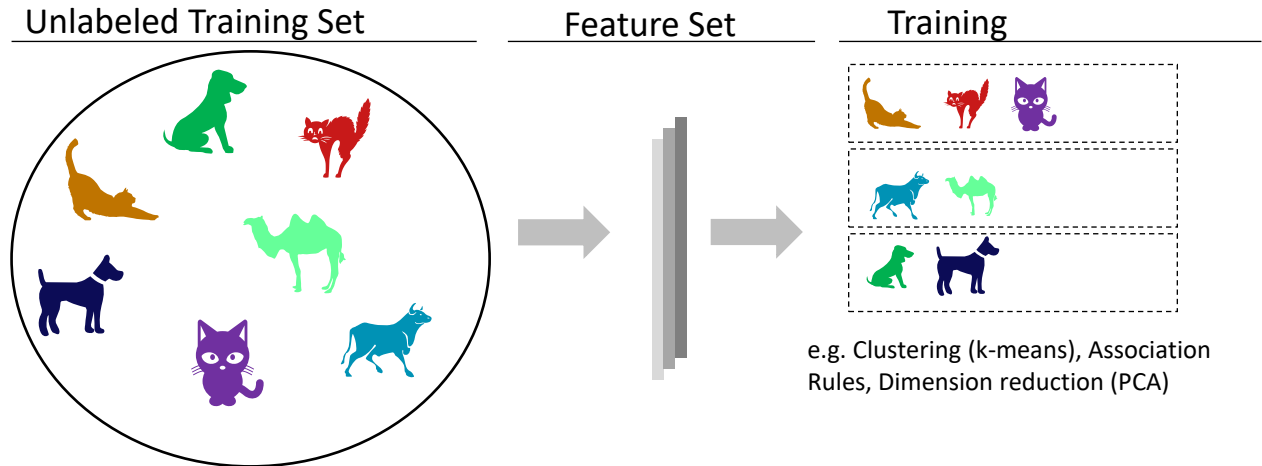
- Predicting Demand & Sales

- Predicting Machine Failures

- Predicting Warehouse Operations

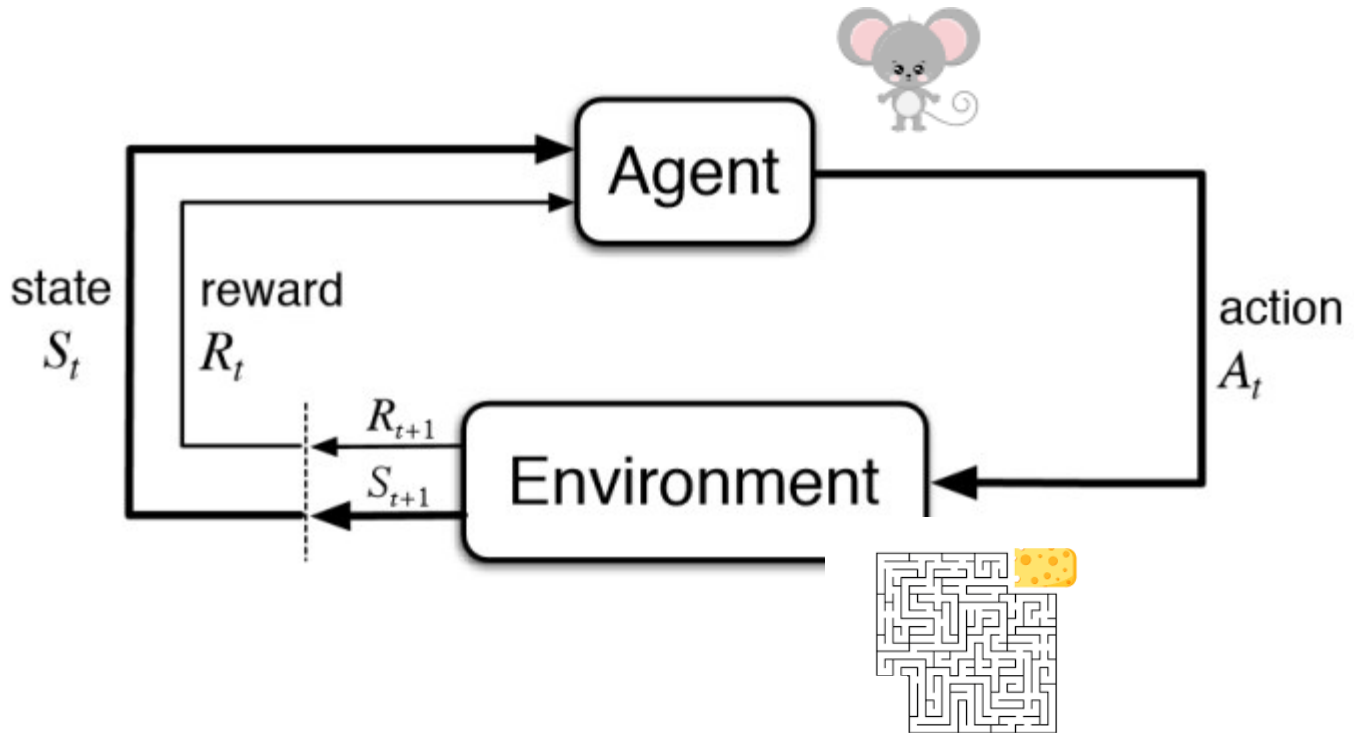- Predicting Decision-Maker Behavior

- Predicting Preferences

**TU/e**

# Unsupervised Learning

Unlabeled Training Set

Feature Set

Training

e.g. Clustering (k-means), Association
Rules, Dimension reduction (PCA)

AI in Logistics | Christina Imdahl

TU/e

# OM Applications

- Market Segementation/Targeted Advertisement

- Fraud Detection

- Basket Analysis

- Recommender Engines (Amazon „Other People Bought", Netflix Movie Recommendation)

**TU/e**

# Reinforcement Learning



state
$S_t$

reward
$R_t$

$R_{t+1}$

$S_{t+1}$

action
$A_t$

AI in Logistics | Christina Imdahl

**TU/e**

# Characteristics of RL

- There is no supervisor, only reward

- Feedback is delayed

- Time matters (sequentiell decision making)
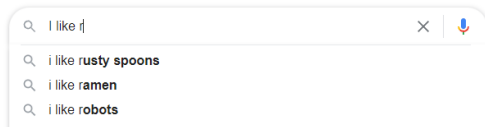
- Actions influence the data the agent receives

