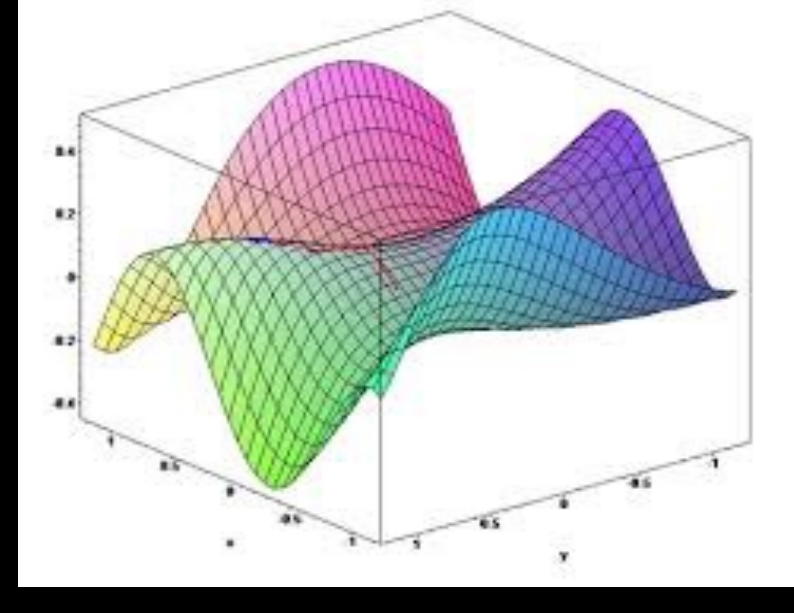# Improving the life of a data scientist
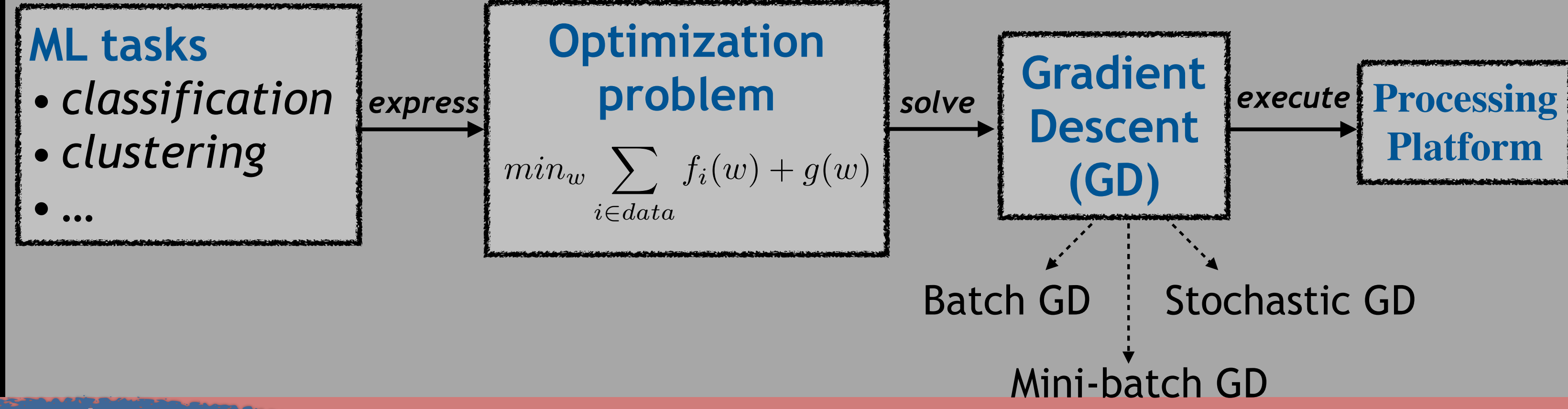
