





# Data Lakehouse: The best of both worlds





**Franco Patano** Lead Product Specialist Databricks

# Why build the Lakehouse ?

#### What we heard from our customers



1

Delta Lake enabled robust data management on the data lake.



Customers increasingly use SQL to directly query data lake data.



All the data is going to the lake. Only a portion gets into the data warehouse.

# "We do not want another data warehouse"

SQL Usage Growth on Databricks (2018-2020)

# **Delivering Analytics on the Lakehouse**

# **Our Vision**

 $\mathbf{O}$ 

- One source of truth for all your data
- World-class performance with data lake economics
- Data, analytics, and AI in one place



## **The Challenges Ahead**

"Data lakes have no structure they're just big data dumps - so it is impossible to turn it into a data warehouse."

"These two worlds are very different - one is based on files, the other on tables. It doesn't make any sense to blend them."

"It takes a very very long time to build a real data warehouse. People have been doing this for decades, **you can't do this overnight**." "It will be impossible to leverage an open format like Parquet and get world-class performance out of it."

"Can you even layout the data automatically as you ingest it? In Spark, you have to specify partitioning, etc... but not in a data warehouse, how will you handle this ?"

"Can you even get concurrency out of a platform like this (for like 1,000's of users)? How would you handle concurrency on a framework built for batch analytics?" "Can you even get existing BI and analytics tools to work with this open format and deliver data warehousing performance?"

"Data Scientists love the Databricks notebooks, but they are not suited for the majority of analysts"

"You can't get great performance out of open source tools"

"The data is not in this data mart, or the data warehouse"

# **Analyst Dilemma**

#### Where is the data we need?



**Business Partner** 

I need data on this new widget we put into production



Analyst

I don't see that data in the warehouse, let me check





#### Data Warehouse Team

That data was not warehoused, it's in the data lake

#### Data Lake Team

We have the data in the data lake, but I can't give you access. How about an Extract?

### **Databricks Lakehouse Platform**



#### Simple

Unify your data warehousing and Al use cases on a single platform

#### Open

Built on open source and open standards

#### **Multi-cloud**

One consistent data platform across clouds

## Databricks thrives within your modern data stack



# **Databricks SQL**

#### Data Warehousing on the Lakehouse

- Analytics on all your data with your tools of choice
- Simplified administration and fine-grained governance
- In-place, lightning-fast analytics on data lake data



### **Analyst Journey**

#### What can we do with the data?



**Business Partner** 



I need data on this new widget we put into production

Here you go, with Lakehouse, we can do analytics on all the data!



Great! Can you classify these actions for me real quick?

Uh?



# **Databricks AutoML**

A glass-box solution that empowers data teams without taking away control



# "Glass-Box" AutoML with a UI

igure AutoML experiment Preview Provide F	Feedback				
nents > Configure AutoML experiment					
Configure (2) Augment (3) Tr	rain —	luate			
AutoML Experiment Configuration					
Compute	✓ AutoML				
dais mlr 8 new	Configure	- 🕑 Train	-		
elect an existing cluster with a Databricks Runtime for ML 8.0-	AutoML Evaluation				
ML problem type	All runs have completed, and have been added to the table below. Click a specific run	to view details or review the data exploration notebook.	Customize		
Classification	Model with best val_f1_score	Generated Trial Notebook (Python)		o ? 🛔	
		▲ • dals_mir_8_new	<ul> <li> عند المحمد ا </li></ul>	C 1 * 1	
defeult veges lans	Register and deploy model Edit model	Random Forest training     6     7     * Characteristics	diction to evolution or	multiple examples for more thorough	
deraun.usage_logs	Showing 16 matching runs	Load Data results.	diccion to exptain, or sample i	matchipte examples for more chorough	
Target column	Grant Compare Delete Download CSV. ST	Yreprocessors     S     example = X_val     Val     Vumerical columns     9	sample(n=25)	Deploy	
isMining		One-hot encoding	P to explain feature importanc x: model.predict(pd.DataFrame	e on example Deploy	
Europiment nome	Columns @ Q metrics.rmse < 1 and params.rr	Feature standardiza 12 explainer = Ken 13 shap_values = e			
infining users loss 2021 05 05 22 22		Train classification mo 14 summary_plot(sh 15 except Exception as	Details Servin	ng	
	Start Time Run Name User Source	Inference 16 print(f"An unex	Status: Beady .	Stop Cluster: mlfrow-model-dais d	amo 🖗
Data directory	Ø 2021-08-05 1 logistic_r kase B Notebook: LogisticRegres	100% 25/25	•		
		in address online	Model Versions	Model Events Cluster Settings	
2	⊘ 2021-05-05 1 logistic_r kase      Notebook: LogisticRegres	account_id_online_feature_look	Model Versions		
	📄 📀 2021-05-05 1 logistic_r kase 🖹 Notebook: LogisticRegres 🋫	account_id_online_feature_los	Version 2	Ready     Model URL:	
	Ø 2021-05-05 1 logistic_r kase      Notebook: LogisticRegres	_ip_address_or	Production	https://dbc-60ef17e8-d99b.dev.databricks.com/model	/dais_demo/2/invocations
	2021-05-05 1 decision kase In Notebook: DecisionTree		Q	https://dbc-ouer1/eo-de9b.dev.databricks.com/model	dais_demovProduction/invocations
_	C 2021 OF OF 1 desition know D Notobeck DesiderTern		A	✓ Call The Model	
	A		~	Browser Curl Python	
			588	Request @	Response @
			≌	I.	[true,true,true,true,true]
			因	"notebookLanguage": "python",	
		Command took 3.08 minutes -	-m	"ip_address": 165225104128,	
		Cad 25	~	"account_id": 1587714725935792, "memberStartTime": 1614074788455.	
			24	Send Request Show Example	

. . .