AI in Logistics | Christina Imdahl

TU/e

# RL breaktrough



AI in Logistics | Christina Imdahl

*„I thought AlphaGo was based on probability calculation and that it was merely a machine. But when I saw this move, I changed my mind. Surely, AlphaGo is creative."*

-- Lee Sedol

TU/e

# OM Applications



**Retail**

- Dynamic Pricing
- Anticipatory Shipping
- Item Descriptions
- Fraud Detection



**Robotics**

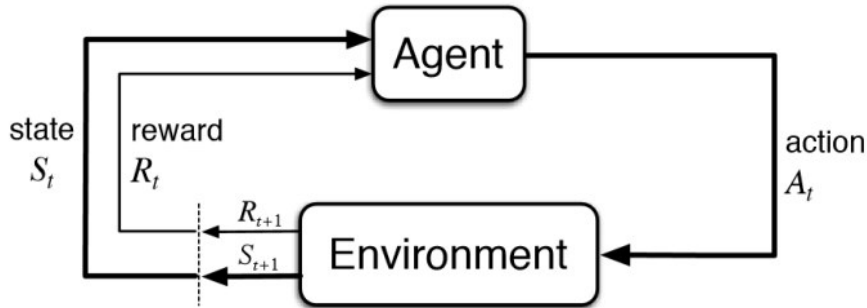- Picking in Warehouses
- Routing Robots
- Manufacturing



**Speech Analysis**

- Sentiment Analysis
- Chat robots
- Smart Home Tools

TU/e

# Agenda

AI in Logistics | Christina Imdahl

TU/e

# Again: Markov Decision Processes



Agent observes state in t: $S_t$
Agent takes action in t: $A_t$
Agent receives reward R in t+1: $R_{t+1}$
Agent moves to the next state: $S_{t+1}$

TU/e

# RL Terminology

Policy: The agent's behavior given a state a

State-Value Function: how good is a state (s)?

