# Type-Driven Automated Learning with LALE

Martin Hirzel, Kiran Kate, Avi Shinnar,
Subhrajit Roy, Pari Ram, and
Guillaume Baudart
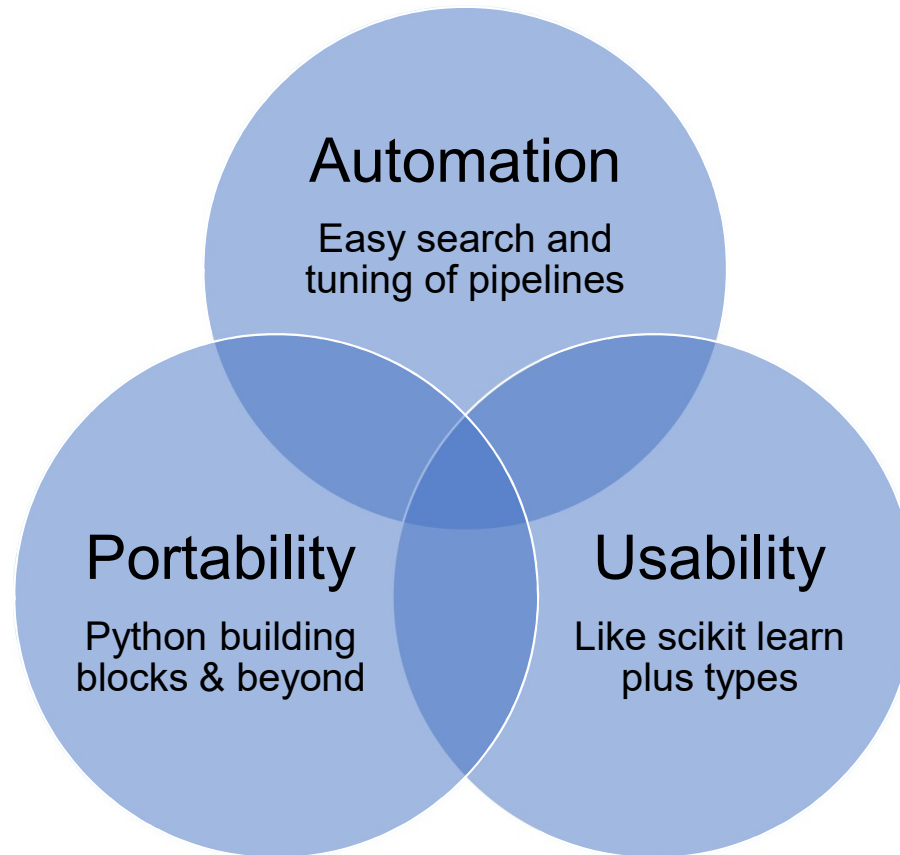
Wednesday 29 May 2019

Global Data Scientist Profession
1/2 Day Conference

# Value Proposition

Augment, but don't replace, the data scientist.



Automation
Easy search and tuning of pipelines

Portability
Python building blocks & beyond

Usability
Like scikit learn plus types

# Manual ML with Sklearn

Prior work: scikit learn, popular machine learning package

```
1   pca_lr = make_pipeline(PCA(svd_solver='full', n_components=0.3),
2                          LR(solver='liblinear', penalty='l1'))
```
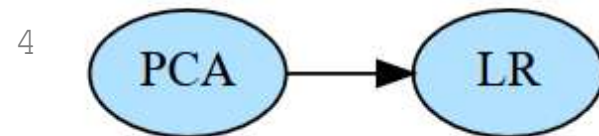
```
3   pca_lr.fit(train_X, train_y)
4   predicted = pca_lr.predict(test_X)
5   print(f'accuracy  {accuracy_score(test_y, predicted):.1%}')
```

```
6   accuracy  70.2%
```
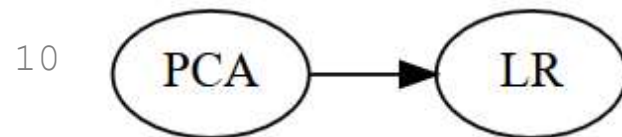
# Manual ML with LALE

Our work: Language for Automated Learning Exploration

```
1  pca_lr = PCA(PCA.svd_solver.full, n_components=0.3) \
2         >> LR(LR.solver.liblinear, LR.penalty.l1)
3  to_graphviz(pca_lr)
```