## Data scientist today
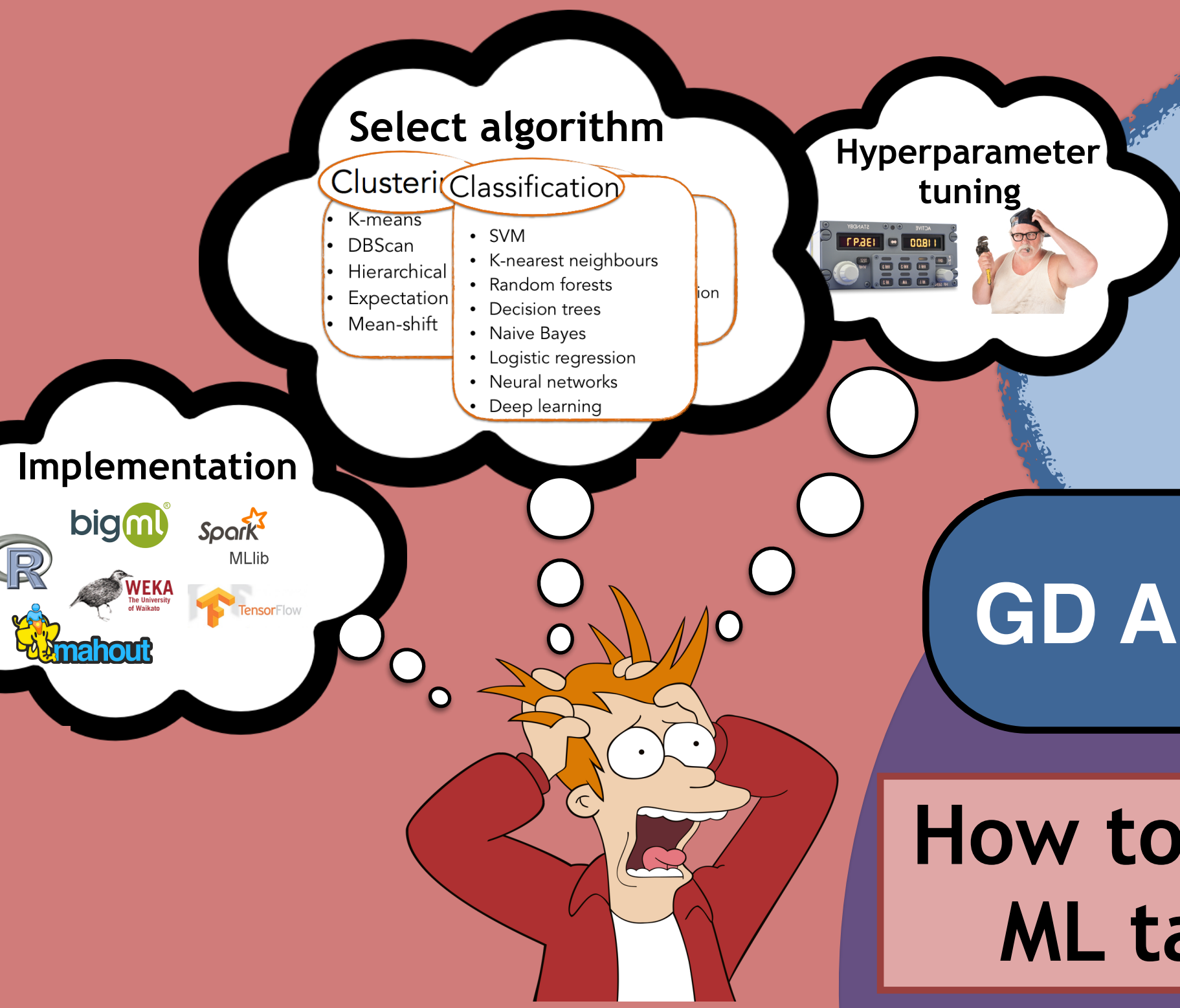
What people think he does   What he thinks he does   What he actually does

