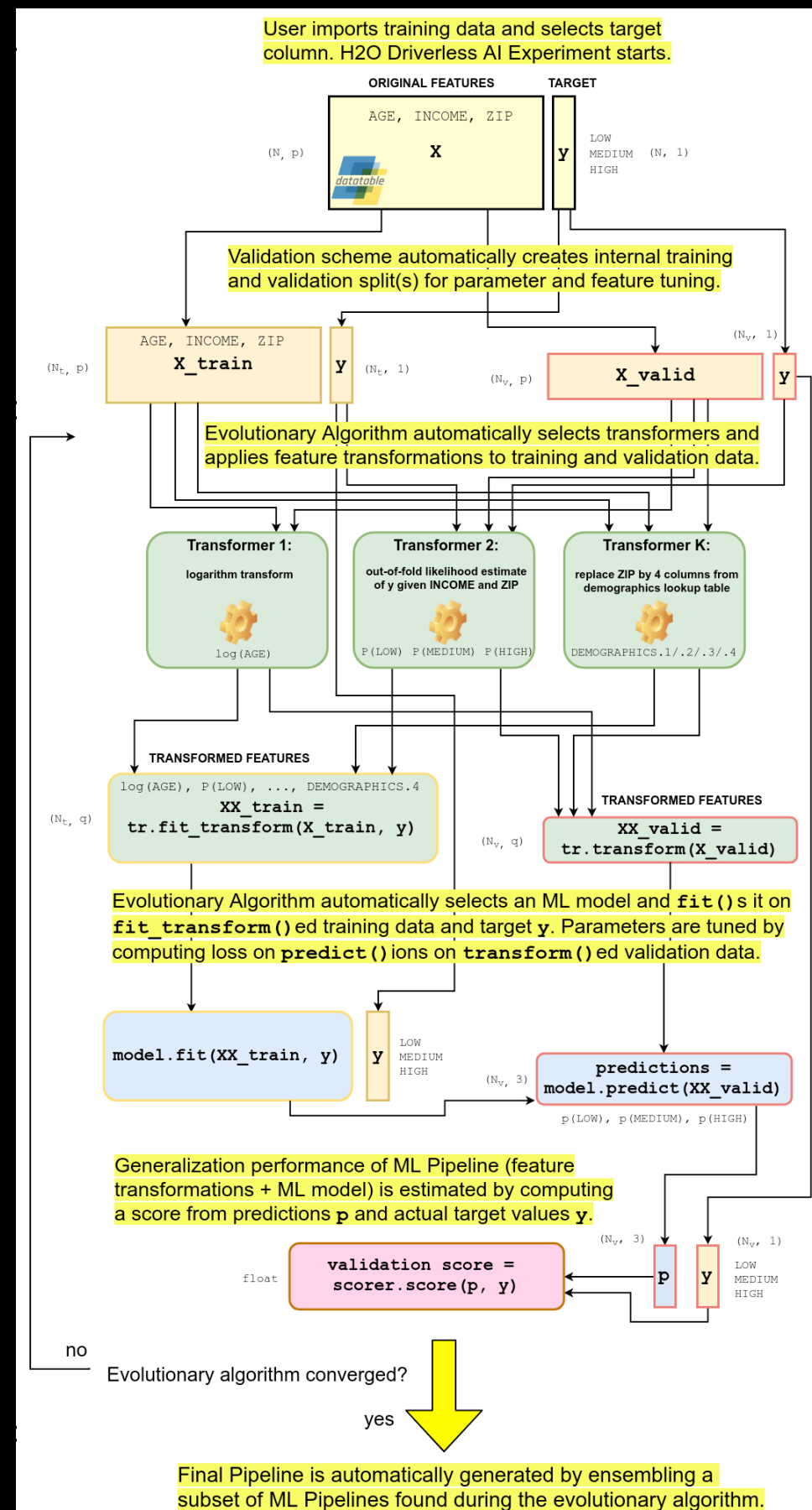


# DAI Expert settings

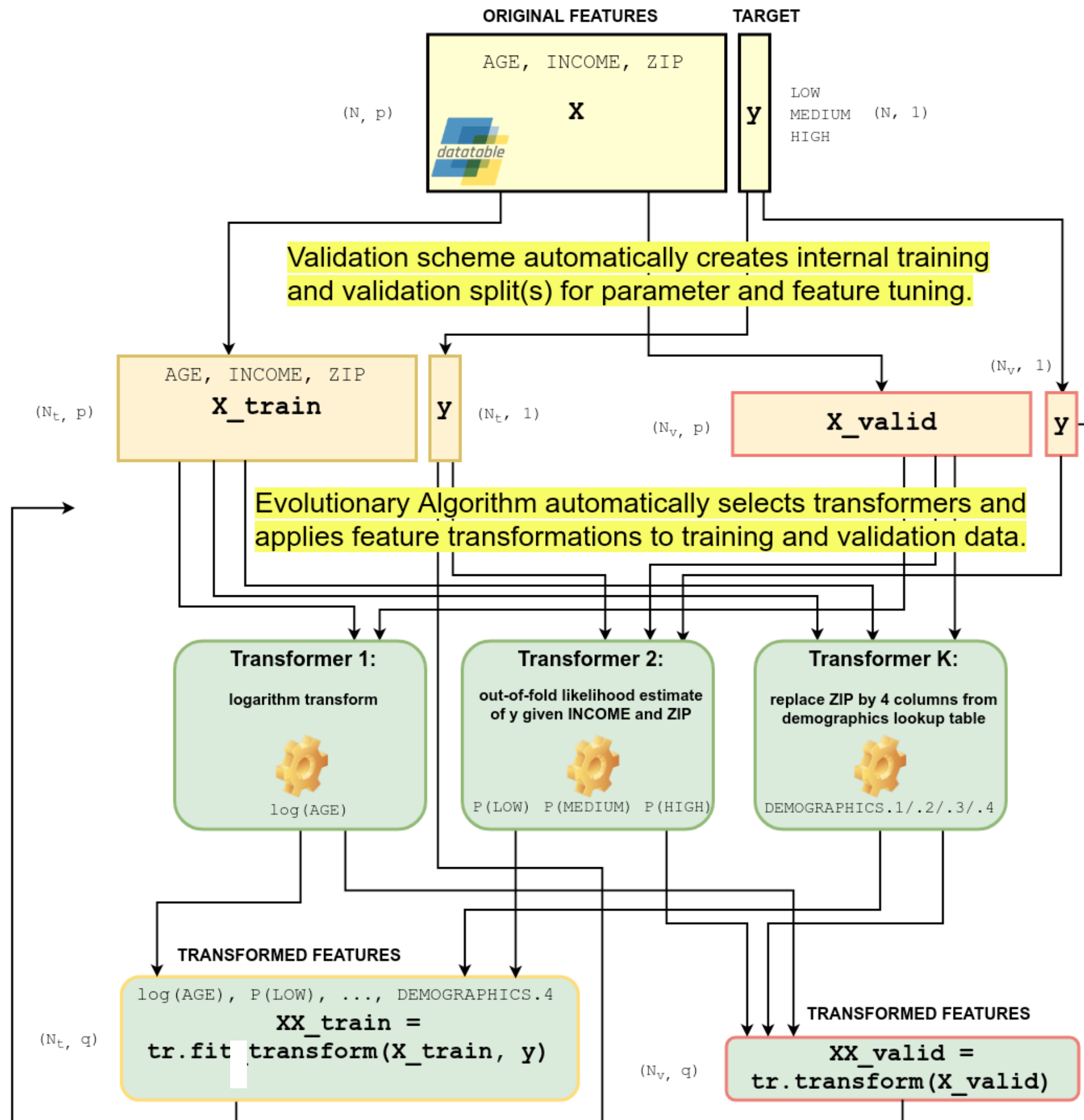
## Transformers

Transformer recipes = feature engineering

# Feature engineering in DAI



User imports training data and selects target column. H2O Driverless AI Experiment starts.



Evolutionary Algorithm automatically selects an ML model and `fit()`s it on `fit_transform()`ed training data and target `y`. Parameters are tuned by computing loss on `predict()`ions on `transform()`ed validation data.

