

The Contextual Appointment Scheduling Problem

Nima Salehi Sadghiani

Department of Industrial and Operations Engineering
University of Michigan
Ann Arbor, MI 48105
nsalehi@umich.edu

Saeid Motiian

Adobe Inc.
San Jose, California, USA
motiian@adobe.com

Abstract

This study is concerned with the determination of optimal appointment times for a sequence of jobs with uncertain duration. We investigate the data-driven Appointment Scheduling Problem (ASP) when one has n observations of p features (covariates) related to the jobs as well as historical data. We formulate ASP as an Integrated Estimation and Optimization problem using a task-based loss function. We justify the use of contexts by showing that not including them yields to inconsistent decisions, which translates to sub-optimal appointments. We validate our approach through two numerical experiments.

Introduction

Appointment Scheduling Problem (ASP) arises in many service industries to improve the service quality, operational efficiency, utilization of resources, match workload to available capacity, and smooth the flow of customers (e.g. health care, legal services, accounting services, and loading/unloading ships/planes). In health care industry, ASP emerges in numerous settings, such as scheduling outpatient appointments in primary care (Cayirli, Veral, and Rosen 2006), surgery scheduling (Denton and Gupta 2003), and call-center/nurse staffing (Gurvich, Luedtke, and Tezcan 2010; Koçağa, Armony, and Ward 2015).

Assuming punctuality, a common problem faced by decision makers is how to determine the starting time of job when its duration is uncertain. The decision makers' goal is to minimize the expected costs associated with jobs' waiting time, server idle time, and overtime. For example, when appointments are set, jobs are not available prior to the scheduled start time, even if the server becomes idle. Thus, at the cost of additional waiting (for jobs), choosing an early start time will lead to better server utilization, whereas a late start time will reduce the waiting time at the cost of additional server idling.

A common assumption in the literature of the ASP is that we can accurately estimate the probability distribution of the job duration. However, in reality, this may not be the case. The scarcity of the historical data, diversity in types of jobs can easily violate the underlying assumptions behind distribution fitting (like independent and identically distributed

(i.i.d.) assumption, the choice of distribution, or the correlation among the job duration distributions). Although the ASP is a well-explored problem, most studies in the literature made these assumptions without considering the consequences.

In this paper, we investigate the ASP under two settings: classic models in which we only have access the previous job duration, and the feature-based models in which we have access to some exogenous information about the job duration as well as the actual value. Rather than a two-step process of first estimating the duration distribution then optimizing for the optimal scheduling, we propose solving the ASP via a single step by training a Deep Neural Network (DNN) in an end-to-end task-based framework. The difficulty of the integrated learning and optimizing task in ASP arises from the fact that the decision about each job duration depends on the actual realization of the previous jobs.

Related Work

ASP literature

ASP studies with uncertain service times often assume full distributional information about the job durations (i.e. full joint distribution of the jobs is known, or jobs are independent and marginal distributions are known). Excellent surveys of ASP models can be found in (Gupta and Denton 2008; Cardoen, Demeulemeester, and Beliën 2010).

Sample Average Approximation (SAA) is a common method used to solve the stochastic ASP. SAA approach can be computationally intensive, and require knowledge of the joint distribution of the job service durations (Denton and Gupta 2003).

Further, when the joint distribution of the uncertain parameters is not known to us, other types of modeling is proposed in the literature. Mittal, Schulz, and Stiller 2014 proposed Robust Optimization (RO) for the ASP where job durations lie within interval uncertainty sets. Mak, Rong, and Zhang 2014 proposed Distributionally Robust Optimization (DRO) for ASP assuming partial distributional information of the uncertain parameters such as support and moments are known. The appointments are made to minimize the worst-case expected objective value among all possible distributions with the specified support and moments information. Finally, Zhang, Shen, and Erdogan 2017 proposed

using DRO with moment uncertainty set and derived an approximate semidefinite programming model.

Prescriptive Analytics

Solving many real world problems requires two main steps of estimation (prediction) and optimization (prescription). Due to the complex nature of each step, researchers typically solve these two steps separately. The major issue with this approach is that the optimization problem is blind to the prediction process and its sensitivity to over/under estimation of the parameter. In contrast to the Separated Estimation and Optimization (SEO) approach, another framework is proposed by Bertsimas and Kallus 2014. They referred to this method as "prescriptive analytics". They combined the ideas from Machine Learning (ML) and Operations Research and Management Science (OR/MS) in developing a framework for using auxiliary data to prescribe optimal decisions.

A general framework for the Integrated Estimation and Optimization (IEO) approach was proposed later by Elmachtoub and Grigas 2017. They called their framework "Smart Predict, then Optimize" (SPO). Training a model with respect to the SPO loss is computationally challenging, and therefore they also proposed a surrogate loss function, called the SPO+ loss, which upper bounds the SPO loss, has desirable convexity properties, and is statistically consistent under mild conditions. Minimizing an application-specific loss function is referred to as "training with direct loss" (Song et al. 2016) or task-based model learning (Donti, Amos, and Kolter 2017) in the ML literature. In this paper, we proposed using DNN to solve the feature-based ASP.

ASP Modeling

We consider an ASP with a single server at which jobs arrive punctually at scheduled appointment times. The set of jobs follows a fixed order of arrivals. We formulate the stochastic ASP with random service duration. Throughout this paper, we use boldface notation to denote vectors.

ASP Classic Models

We consider n jobs arriving at a single server given a predetermined sequence $1, 2, 3, \dots, n$. Job i has a random service duration p_i . The planner determines time allowances $S_i \geq 0$ for each job i . If job i cannot be started at its planned start time due to a delay of completion of the previous job, a waiting time cost will be incurred per unit time of delay. Let W_i denote the waiting time of the job i . On the other hand, if job $i - 1$ completed sooner than its planned start time, an idling time cost will be incurred per time. Let I_i denote the idling time of the server between job $i - 1$ and job i . The goal is to optimize the decisions of scheduling for arrival times of each job, or equivalently, assigning time intervals between each two jobs.