Action-Value Function: how good is a state-action pair (s,a)?

Model: *agent*'s representation of the environment.

AI in Logistics | Christina Imdahl

TU/e

# Policy

- Determines the agent's behavior given a state

- The policy maps from state to actions

- Deterministc policy: $a = \pi(s)$

- Stochastic policy: $\pi(a|s) = P(A_t = a|S_t = s)$

AI in Logistics | Christina Imdahl

TU/e

# State-Value function / Value function

- Value of a  state under $\pi$

- Used to evaluate how good/bad a certain state is

- Used to select between actions

- E.g.:

$$v_\pi(s) = \mathbb{E}_\pi(\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}(s_{t+k}, \pi(a|s_{t+k}), s'_{t+k}))$$

Disc. reward of following policy $\pi$

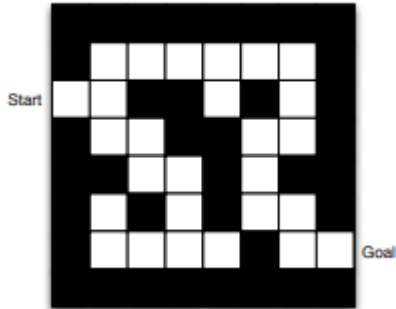AI in Logistics | Christina Imdahl

**TU/e**

# Action-Value function / Q-Function

- Value of an action under $\pi$

- Used to evaluate how good/bad a certain state-action pair is

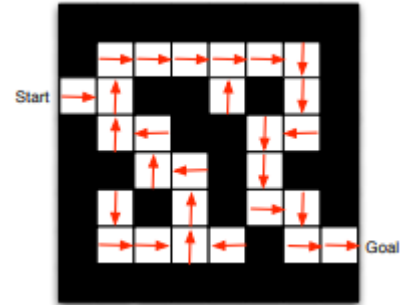- Used to select between actions

- E.g.:

$$q_\pi(a,s) = \mathbb{E}_\pi(\underbrace{R_{t+1}(s_{t+k}, a, s'_{t+k})}_{\text{Immediate reward taking a}} + \underbrace{\sum_{k=1}^{\infty} \gamma^k R_{t+k+1}(s_{t+k}, \pi(a|s_{t+k}), s'_{t+k}))}_{\text{Disc. reward of following policy } \pi}$$

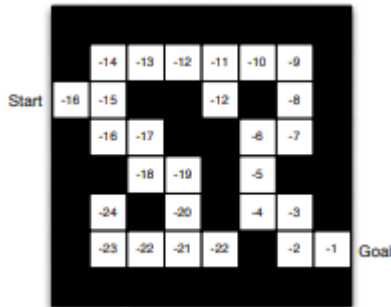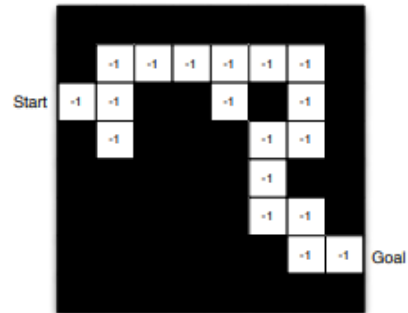AI in Logistics | Christina Imdahl

TU/e

# Example: MAZE

## The game



## The policy



## The value function



## The model



AI in Logistics | Christina Imdahl

**TU/e**

# Key goal of RL

> ## Determine the policy that maximizes cumulative reward

- Determine the policy directly (policy-based methods)
  - Policy gradient methods, e.g. REINFORCE

