# AI-POWERED FINTECH:
# AR FORECASTING WITH DATABRICKS & MLFLOW

**Floriant Sturm & Julie Vanackere**
**12/06/2024**

# Nice to meet you!

Julie Vanackere
Data scientist

Floriant Sturm
Co-founder

DATA AI SUMMIT

# Outline

## 2 main topics

| Business perspective | Technical perspective |
|---|---|
| 1. Why Accounts receivable (AR) forecasting?<br><br>2. How do we approach this?<br><br>3. How can AR forecasting become a plug & play solution and what are the **technical requirements**? | 1. Automated infra deployments<br><br>2. Standardized feature engineering<br><br>3. Standardized ML training<br><br>4. An automated ML lifecycle<br><br>5. Monitoring for customer confidence |

# Business perspective

# Why did we focus on AR forecasting?

## What is AR forecasting?

### Accounts-Receivable
~ Outstanding invoices

- Currently companies use a <u>reactive approach</u> to chase "late-payers"
  - They contact them after it is too late
  - They know their outstanding invoices

- Big issue
  - "Cash flow is the pulse of the company"

### Forecasting

- Predicting <u>when</u> a customer will pay their invoice in the future

- This will help you to anticipate
  - Who will pay late
  - How much cash is to be expected in the next x period

- Contact strategy: incentivize customers that are likely to pay late = <u>proactive approach</u>

# Why did we focus on AR forecasting?

## Why is it relevant for all companies?

- Quick win!
  We only need the historical invoice data to get started

- Predicting future cash flows reduces:
  - Credit risk
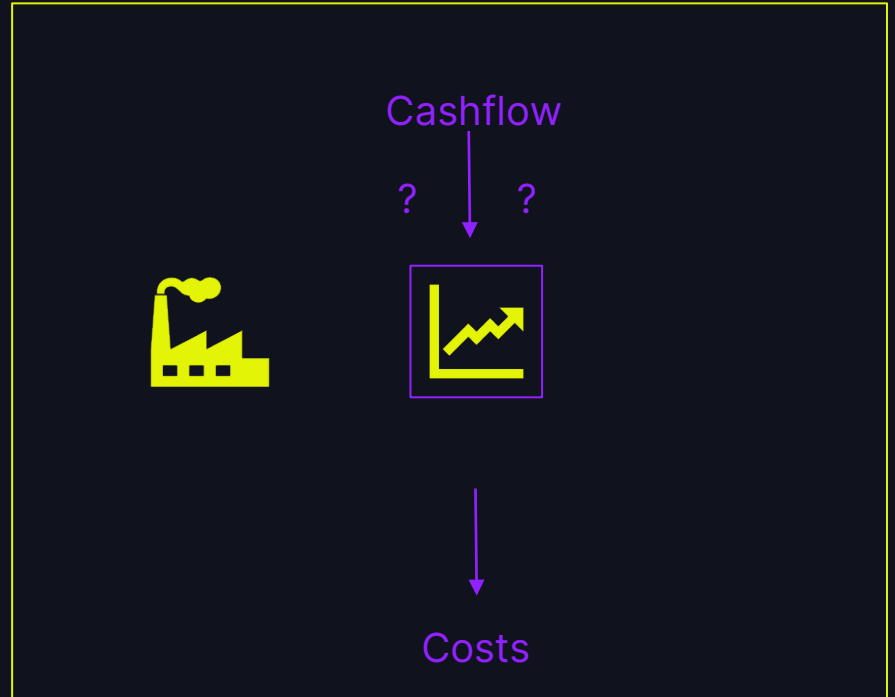  - Plan expenses, investments & potential savings

So... the proactive approach is the clear way to go!

# Why did we focus on AR forecasting?

## A tangible example

Context:

- **Production** company of pharmaceuticals
- Need to expand the production plant and invest in machines
- Do we have the cash flow to cover the costs?

Cashflow

? ?