`model.fit(XX_train, y)`

`y`  
LOW  
MEDIUM  
HIGH

$(N_v, 3)$

`predictions = model.predict(XX_valid)`

`p(LOW), p(MEDIUM), p(HIGH)`

Generalization performance of ML Pipeline (feature transformations + ML model) is estimated by computing a score from predictions `p` and actual target values `y`.

float

`validation score = scorer.score(p, y)`

$(N_v, 3)$

`p`

$(N_v, 1)$

`y`

LOW  
MEDIUM  
HIGH

no

Evolutionary algorithm converged?

yes

Final Pipeline is automatically generated by ensembling a subset of ML Pipelines found during the evolutionary algorithm.

Demo : Expert settings

# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

1. Numeric (int, real, binary)
2. Time series
3. Categorical (string)
4. Text (string)
5. Time (date, time)

# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

## 1. Numeric (int, real, binary)

## 2. Time series

## 3. Categorical (string)

## 4. Text (string)

## 5. Time (date, time)

- Original
- Interactions
- ClusterDist
- NumCatTE
- NumToCatTE
- ClusterTE
- NumToCatWoE
- NumToCatWoEMonotonic
- TruncSVDNum



# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

## 1. Numeric (int, real, binary)

2. Time series

3. Categorical (string)

4. Text (string)

5. Time (date, time)

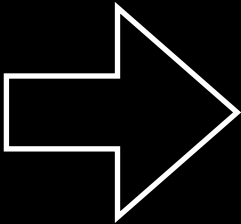
- Original
- Interactions
- ClusterDist
- NumCatTE
- NumToCatTE
- ClusterTE
- NumToCatWoE
- NumToCatWoEMonotonic
- TruncSVDNum

• Interactions

$f : +, -, \times, /$



y	x1	x2	x3
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>



y	x1	x2	x3	<i>f</i> (2,3)
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>

Name of variable *f*(2,3) : Interaction\_x2#subtract#x3

# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

## 1. Numeric (int, real, binary)

## 2. Time series

## 3. Categorical (string)

## 4. Text (string)

## 5. Time (date, time)

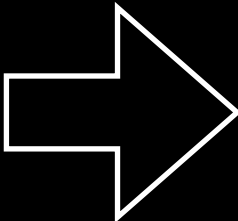
- Original
- Interactions
- **ClusterDist**
- NumCatTE
- NumToCatTE
- ClusterTE
- NumToCatWoE
- NumToCatWoEMonotonic
- TruncSVDNum

- ClusterDist

$f$  : 2,3 distance from cluster n (total k clusters, k-Means)

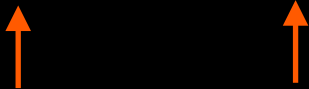


y	x1	x2	x3
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>



y	x1	x2	x3	f(2,3)
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>

Name of variable  $f(2,3)$  : ClusterDist\_k\_x2\_x3\_n



# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

## 1. Numeric (int, real, binary)

## 2. Time series

## 3. Categorical (string)

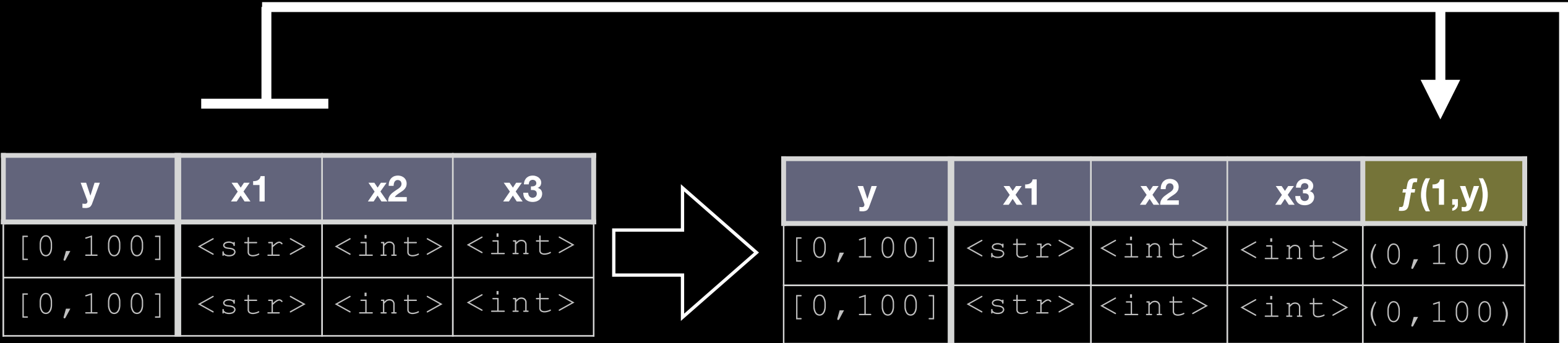
## 4. Text (string)