- Determine the value function (value-based methods)
  - E.g. Q-Learning, SARSA

TU/e

# Dynamic Programming

- Two methods of dynamic programming can be used to solve MDPs
  - Value Iteration and Policy Iteration


- Let's revisit them…

AI in Logistics | Christina Imdahl

**TU/e**

# Value Iteration

Update of values based on Bellman Eq.

$$v_{k+1}(s) = \max_a \sum_{s' \in S} P(s'|s,a) \underbrace{(R(s,a,s')}_{\text{Immediate Reward}} + \underbrace{\gamma v_{\pi_i,k}(s'))}_{\text{Disc. Value of Sucessor State}}$$

1. Randomly initiate value function
2. Update values by above equation until convergence
3. Optimal policy is the action ending in the state with the maximum value

TU/e

# Policy Iteration

1. Randomly initiate policy
2. Policy evaluation: evaluate value function for current policy until convergence

$$v_{\pi_i,k+1}(s) = \sum_{s' \in S} P(s'|s, \pi_i(s))(R(s, \pi_i(s), s') + \gamma v_{\pi_i,k}(s'))$$

3. Policy improvement: improve policy by

$$\pi_{i+1}(s) = \text{argmax}_a \sum_{s' \in S} P(s'|s, a)(R(s, \pi_i(s), s') + \gamma v_{\pi_i}(s')$$

AI in Logistics | Christina Imdahl

TU/e

# Dynamic Programming Differences and Similarities

- Key methods from dynamic programing can be used to solve MDPs
  - Value iteration and policy iteration

- They assume full knowledge of the underlying MDP (transition probabilities, rewards)

- The *agent* often does not have these (we as programmers may have it)

- The agent shall be able of entering new situations that we don't know

TU/e

# Next Lecture ….

# How do we learn when transition probabilities and values are unkown to the agent.

TU/e

# Main idea: Try & Fail & Try….



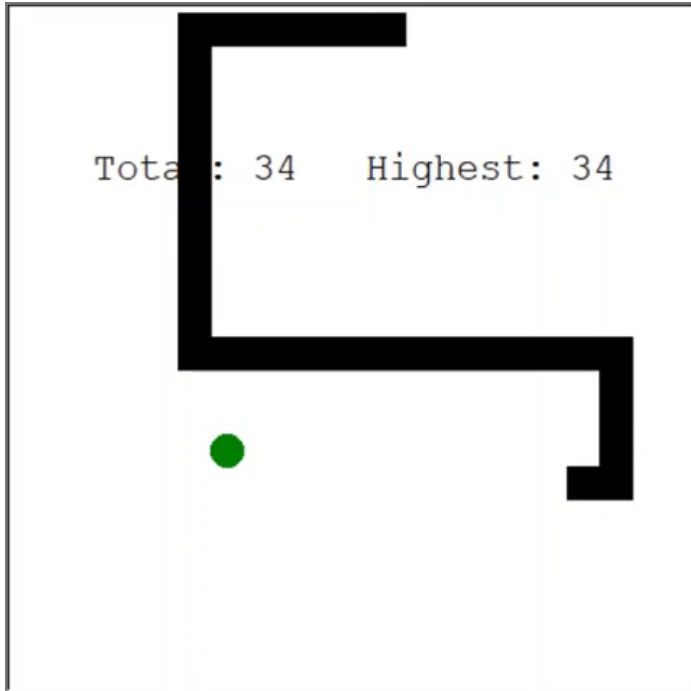Mom : If your friends jumped off a bridge, would you jump too?