The random job durations \mathbf{p} follow a joint probability distribution $F(\mathbf{p})$. For any job i , the support of p_i is independent from other jobs and is denoted by \mathbb{D}_i ; therefore, $\mathbb{D} = \mathbb{D}_1 \times \dots \times \mathbb{D}_n$. Given time allowances $\mathbf{S} \in \mathbb{R}^+$ and a realization of the random service times, denoted by \mathbf{p} , the

waiting and idling times should satisfy the following constraints:

$$W_i - I_i = W_{i-1} + p_{i-1} - S_{i-1} \quad \forall 2 \leq i \leq n \quad (1)$$

Moreover, $W_i, I_i \in \mathbb{R}^+$. It is also worth noting that because of the punctually arrival assumption $W_1 = I_1 = 0$. Assuming linear costs for waiting and idling, the total waiting and idling cost is denoted by

$$f(\mathbf{S}, \mathbf{p}) = c^W \sum_{i=1}^n W_i + c^I \sum_{i=1}^n I_i \quad (2)$$

where c^W and c^I (non-negative) represent the costs of waiting and idling, respectively.

ASP Model with Distributional Information. Assuming the joint probability distribution of job durations $F(\mathbf{p})$ is known, we can formulate the stochastic ASP model as follow

$$\min_{\mathbf{S} \in \mathbb{R}^+} \mathbb{E}_{\mathbf{p} \sim F(\cdot)} \mathcal{F}(\mathbf{S}, \mathbf{p}) \quad (3)$$

where \mathbb{E} denotes the expected value and $\mathcal{F}(\cdot) := \min_{\mathbf{W}, \mathbf{I} \in \mathbb{R}^+ \cup \{eq.(1)\}} f(\cdot)$. The objective function of the stochastic ASP is to minimize the expected total waiting and idling costs under the assumption that random job durations come from the known joint probability distribution $F(\cdot)$ with support \mathbb{D} .

ASP Model with Limited Distributional Information.

A key assumption of the previous model is that the exact joint distribution $F(\cdot)$ is known to the planner at the beginning. However, this assumption may not be held under different practical settings. To address the distributional ambiguity issue, in this section, we assume only some limited distributional information about $F(\cdot)$ is available to the planner as well as the support \mathbb{D} . We consider a moment-based ambiguity set for the distributions of the uncertain parameters.

For moment-based ambiguity set, the q th moment of p_i is denoted by M_{iq} . We also use μ_i and σ_i to denote the mean and standard deviation of p_i , respectively; that is, $M_{i1} = \mu_i$ and $M_{i2} = \mu_i^2 + \sigma_i^2$. For any $q \in Q$ where Q is a finite set of positive integers, the distribution $F(\cdot)$ belongs to the ambiguity set $\mathcal{M}(\mathbb{D}, Q)$

$$:= \left\{ \begin{array}{l} \int_{\mathbb{D}} dF(\mathbf{p}) = 1 \\ \int_{\mathbb{D}} p_i^q dF(\mathbf{p}) = M_{iq} \quad \forall q \in Q, \quad 1 \leq i \leq n \\ dF(\mathbf{p}) \geq 0 \quad \forall \mathbf{p} \in \mathbb{D} \end{array} \right\} \quad (4)$$

Any set of \mathbf{p} belong to $\mathcal{M}(\cdot)$ is a distribution, match the moments, and it is non-negative. Under this assumption, we can formulate the stochastic ASP as a min-max problem robust to the distribution ambiguity

$$\min_{\mathbf{S} \in \mathbb{R}^+} \max_{F \in \mathcal{M}} \mathbb{E}_{\mathbf{p} \sim F(\cdot)} \mathcal{F}(\mathbf{S}, \mathbf{p}) \quad (5)$$

We shall refer to this framework as Distributionally Robust Optimization (DRO).

Data-driven ASP without Distributional Assumptions.

In practice, the planner does not know the true joint distribution of the job durations. If one has access to historical job durations for T periods: $(\mathbf{p}^t) = (p_1^t, \dots, p_n^t)$, $1 \leq t \leq T$, but no other information (no context), then the sensible approach is to substitute the true expectation with a sample average expectation and solve the resulting problem.

$$\min_{\mathbf{S} \in \mathbb{R}^+} \frac{1}{T} \sum_{t \in T} \mathcal{F}(\mathbf{S}, \mathbf{p}^t) \quad (6)$$

This approach is called the Sample Average Approximation (SAA) approach in stochastic optimization (Kleywegt, Shapiro, and Homem-de Mello 2002). The consistency of the SAA estimator for the expectation of the objective function of the stochastic models has been discussed thoroughly in (Shapiro, Dentcheva, and Ruszczyński 2009).

Drawbacks of the Classic Models

In the previous section, we reviewed classic approaches to model ASP. In a realistic situation, the classic models are too simplistic to represent many real scenarios because one can collect data on exogenous information about the jobs' characteristics as well as the job durations. In other words, the planner has access to a richer data set for decision making. Further, most of the models fail to consider context in estimating the job durations. The classic models obtain an estimate of the probability distribution and solve the optimization problem separately. The approximation process of the probability distributions involves errors, especially when historical data is scarce or when durations are not i.i.d.

The Feature-based ASP Models

In practice, the job durations depend on many observable features (equivalently, covariates or context). For example, in patient appointment scheduling, the characteristics of the patients can be used as features. If the jobs are flight arrivals, then the attributes of the flights can be used (i.e. day, month, season, weather, origin and destination airports).

