

# **Tabular data and Deep learning - state of affairs**

Ivan Rubachev - 2022

# Who am I



# Outline

## What we'll talk about

- what's the problem (tabular data)
- Survey:
  - History up until very recently (just good old ml)
  - Recent times (a bunch of tabular deep learning papers)
  - Better times (this bunch of tabular deep learning doesn't really *work*)
  - Our new stuff (but is it?)
- Closing remarks

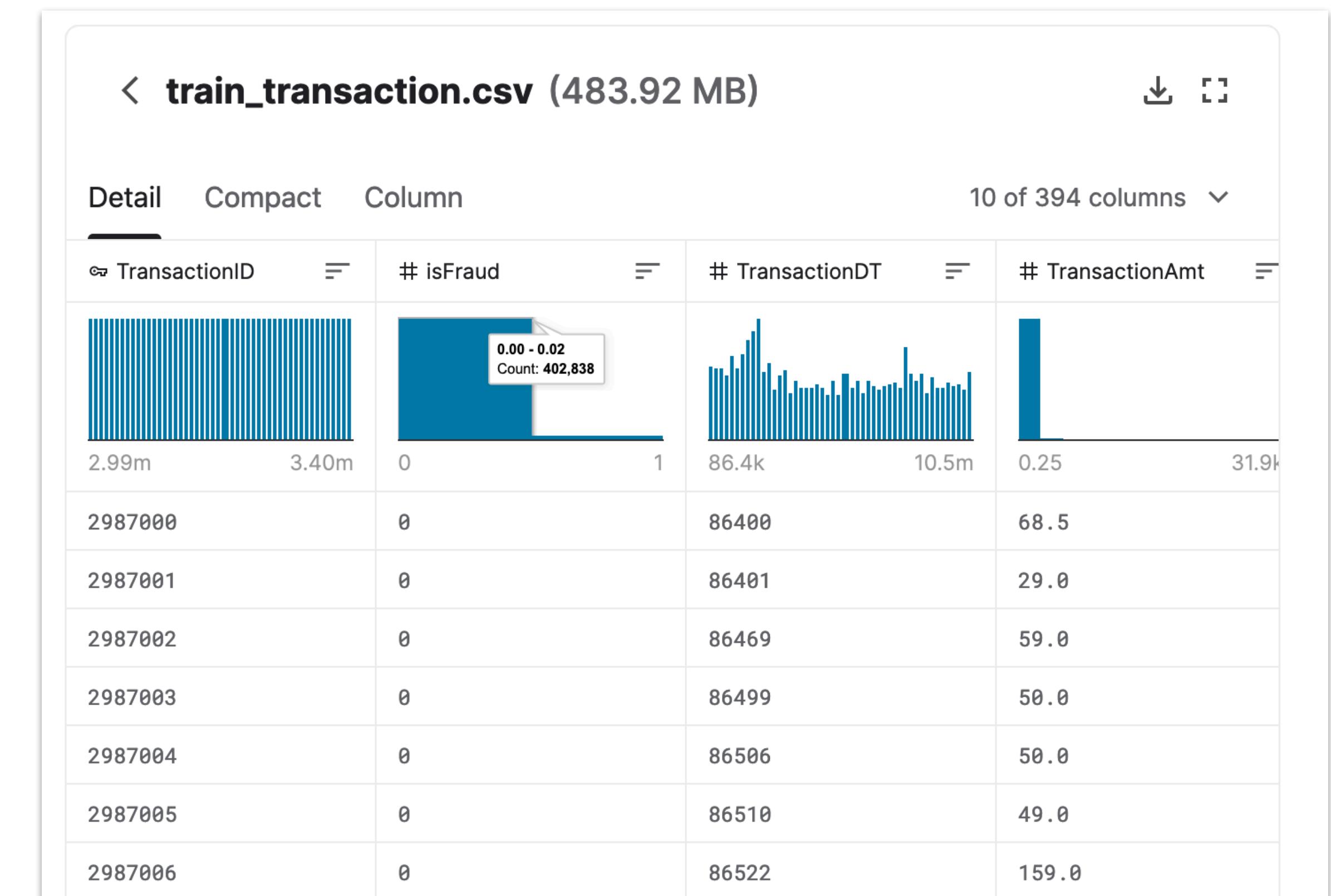
# The field

## Tabular data

Tabular data – two dimensional tables

rows ~ objects

columns ~ features



# The field

## Deep learning

Deep learning

- ~ more like a collection of methods
- ~ approach to building models
- ~ **very** successful in some fields



Transformers

# The field

Tabular data + Deep learning =



Bojan Tunguz  
@tunguz

...