Machine learning algorithm :

AI in Logistics | Christina Imdahl

TU/e

# Agenda

TU/e

# Snake



Tota : 34    Highest: 34

Discuss with your neighbour:

Actions

Rewards



State space
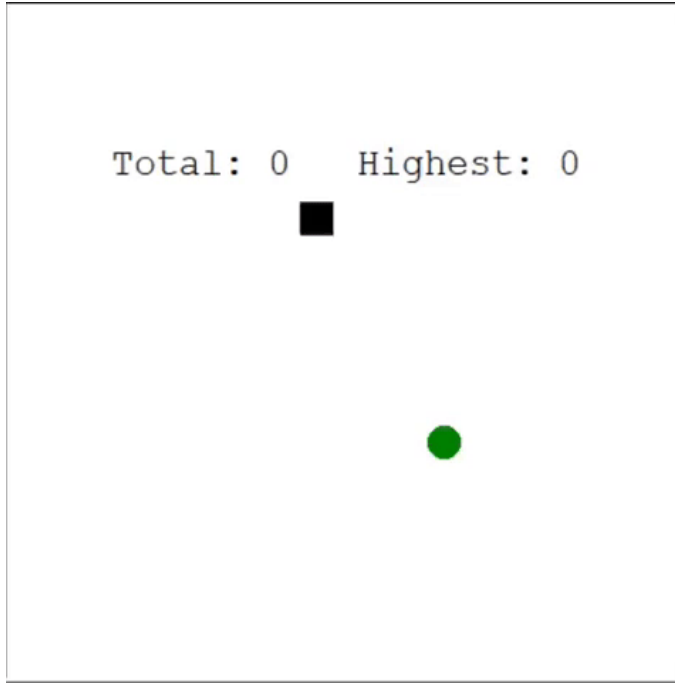
TU/e

# Snake

## Rewards



## Actions

Snake moves up
Snake moves right
Snake moves down
Snake moves left

## State

Apple is above the snake
Apple is on the right of the snake
Apple is below the snake
Apple is on the left of the snake
Obstacle directly above the snake
Obstacle directly on the right
Obstacle directly below the snake
Obstacle directly on the left
Snake direction == up
Snake direction == right
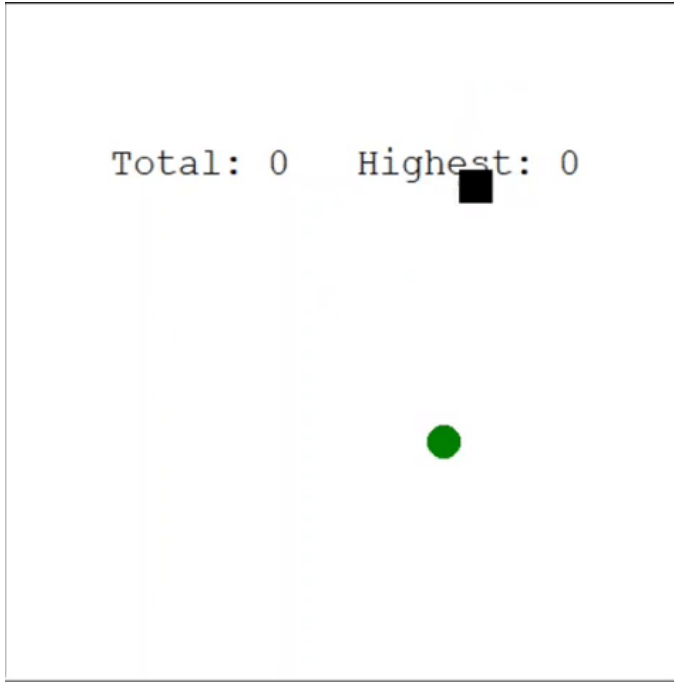Snake direction == down
Snake direction == left

