## Scaling Machine Learning @ Careem

Careem

VITALII DUK October, 2019 Careem is the Middle East's leading platform in transportation, delivery & payments.





What's hard about scaling Machine Learning?

### **Problem Statement**









## **Production ML workflow**









At Careem we've faced several issues while productizing ML models in the past:

- How to train & deploy a lot of **city-specific ML models** fast?
- How to **A/B test** the performance of several ML models running in parallel?
- How to perform **real-time monitoring** of hundreds of models running in production?





**CMLP** is an in-house built **Machine Learning Platform** designed to:

- Automate and accelerate the adoption of Machine Learning to support Careem's Platform vision
- **Reduce** the amount of **time** Data Scientists and Engineers need to deploy ML models into production
- Make Data Scientists **spend time more effectively**, working on more exciting tasks
- Ensure that ML models perform the same in **R&D and production** environments



Predictions provided per day

**Step 1: Data Preparation** 







Dalasels			
Project cancellation_prediction	✓ State	us	~
SQL File Name (i)	Last Change 🛈	Status (i)	
cancellation_training_dataset.sql	26/08/2019 01:24 PM	Approved	EDIT VIEW

### **Step 2: Model Design & Training**





## Model Training - code example



```
# Read generated dataset
training_input = TrainingInput('dataset.pickle')
# Define classification model parameters
model parameters = {
    'algorithm':
                        CatBoostClassifier,
    'loss function':
                        'Logloss',
    'iterations':
                        JobOptions.get_as_int('iterations', 1500),
                        JobOptions.get_as_float('learning_rate', 0.04),
    'learning_rate':
    'depth':
    'use_best_model':
    'one hot max size': 10
job info = JobInfo(
   name=config['job_name'],
   project_name=config['project_name'],
   dimension_type=DimensionType.SERVICE_AREA,
   dimension=config['dimension']
```

```
# Model wrapper
catboost_model = CatBoostTrainer(**model_parameters)
```

```
trainer = Model(
    input_data=training_input,
    version='0.0.3',
    model_instance=catboost_model,
    job_info=job_info
```

```
# Run model training
trainer.run()
```

```
# Report metrics
reporter = trainer.get_reporter()
reporter.evaluate_metrics(
    metric_calculator_list=[roc_auc_score],
    target='fraud_probability'
```





< Trainin	ıg job detail	s #3460			
Model Name xgb_model	Version 0.0.3	Project booking_decline	Dimension Type Country	Dimension UAE	
Training Time 15th October 201	19, 07:34 AM				¥
Mode production					
Job Configs					~
	RFLOW 📝 EDIT PIP	ELINE			
O DOWNLOAD	MODEL FILES				

## **Model Training - metrics**



### **Step 3: Scheduling Model Training**

## **Training Scheduler - UI**

### < Create a new pipeline

~
~
~
~

Dataset Generation Options	5	
e.g.: trainingDays	30	+
Start Time 25/10/2019 07:24 PM		
Schedule Interval Daily		~
Run Job Options		
Кеу	Value	
Pipeline Options		
Key	Value	
CREATE PIPELINE	CANCEL	

### **Step 4: Deploying Trained Model**



## **Model Serving - canary deployment**











### **Model Serving - UI**



Contract Continuation for project "fraud\_detection"

Configuration Version		372	
		<u> </u>	
Version 4			
Environment			
Production		~	
Servo Version			
3.0.0		~	
MODELS	DEPLOYMENT		
Dimension			
Islamabad (58)		~	
Models			
fraud_detection-58-catboost	_model-latest 🛞		
fraud_detection-58-sklearn_l	ogreg_model-latest 🔞		

### Metadata Governance

## **Metadata Governance**



We are using Apache Atlas as the metadata management and governance tool to maintain a catalog of our data assets.

trip_details_	by_zone (h	ive_table	e)	
Classifications: +				
Term: +				
Properties Lineage	Relationships Classi	ifications Audits	Schema	
/atlas/trip_detai	create external t	trip_details		
			create table trip	trip_details_by_z
/atlas/zone_lookup	create external t	trip_zone_lookup	0	
	<b>►</b> Ø			

### **Hyper-Parameter Optimization (HPO)**

## **Hyper-parameter optimization**







## **HPO - code examples**

# Read generated dataset
training\_input = TrainingInput('dataset.pickle')

# Job information
job\_info = JobInfo(
 name=config['job\_name'],
 project\_name=config['project\_name'],
 dimension\_type=DimensionType.SERVICE\_AREA,
 dimension=config['dimension']

# Model wrapper
catboost\_hpo\_model = HpoCatBoostTrainer(
 problem\_type='classification'

trainer = Model( input\_data=training\_input, version='0.0.4', model\_instance=catboost\_hpo\_model, job\_info=job\_info

# Run model training
trainer.run()

# Report metrics
reporter = trainer.get\_reporter()
reporter.evaluate\_metrics(
 metric\_calculator\_list=[roc\_auc\_score],
 target='fraud\_probability'



## **HPO - code examples**

# Read generated dataset
training\_input = TrainingInput('dataset.pickle')

# Job information
job\_info = JobInfo(
 name=config['job\_name'],
 project\_name=config['project\_name'],
 dimension\_type=DimensionType.SERVICE\_AREA,
 dimension=config['dimension']