## 5. Time (date, time)

- Original
- Interactions
- ClusterDist
- NumCatTE
- NumToCatTE
- ClusterTE
- NumToCatWoE
- NumToCatWoEMonotonic
- TruncSVDNum

- NumCatTE

Categorical variable **x1** holds m categories



Training dataset

$f$  :Mean of out-of-sample target for category n of variable **x1**

y	x1	x2	x3
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>

Validation dataset (k-folds)

Name of variable  $f(1,y)$  : CV\_TE\_x1

## Target Encoder Transformers in a nutshell

Encoder name	Variables in training dataset...	Response in validation dataset...	Name of output variable
NumCatTE	...are identified as categorical	...is averaged for each category	CV_TE_<str>
NumToCatTE	...are binned to become categorical	...is averaged for each category	CV_TE_<num>
ClusterTE	...are clustered	...is averaged for each cluster	CV_TE_<num>_<num>_...

# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

## 1. Numeric (int, real, binary)

## 2. Time series

## 3. Categorical (string)

## 4. Text (string)

## 5. Time (date, time)

- Original
- Interactions
- ClusterDist
- NumCatTE
- NumToCatTE
- ClusterTE
- NumToCatWoE
- NumToCatWoEMonotonic
- TruncSVDNum



## Weight of Evidence (WoE) in a nutshell

High cardinality categorical variables mapped to binary outcomes

Category	Number of observations	Distribution of outcome Y	Distribution of outcome N	WoE
A	$Y_A + N_A$	$D_A^Y = Y_A / Y$	$D_A^N = N_A / N$	$\ln(D_A^Y / D_A^N)$
B	$Y_B + N_B$	$D_B^Y = Y_B / Y$	$D_B^N = N_B / N$	$\ln(D_B^Y / D_B^N)$
C	$Y_C + N_C$	$D_C^Y = Y_C / Y$	$D_C^N = N_C / N$	$\ln(D_C^Y / D_C^N)$
D	$Y_D + N_D$	$D_D^Y = Y_D / Y$	$D_D^N = N_D / N$	$\ln(D_D^Y / D_D^N)$
All	$Y + N$	Y	N	

### Example:

Categories: Education level bins

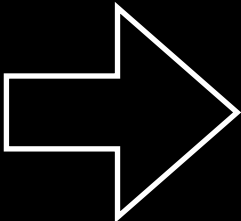
Outcomes: Good loan status (Y), Defaulted (N)

- NumToCatWoE

$f : (1) \text{ binning}(x) \rightarrow \text{categorical\_x}$   
 $(2) \text{ WoE}(\text{categorical\_x})$



y	x1	x2	x3
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>
[0,100]	<str>	<int>	<int>



y	x1	x2	x3	$f(3)$
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>
[0,100]	<str>	<int>	<int>	<int>

# Transformers included in DAI

<http://docs.h2o.ai/driverless-ai/latest-stable/docs/userguide/transformations.html>