Assuming the planner schedules n jobs with m features at each time period and s/he has access to historical data for T periods, we have $\{(\mathbf{x}^t, \mathbf{p}^t)\}_{t=1}^T$ where $\mathbf{x}^t = (x_1^t, \dots, x_n^t)$ and $\mathbf{p}^t = (p_1^t, \dots, p_n^t)$. For notation brevity, we denote all the m features for the job i at time t with x_i^t , $1 \leq i \leq n$, i.e. $x_i^t = (x_{i1}^t, \dots, x_{im}^t)$.

In the feature-based ASP model, the features are available to the planner prior to the decision making process. In other words, the decision maker optimizes the conditional expected cost function:

$$\min_{\mathbf{S}(\cdot) \in \mathcal{S}, \{\mathbf{S}: \mathcal{X} \rightarrow \mathbb{R}^+\}} \mathbb{E}_t f(\mathbf{S}(\mathbf{x}^t), \mathbf{p}^t | \mathbf{x}^t) \quad (7)$$

where the scheduling decisions $\mathbf{S}(\cdot)$ are now functions mapping the feature space $\mathcal{X} \subset \mathbb{R}^m$ to the positive real numbers (\mathbb{R}^+) and the expected waiting and idling costs that we minimize is now conditional on the feature vector $\mathbf{x}^t \in \mathcal{X}$.

The planner needs to be able to solve the described problem in equation (7) efficiently given the observation \mathbf{x}^t . The main issue in this model is that the job durations might be

correlated and allowance decisions have to be made sequentially. For example, when the first job is being processed on the server, we do not know whether the server will be idle until the next job arrives or the next job will be delayed. The uncertainty about the starting time of each job propagates through the sample of n jobs at each time unit t . The planner needs to know the actual realization of the previous job and its starting time to optimally plan for the next one.

2-Job ASP Model. The special case of this problem is when we have $n = 2$. Because of punctual arrivals the first job arrives at time zero. For this case, it is only necessary to predict the duration of the first job, and the second job starts its process right when the first job is done. Weiss (1990) showed this special case corresponds to the *News vendor* problem. Therefore, under the assumption that no data features is available but the underlying distribution of the job durations is known, a closed form expression for the optimal allowance of the first job can be obtained. The closed form solution for the classic *News vendor* problem described in Gallego and Moon (1993) is as follow

$$\hat{p}_1^* = S_1^* = F^{-1} \left(\frac{c^{\mathcal{I}}}{c^{\mathcal{I}} + c^{\mathcal{W}}} \right) \quad (8)$$

where \hat{p}_1^* is the optimal estimation for the duration of the first job and also its time allowance. The $F(\cdot)^{-1}$ refers to the inverse distribution of $F(\cdot)$. Recently, some scholars used ML techniques to solve the data-driven *News vendor* problem. Here, we discuss the ones that are closely related to our model.

Ban and Rudin (2018) has shown analytically the superiority of using Linear Model (LM) to describe the optimal solution than other data-driven approaches such as SAA. Basically, in their models they defined the mapping function $\mathbf{S}(\mathbf{x}^t)$ as:

$$\mathcal{S}_{LM} = \left\{ \mathbf{S} : \mathcal{X} \rightarrow \mathbb{R}, \mathbf{S}(\mathbf{x}^t) = \sum_{k=0}^m \beta_k x_k^t \right\} \quad (9)$$

where β_k s are the coefficients of the linear model and the $x_0^t = 1$ accommodates for a feature-independent term. Further, they have shown that a complex-valued $\mathbf{S}(\mathbf{x}^t)$ mapping functions (but infinitely differentiable at 0) can be approximated by polynomial terms by Taylor series. In their setting, the objective function becomes:

$$\min_{\mathbf{S}(\cdot) \in \mathcal{S}_{LM}, \{\mathbf{S}: \mathcal{X} \rightarrow \mathbb{R}\}} \frac{1}{T} \sum_{t=1}^T f(S_1(\mathbf{x}^t), p_1^t) = \frac{c^{\mathcal{W}}}{T} \|\mathbf{p}_1 - S_1(\mathbf{x})\|_+ + \frac{c^{\mathcal{I}}}{T} \|(S_1(\mathbf{x}) - \mathbf{p}_1)_+\|_1 \quad (10)$$

In equation (10), the objective function has the form of \mathcal{L}_1 asymmetric loss function where $c^{\mathcal{W}}$ and $c^{\mathcal{I}}$ represent the unit cost of overestimation and underestimation, respectively. This special structure is due to that fact that $W_1^t \cdot I_1^t = 0$, $1 \leq t \leq T$. Further, they discuss the conditions for the optimality of the decisions from the feature-based model and the asymptotics of the optimal solutions. Finally, they have

shown the decisions from SAA (with finite sample sizes) has bias of $\mathcal{O}(1)$; hence the SAA decisions may not be asymptotically optimal.

Oroojlooyjadid, Snyder, and Takac (2016) have shown the benefits of using non-linear models to estimate the mapping function $\mathcal{S}(\mathbf{x}^t)$ in comparison with linear models. In their model $\mathcal{S}(\mathbf{x}^t)$ is a complex and non-linear function (dense layer architecture).

$$\mathcal{S}_{DNN} = \{\mathcal{S} : \mathcal{X} \rightarrow \mathfrak{R}, \mathcal{S}(\mathbf{x}^t, \boldsymbol{\omega}, \mathbf{b})\} \quad (11)$$

where $\boldsymbol{\omega}$, \mathbf{b} are the weights and biases for a fully-connected network architecture. They used sigmoid functions to create non-linearity in the model. They also suggested using a surrogate \mathcal{L}_2 loss function for faster convergence.