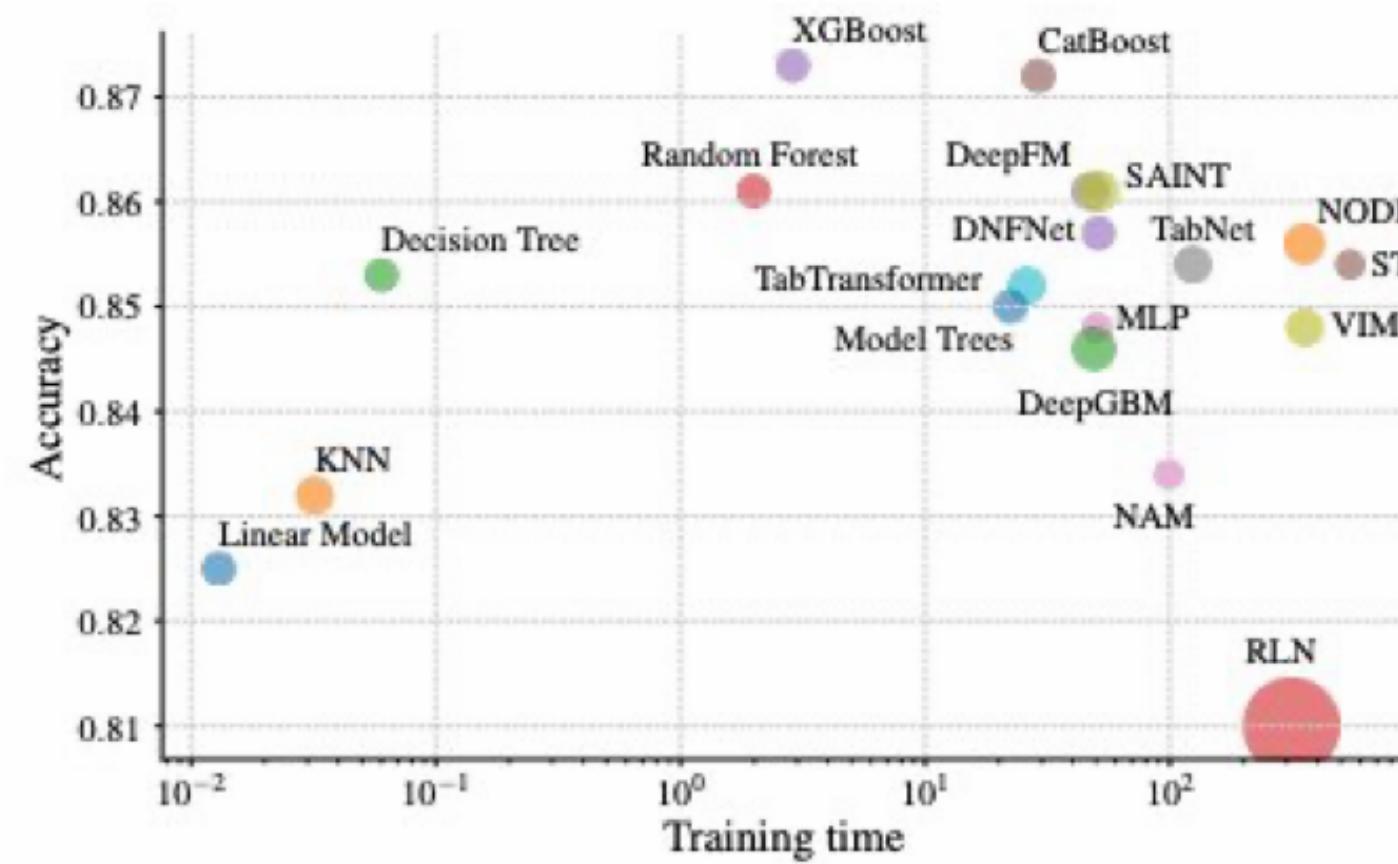
Mark Tenenholz  
@marktenenholz

...

XGBoost Is All You Need

Deep Neural Networks and Tabular  
Data: A Survey

[arxiv.org/pdf/2110.01889...](https://arxiv.org/pdf/2110.01889.pdf)



Adult dataset

Why is XGBoost so incredible?