## Observation
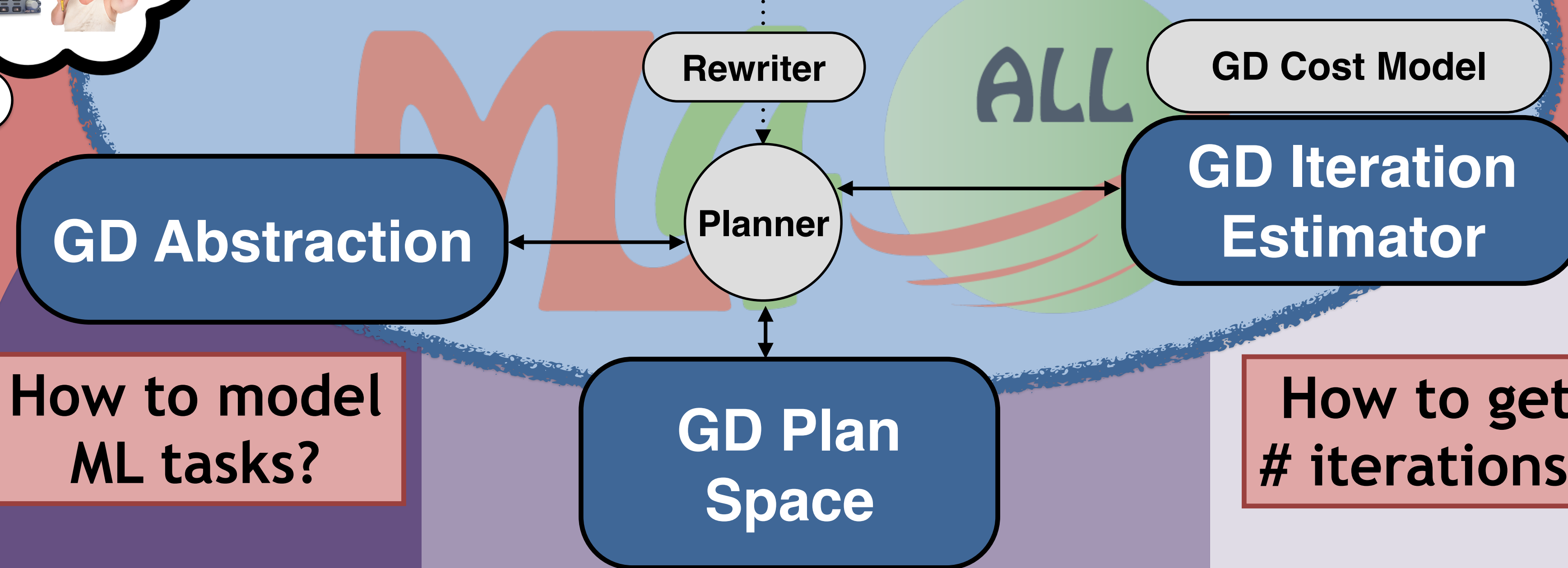
**ML tasks**
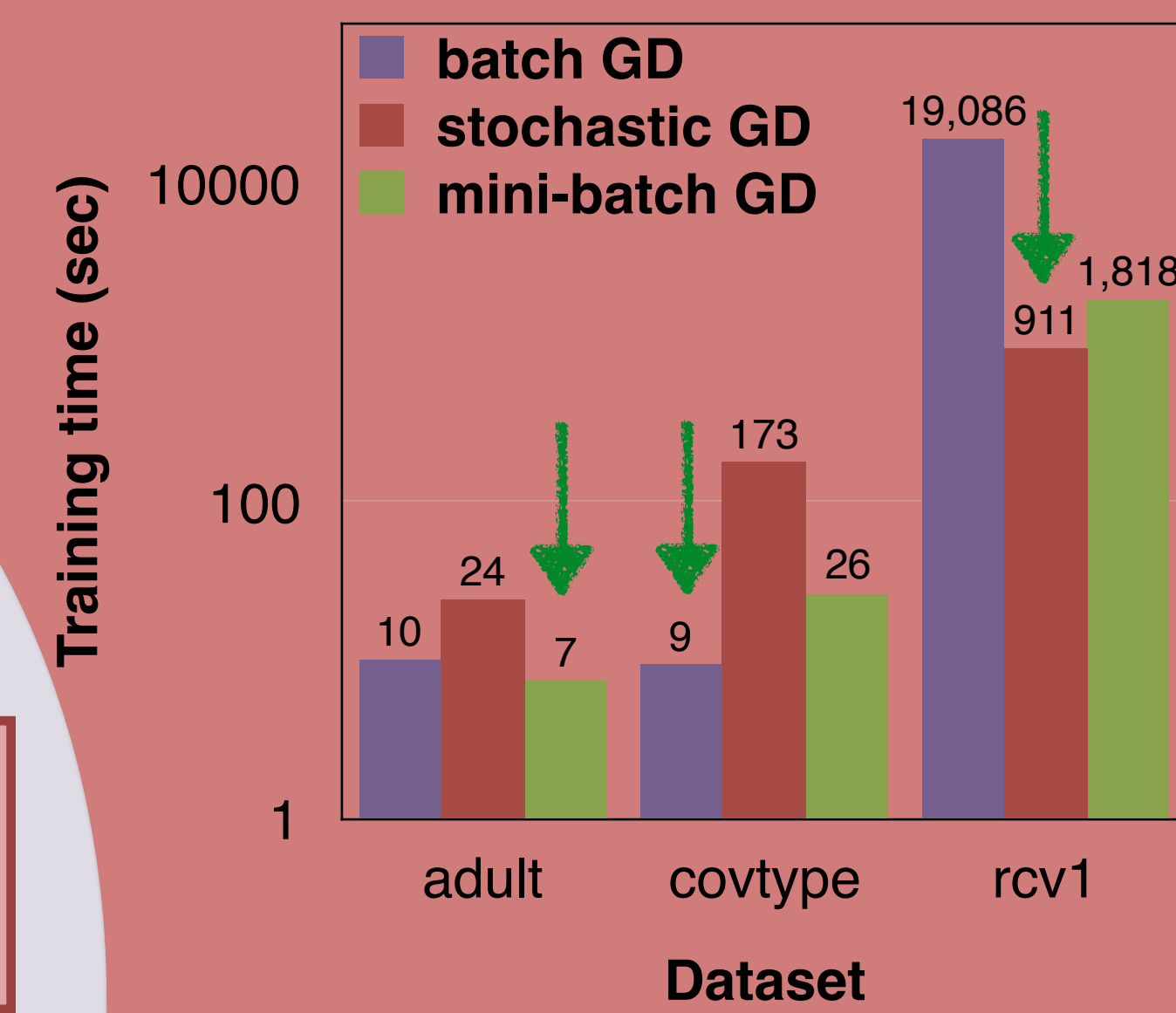- classification
- clustering
- ...

