
Benchmarks for Online Stochastic Operations Research Problems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Benchmarks to compare reinforcement learning algorithms with state of the art for
2 solving online stochastic operations research problems

3 1 Introduction

4 Operations Research is an awesome field. Many online stochastic problems exist in practice but
5 solutions make simplistic assumption. Reinforcement learning is a good fit to these problems and
6 there have been some initial results. However, no standard set of problems or algorithms to compare
7 against exist. We prepare these benchmarks to compare different algorithms and push the community
8 towards developing better ones.

9 2 Related Work

10 3 Bin Packing

11 In the classic bin packing problem, we fit given items of varying sizes into as few fixed size bins as
12 possible. In the online stochastic version of this problem, each item arrives one at a time and its size
13 varies as per an unknown distribution. We formulate the problem as a Markov Decision Process and
14 compare reinforcement learning algorithms against a well known baseline called Sum of Squares that
15 asymptotically converges to a good solution regardless of the item size distribution.

16 3.1 Related Work

17 Online stochastic bin packing [?]

18 3.2 Problem Formulation

19 Items can be of different types $j \in \{1, \dots, J\}$. The size of type j is s_j and the probability that an item
20 is of type j is p_j . Without loss of generality, we assume item types are enumerated in the increasing
21 order of their size: $s_1 < s_2 < \dots < s_J$. Items arrive one at a time and are packed into bins of size
22 B , where $s_J < B < \infty$. We assume the item sizes s_j and bin size B are integers. We assume the
23 number of bins are unlimited and denote the sum of item sizes in a bin as *level* h . After n items have
24 been packed, we denote the number of bins at level h as $N_h(n)$, where $h = 1, \dots, B$.

Our objective is to reduce the number of non-empty bins. We can achieve this by minimizing the total waste (i.e. empty space) in the partially filled bins. Hence our objective is to minimize waste at any point in time:

$$W(n) \triangleq \sum_{h=1}^{B-1} N_h(n)(B-h) \quad (1)$$

3.3 Baseline Algorithm

3.4 Reinforcement Learning Algorithm

3.5 Results

4 Newsvendor

5 Vehicle Routing

5.1 Literature Review

5.2 Problem Formulation

We consider a version of VRP that is of a food delivery driver (e.g. using Amazon Restaurants). Orders arrive at the driver's phone app over time in a dynamic manner over time. Each order has a reward (e.g. delivery fee and tip) associated with it and it is assigned to a specific restaurant in the city. "City" here means the whole Euclidean space in which the VRP problem lives. Orders arrive according to a Poisson process and the rate depends on the region of the city that the order is created in. Also, the reward of an order comes from a Gamma distribution, again, with parameters specific to the region. The driver has to accept an order and pick the items up from the restaurant the order is mapped to. The order needs to be delivered within a time window since its creation, which is imposed as a hard constraint. If the driver does not accept a specific order, it remains open for a while and then disappears according to a probability distribution, meaning that either it has expired or accepted by some other driver. There is a capacity on the number of orders a driver can carry in the vehicle, however there is no limit on the number of orders that are accepted (but not delivered) by the driver at the same time. Finally, there is a cost associated with travel. The driver's goal is to maximize the average reward over an infinite horizon.

This problem is known as stochastic and dynamic capacitated vehicle routing problem with pick up and delivery, time windows and service guarantee. (SDPDPTW with service guarantee).

5.3 Baseline Solution

The problem for a given set of orders can be expressed and solved using the following Mixed Integer Programming formulation.

Sets

V : Current vehicle location, $V = \{0\}$

P : Pick up nodes (copies of the restaurant nodes, associated with the orders that are not in transit)

D : Delivery nodes representing the orders that are not in transit

A : Nodes representing the orders that are accepted by the driver; $A \subset D$

T : Delivery nodes representing the orders that are in transit

R : Nodes representing the restaurants, used for final return)

N : Set of all nodes in the graph, $N = V \cup P \cup D \cup T \cup R$

E : Set of all edges, $E = \{(i, j), \forall i, j \in N\}$

55 Decision variables

- x_{ij} : Binary variable, 1 if the vehicle uses the arc from node i to j , 0 otherwise; $i, j \in N$
 y_i : Binary variable, 1 if the order i is accepted, 0 otherwise; $i \in P$
 Q_i : Auxiliary variable to track the capacity usage as of node i ; $i \in N$
 B_i : Auxiliary variable to track the time as of node i ; $i \in N$

56 Parameters

- n : Number of orders available to pick up, $n = |P|$
 C_{ij} : Symmetric Manhattan distance (in miles) matrix between node i and j , $(i, j) \in E$
 q_i : Supply (demand) at node i , $q_0 = |T|$; $q_i = 1, \forall i \in P$; $q_i = -1, \forall i \in D \cup T$; $q_i = 0 \in R$
 m : Travel cost per mile
 r_i : Revenue for order associated with pick up node i , $i \in P$
 U : Vehicle capacity
 M : A very big number
 t : Time to travel one mile
 d : A constant positive service time per stop

57 Model

$$\begin{aligned} & \text{maximize} && \sum_i r_i y_i - m \sum_{(i,j) \in E} C_{ij} x_{ij} && (2a) \\ & x, y, Q, B \end{aligned}$$

$$\text{subject to} \quad \sum_{j \in N} x_{ij} = y_i \quad \forall i \in P, \quad (2b)$$

$$\sum_{j \in N} x_{ij} - \sum_{j \in N} x_{i+n,j} = 0 \quad \forall i \in P, \quad (2c)$$

$$y_i = 1 \quad \forall i \in A, \quad (2d)$$

$$\sum_{j \in N} x_{ij} = 1 \quad \forall i \in T, \quad (2e)$$

$$\sum_{j \in N} x_{0j} = 1, \quad (2f)$$

$$\sum_{j \in N \setminus R} x_{ji} = 1 \quad \forall i \in R, \quad (2g)$$

$$\sum_{j \in N \setminus R} x_{ji} - \sum_{j \in N} x_{ij} = 0 \quad \forall i \in P \cup D \cup T, \quad (2h)$$

$$Q_i + q_j - M(1 - x_{ij}) \leq Q_j \quad \forall i, j \in N, \quad (2i)$$

$$\max(0, q_i) \leq Q_i \quad \forall i \in N, \quad (2j)$$

$$\min(U, U + q_i) \geq Q_i \quad \forall i \in N, \quad (2k)$$

$$B_i + d + C_{ij}t - M(1 - x_{ij}) \leq B_j \quad \forall i, j \in N, \quad (2l)$$

$$B_i + C_{i,j+n}t \leq B_{i+n} \quad \forall i \in P, \quad (2m)$$

$$x_{ij} \in \{0, 1\} \quad \forall i, j \in N, \quad (2n)$$

$$y_i \in \{0, 1\} \quad \forall i \in P \quad (2o)$$

58 6 Conclusion

59 Amazing! Awesome results! Accept!

60 **Acknowledgments**

61 Do not include acknowledgments in the anonymized submission, only in the final paper.