Costs

# How do we approach this?

## How do we provide a sustainable approach?

### Where do we generally focus on?

- Identifying the business problem

- Strategy - focus on a **sustainable** solution
  - Provides direct impact
  - Efficient implementation
  - Easy maintainable by the client

- Business validation

- Coaching & development

### How do we look at Data Science?

- We try to go beyond, but how?
  - We keep the baseline structure (gathering data, etc...)
  - ...but the core business use case is tackled by Data Science

- Data Science means ML, AI,... whatever suits the business case best

# How can AR forecasting become a plug and play solution?

## What were our initial requirements?

1. Automated infra deployments

2. Standardized feature engineering

3. Standardized ML training

4. Automated ML lifecycle

5. Monitoring for customer confidence

# Technical perspective

# Automated infra deployments

## Deploy Azure infrastructure quickly through Terraform

Infrastructure-as-code that deploys:

- Resources
  - Databricks
  - ADF
  - Storage Account
  - Keyvault

- Networking

- Roles and responsibilities

- 2 environments

All automated through scripting

### What does terraform look like?

```
module "e61-tff" {
  source = "../e61-tif"

  tags            = var.tags
  global_settings = var.global_settings

  resource_groups = var.resource_groups

  networking = {
    vnets                           = var.vnets
    route_tables                    = var.route_tables
    routes                          = var.routes
    network_security_group_definition = var.network_security_group_definition
  }

  security = {
    keyvaults               = var.keyvaults
    keyvault_access_policies = var.keyvault_access_policies
  }

  storage_accounts = var.storage_accounts

  analytics = {
    databricks_workspaces = var.databricks_workspaces
  }

  role_mapping = var.role_mapping
```
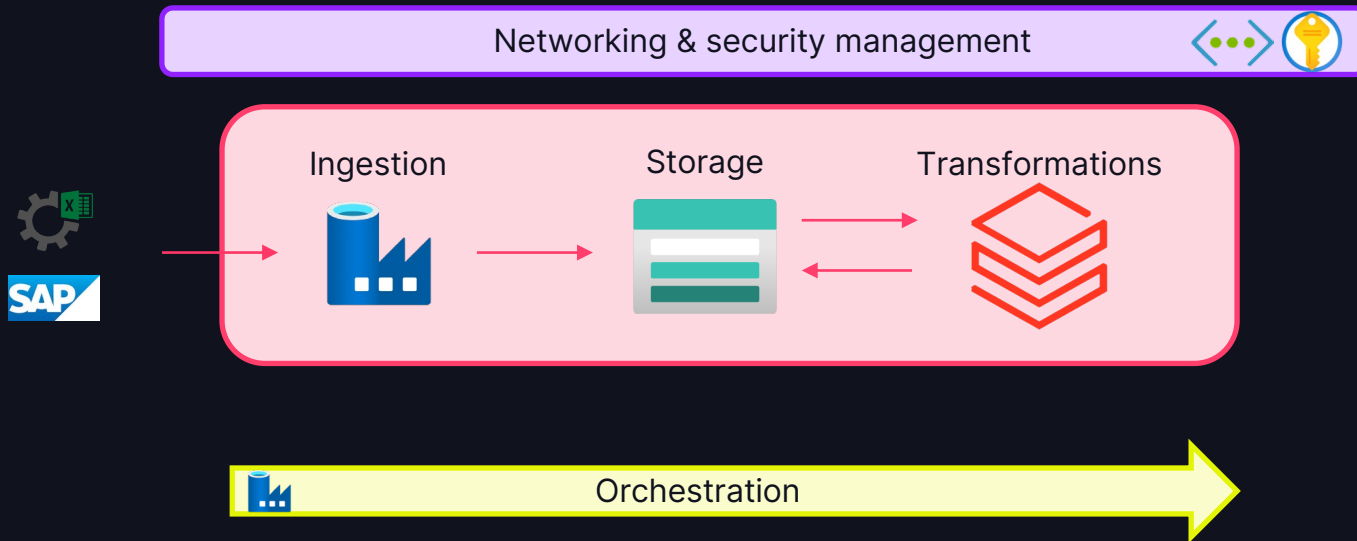
# Automated infra deployments
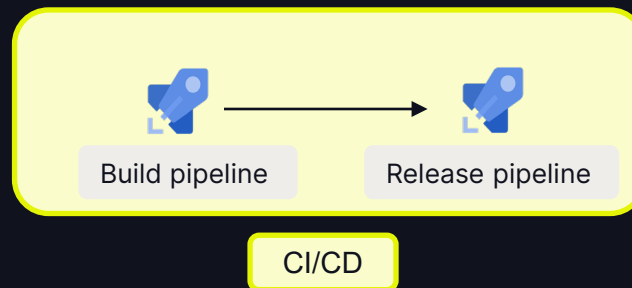
We start with a Modern Data Platform in Azure