$$\begin{aligned} \min_{\mathcal{S}(\cdot) \in \mathcal{S}_{DNN}, \{\mathcal{S} : \mathcal{X} \rightarrow \mathfrak{R}\}} \frac{1}{2T} \sum_{t=1}^T f^2(\mathcal{S}_1(\mathbf{x}^t, \boldsymbol{\omega}, \mathbf{b}), p_1^t) = \\ \frac{c^{\mathcal{W}}}{2T} \left\| (\mathbf{p}_1 - \mathcal{S}_1(\mathbf{x}, \boldsymbol{\omega}, \mathbf{b}))_+ \right\|_2^2 + \frac{c^{\mathcal{I}}}{2T} \left\| (\mathcal{S}_1(\mathbf{x}, \boldsymbol{\omega}, \mathbf{b}) - \mathbf{p}_1)_+ \right\|_2^2 \end{aligned} \quad (12)$$

Zhang and Gao (2017) discuss the benefits of using an extra layer of Rectified Linear Units (ReLUs) to decompose the original \mathcal{L}_1 loss function into a constant (p_1^t) and a differentiable function ($\mathcal{S}_1(\mathbf{x}^t, \boldsymbol{\omega}, \mathbf{b})$). They used $c^{\mathcal{W}}$ and $c^{\mathcal{I}}$ as the weights of the new layer to show the original \mathcal{L}_1 loss function can be reconstructed from these two parts.

n -Job ASP Model. Assuming n jobs are needed to be scheduled, for scheduling each job, the planner needs to take into account the actual realizations of the previous jobs. One can clearly see how the previous methods in the literature fail to consider this recursive relationship among the predictors (decision variables of the optimization model). To address this shortcoming, we propose a different loss function to learn the complex and recursive behavior of the job durations from the context while it accounts for the sequential relationships among jobs that are required to be scheduled in each time unit.

To estimate the mapping function $\mathcal{S}(\mathbf{x}^t)$, we propose adding a sigmoid layer on the outputs of the dense layer to satisfy the non-negativity of the decision variables. We will refer to this mapping as DNN^+ ,

$$\mathcal{S}_{DNN^+} = \{\mathcal{S} : \mathcal{X} \rightarrow \mathfrak{R}^+, \mathcal{S}(\mathbf{x}^t, \boldsymbol{\omega}, \mathbf{b})\} \quad (13)$$

Further, we suggest using Quantile Regression (QR) Loss Function for the neural network. White (1992) presents theoretical support for the use of QR for an ANN to estimate the potentially non-linear quantile models. Taylor (2000) has shown the application of QR-NN to estimate the conditional density of the uncertain parameters in a multi-period forecasting framework. Minimizing the Mean Absolute Deviation (MAD) leads to an estimate of the conditional median of the predict and data. By applying asymmetric weights to positive/negative errors by using a tilted form of the absolute value function, one can instead compute conditional quantiles of the predictive distribution (Koenker and Bassett Jr 1978). The tilted absolute value function (also known as the

pinball loss function) for each individual prediction is given by

$$\mathcal{L}(p_i - \hat{p}_i | q) = \begin{cases} q(p_i - \hat{p}_i), & \text{if } p_i - \hat{p}_i \geq 0, \\ (1 - q)(\hat{p}_i - p_i), & \text{otherwise.} \end{cases} \quad (14)$$

where q is the quantile of interest. We can show the loss function for the ASP can be represented as follows:

Proposition 1 *The proposed surrogate loss function for the ASP can be represented by*

$$\begin{aligned} \min_{\mathcal{S}(\cdot) \in \mathcal{S}_{DNN^+}, \{\mathcal{S} : \mathcal{X} \rightarrow \mathfrak{R}^+\}} \frac{1}{2T(c^{\mathcal{W}} + c^{\mathcal{I}})} \sum_{t=1}^T f(\hat{p}^t, p^t) = \\ \frac{1}{2T(c^{\mathcal{W}} + c^{\mathcal{I}})} \left(c^{\mathcal{W}} \left\| \sum_{i=2}^n (\mathbf{W}_{i-1} + \mathbf{p}_{i-1} - \hat{\mathbf{p}}_{i-1})_+ \right\|_1 \right) + \\ \left(c^{\mathcal{I}} \left\| \sum_{i=2}^n (-\mathbf{W}_{i-1} + \hat{\mathbf{p}}_i - \mathbf{p}_i)_+ \right\|_1 \right) \end{aligned} \quad (15)$$

where $q = \frac{c^{\mathcal{W}}}{c^{\mathcal{W}} + c^{\mathcal{I}}}$, $\hat{p}_i = \mathcal{S}(\mathbf{x}_i, \boldsymbol{\omega}, \mathbf{b})$, and W_i is calculated from a modified version of Lindley's recursion (1952). Note that we use $(\cdot)_+$ indicating $\max(\cdot, 0)$.

Optimization Methods for ASP Models

Approaches for Classic ASP

For the classic ASP models, we propose using SAA for the stochastic optimization models. This approach is applicable when the joint distribution of the job durations is known or when we have enough historical data. In the former case, we sample from the distribution and generate scenarios. In the latter, we treat each historical realization as a scenario. If limited distributional information is available, we propose using a cutting-plane approach to solve the DRO models.

SAA for Two-stage Stochastic Optimization. Having known the distribution of the job durations, we can sample and use the SAA scheme to solve the two-stage stochastic optimization model. Similar to Denton and Gupta (2003), we can rewrite the ASP problem as follow into a two-stage optimization model:

$$\min_{\mathcal{S}, \mathbf{W}, \mathcal{I} \in \mathfrak{R}^+} \frac{1}{N} \sum_{j=1}^N \left(c^{\mathcal{W}} \sum_{i=1}^n W_i(\xi_j) + c^{\mathcal{I}} \sum_{i=1}^N I_i(\xi_j) \right) \quad (16a)$$

$$W_i(\xi_j) \geq W_{i-1}(\xi_j) + p_{i-1}(\xi_j) - S_{i-1}, \forall \xi_j \in \Omega, 2 \leq i \leq n \quad (16b)$$

$$I_i(\xi_j) \geq -W_{i-1}(\xi_j) - p_{i-1}(\xi_j) + S_{i-1}, \forall \xi_j \in \Omega, 2 \leq i \leq n \quad (16c)$$

where ξ_j is a scenario in scenario set Ω . Each scenario is a set of realizations for job durations $1 \leq i \leq n$.