4


```
5  trained = pca_lr.fit(train_X, train_y)
6  predicted = trained.predict(test_X)
7  print(f'accuracy  {accuracy_score(test_y, predicted):.1%}')
8  to_graphviz(trained)
```

9  accuracy  70.2%

10

# Automated ML with LALE

Combined algorithm selection and hyperparameter tuning

```
1  planned = PCA >> (J48 | LR)
2  to_graphviz(planned)
```

3



```
4  hyperopt_classifier = HyperoptClassifier(planned, max_evals=5)
5  best_found = hyperopt_classifier.fit(train_X, train_y)
6  predicted = best_found.predict(test_X)
7  print(f'accuracy  {accuracy_score(test_y, predicted):.1%}')
8  to_graphviz(best_found)
```
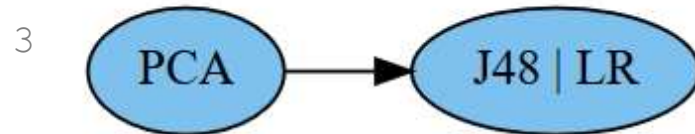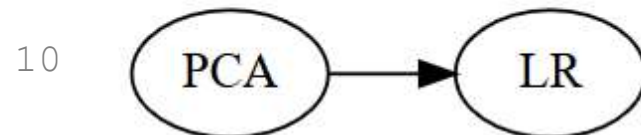
9
```
accuracy  96.4%
```

10

# Constraints in Manual ML

Conditional hyperparameters

```
1   pca_lr = make_pipeline(PCA(svd_solver='full', n_components=0.3),
2                           LR(solver='sag', penalty='l1'))
```

```
3   pca_lr.fit(train_X, train_y)
```

```
---------------------------------------------------------------------
ValueError                              Traceback (most recent call last)
<ipython-input-7-de82d92d1962> in <module>
----> 1 pca_lr.fit(train_X, train_y)

~/python3.7venv/lib/python3.7/site-packages/sklearn/pipeline.py in fit(self, X, y, **fit_params)
    265         Xt, fit_params = self._fit(X, y, **fit_params)
    266         if self._final_estimator is not None:
--> 267             self._final_estimator.fit(Xt, y, **fit_params)
    268         return self
    269

~/python3.7venv/lib/python3.7/site-packages/sklearn/linear_model/logistic.py in fit(self, X, y, sample_weight)
   1275                         "positive; got (tol=%r)" % self.tol)
   1276
-> 1277         solver = _check_solver(self.solver, self.penalty, self.dual)
   1278
   1279         if solver in ['newton-cg']:

~/python3.7venv/lib/python3.7/site-packages/sklearn/linear_model/logistic.py in _check_solver(solver, penalty, dual)
    445     if solver not in ['liblinear', 'saga'] and penalty != 'l2':
    446         raise ValueError("Solver %s supports only l2 penalties, "
--> 447                          "got %s penalty." % (solver, penalty))
    448     if solver != 'liblinear' and dual:
    449         raise ValueError("Solver %s supports only "
```

28  `ValueError: Solver sag supports only l2 penalties, got l1 penalty.`

# Constraints in AutoML

**Problem:** Some automated iterations raise exceptions

**Solution 1:** Unconstrained search space
- $\{S:[linear,sag,lbfgs], P: [l1,l2]\}$
- Catch exception
- Return made-up loss `np.float.max`

**Solution 2:** Constrained search space
- $\{S:[linear,sag,lbfgs], P: [l1,l2]\}$ **and** (**if** $S: [sag,lbfgs]$ **then** $P: [l2]$)
- No exceptions
- No made-up loss

# GridSearchCV Search Space

AutoML included with Sklearn



$$
\begin{array}{llll}
\text{dict}\{N: [0.21, 0.65, 0.84], & D: [J48], & R: [false], & C: [0.07, 0.30, 0.89]\} \\
\vee \text{dict}\{N: [0.21, 0.65, 0.84], & D: [J48], & R: [true, false], & C: [0.25]\} \\
\vee \text{dict}\{N: [mle], & D: [J48], & R: [false], & C: [0.07, 0.30, 0.89]\} \\
\vee \text{dict}\{N: [mle], & D: [J48], & R: [true, false], & C: [0.25]\} \\
\vee \text{dict}\{N: [0.21, 0.65, 0.84], & D: [LR], & S: [linear], & P: [l1, l2]\} \\
\vee \text{dict}\{N: [0.21, 0.65, 0.84], & D: [LR], & S: [linear, sag, lbfgs], & P: [l2]\} \\
\vee \text{dict}\{N: [mle], & D: [LR], & S: [linear], & P: [l1, l2]\} \\
\vee \text{dict}\{N: [mle], & D: [LR], & S: [linear, sag, lbfgs], & P: [l2]\}
\end{array}
$$