AI in Logistics | Christina Imdahl

TU/e

# Snake - Learning

Game 1-4



Total: 0    Highest: 0

TU/e

# Snake - Learning

Game 7 – First Apple



AI in Logistics | Christina Imdahl

TU/e

# Snake - Learning

Game 13



Total: 22   Highest: 29

TU/e

# Snake - Learning

Game 30



Total: 5   Highest: 5

TU/e

# Agenda

TU/e

# Reward 1 – MOVE!

## Rewards



## Result



Total: 0     Highest: 0

# Reward 2 – RUN!

## Rewards

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 |
| -100 | | | | | | | | | | | -100 |
| -100 | | | 10 | | | | | | | | -100 |
| -100 | | | | | | -1 | | | | | -100 |
| -100 | | | | -1 | | -1 | | | | | -100 |
| -100 | | | | | | -1 | | | | | -100 |
| -100 | | | | | | | | | | | -100 |
| -100 | | | | | | | | | | | -100 |
| -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 | -100 |

## Result

Total: 0    Highest: 1

TU/e

# Agenda

TU/e

# Different State Space Definitions

**1** State (Obstactles = Body/Wall)

Apple is above the snake
Apple is on the right of the snake
Apple is below the snake
Apple is on the left of the snake
Obstacle directly above the snake
Obstacle directly on the right
Obstacle directly below the snake
Obstacle directly on the left
Snake direction == up
Snake direction == right
Snake direction == down
Snake direction == left

**2** State - only walls (Obstacles = Wall)

Same as **1** but body location not included

**3** State – coordinates

(x,y) Apple
(x,y) Snake

Obstacle directly above the snake
Obstacle directly on the right
Obstacle directly below the snake
Obstacle directly on the left
Snake direction == up
Snake direction == right
Snake direction == down
Snake direction == left

**4** State – no direction

Apple is above the snake
Apple is on the right of the snake
Apple is below the snake
Apple is on the left of the snake
Obstacle directly above the snake
Obstacle directly on the right
Obstacle directly below the snake
Obstacle directly on the left

AI in Logistics | Christina Imdahl

TU/e

# Learning Curve



Learning Curve

State : only walls
State : direction 0 or 1
State : coordinates
State : no direction

Sum of reward during episode

Episodes

AI in Logistics | Christina Imdahl

TU/e

# Agenda

TU/e

# Set-up

- Provider of bike spare parts (brakes, tires …)

- You start managing the inventory of one item based on part demands

- There is a penalty if you have to much inventory
- There is a penalty if you have to little inventory

- Your objective is to maximize your sales taking into account the costs.

AI in Logistics | Christina Imdahl

TU/e

# Sequence of events



**Period t**

**Place order to arrive in t+L**
*env.takeAction() updates pipeline inv. and action_log*

**Observe inventory & compute reward**
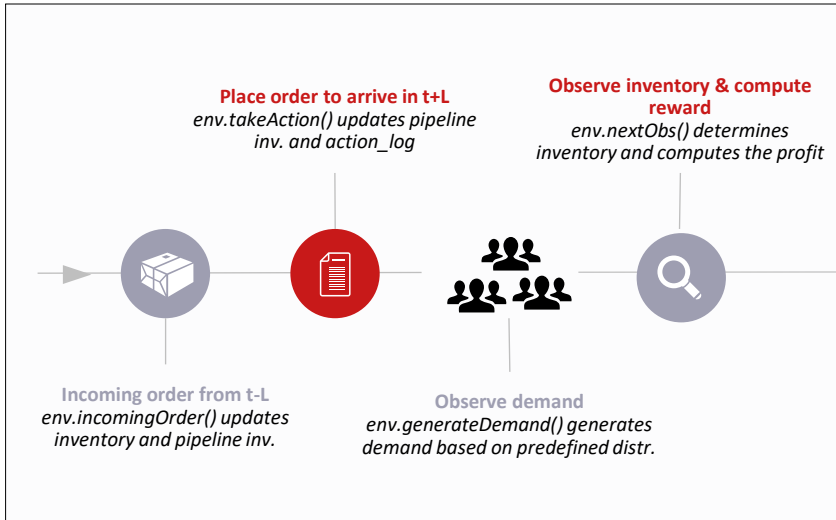*env.nextObs() determines inventory and computes the profit*

**Incoming order from t-L**
*env.incomingOrder() updates inventory and pipeline inv.*

**Observe demand**
*env.generateDemand() generates demand based on predefined distr.*

**Period t+1**

**Place order to arrive in t+L+1**
*env.takeAction() updates pipeline inv. and action_log*

**Incoming order from t-L+1**
*env.incomingOrder() updates inventory and pipeline inv.*

TU/e

# Two systems to study

## Backorder System

- Demand not served at the end of the period can be served the next period

- Optimal policy is a base-stock policy with $S = F^{-1}(\frac{h}{h+b})$

- Can you be equally good?
- How do different parameter/model choices affect your learning?