A Cutting-plane algorithm for DRO. When we have limited distributional information, we propose modeling ASP as a DRO model. Mak, Rong, and Zhang (2014) formulated the equivalent models for the min-max problem with moment-based ambiguity set under some specific assumptions. In here, we develop a generic cutting-plane algorithm

to solve the min-max model for the moment-based ambiguity set. The generic form of the ASP DRO model is as follow

$$\begin{aligned} & \min_{\mathbf{S} \in \mathbb{R}^+} \max_{F \in \mathcal{M}} \mathbb{E}_{\mathbf{p} \sim F(\cdot)} \mathcal{F}(\mathbf{S}, \mathbf{p}) \\ & \equiv \min_{\mathbf{S} \in \mathbb{R}^+} \max_{F \in \mathcal{M}} \int_{\mathcal{D}} \mathcal{F}(\mathbf{S}, \mathbf{p}) dF(\mathbf{p}) \end{aligned} \quad (17)$$

Although this model is computationally intractable, an approximation of this problem with discrete job durations can be solved with the following cutting-plane algorithm

Algorithm 1 A cutting-plane algorithm for DRO ASP

- 1: Set tolerance level $\epsilon \leftarrow (0, 1)$
- 2: Initialize iteration counter $k \leftarrow 0$
- 3: Initialize cuts list $cuts \leftarrow \emptyset$
- 4: Set $g(\mathbf{S}^0) \leftarrow \infty$
- 5: Set stop condition $newCuts \leftarrow \text{True}$
- 6: **while** $newCuts$ **do**
- 7: $k \leftarrow k + 1$
- 8: Solve (*Master*):

$$\min_{\mathbf{S}, \alpha, \delta} \sum_{i=1}^n \sum_{q \in Q} \alpha_{iq} M_{iq} + \delta$$

subject to:

$$\delta \geq g(\mathbf{S}^k), \quad k = 0, 1, \dots, K$$

- 9: Let $\bar{\mathbf{S}}^k$, $\bar{\alpha}^k$, and $\bar{\delta}^k$ be the optimal solutions for the (*Master*), then solve (*Subproblem*):

$$\begin{aligned} g(\bar{\mathbf{S}}^k) & \equiv \max_{\hat{\mathbf{p}} \in \{0,1\}^{n \times L}, \lambda \in \mathbb{R}^+, \mathbf{Y} \in F^D(\mathbf{Y})} \\ & \sum_{i=1}^n \sum_{l \in L_i} l \lambda_{il} - \sum_{i=1}^n \bar{\mathbf{S}}_{i-1}^k Y_i - \sum_{i=1}^n \sum_{q \in Q} \sum_{l \in L_i} \bar{\alpha}_{iq}^k l^q \hat{p}_{il} \end{aligned}$$

subject to:

$$\sum_{l \in L_i} \hat{p}_{il} = 1, \forall i$$

$$\lambda_{il} \geq -c^{\mathcal{I}} \hat{p}_{(i-1)l}$$

$$\lambda_{il} \geq Y_i - M(1 - \hat{p}_{(i-1)l})$$

$$\lambda_{il} \leq M \hat{p}_{(i-1)l}$$

$$\lambda_{il} \leq Y_i + c^{\mathcal{I}}(1 - \hat{p}_{(i-1)l})$$

- 10: **if** $\bar{\delta}^k < (1 - \epsilon)g(\bar{\mathbf{S}}^k)$ **then**
 - 11: $cuts \leftarrow cuts \cup \{\bar{\delta}^k \geq g(\bar{\mathbf{S}}^k)\}$
 - 12: **else**
 - 13: $newCuts \leftarrow \text{False}$
-

Approaches for Feature-based ASP

In the feature-based ASP models, we incorporate the features to predict the job durations. This can be done in two settings: (1) SEO: Separate Estimation and Optimization, (2) IEO: Integrate Estimation and Optimization.

Separate Estimation and Optimization (SEO). In this approach, we incorporate the data feature information in the

ASP by first mapping the job durations on the feature space assuming a normally distributed error term; then, we treat the estimates as the deterministic values in the ASP model. Ban and Rudin (2018) refer to this approach as SEO. In this paper, we use neural nets to estimate the job durations from the feature space. Then, we solve the optimization problem given the predicts from the network.

Integrate Estimation and Optimization (IEO) with QR-NN. In this approach, we try to find a parameterization of the features that optimizes the objective function of the ASP. In other words, we try to find a mapping function from the feature space to the optimal appointment decisions. To do so, we minimize the surrogate loss function in equation (15) (direct loss for the neural network). The estimates of this loss function are the optimal decisions of the stochastic ASP. In general, computing the gradients of an objective that depends upon the argmin/argmax operation is challenging. Recently, several authors developed different techniques for argmin/argmax differentiating (Gould et al. 2016; Donti, Amos, and Kolter 2017). For the ASP, due to the special structure of the optimization problem, we do not need to worry about this issue. However, still the gradients of the loss (15) are convolutional and hard to derive. Specifically, for the loss function $\mathcal{L}(\boldsymbol{\theta})$:

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{i=2}^n \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \hat{p}_i} \frac{\partial \hat{p}_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \quad (18)$$

The first term in equation (18) is the partial derivative of the loss function to the optimal job duration i and the second term is the partial derivative of each estimate with respect to the network parameters $\boldsymbol{\theta}$ (weights and biases). Calculating both terms in equation (18) is a complicated task since it requires the estimate and actual job durations of the previous jobs. Therefore, to facilitate the gradient computation in equation (18), we use Tensorflow 1.8 Python API to benefit from automatic gradient computation (automatic differentiation) (Abadi et al. 2016).