It's not because:

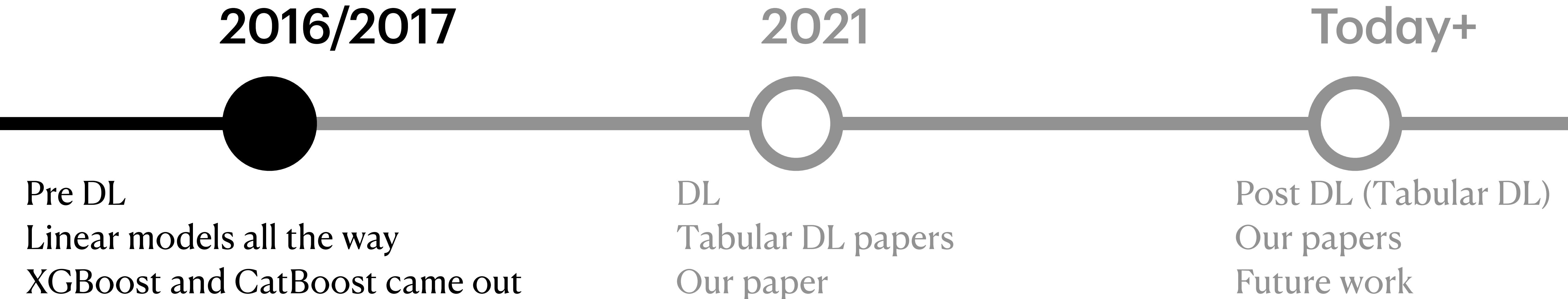
- It's always the most accurate
- It's always the fastest
- It's always the most interpretable

It's because:

- It's always the fastest to a good (tabular) solution.

Yeah, I can beat it with deep learning.  
But time is money.

# Timeline



# “Classical” ML algorithms



## Scaling Up the Accuracy of Naive-Bayes Classifiers: a Decision-Tree Hybrid

Ron Kohavi

Data Mining and Visualization

Silicon Graphics, Inc.