8

# SMAC Search Space

Sequential Model-based Algorithm Configuration



$$
\begin{aligned}
& \text{dict}\{N\colon (0..1),\ D\colon [J48],\ R\colon [false], && C\colon (0..1)\ \} \\
\vee\ & \text{dict}\{N\colon (0..1),\ D\colon [J48],\ R\colon [true, false], && C\colon [0.25]\ \} \\
\vee\ & \text{dict}\{N\colon [mle],\ \ D\colon [J48],\ R\colon [false], && C\colon (0..1)\ \} \\
\vee\ & \text{dict}\{N\colon [mle],\ \ D\colon [J48],\ R\colon [true, false], && C\colon [0.25]\ \} \\
\vee\ & \text{dict}\{N\colon (0..1),\ D\colon [LR],\ \ S\colon [linear], && P\colon [l1, l2]\} \\
\vee\ & \text{dict}\{N\colon (0..1),\ D\colon [LR],\ \ S\colon [linear, sag, lbfgs],\ P\colon [l2]\ \ \ \} \\
\vee\ & \text{dict}\{N\colon [mle],\ \ D\colon [LR],\ \ S\colon [linear], && P\colon [l1, l2]\} \\
\vee\ & \text{dict}\{N\colon [mle],\ \ D\colon [LR],\ \ S\colon [linear, sag, lbfgs],\ P\colon [l2]\ \ \ \}
\end{aligned}
$$

# Hyperopt Search Space

Supports parallel search



$$\text{dict}\left\{\begin{array}{l} 0 : \text{dict}\{N\!:\!(0..1)\} \vee \text{dict}\{N\!:\![mle]\} \\ 1 : \left(\begin{array}{c}\left(\begin{array}{l}\text{dict}\{D\!:\![J48], R\!:\![false], \quad C\!:\!(0..1)\} \\ \vee\, \text{dict}\{D\!:\![J48], R\!:\![true,false], C\!:\![0.25]\}\end{array}\right) \\ \vee\left(\begin{array}{l}\text{dict}\{D\!:\![LR], S\!:\![linear], \qquad\qquad P\!:\![l1,l2]\} \\ \vee\, \text{dict}\{D\!:\![LR], S\!:\![linear, sag, lbfgs], P\!:\![l2]\quad\}\end{array}\right)\end{array}\right)\end{array}\right\}$$

# Types as Search Spaces

Lᴀʟᴇ auto-generates search spaces for AutoML tools

**Planned pipeline**

PCA → J48 | LR

**Operator schemas**

*PCA*: …
*J48*: …
*LR*: …

**Compiler**

**GridSearchCV**

$$
\begin{array}{l}
\text{dict}\{N: [0.21, 0.65, 0.84], D: [J48], R: [false], \quad C: [0.07, 0.30, 0.89]\} \\
\lor \text{dict}\{N: [0.21, 0.65, 0.84], D: [J48], R: [true, false], \quad C: [0.25]\} \\
\lor \text{dict}\{N: [mle], \qquad D: [J48], R: [false], \quad C: [0.07, 0.30, 0.89]\} \\
\lor \text{dict}\{N: [mle], \qquad D: [J48], R: [true, false], \quad C: [0.25]\} \\
\lor \text{dict}\{N: [0.21, 0.65, 0.84], D: [LR], S: [linear], \qquad P: [l1, l2]\} \\
\lor \text{dict}\{N: [0.21, 0.65, 0.84], D: [LR], S: [linear, sag, lbfgs], P: [l2]\} \\
\lor \text{dict}\{N: [mle], \qquad D: [LR], S: [linear], \qquad P: [l1, l2]\} \\
\lor \text{dict}\{N: [mle], \qquad D: [LR], S: [linear, sag, lbfgs], P: [l2]\}
\end{array}
$$

