## Almosüfer تجوّل:tajawal

How Data Science & Analytics Tools are used by the Al Tayyar Group to Support & Enable the Data-Driven Organisation

**DataCon Africa** 

### Who are we?

**Online Business Unit** 

NICH

AMSTERDAM

HONG KONG

NEW

LONDON





Vision: To be the **dominant** online travel solution in the **MENA** region, powered by **innovative technology**.

## Alussier's exceptional Journey of growth



Largest Online Travel Agency in Saudi Arabia

### The data strategy follows the customers





## Who do we support?



## **Different tools for different purposes**



Business requirements define the level of reporting and data integration needed.

Tools can be used in isolation or combined for cross-platform analysis where needed.

## Implementation of self-service tools are crucial to **empower a data driven organization**

→ Simple requests/tasks can be handled by anyone, leaving more time & resources for deeper business analytics to the Data team.

### How we got started?



## I. Cancellation Prediction model v1

## Almosüfer تجوّل:tajawal

### Business Background Problem statement

A customer reserves a hotel and chooses the payment option 'Pay Later' Customer service team sends several reminders to the customer to check if they'd like to make the payment

The customer has until 2 days before the check-in date to make the payment If no payment is made and no feedback from the customer CS team manually cancels the booking 2 days prior to the check-in date

A large share of Hotel bookings made are "Book Now / Pay Later" (~30-50%)

- Depending on seasonality, ~ 70-90% of these bookings are cancelled.
- For 'Pay Now' payment option only ~3% of bookings are cancelled
- Customers are reminded via SMS, in-App, Emails, Phone calls to complete the payment

"Our <u>revenue estimates</u> for the month <u>are inaccurate</u> due to the large amount of last minute cancelled pay later bookings."

"The several follow up attempts require <u>time and effort from our customer</u> <u>service</u> team and <u>cost money</u>, while our <u>customers</u> eventually are <u>inconvenienced</u>."

## Almosüfer تجۆل:tajawal

### **AIM OF THE CANCELLATION PREDICTION MODEL**

### TO PREDICT THE PROBABILITY OF CANCELLATION OF PAY LATER HOTEL BOOKINGS AT TRANSACTION LEVEL



## **Model Preparation**

**Model Selection & Feature Engineering** 



## **Model Selection**

We used **Dataiku**, a Data Science platform to run preliminary models on the data and see which algorithm performs best and built our model in Python to implement it on a transaction level

			dat ku
Random fo 2018-10-24 - 16:1	1:07		
Trees	100		
Depth	20		
Min samples	1		
Was grid of size	12		
		Log Loss: 0.09	93

XGBoost 2018-10-24 – 11:01:32

Trees	300
Max depth	5
Was grid of size	100

#### *Log Loss:* 0.118



Log Loss: 0.134

#### tajawal: تجوّل Almos Ter 12

## **Feature Engineering**

#### Dataset:

- transactional data for Hotel and Flight Bookings from Almosafer since 2017 [>2,000,000 data points]
- Created **34 explanatory variables**, that were potentially expected to influence a customer's cancellation

#### • Examples:

- Count of previous cancellations of the customer
- Customer also booked flight to the same city
- If customer has any other Hotel booking with the same check-in date
- Ratio of pay later bookings of each customers
- Ratio of bookings with promo code per customer
- CLV vs average historical CLV
- Average stay duration vs current stay duration per customer

#### Distribution of cancelled (1) and not cancelled (0) PAY LATER transactions



## **Model Outcome**

**Model Evaluation & Next Steps** 



## Model evaluation & Variable Importance

Ongoing optimization of input variables.

n = 225,176	Predicted					
Actual	FALSE	TRUE				
FALSE	TN = 28,347	FP = 16,049				
TRUE	FN = 4,239	TP = 176,541				





 To be implemented at a transaction level in real time

Model to be enhanced further to improve results even more



## **Cumulative Gain Chart**



The goal of this curve is to visualize the benefits of using a model for targeting a subset of the population.

On the horizontal axis, we show the percentage of the population which is targeted and on the vertical axis the percentage of found positive records.

 $\rightarrow$ This basically represents the benefit of using the model vs without the model.

The **dotted diagonal** illustrates a random model

![](_page_15_Picture_6.jpeg)

# II. BI Chatbot

Kick-Off v1.0

## Almosöfer tajawal:تجوّل

#### **Problem Statement:**

There a several dashboards & recurring reports that provide rich data for specific purposes.

To support the correct interpretation of this data, a chatbot can be installed on slack to answer common data enquires and enable fast & accurate query resolution.

### **Opportunities:**

- Make data more accessible & interpretable
- Provide on instant results for regular requests
- Direct people to appropriate dashboards for further analysis.