Experiments

In this section, we present our results for a randomly generated dataset, and CT scan images from the Lung Image Database Consortium image collection (LIDC-IDRI) (Armato III et al. 2015) with randomly generated appointment durations.

Randomly Generated Data

Inspired by medical appointment dataset (Hoppen 2016) on Kaggle, we generated a random dataset with 4 categorical covariates (gender, day of month, time of day, intensity) and the response variable shows the duration of the appointment (non-negative). We assume the ground truth distribution of the duration times can be from multiple different distributions, but in all of them mean and standard deviation is fixed:

$$\mu_{(x_2, x_3, x_4)} = 1 + 0.1x_2 + 0.4x_3 + 1.5x_4$$

$$\sigma_{(x_1, x_4)} = 0.1 + 0.2x_1 + 0.2x_4$$

Let us consider 4 possible distributions for the response in here:

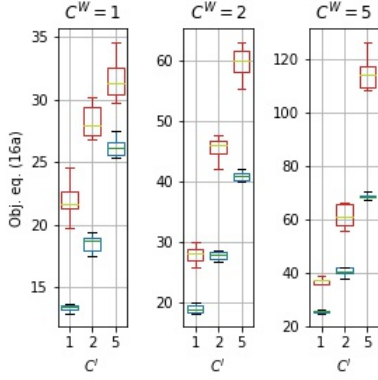


Figure 1: The Box-plots of the objective function (16a) for different combinations of (C^W, C^I)

Distribution	Description
Normal (trunc.)	$\mathcal{N}(\mu_{(\cdot)}, \sigma_{(\cdot)})$
Logistic (trunc.)	$LOG(\mu_{(\cdot)}, \sigma_{(\cdot)})$
Beta (scaled)	$B(\alpha(\mu_{(\cdot)}, \sigma_{(\cdot)}), \beta(\mu_{(\cdot)}, \sigma_{(\cdot)}))$
Uniform	$\mathcal{U}(a(\mu_{(\cdot)}, \sigma_{(\cdot)}), b(\mu_{(\cdot)}, \sigma_{(\cdot)}))$

Later, we add some noise to the labels from $\mathcal{U}(-1, 1)$. The appointment durations are not i.i.d, however, i.i.d is a conventional assumption for the inputs of the classic models in stochastic optimization.

To make sure predictions are non-negative, we scale the appointment durations between 0 and 1. This transformation is necessary since our predictions are outputs of a sigmoid function. Since the covariates are categorical, we consider One-Hot encoding transformations of them for the training.

Architecture: We considered a Fully-Connected (FC) architecture with two dense layers for this experiment. The input layer has $2 + 30 + 5 + 4 = 41$ nodes. The number of hidden nodes are chosen among powers of 2 (i.e. $2^n, 2 \leq n \leq 8$). Using grid search on the hyperparameters, we used the best performance among all of our runs. The weights are initialized using the Xavier initializer (Glorot and Bengio 2010). The output layer has only one node and returns the scaled predictions. For the choice of optimizer, we compared the results between SGD, ADAM (Kingma and Ba 2014), L-BFGS (Nocedal 1980) and chose the best one in each run. The learning rate is also chosen from $[0.001, 0.1]$ range using a grid search. We used different strategies such as constant, time inverse decaying, and exponential decaying to change the learning rates in each epoch and for each run.

Solutions of the Classic ASP

The Maximum Likelihood Estimates (MLEs) of the mean and standard deviation parameters from the generated data are computed and used in this section. From the dataset, we have estimated $\hat{\mu}_{MLE} = 5.9254$, and $\hat{\sigma}_{MLE} = 2.2151$ for the truncated Normal distribution.

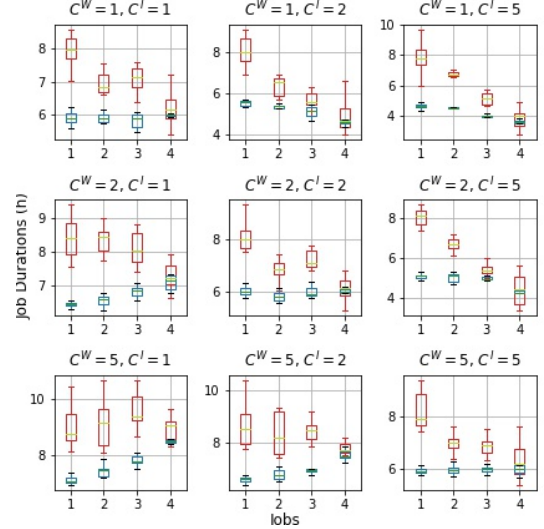


Figure 2: The Box-plots of the ASP solutions for different combinations of (C^W, C^I)

We use the Gurobi 7.0 solver API for Python to solve the models. The solver parameters for tuning are all set to their default values. We use SAA with 10 reps and equal number of samples to the original dataset for a model with $n = 5$ (scheduling 5 jobs on a single server). The Box-plots of the objective function (16a) for different combinations of C^W and C^I parameters from the list of $[1, 2, 5]$ are depicted in Figure 1. The blue plots are related to the SAA model and the red ones are related to the DRO model.

As we can see in Figure 1, the base cost of the planning for $C^W = 1$ and $C^I = 1$ is estimated to be 13.40 for SAA and 22.02 for DRO and it is increasing proportionally by C^W and C^I . Since DRO considers the worst-case distribution given the fixed mean and standard deviation, in all cases the total objective of DRO is higher than SAA. The actual decisions and their Box-plots are shown in Figure 2.

The interesting trends in Figure 2 are related to the optimal scheduled appointments when changing C^W and C^I . As we can see, in settings with high ratio of $\frac{C^W}{C^I}$, the length of the allowances are increasing by the order of the jobs and in settings with high ratio of $\frac{C^I}{C^W}$, the allowances are decreasing by the order of the jobs. Job duration estimates from DRO are generally higher than SAA. This is due to the fact that DRO considers the worst-case distribution of the job durations whereas SAA assumes the exact distribution. Further, DRO is more sensitive to the ratio of the wait/idle costs.