Networking & security management

Ingestion → Storage → Transformations

Orchestration

DATA AI SUMMIT

# Automated `infra` deployments

## This facilitates a development – (acceptance) – production set up
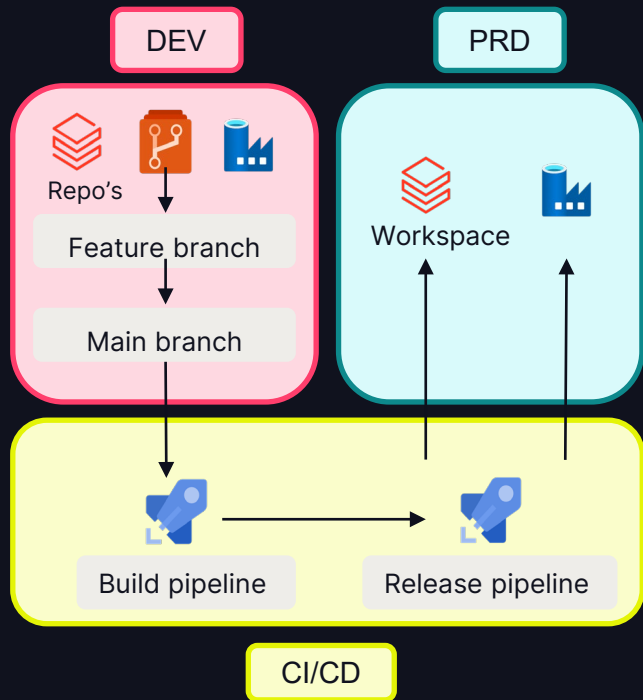
- Because of **terraform** different <u>environments</u> with the same resources can be easily setup

- Because of the **CICD pipelines**, <u>code</u> can be reproduced in these environments

- But how do we do this <u>practically</u>?



Build pipeline → Release pipeline

CI/CD

# Automated infra deployments

## Afterwards, we deploy our code to Databricks and ADF using Devops CI/CD



- Development environment
  - Code is stored in <u>Azure devops</u>
  - <u>Databricks repo</u> code is deployed in Databricks workspace (ML models)
  - <u>ADF GIT</u> is deployed to ADF live mode

- Production environment

  <u>Stable code</u> (finished ML models) runs in prd environment if used for critical business processes = extra layer

  ☰ Fully managed through scripts

DATA AI SUMMIT

# How can AR forecasting become a plug and play solution?

## What were our initial requirements?

1. Automated infra deployments ✅

2. Standardized feature engineering

3. Standardized ML training

4. Automated ML lifecycle

5. Monitoring for customer confidence

# Standardized feature engineering

## Feature table for AR forecasting

### Invoice-level features (mandatory)

- Year invoice was created
- Month in which the invoice is due
- Document type
- # Line items in invoice
- The invoice amount

### Customer-aggregated features (optional)

- % previous invoices late
- # of previous invoices
- Whether the last invoice was late (0/1)
- Preferred payment date

### Data collections (optional)

- When and with what action did we contact the customer?
- At which dunning level?
  - 1: sending reminder
  - 2: calling
  - 3: giving a fee

DATA AI SUMMIT

# Standardized ML training

## We use the specified features, to make predictions

### Features

**Invoice-level features**

- Year invoice was created
- Month in which the invoice is due
- Document type
- # Line items in invoice
- The invoice amount

**Customer-aggregated features**

- % previous invoices late
- # of previous invoices
- Whether the last invoice was late (0/1)
- Preferred payment date

**Data collections**

- When and with what action did we contact the customer?
- At which dunning level?
  - 1: sending reminder
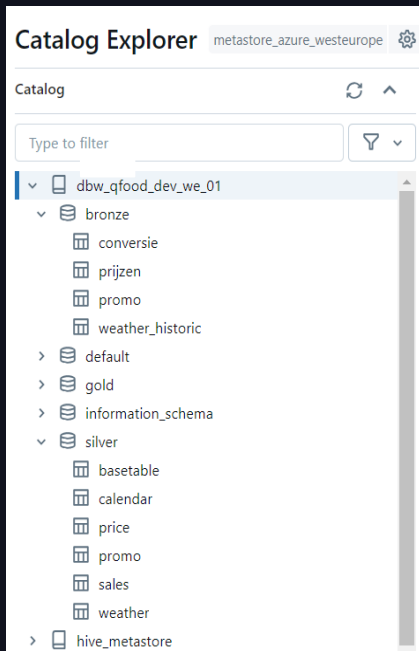  - 2: calling
  - 3: giving a fee