# **Databricks AutoML Supported Features**

Time Series Forecasting

**Problem Types** Models / Tuning Tracking Features Evaluation Deployment ml*flow* 👫 Shap learn Batch Scoring Classification Numerical mlflow N N N  $( \mathcal{C} )$ ЛП XGBoost Regression Categorical Model Serving Metrics Parameters <u>جابی</u> جربی Η HYPEROPT Artifacts Models Timestamp **TensorFlow** 

Text

...

#### **Analyst Saves the Day**

Serving business solutions every day



**Business Partner** 



I need data on this new widget we put into production

Here you go, with Lakehouse, we can do analytics on all the data!





Great! Can you classify these actions for me real quick?

Yeah, with Databricks AutoML I can do that with a couple of clicks!



•••





# **Time to Science**

ThoughtSpot for Pharmaceutical Research





#### Gregory M. Troup PhD

Sr. Principal Scientist Merck Research Labs Pharmaceutical Sciences & Clinical Supplies Merck & Company Inc. Rahway NJ USA

# So, what is Merck Pharmaceutical Sciences doing?





#### **Drug Product Development/Clinical Plant Modernization**

#### State of the art product development facilities

- IT systems
- Operational model •
- Small scale, portable equipment



#### What does it feature?

- Designed for flexibility
- Modular Wall Systems/ POD suites •
- Product facing staff model, self service oriented
- New mindsets around data.

### **Current State of Scientific Data: Unstructured & Essentially Single Use**



Spreadsheet that the scientists use to use pool data for analysis and plotting

Slide deck that the scientist uses to share results with the team

Business doc used as reference in another business document There is a loss of context as our traverses the primary use arc

#### Where we are moving to

#### **Current State**

- Assembly of data sets for analysis is a manual process relies on , personal connections, emailing files, extensive file gymnastics
- Loss of context as our data traverses the "single use" arc.
- We have some desktop tools for shaping and analyzing data sets, but they are static
- No way to operationalize analytics, analysis is static

#### **Future State**

- Easy, intuitive searching, and assembling data sets for analysis.
- We have full data lineage to accompany our data collections, our analytics, and our models.
- In addition to our desktop tools, data science workbench and cloud-based data science tools for query, shaping, analytics, and modeling.
- Our analytics and models are now dynamic and easily deployed into routine use.



Lens User

Adaptability to manage the complexity of our field instruments



Flexibility & User Experience to meet the business needs



GMP Store of official records: native and cloud operable, fed with complete data

#### A new approach: Separate documentation of experimental exaction and analysis

Automate the collection of the complete record of the experiment, the who, what, how, where and when,

Make the user experience intuitive, think Panera Bread Kiosk intuitive

- Internal solution (a combination of technology that we already 1. had in house) to address the gap in the LAB execution space.
- 2. IoT platform that manages our equipment, instruments and sensors as IoT assets.
- 3. MES project to bring electronic batch records to development and clinical supplies.



v24.00

- Data between labs can be connected
- Sample information is seamlessly passed downstream
- Data Acquisition integration



# **Bringing Together Data Platforms: Accelerating Data Use**



Google Cloud

ThoughtSpot gives us focus and a target to drive our data model, and

#### Value Perspectives:

"I need real time aggregated data for the process I am currently running in one easy to set up view that I can see from anywhere in the room."

"I can easily set up any combination of room, equipment, data that I need to displayed for run time viewing on our devices."

" I have easy access to historical data, so makes it easier for me to plan and execute new work."

" I can instantly see the results generated by my product team members; we now can focus our energy on planning the next steps."







"I can quickly configure a data visualization that shows all of the relevant data for my equipment/process subject matter expert group." "Tech transfer risks are reduced because we can now quickly find development data related to the product we are working on and similar products." "As a data scientist, I now have easy access to training and validation data, and I have a fast path to deployment."

"Our ML model outputs are now easily accessible right along with our experimental data."





# A New Mindset is needed

#### **Everyone is an Analyst!**

- Using the best tools for the job, breaking the monolithic system trap
- Real time scientific inquiry, the end of slide decks as primary tool to sharing results.
- Ask and answer questions on the fly
- The second use, easy to find data sets for cross "fill in the blank" analysis (blank = program/product/material/equipment/site/...)
- Creating an environment where our scientists can learn faster and execute efficiently. More knowledge = less risk.



