# Dataiku DSS

## A Collaborative Data Science Platform



**PRÉPARER**
Chargez et préparez vos données

**ANALYSER**
Visualisez et partagez vos découvertes

**MODÉLISER**
Construisez vos modèles

**AUTOMATISER**
Ré-exécutez tous les jours…

**EXECUTER**
… ou en Temps Réel

**MONITORER**
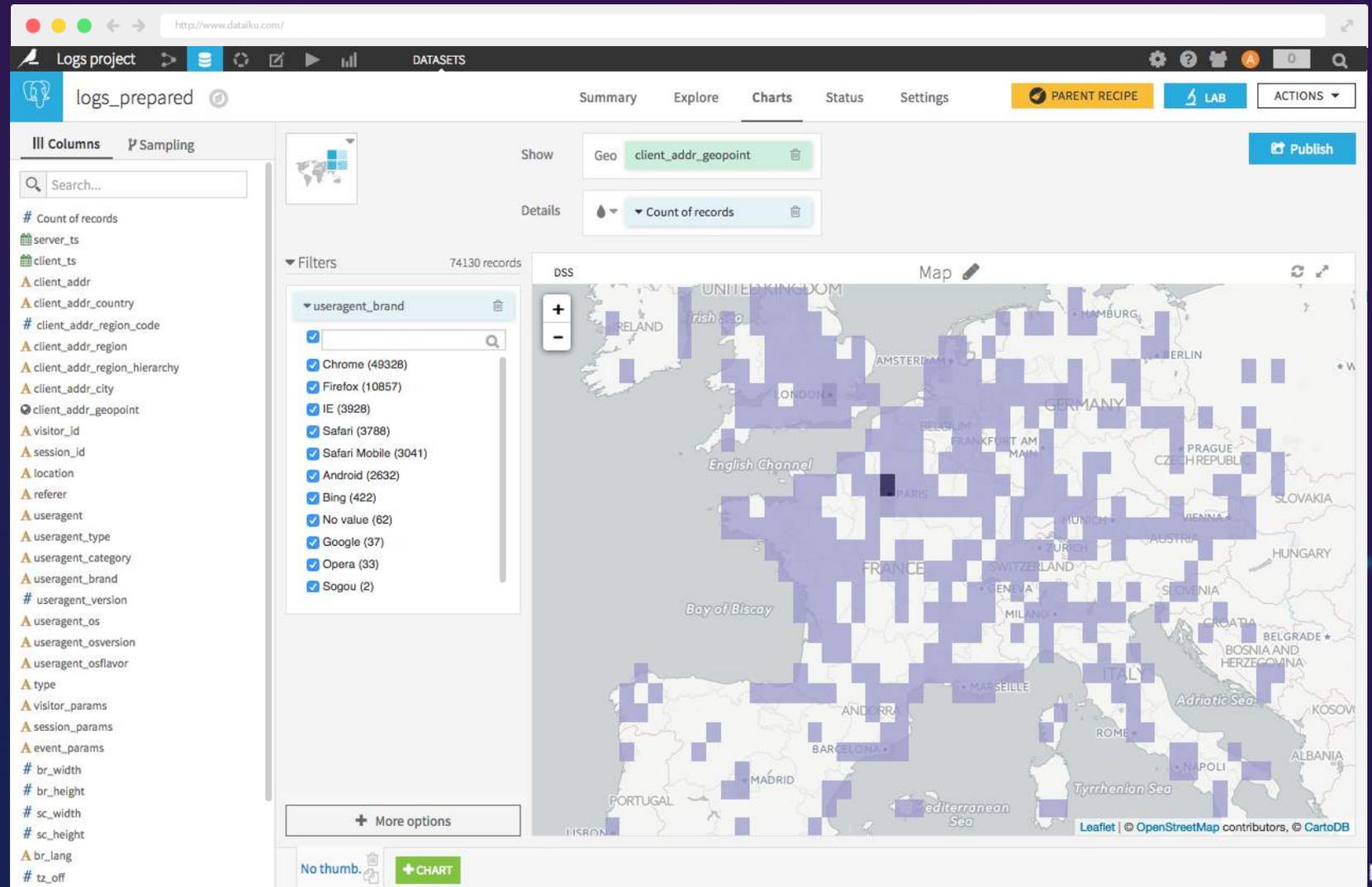Suivez votre production

OVH.com

# Logs webs

- Une ligne par action sur un site web
- Des informations typiques:
  - Date (timezone?)
  - IP
  - Page/ressource
  - Identifiant?
  - …
- Deux "catégories":
  - Logs serveurs
  - Logs Javascript