### Target

- Regression: # Days late

- Classification: in buckets
  - On time
  - 0-30 days late
  - 30-60 days late
  - > 60 days late

DATA AI SUMMIT

# Standardized feature engineering

## Where do we store these features & labels?



Delta tables

- Delta files stored on data lake

- ACID

- Natively integrated with Unity Catalog

- Upserts & truncate insert

DATA+AI SUMMIT

# How can AR forecasting become a plug and play solution?

## What were our initial requirements?

1 Automated infra deployments ✅

2 Standardized feature engineering ✅

3 Standardized ML training

4 Automated ML lifecycle

5 Monitoring for customer confidence

# Standardized ML training

## What is AutoML and how do we use it?



- AutoML: interface & code based model training in databricks

- Can be used for exploration, but we use it for model training as a whole

- The best model (according to R2) is automatically stored in Mlflow registry

- We track the feature importance to iterate on

DATA+AI SUMMIT

# How can AR forecasting become a plug and play solution?

## What were our initial requirements?

1 Automated infra deployments ✅

2 Standardized feature engineering ✅

3 Standardized ML training ✅

4 Automated ML lifecycle

5 Monitoring for customer confidence

# An automated ML lifecycle

## How do we manage the model lifecycle?



**Train a model**

Experiments

**Model in production**

Model registry

**Model in staging**

Model registry

- We use MLFlow – natively integrated in Databricks

- Everything we need:
  - Experiments
  - Model registry with lifecycle mgmt.
  - Python SDK (automation ☺)

# An automated ML lifecycle

## How do we choose to update the production model?

- Difference between technical KPI (R2) and business KPI (# Days late)
- Technical KPI as benchmark
  - <u>No actions</u> towards customers
  - <u>Gradually improves</u> when retraining the model
- Business KPI as benchmark
  - Switch to Business KPI when <u>technical KPI gradually worsens</u>
  - <u>Actions</u> have been taken towards customers
  - The ML model is used in the business for a while now
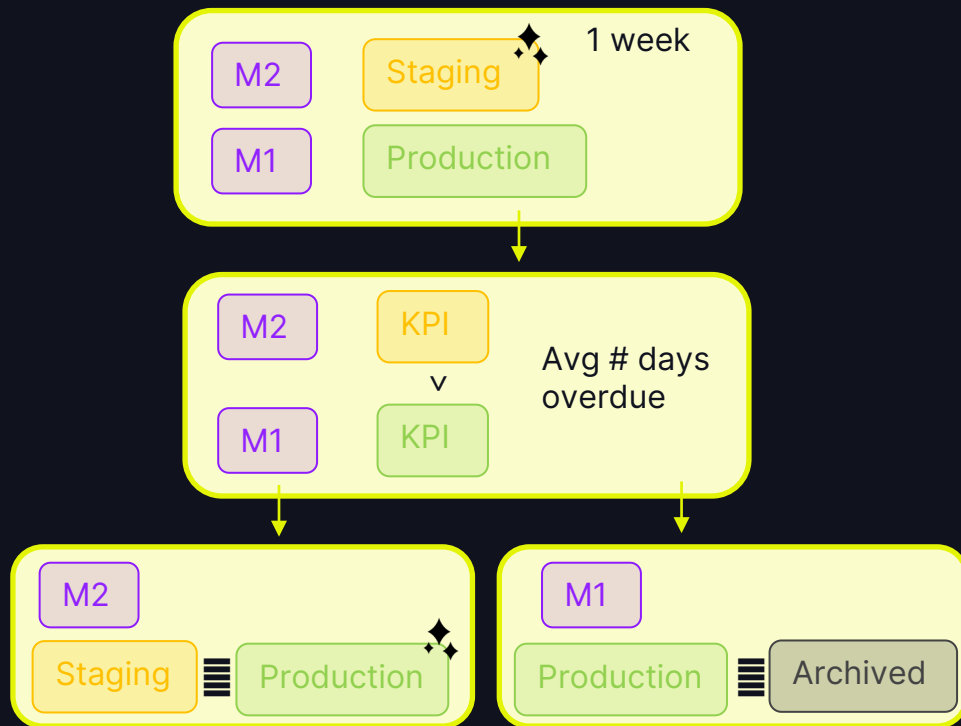- Let's say we focus on the second scenario