## 1. Numeric (int, real, binary)

## 2. Time series

## 3. Categorical (string)

## 4. Text (string)

## 5. Time (date, time)

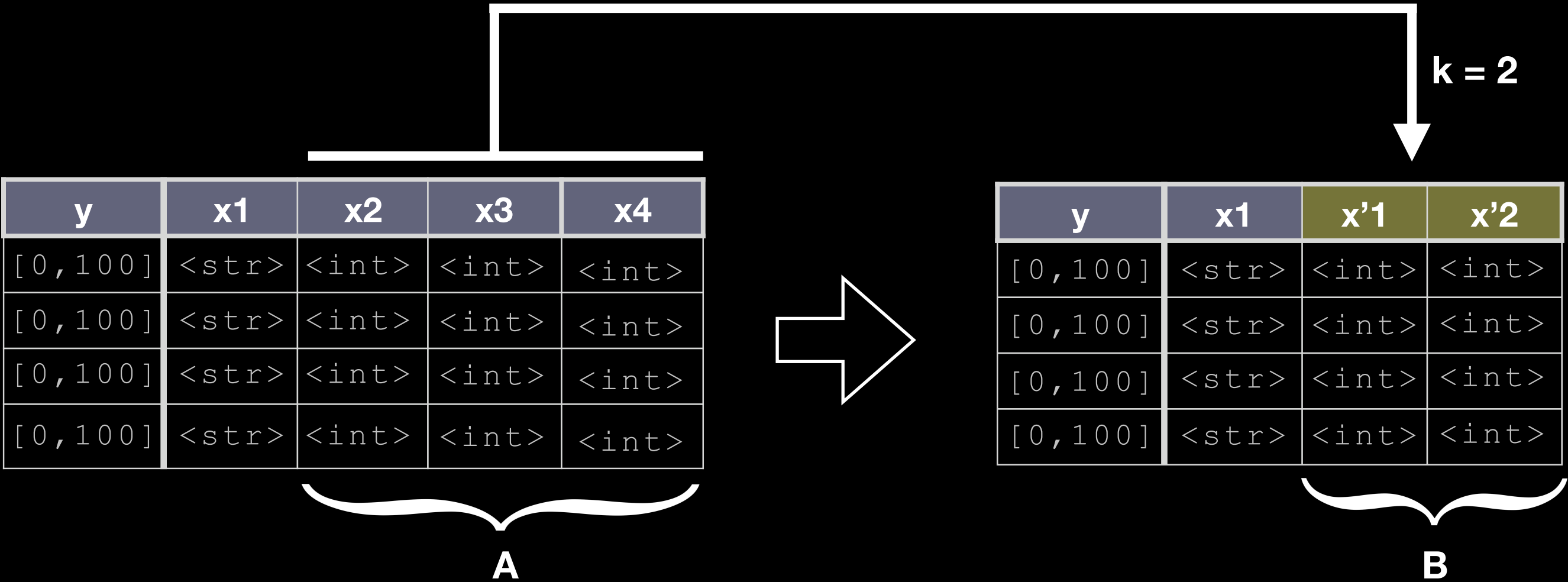
- Original
- Interactions
- ClusterDist
- NumCatTE
- NumToCatTE
- ClusterTE
- NumToCatWoE
- NumToCatWoEMonotonic
- TruncSVDNum

- TruncSVDNum

$$A = U \cdot \Sigma \cdot V^T$$

$\Sigma$  (=singular values): diagonal matrix with ordered values

- f* : (1) calculate  $U, \Sigma, V^T$  from original data (A)  
(2) Truncate  $\Sigma$  to the top  $k$  singular values ( $\Sigma_k$ )  
(3) Use  $\Sigma_k$  and  $V_k^T$  to calculate new matrix B



Name of variable  $f(1,y)$  : TruncSVD\_x2\_x3\_x4\_k( $n^{\text{th}}$ )

# DAI, custom recipes

```
1  """Adds together 3 or more numeric features"""
2  from h2oai.core.transformer_utils import CustomTransformer
3  import datatable as dt
4  import numpy as np
5
6
7  class SumTransformer(CustomTransformer):
8
9      _regression = True
10     _binary = True
11     _multiclass = True
12     _numeric_output = True
13     _is_reproducible = True
14     _included_model_classes = None # List[str]
15     _excluded_model_classes = None # List[str]
16
17     @staticmethod
18     def is_enabled():
19         return True
20
21     @staticmethod
22     def do_acceptance_test():
23         return True
24
25     @staticmethod
26     def get_default_properties():
27         return dict(col_type="numeric", min_cols=3, max_cols="all", relative_importance=1)
28
29     def fit_transform(self, X: dt.Frame, y: np.array = None):
30         return self.transform(X)
31
32     def transform(self, X: dt.Frame):
33         return X[:, dt.sum([dt.f[x] for x in range(X.ncols)])]
```

Demo : Custom recipes