## Lost Sales System

- Demand not served at the end of the period is lost

- Optimal policy is unknown, approximative policies exist

- What is a suitable benchmark?
- Can RL outperform your benchmark?

AI in Logistics | Christina Imdahl

TU/e

# Acknowledgements

Snake: https://towardsdatascience.com/snake-played-by-a-deep-reinforcement-learning-agent-53f2c4331d36

Snake Code: https://github.com/henniedeharder/snake

RL course by David Silver:
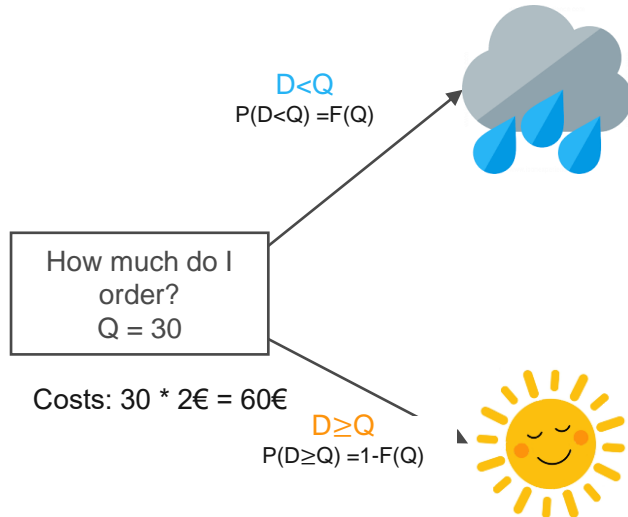https://www.youtube.com/watch?v=2pWv7GOvuf0&list=PLqYmG7hTraZDM-OYHWgPebj2MfCFzFObQ&index=1&t=74s

TU/e

# Back-Up Newsvendor

AI in Logistics | Christina Imdahl

TU/e

## EXAMPLE OF TODAY
## THE NEWSSTAND

- YOU are a running a small kiosk and need to decide on the optimal amount of papers at the **beginning of the day**

- You buy newspaper at a **costs of 2€**

- You sell newspaper at a **price of 10€**

- At the end of a day all newspapers are shredded

- You expect 30 customers (with a standard deviation of 10, normally distributed).

# NEWSVENDOR MISMATCH COSTS

How much do I order?
Q = 30

Costs: 30 * 2€ = 60€

D<Q
$P(D<Q) = F(Q)$

D≥Q
$P(D≥Q) = 1-F(Q)$

Demand = 10

Revenue: 10 * 10€ = 100€
Profit: 100€ - 60€ = 40€

Overstock quantity: 30-10 = 20
Costs of overstocking $c_o$: 20 * 2€ = 40€

Demand = 50

Revenue: 30 * 10€ = 300€
Profit: 300€ - 60€ = 240€

Understock quantity: 50-30 = 20
Costs of understocking $c_u$: 20 * 8€ = 160€

$$C(Q;D) = c_o \underbrace{\max(0, Q - D)}_{\text{Overstock quantity}} + c_u \underbrace{\max(0, D - Q)}_{\text{Understock quantity}}$$

# NEWSVENDOR MISMATCH COSTS

D<Q
P(D<Q) =F(Q)

How much do I order?