# An automated ML lifecycle

## How do we choose to update the production model?

We are solving a business problem, so

- We use a business KPI
  – RMSE on actuals

- We retrain the model with new data every week.

- The new model becomes a staged model that « shadows » the production model

- We only update the production model when our business KPI improved

DATA AI SUMMIT

# How can AR forecasting become a plug and play solution?

## What were our initial requirements?

1. Automated infra deployments ✅

2. Standardized feature engineering ✅

3. Standardized ML training ✅

4. Automated ML lifecycle ✅

5. Monitoring for customer confidence

DATA⁺AI SUMMIT

# Monitoring for customer confidence

## How do we enable trust? – the most important 'KPI' of an AI-solution

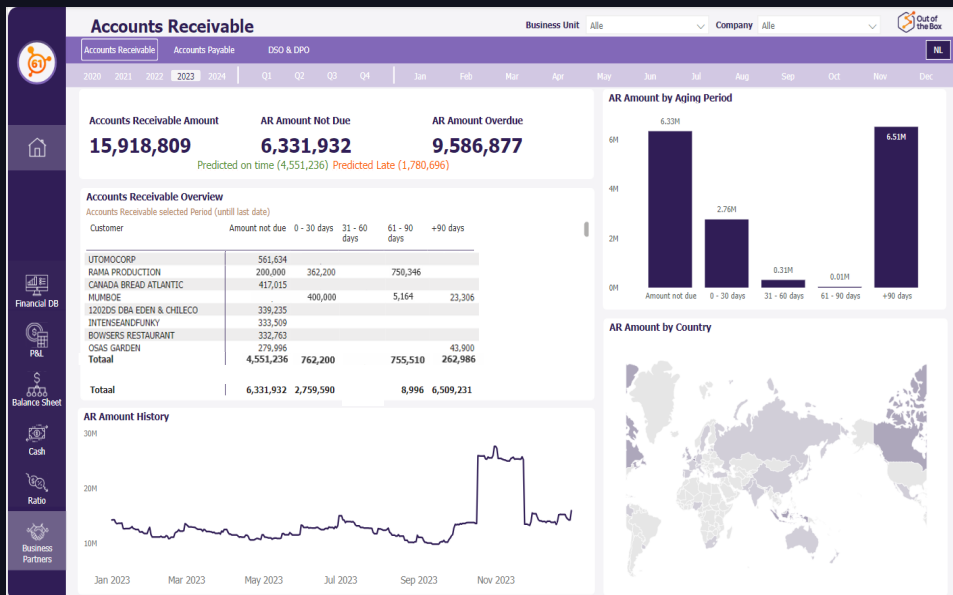| Dashboard for customers: | Dashboard for data scientists: |
|---|---|
| Build insights in the actuals vs predictions<br>• Amount/invoice that will be overdue per customer<br>• How many days this will be overdue<br>• Action list: which customers to target?<br>• Comparison of cash flows to future investments | Monitor data & models over time:<br>• Model versions - keep track of historic versions<br>• Model performance - technical KPIs<br>• Model performance - business KPIs<br>• Data drift |