**SMAC**

$$
\begin{array}{l}
\text{dict}\{N: (0..1), D: [J48], R: [false], \qquad C: (0..1)\} \\
\lor \text{dict}\{N: (0..1), D: [J48], R: [true, false], \qquad C: [0.25]\} \\
\lor \text{dict}\{N: [mle], D: [J48], R: [false], \qquad C: (0..1)\} \\
\lor \text{dict}\{N: [mle], D: [J48], R: [true, false], \qquad C: [0.25]\} \\
\lor \text{dict}\{N: (0..1), D: [LR], S: [linear], \qquad P: [l1, l2]\} \\
\lor \text{dict}\{N: (0..1), D: [LR], S: [linear, sag, lbfgs], P: [l2]\} \\
\lor \text{dict}\{N: [mle], D: [LR], S: [linear], \qquad P: [l1, l2]\} \\
\lor \text{dict}\{N: [mle], D: [LR], S: [linear, sag, lbfgs], P: [l2]\}
\end{array}
$$

**hyperopt**
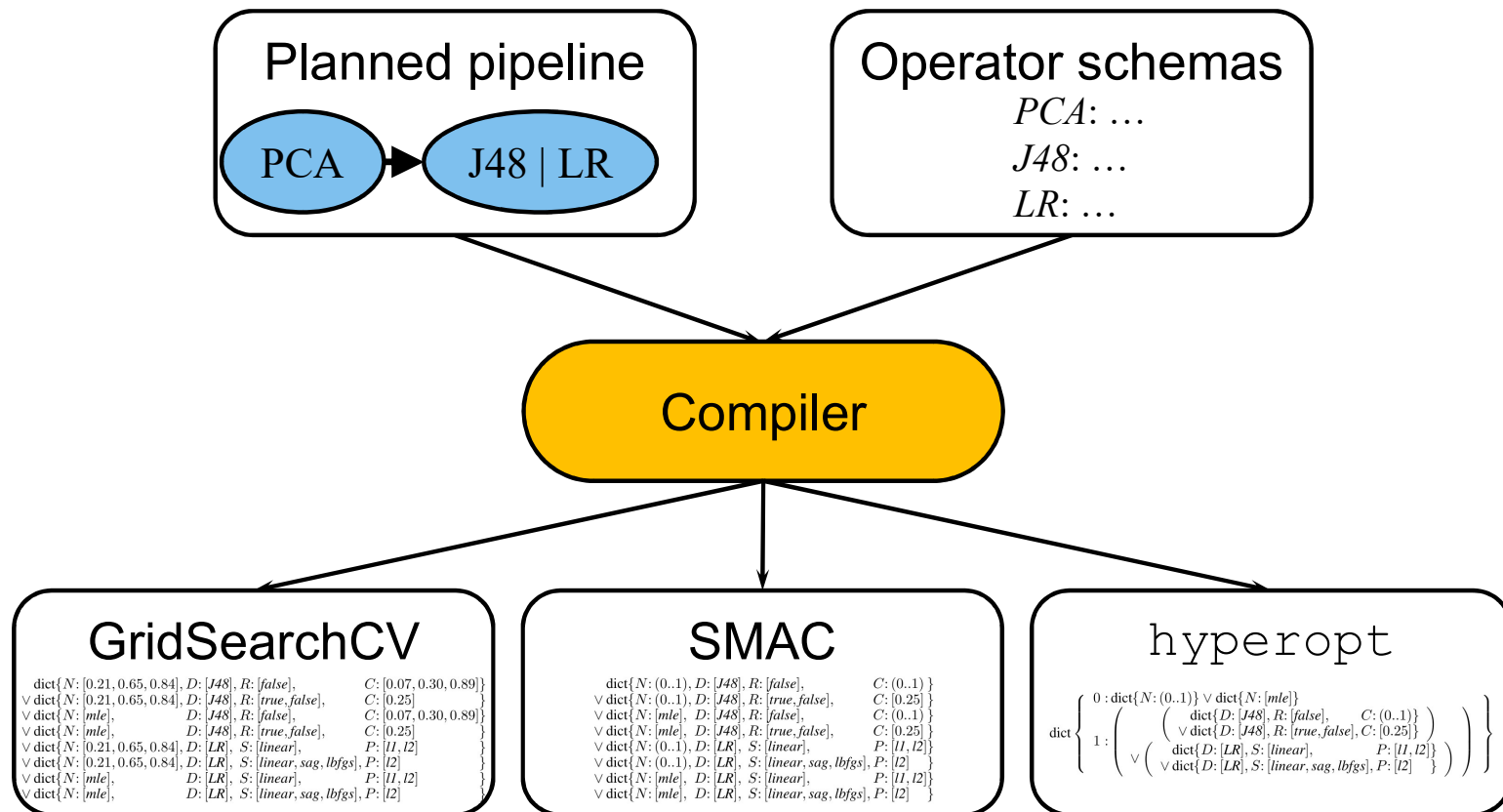
$$
\text{dict}\left\{
\begin{array}{l}
0 : \text{dict}\{N: (0..1)\} \lor \text{dict}\{N: [mle]\} \\
1 : \left(
\begin{array}{l}
\left(\begin{array}{l} \text{dict}\{D: [J48], R: [false], \quad C: (0..1)\} \\ \lor \text{dict}\{D: [J48], R: [true, false], C: [0.25]\}\end{array}\right) \\
\lor \left(\begin{array}{l} \text{dict}\{D: [LR], S: [linear], \qquad P: [l1, l2]\} \\ \lor \text{dict}\{D: [LR], S: [linear, sag, lbfgs], P: [l2]\}\end{array}\right)
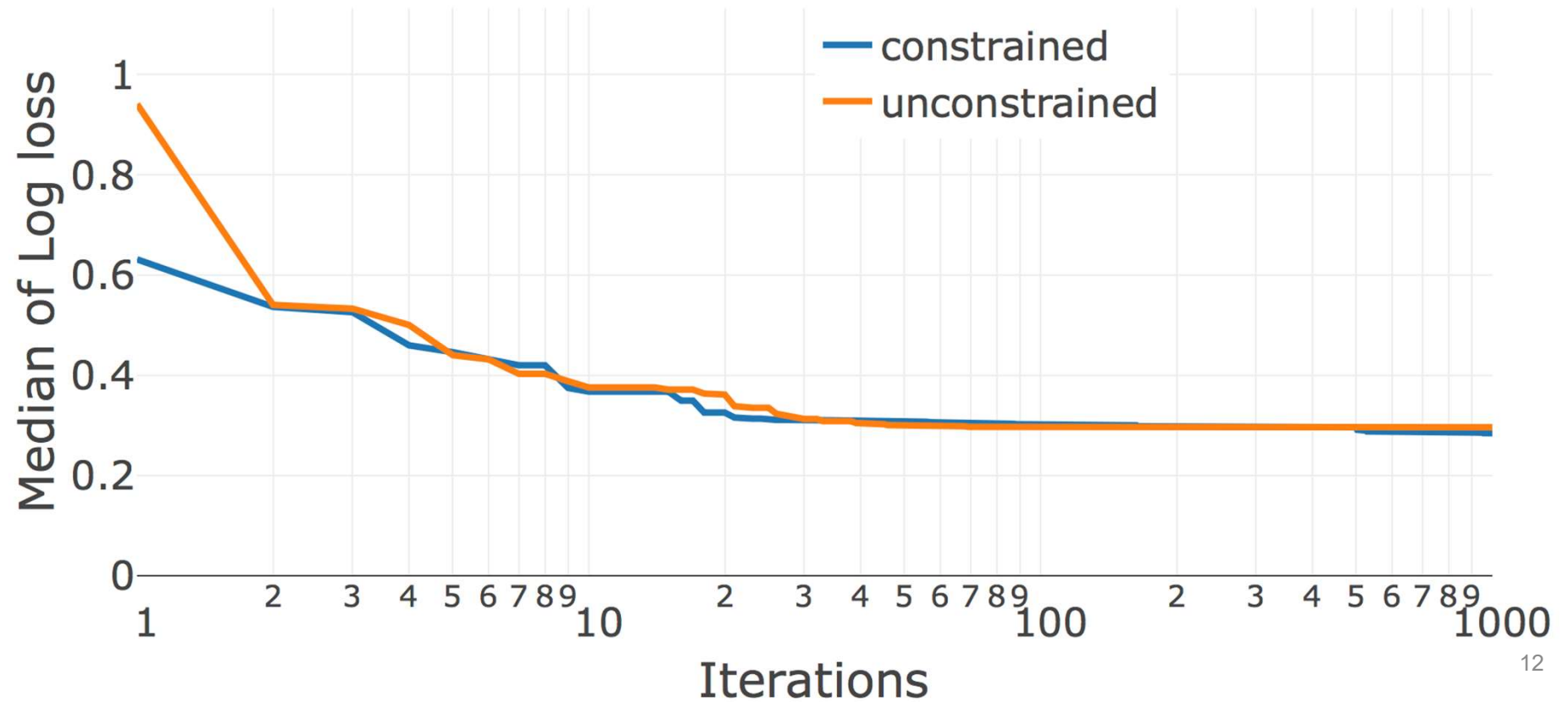\end{array}
\right)
\end{array}
\right\}
$$