D≥Q
P(D≥Q) =1-F(Q)

## Overstocking

We need to pay *c* (*2€*) for every newspaper bought, but not sold.

$$c_o = c$$

Marginal expected costs of one additional unit:
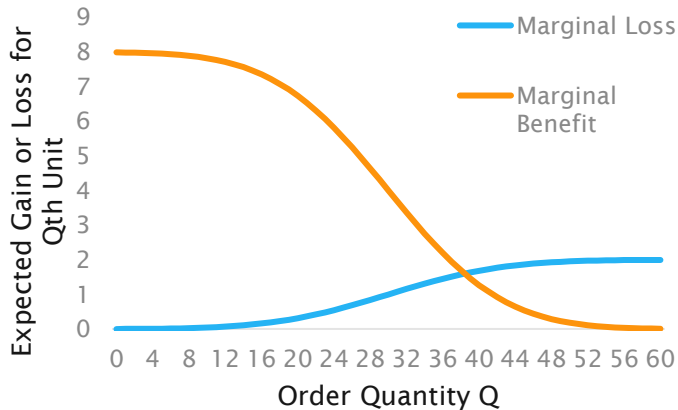
$$c_o \cdot F(Q)$$

## Understocking

We are losing the margin *p-c (10€-2€)* for every newspaper that we did not sell, because we had too little.

$$c_u = p - c$$

Marginal expected costs of one additional unit:

$$c_u \cdot (1 - F(Q))$$

# MARGINAL ANALYSIS: IS IT BENEFICIAL TO ORDER ONE MORE UNIT?



The optimal order quantity is the point where the expected marginal benefit through an order quantity equals the expected marginal loss through an additional order quantity.

$$c_o \cdot F(Q) = c_u \cdot (1 - F(Q))$$
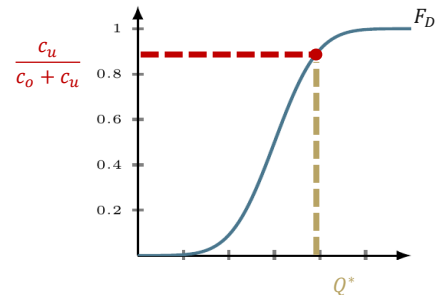
## NEWSVENDOR MISMATCH COSTS
## OPTIMAL ORDER QUANTITY

$$C(Q; D) = c_o \underbrace{\max(0, Q - D)}_{\text{Overstock quantity}} + c_u \underbrace{\max(0, D - Q)}_{\text{Understock quantity}}$$
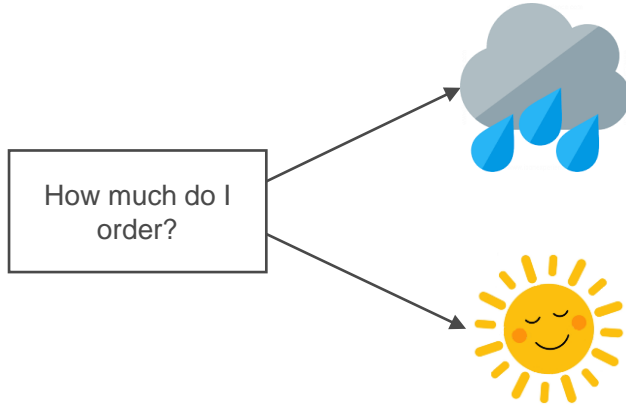
Minimizing the expected mismatch costs via optimization yields a simple formula for the optimal order quantity:

$$F(Q^*) = \frac{c_u}{c_o + c_u}$$

**Critical Ratio**
**= Service level**

# NEWSVENDOR MISMATCH COSTS



### Overstocking

We need to pay $c$ (*2€*) for every newspaper bought, but not sold.

$$c_o = c$$

### Understocking

We are losing the margin *p-c (10€-2€)* for every newspaper that we did not sell, because we had too little.

$$c_u = p - c$$

$$F(Q^*) = \frac{c_u}{c_u + c_o} = \frac{p-c}{p-c+c} = \frac{p-c}{p}$$

$$= \frac{10-2}{10} = 80\%$$