# Monitoring for customer confidence

## Example of a customer insights dashboard



- Periodic buckets with amounts

- Actionable dashboard

▤ Contact those with large amounts with 90+ days predicted

DATA+AI SUMMIT

# How can AR forecasting become a plug and play solution?

## What were our initial requirements?

1 Automated infra deployments ✅

2 Standardized feature engineering ✅

3 Standardized ML training ✅

4 Automated ML lifecycle ✅

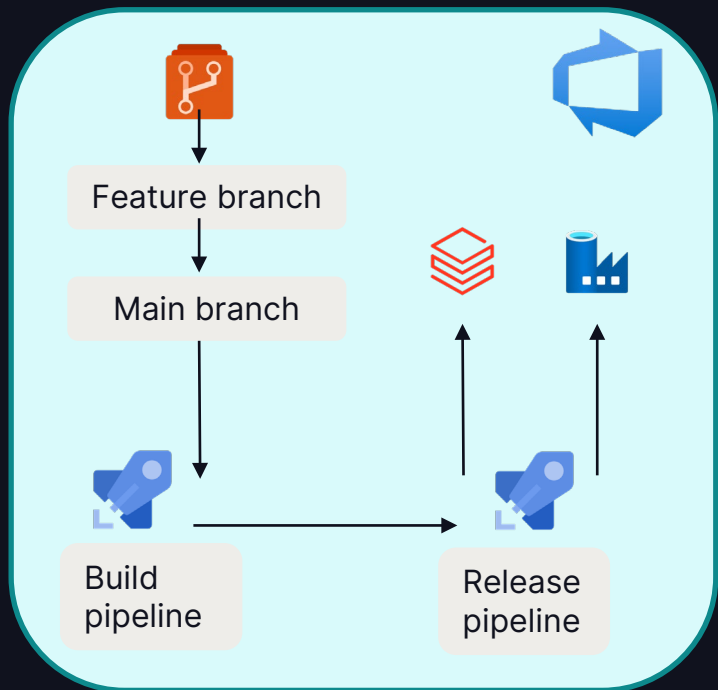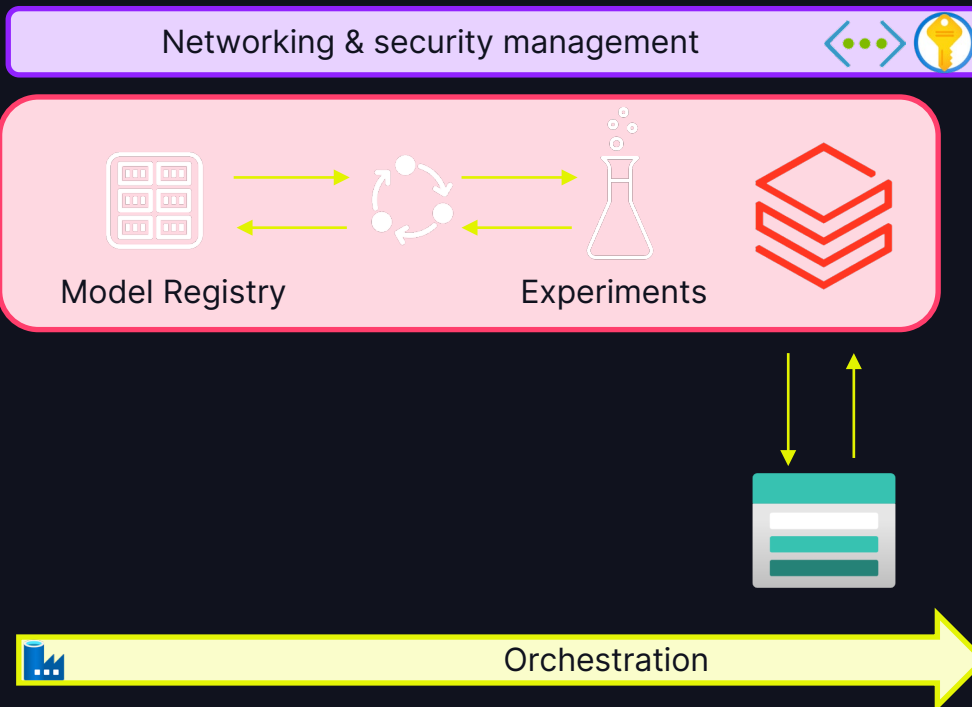5 Monitoring for customer confidence ✅

# Great! All requirements fulfilled!

We managed to build a plug -and-play ML solution!



DevOps

Azure

Feature branch

Main branch

Build pipeline

Release pipeline

Networking & security management

Model Registry

Experiments

Orchestration

©2024 Databricks Inc. — All rights reserved

# Let's have a chat!

**Julie Vanackere**
**Data scientist**

**Floriant Sturm**
**Co-founder**

https://www.linkedin.com/in/julie-vanackere/

https://www.linkedin.com/in/floriantsturm/

DATA'AI SUMMIT

31