# Data cleaning (dates, IPs…)

Exemples:

- Parsing dates
- Géolocalisation IP
- Nettoyage des valeurs
- Extraire composants d'URL
- Enrichissement user-agent
- Filtrage
- Catégorisation

# Réduction de dimension (*group by*)

Agrégation au niveau de l'utilisateur:

- Cookie

- Identifiant

- IP

- Hash IP + user-agent

Recette visuelle SQL de *group by*.

Exécution *in-database* possible: MySQL, PostgreSQL, Vertica, Hive (Hadoop ou Spark)...

# Enrichissement en SQL

Objectif: construire des variables au niveau agrégé (visiteur)

Exemples:

- Récupération du premier *referer* avec une *Window Function*

- Enrichissement avec des données à partir du CRM avec un *Left Join*

# Prêt pour le Machine Learning ?!

Prérequis: des données agrégées, propres et enrichies

Machine Learning:

- Supervisé: pour prédire une variable cible

- Non-supervisé: pour trouver des groupes similaires (*clustering*)

# Prédiction

Exemple: prédire une conversion

Etapes:

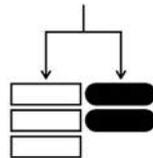- Choix des variables et types (numériques, catégorielles, …)

- Choix des algorithmes (linéaires, arbres…)

- Phase d'entrainement

- Etude des résultats

# Un peu de Python ou R?

- Pour aller plus loin que l'interface visuelle dans la préparation de données

- Pour intégrer des éléments extérieurs (API, …)

- Pour utiliser des librairies de Machine Learning

# Le « flow » pour l'industrialisation

- Exécuter l'ensemble de la préparation des données et la modélisation en quelques clics

- Possible de déployer sur une infrastructure de production dédiée (*batch & real time scoring*)

# Les cas d'usages classiques avec les logs

- Optimisation des conversions

- Segmentation des comportements

- Travail sur la recommandation

- Calculs de score de satisfaction client

- Détection comportements suspects

- …

OVH.com

# Nos clients (80+)

Web

Banques

Assurances

Industrie

Infrastructure

Santé

Médias

Jeux vidéo

# Essayez Dataiku DSS

www.dataiku.com/try
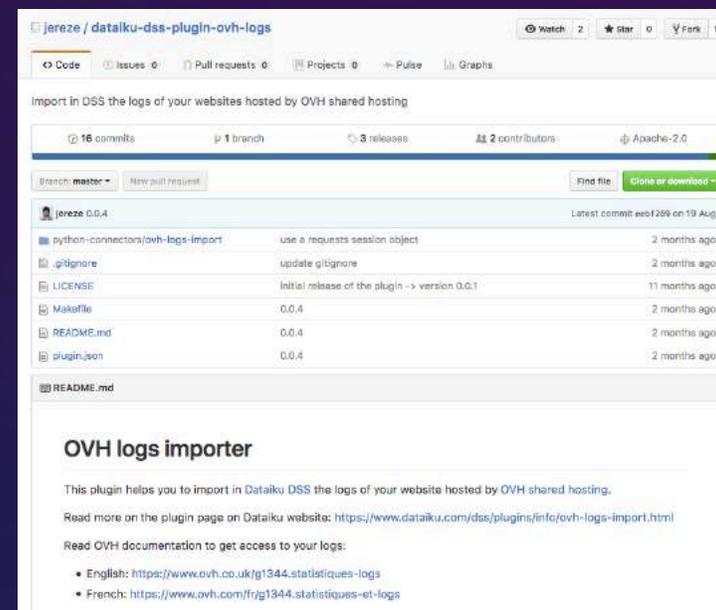
Free Edition

Enterprise Edition

Bonus:

- Un plugin pour récupérer les logs OVH mutu
  github.com/jereze/dataiku-dss-plugin-ovh-logs

- Un tracker web open-source qui scale
  github.com/dataiku/wt1

OVH.com