## **Meet Andy**

![](_page_17_Figure_8.jpeg)

### Almosüfer

تجوّل:tajawal

### **Design Overview**

	Requirements	Possible Approaches			
Data Layer	<ul> <li>Access database to be able to respond to requests.</li> <li>A data-structure to help the bot recommend appropriate tableau dashboards based on the request.</li> </ul>	<ul> <li>Specific SQL queries for specific requests.</li> <li>Read in and store locally the data for specific requests to allow for fast responses.</li> </ul>			
Business Logic Layer	<ul> <li>Allow the bot to interpret the message and present the appropriate result for each response.</li> </ul>	<ul> <li>Text based response.</li> <li>Create a small tableau dashboard for each metric and direct people to this for further analysis</li> </ul>			
Natural Language Processing/ Generative Speech	<ul> <li>Understand conversational requests from users</li> <li>Interpret them into an appropriate action/response.</li> </ul>	<ul> <li>Procedural/pre-defined responses to key questions.</li> <li>NLP using deep learning Neural Networks for a more generalised response.</li> </ul>			
API/ Presentation Layer	<ul><li>Send and receive requests/responses to Slack.</li><li>Push notifications for recurring updates or "mini-reports"</li></ul>	The API to communicate with Slack			

**Note:** Tableau provides an API for python. The API only allows you to access the fields/views that are available in tableau. It is not possible to access the data values or filters from the API and therefore this cant be data source.

### NLP – What question is being asked?

This can be used to **determine the context of the question** and therefore the information needed to **derive an appropriate response**.

#### 1. Procedural keyword identification and pre-defined responses.

- Straight forward procedural implementation.
- Predictable

### 2. NLP

- A. NLTK [http://www.nltk.org/book/]
  - The core of most NLP packages in python.
  - Very low level and allows you to create very bespoke implementations.
  - Can be optimized to be very fast.

#### B. SpaCy

- Wraps the functionality of NLTK into an easier implementation.
- Simplifies many tasks, but doesn't allow for as much freedom as NLTK

### **Business Logic – How to respond?**

After a question has been identified, this will be used to generate the appropriate response.

- 1. **Procedural/ structured response** based on the NLP generated keywords.
  - The simplest to implement.
  - Would become very difficult to manage once the business logic become large.
- 2. More structured/tree based model.
  - If used in conjunction with a good classification solution > more manageable than procedural.
  - Would allow for the traversal of the tree for more depth of detail. E.g. revenue -> tajawal -> Organic.
- 3. NLP: Use machine learning to automatically build a response based on previous enquiries
  - Can generate responses based on previous interactions
  - Much more difficult to implement but more scalable.
  - Prone to mistakes if badly trained. Requires good training data (example questions).

→ Need to evaluate Value of each enhancement vs. Effort/Resources

### **Putting it together**

#### Question:

"What was the IBV for Almosafer in Q1 2017?"

#### NLP:

Extract keywords: "IBV", "Almosafer, Q1, 2017.

#### Business Logic:

- METRIC [IBV] = Revenue (dictionary Revenue, IBV, Sales, etc..)
- SUBSET [dim Brand] = Almosafer
- PERIOD [Year, Quarter] = Year=2017, Quarter=Q1

#### Data Layer:

Locate the most appropriate response  $\rightarrow$  Revenue\_request()

- Response: The <METRIC> in <PERIOD> for <SUBSET> was <VALUE><FORMAT>
- Extract Revenue from transactional database for the period and subset and populate the response

![](_page_21_Picture_13.jpeg)

#### Response/API:

Response: "The Revenue in Q1 for Almosafer was \$1.57 million"

### Scope

- This is an ongoing project that will be completed in phases.
- Each phase will add additional functionality and intellect to the bot.

	Scope
Phase 1	<ul> <li>Complete the NLP and API to allow for communication.</li> <li>Implement basic/common responses in the business logic and data layers.</li> </ul>
Phase 2	<ul> <li>Develop a more saleable/maintainable business logic layer and implement further responses.</li> <li>Fine tune the NLP layer.</li> </ul>
Phase 3	<ul> <li>Create a business structure for staff:</li> <li>Automatically message people with relevant KPI's on a regular basis ("Mini reports")</li> <li>Message relevant parties when KPI's are beyond a set limit. (Alerts)</li> </ul>
Phase 4	<ul> <li>Add a more advanced business logic to allow the bot to add helpful details. E.g. revenue is down this week. What might be the cause?</li> </ul>
Phase 5	• Implement a text-to-speech component for audio questions and responses (e.g. alexa/siri integration)

### Examples

### Andy's responses & reports

	lexander	Metlewic	<b>7</b> 10.52 AM						-	Alexand	er Met	lewicz 1:36 PM
Strictly Co	Arexander V So Ar C C C C	Alexand Ik A T 2 Ir ** fl h p ((	der Metlewicz Alexander M What was th A A A A Andy The t 	11:32 AM letlewicz 2:10 P e revenue by bu APP 10:18 AM Andy APP 1 The top airl saudi arab flynas: Flyadeal: air arabia: emirates: { saudi gulf egyptair: { indigo: flydubai: { etihad airv You can ge	In and and p In and and p In a second seco	oroduct onth are: op campa s: S: S: S: S: S: S: S: S: S: S: S: S: S:	Iast year?         Image: Constraint of the second of th	J ⊆ ⊆ 云 ☆         9-02-24 - 2019-03-02         WAL************************************	Andy A Questic @Mam week top p top p top p top p top f DIAG DIAG @Saud Hi help What What	Alexand hi hi PP 10:46 AM ons asked by all use douh Ahmed ly kpi romotions almosafer romos this week romo irlines this month ight destinations to i FAQ i UQ Altamimi	er Met	Andy APP 1:36 PM Bonjour!

![](_page_24_Picture_1.jpeg)

Andy APP 2:09 PM Hasta la vista!

## Almosäfer تجوّل:tajawal

Thank you