Comparison of IEO and SEO

To make a fair comparison, we use Maximum Absolute Deviation (MAD) for the loss function of the SEO and the proposed loss function (15) for the IEO. Moreover, we evaluate the actual objective value (16a) for the solutions from both method. We also make sure all other parameters and hyperparameters of the networks are the same in both settings.

Table 1: The objective function (16a) for different combinations of (C^W, C^I)

Dist.	C^w/C^I	SEO			IEO		
		1	2	5	1	2	5
Norm.	1	10.56	15.63	30.86	10.86	14.92	19.93
	2	16.04	21.11	36.34	14.85	21.73	32.49
	5	32.48	37.55	52.78	20.64	32.40	54.44
Log.	1	11.09	16.81	33.98	11.25	16.56	23.10
	2	16.46	22.18	39.35	15.81	23.79	36.54
	5	32.55	38.28	55.45	22.34	36.69	59.73
Beta	1	10.61	15.99	31.53	10.94	14.87	20.23
	2	15.84	21.22	37.35	14.90	21.78	32.15
	5	31.53	36.91	53.04	20.40	32.48	54.70
Uni.	1	10.50	15.81	31.72	10.91	14.75	20.16
	2	15.70	21.01	36.92	14.99	21.79	32.28
	5	31.30	36.61	52.52	20.42	32.28	54.49

As we can see in Table 1, the performance of the IEO is better than SEO in most of the settings. The most interesting part of this table is the objective function for the base parameters $C^W = 1$ and $C^I = 1$. In here, the solutions from the IEO is 10.56, and from SEO to 10.86 and both of them are better than the objective from classic model (13.40). This fact indicates that the i.i.d assumption is not valid.

From Table1, we can also conclude that the choice of the distribution does not really affect the performance of the both models. As we expected, the performance of the SEO drops for the cases with high values of C^W or C^I . This is due to the fact that in SEO we ignore these costs when we are estimating the job durations. The results clearly show how SEO can lead to sub-optimal solution.

LIDC-IDRI Dataset

The benefits of using DNNs for the contextual ASP do not end with bringing in context in decision making and creating a task-based loss function. In some cases, other ML algorithms such as SVM can be used, but they are often constrained by the assumption that we can define all the covariates. However, this is not necessarily true in many cases and making assumptions about non-numeric features can lead to low sensitivity. Deep learning is a fast-growing field that could be an ideal solution for these cases. This is due to the fact that we can learn features from raw images. In screening process for lung cancer, millions of CT scans will have to be analyzed. This is a controversial task even among radiologists. Further, radiologists often have to look through large volumes of these images that can lead to mistakes. Here, we skip the details about the screening process, and the dataset. We refer the interested readers to (Armato III et al. 2011) for more details about lung nodule analysis.

Consider a scenario in which we try to schedule patients for a second screening based on their initial phase results. Therefore, if the radiologist in the initial phase detects a nodule, the patient needs more time for the second phase. The process of scheduling patients in this manner can be an

Table 2: The objective function (16a) for different combinations of (C^W, C^I)

Dist.	C^w/C^I	IEO			SEO		
		1	2	5	1	2	5
Norm.	1	11.88	16.01	21.34	11.90	19.15	40.92
	2	16.84	23.65	34.19	16.54	23.79	45.56
	5	23.89	24.10	59.37	30.47	37.72	59.48

overwhelming burden for the radiologists. In this paper, we use the LIDC-IDRI dataset with 8723 CT scans. Similar to the previous example, we add more covariates (gender, day, hour) and we sample the response variable (appointment duration) from a Normal distribution using similar means and standard deviations to the previous experiment. The only difference in this experiment with the previous one is that the intensity is not a categorical column anymore and it is replaced by the image.

Architecture: We used 3 convolutional layers with 32, 64, and 128 filters of 2×2 kernels followed by 2 dense layers. The last dense layer has one output node with a sigmoid activation corresponding to positive and negative cases. Using small kernel sizes and not using max-pooling layers allow us to detect small areas of interest in the input images. We first resized all images to 55×55 pixels and then used several data augmentation methods including random cropping of 50×50 , random horizontal flipping, color jittering. In addition to our surrogate loss function, we used a binary crossentropy loss together with Adam optimizer for the classification task.

The training process with the convolutional layers is significantly more expensive, therefore, we allow for more training epochs in this experiment. The accuracy of the convolutional layers for the classification task increased up to %96 after the augmentation. The results in Table 2 suggests that the IEO outperforms SEO in all cases, specifically for higher values of C^W and C^I .

Conclusions

In this paper, we consider the contextual Appointment Scheduling Problem (ASP). If the probability distribution of the job durations is known for every possible combination of the data features, there is an exact solution for this problem. However, approximating a probability distribution is inaccurate, and impossible when data is scarce. The SEO approach to solve ASP is also shown to yield sub-optimal decisions. In particular, in cases that the idling or waiting cost is significantly higher than the other cost. To address this issue, we proposed the IEO approach with a surrogate loss function in which we estimate the durations using non-linear models. This approach does not require the knowledge of the job durations probability distribution and uses only historical data. We have shown the validity of our approach though the comparisons of the solutions from IEO with other approaches.