# Search Convergence (1/3)

LR | KNN

Car dataset

`hyperopt`

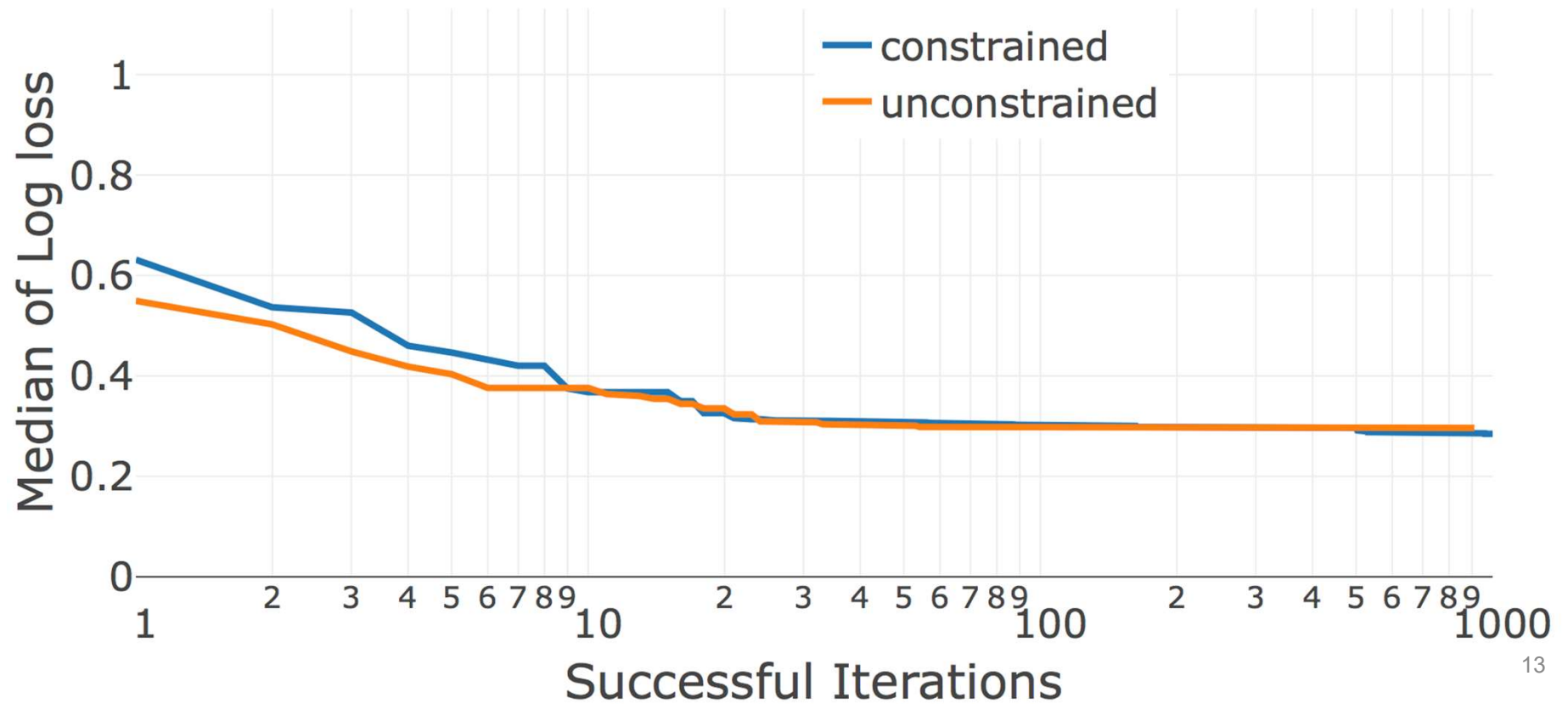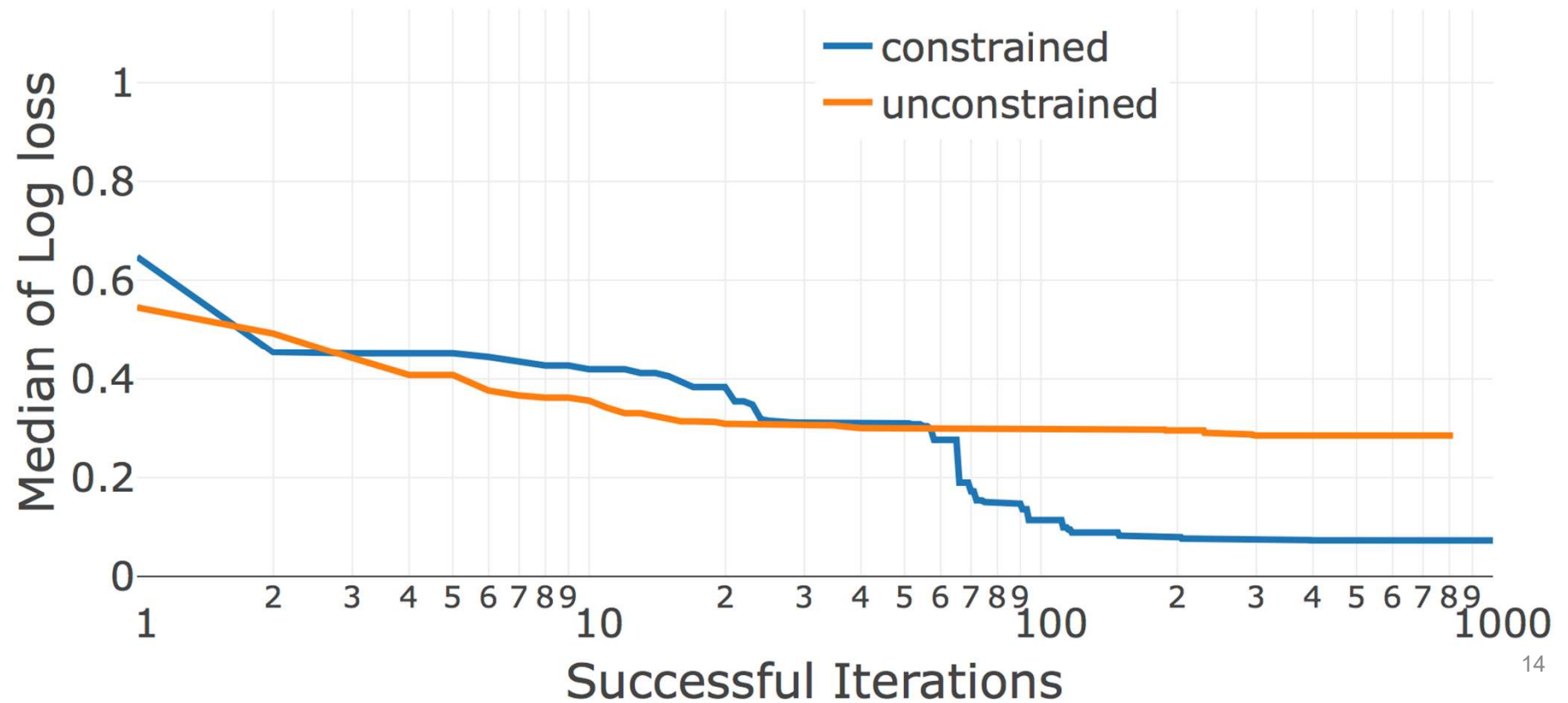# Search Convergence (2/3)

LR | KNN

Car dataset

`hyperopt`

# Search Convergence (3/3)



J48 | LR | KNN

Car dataset

`hyperopt`

# Portability

| Modality | Dataset | Pipeline (**bold**: best found choice) |
|---|---|---|
| Text | Movie reviews (sentiment analysis) | (**BERT** \| TFIDF) >> (**LR** \| MLP \| KNN \| SVC \| PAC) |
| Table | Car (structured with categorical features) | **J48** \| ArulesCBA \| LR \| KNN |
| Images | CIFAR-10 (image classification) | **ResNet50** |
| Time-series | Epilepsy (seizure classification) | **WindowTransformer** >> (**KNN** \| XGBoost \| LR) >> **Voting** |

# Status

https://github.ibm.com/aimodels/lale