# Model wrapper catboost\_hpo\_model = HpoCatBoostTrainer( problem\_type='classification'

trainer = Model(
 input\_data=training\_input,
 version='0.0.4',
 model\_instance=catboost\_hpo\_model,
 job\_info=job\_info

# Run model training
trainer.run()

# Report metrics
reporter = trainer.get\_reporter()
reporter.evaluate\_metrics(
 metric\_calculator\_list=[roc\_auc\_score],
 target='fraud\_probability'

# Read generated dataset
training input = TrainingInput('dataset.pickle')

# Define classification model parameters

model\_parameters = {
 'algorithm': CatBoostClassifier,
 'loss\_function': 'Logloss',
 'iterations': JobOptions.get\_as\_int('iterations', 1500),
 'learning\_rate': JobOptions.get\_as\_float('learning\_rate', 0.04),
 'depth': 8,
 'use\_best\_model': True,
 'one\_hot\_max\_size': 10

# Job information
job\_info = JobInfo(
 name=config['job\_name'],
 project\_name=config['project\_name'],
 dimension\_type=DimensionType.SERVICE\_AREA,
 dimension=config['dimension']

# Model wrapper
catboost\_model = CatBoostTrainer(++model\_parameters)

trainer = Model(
 input\_data=training\_input,
 version='0.0.3',
 model\_instanc==catboost\_model,
 job\_info=job\_info

# Run model training
trainer.run()

# Report metrics
reporter = trainer.get\_reporter()
reporter.evaluate\_metrics(
 metric\_calculator\_list=[roc\_auc\_score],
 target='fraud\_probability'

## **HPO - code examples**

# Read generated dataset
training\_input = TrainingInput('dataset.pickle')

# Job information
job\_info = JobInfo(
 name=config['job\_name'],
 project\_name=config['project\_name'],
 dimension\_type=DimensionType.SERVICE\_AREA,
 dimension=config['dimension']

# Model wrapper catboost\_hpo\_model = HpoCatBoostTrainer( problem\_type='classification'

trainer = Model( input\_data=training\_input, version='0.0.4', model\_instance=catboost\_hpo\_model, job\_info=job\_info

# Run model training
trainer.run()

# Report metrics
reporter = trainer.get\_reporter()
reporter.evaluate\_metrics(
 metric\_calculator\_list=[roc\_auc\_score],
 target='fraud\_probability'

# Read generated dataset
training input = TrainingInput('dataset.pickle')

Define classification model parameters

model\_parameters = {
 'algoritmm': CatBoostClassifie,
 'loss\_function 'Logloss',
 'iterations': Juc90n\*Lons.get\_as\_int('iterations', 1500),
 'learning\_rate': Juc00ptices.get\_as\_float('learning\_rate', 0.04),
 'depth': 8,
 'use\_best\_model': True,
 'ost\_not\_max\_size': 10

# Job information
job\_info = JobInfo(
 name=config['job\_name'],
 project\_name=config['project\_name'],
 dimension\_type=DimensionType.SERVICE\_AREA,
 dimension=config['dimension']

# Model wrapper
catboost\_model = CatBoostTrainer(\*\*model\_parameters)

trainer = Model(
 input\_data=training\_input,
 version='0.0.3',
 model\_instanc=catboost\_model,
 job\_info=job\_info

# Run model training
trainer.run()

# Report metrics
reporter = trainer.get\_reporter()
reporter.evaluate\_metrics(
 metric\_calculator\_list=[roc\_auc\_score],
 target='fraud\_probability'

### AutoML

**AutoML** 





Al use-cases @ Careem

# **Suggested drop-off locations**



#### Why?

- better user experience
- reduced time to booking
- fewer cancellations

#### How?

- supervised ML algorithms
- trained on historical trips data
- >200 user, location and time-related factors used



# **Meeting points**



#### Why?

Customer and Captain can't find each other,

#### which has a direct impact on:

- cancellation rate
- contact rate
- the actual time of arrival (ATA)

#### and indirect impact on:

- retention
- billable availability (BA)
- customer throughput rate (CTP)

#### How?

- unsupervised ML algorithms
- based on historical data about "successful" pickups / drop-offs



## **Careem NOW - recommendation engine**









Man'oushe Street 4.7 Lebanese, Arabic \$\$\$ 15% Off - Limited Time Offer



Shaikh A 4.6 Ar 
⇒ Filter ↓↑ Sort 20% Off - Limited Time Offer

#### Why?

- provide personalized experience: no burgers for vegetarians
- improve conversion, or secondary metrics keeping the conversion the same

#### How?

 collaborative filtering algorithms trained on historical user's behavior







- Continue working on more sophisticated model **auto rollout** features
- Enhance AutoML capabilities by adding automated feature engineering (e.g. featuretools)
- Supercharge model **drift detection**
- Improve the latency of tree-based ML models using binary compilers (e.g. **treelite**)
- Move more features to the **UI**, to improve the UX
- Add **permutation importance** as a default tracked metric



# Shukran! Thank you! Shukriya! Merci!