express → **Optimization problem** $min_w \sum_{i \in data} f_i(w) + g(w)$ → solve → **Gradient Descent (GD)** → execute → **Processing Platform**

Batch GD    Stochastic GD

Mini-batch GD

## General Problem

**Select algorithm**

Clustering
- K-means
- DBScan
- Hierarchical
- Expectation
- Mean-shift

Classification
- SVM
- K-nearest neighbours
- Random forests
- Decision trees
- Naive Bayes
- Logistic regression
- Neural networks
- Deep learning

**Hyperparameter tuning**

**Implementation**
bigml   Spark   MLlib   R   WEKA   mahout   TensorFlow

**GD Abstraction**

**How to model ML tasks?**

1. Preparation phase
2. Processing phase
3. Convergence phase

Data units
+1 2:0.1 4:0.4 10:0.3

## Results

1. Error sequence follows known distribution

2. Shape of error sequence on sample D' << D

Shape of error sequence over D

### Speculative approach

1. Take a sample D' < D
2. Run GD for a larger error
3. Fit the distribution

Time estimates

adult    covtype    rcv1

Legend: BGD-real, BGD-estim, MGD-real, MGD-estim, SGD-real, SGD-estim

Real □ Estimated

Legend: Min, Max, Plan execution, Speculation

adult   covtype   yearpred   rcv1   higgs   svm1   svm2   svm3

---

### Results table

| | GD time | GD iterations | SGD-eager-random time | SGD-eager-random iterations | SGD-eager-shuffle-partition time |
|---|---|---|---|---|---|
| adult 0.1 | 10.082 | 22 | 11.895 | 147.6666666667 | 2.7973333333333 |
| adult 0.01 | 25.9373333333333 | 227 | 54.0063333333333 | 1908.333333333 | 57.588 |
| adult 0.001 | 139.251 | 2586 | 507.582 | 26044.6666666667 | 721.641333333333 |
| covtype 0.1 | 17.4193333333333 | 26 | 41.3933333333333 | 3085.6666666667 | 2.7056666666667 |
| covtype 0.01 | 38.272 | 134 | 172.2773333333333 | 9481.3333333333 | 55.765666666667 |
| covtype 0.001 | 252.025 | 1856 | 1437.8356666666 | 91121.3333333333 | 404.688333333333 |
| yearpred 0.1 | 15.814 | 44 | 12.3663333333333 | 60 | 12.3663333333333 |
| rcv1 0.1 | 275.9826666666667 | 799 | 61.2373333333333 | 229.6666666667 | 52.495 |
| rcv1 0.01 | 19086.1115 | 21942 | 994.4823333333 | 30743.6666666667 | 910.623 |
| SVM synthetic | 220.4846666666667 | 145 | 104.8505666666 | 245.3333333333 | 22.4806666666667 |
|  | 0.073 | 1901.2606666667 |  |  |  |
|  | 0.757 | 930.8576666667 |  |  |  |
|  | 0.063 |  |  |  |  |

---

batch GD
stochastic GD
mini-batch GD

19,086    1,818    911    173    26

adult    covtype    rcv1

Dataset

## A Cost-based Optimizer for Gradient Descent Optimization

Zoi Kaoudi    Jorge Quiané-Ruiz    Saravanan Thirumuruganathan    Sanjay Chawla    Divy Agrawal

UCSB

QCRI
جامعة حمد بن خليفة
Qatar Computing Research Institute