References

- [Abadi et al. 2016] Abadi, M.; Barham, P.; Chen, J.; Chen, Z.; Davis, A.; Dean, J.; Devin, M.; Ghemawat, S.; Irving, G.; Isard, M.; et al. 2016. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, 265–283.
- [Armato III et al. 2011] Armato III, S. G.; McLennan, G.; Bidaut, L.; McNitt-Gray, M. F.; Meyer, C. R.; Reeves, A. P.; Zhao, B.; Aberle, D. R.; Henschke, C. I.; Hoffman, E. A.; et al. 2011. The lung image database consortium (lidc) and image database resource initiative (idri): a completed reference database of lung nodules on ct scans. *Medical physics* 38(2):915–931.
- [Armato III et al. 2015] Armato III, S. G.; McLennan, G.; Bidaut, L.; McNitt-Gray, M. F.; Meyer, C. R.; Reeves, A. P.; and Clarke, L. P. 2015. Data from lidc-idri. the cancer imaging archive.
- [Ban and Rudin 2018] Ban, G.-Y., and Rudin, C. 2018. The big data newsvendor: Practical insights from machine learning. *Forthcoming in Operations Research*.
- [Bertsimas and Kallus 2014] Bertsimas, D., and Kallus, N. 2014. From predictive to prescriptive analytics. *arXiv preprint arXiv:1402.5481*.
- [Cardoen, Demeulemeester, and Beliën 2010] Cardoen, B.; Demeulemeester, E.; and Beliën, J. 2010. Operating room planning and scheduling: A literature review. *European journal of operational research* 201(3):921–932.
- [Cayirli, Veral, and Rosen 2006] Cayirli, T.; Veral, E.; and Rosen, H. 2006. Designing appointment scheduling systems for ambulatory care services. *Health care management science* 9(1):47–58.
- [Denton and Gupta 2003] Denton, B., and Gupta, D. 2003. A sequential bounding approach for optimal appointment scheduling. *IIE Transactions* 35(11):1003–1016.
- [Donti, Amos, and Kolter 2017] Donti, P. L.; Amos, B.; and Kolter, J. Z. 2017. Task-based end-to-end model learning in stochastic optimization. *arXiv preprint arXiv:1703.04529*.
- [Elmachtoub and Grigas 2017] Elmachtoub, A. N., and Grigas, P. 2017. Smart” predict, then optimize”. *arXiv preprint arXiv:1710.08005*.
- [Gallego and Moon 1993] Gallego, G., and Moon, I. 1993. The distribution free newsboy problem: review and extensions. *Journal of the Operational Research Society* 44(8):825–834.
- [Glorot and Bengio 2010] Glorot, X., and Bengio, Y. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 249–256.
- [Gould et al. 2016] Gould, S.; Fernando, B.; Cherian, A.; Anderson, P.; Cruz, R. S.; and Guo, E. 2016. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*.
- [Gupta and Denton 2008] Gupta, D., and Denton, B. 2008. Appointment scheduling in health care: Challenges and opportunities. *IIE transactions* 40(9):800–819.
- [Gurvich, Luedtke, and Tezcan 2010] Gurvich, I.; Luedtke, J.; and Tezcan, T. 2010. Staffing call centers with uncertain demand forecasts: A chance-constrained optimization approach. *Management Science* 56(7):1093–1115.
- [Hoppen 2016] Hoppen, J. 2016. Data from kaggle datasets. the medical appointment no shows.
- [Kingma and Ba 2014] Kingma, D. P., and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- [Kleywegt, Shapiro, and Homem-de Mello 2002] Kleywegt, A. J.; Shapiro, A.; and Homem-de Mello, T. 2002. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on Optimization* 12(2):479–502.
- [Koçağa, Armony, and Ward 2015] Koçağa, Y. L.; Armony, M.; and Ward, A. R. 2015. Staffing call centers with uncertain arrival rates and co-sourcing. *Production and Operations Management* 24(7):1101–1117.
- [Koenker and Bassett Jr 1978] Koenker, R., and Bassett Jr, G. 1978. Regression quantiles. *Econometrica: journal of the Econometric Society* 33–50.
- [Lindley 1952] Lindley, D. V. 1952. The theory of queues with a single server. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 48, 277–289. Cambridge University Press.
- [Mak, Rong, and Zhang 2014] Mak, H.-Y.; Rong, Y.; and Zhang, J. 2014. Appointment scheduling with limited distributional information. *Management Science* 61(2):316–334.
- [Mittal, Schulz, and Stiller 2014] Mittal, S.; Schulz, A. S.; and Stiller, S. 2014. Robust appointment scheduling. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 28. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- [Nocedal 1980] Nocedal, J. 1980. Updating quasi-newton matrices with limited storage. *Mathematics of computation* 35(151):773–782.
- [Oroojlooyjadid, Snyder, and Takac 2016] Oroojlooyjadid, A.; Snyder, L. V.; and Takac, M. 2016. Applying deep learning to the newsvendor problem. *CoRR* abs/1607.02177.
- [Shapiro, Dentcheva, and Ruszczyński 2009] Shapiro, A.; Dentcheva, D.; and Ruszczyński, A. 2009. *Lectures on stochastic programming: modeling and theory*. SIAM.
- [Song et al. 2016] Song, Y.; Schwing, A.; Urtasun, R.; et al. 2016. Training deep neural networks via direct loss minimization. In *International Conference on Machine Learning*, 2169–2177.
- [Taylor 2000] Taylor, J. W. 2000. A quantile regression neural network approach to estimating the conditional density of multiperiod returns. *Journal of Forecasting* 19(4):299–311.
- [Weiss 1990] Weiss, E. N. 1990. Models for determining estimated start times and case orderings in hospital operating rooms. *IIE transactions* 22(2):143–150.
- [White 1992] White, H. 1992. Nonparametric estimation of conditional quantiles using neural networks. In *Computing Science and Statistics*. Springer. 190–199.

[Zhang and Gao 2017] Zhang, Y., and Gao, J. 2017. Assessing the performance of deep learning algorithms for newsvendor problem. In *International Conference on Neural Information Processing*, 912–921. Springer.

[Zhang, Shen, and Erdogan 2017] Zhang, Y.; Shen, S.; and Erdogan, S. A. 2017. Distributionally robust appointment scheduling with moment-based ambiguity set. *Operations Research Letters* 45(2):139–144.