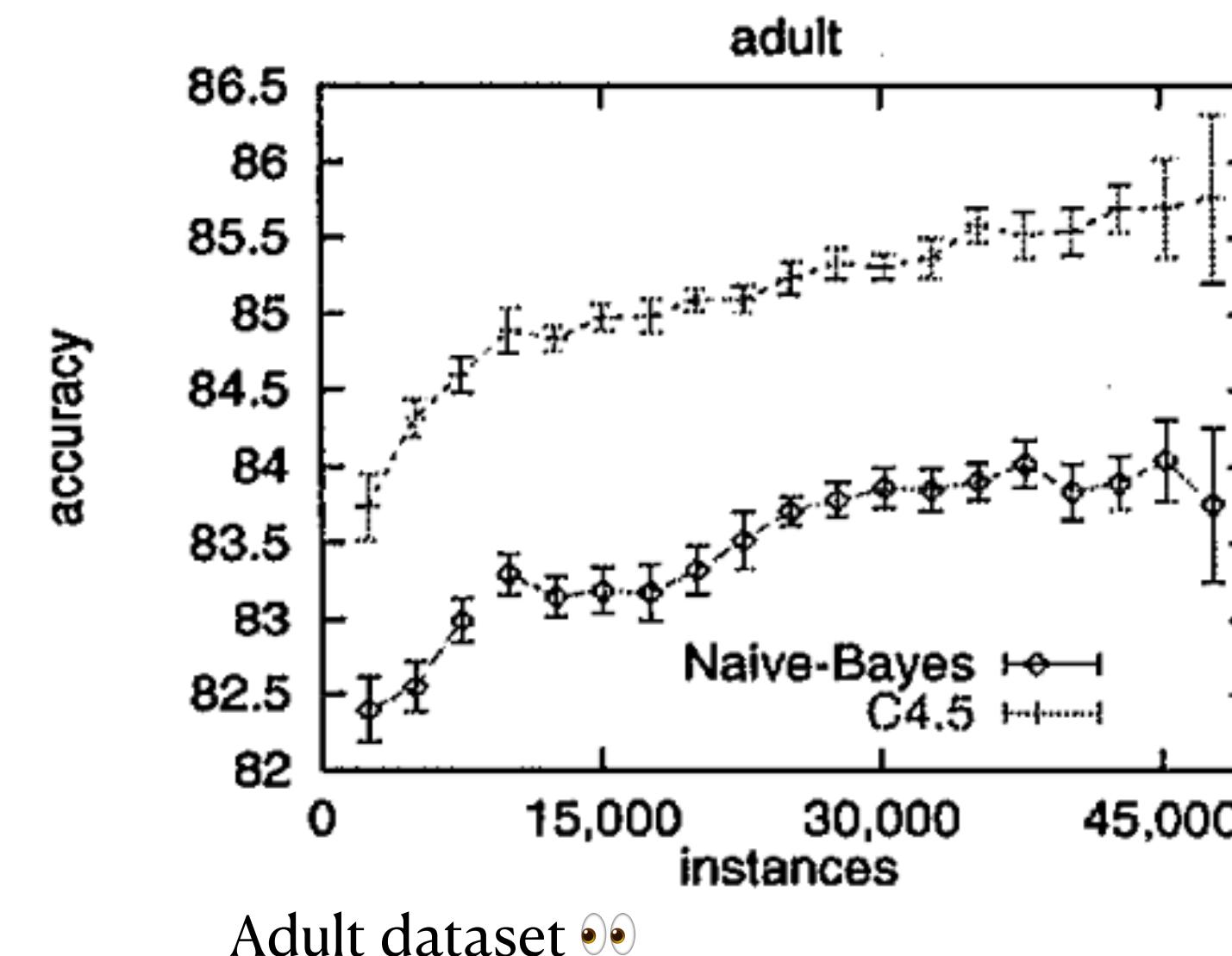
2011 N. Shoreline Blvd

Mountain View, CA 94043-1389

[ronnyk@sgi.com](mailto:ronnyk@sgi.com)

### 3.5. Experimental results

To verify the effects of sample selection bias experimentally, we apply Naive Bayes, logistic regression, C4.5 and SVMLight (soft margin) (Joachims, 2000b) to the Adult dataset, available from the UCI Machine Learning repository (Blake & Merz, 1998). We assume that the original dataset is not biased and artificially simulate biasedness by generating a value for  $s$  for each



# Deep learning is here

## Tabular, but not intentionally

### Self-Normalizing Neural Networks

Günter Klambauer

Thomas Unterthiner

Andreas Mayr

Sepp Hochreiter

LIT AI Lab & Institute of Bioinformatics,  
Johannes Kepler University Linz  
A-4040 Linz, Austria

{klambauer, unterthiner, mayr, hochreit}@bioinf.jku.at

$$\text{selu}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}.$$

**121 UCI Machine Learning Repository datasets.** The benchmark comprises 121 classification datasets from the UCI Machine Learning repository [10] from diverse application areas, such as physics, geology, or biology. The size of the datasets ranges between 10 and 130,000 data points and the number of features from 4 to 250. In abovementioned work [10], there were methodological mistakes [37] which we avoided here. Each compared FNN method was optimized with respect to its architecture and hyperparameters on a validation set that was then removed from the subsequent

# Deep learning is here

## More intentional

---

### Regularization Learning Networks: Deep Learning for Tabular Datasets

---

Ira Shavitt

Weizmann Institute of Science  
[irashavitt@gmail.com](mailto:irashavitt@gmail.com)

Eran Segal

Weizmann Institute of Science  
[eran.segal@weizmann.ac.il](mailto:eran.segal@weizmann.ac.il)

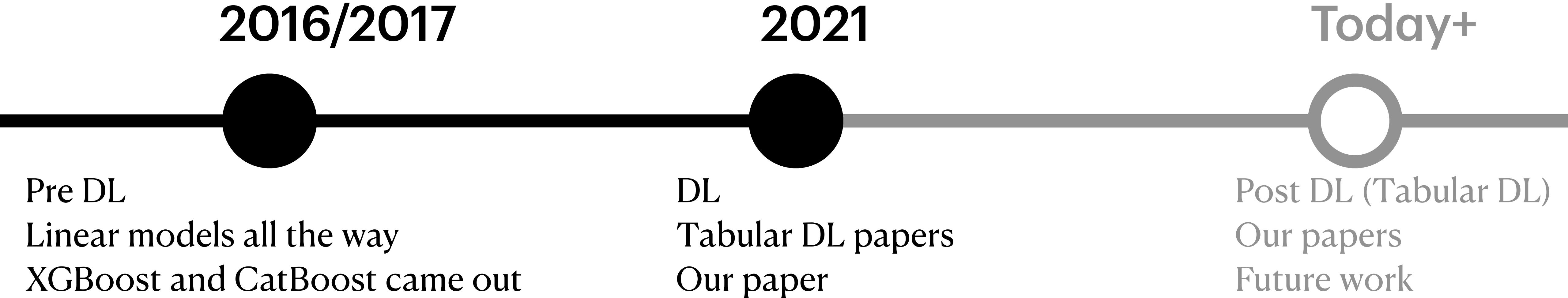
They introduce regularization hyperparameter for each weight in a NN

But optimize via SGD

Test on tabular DL problems

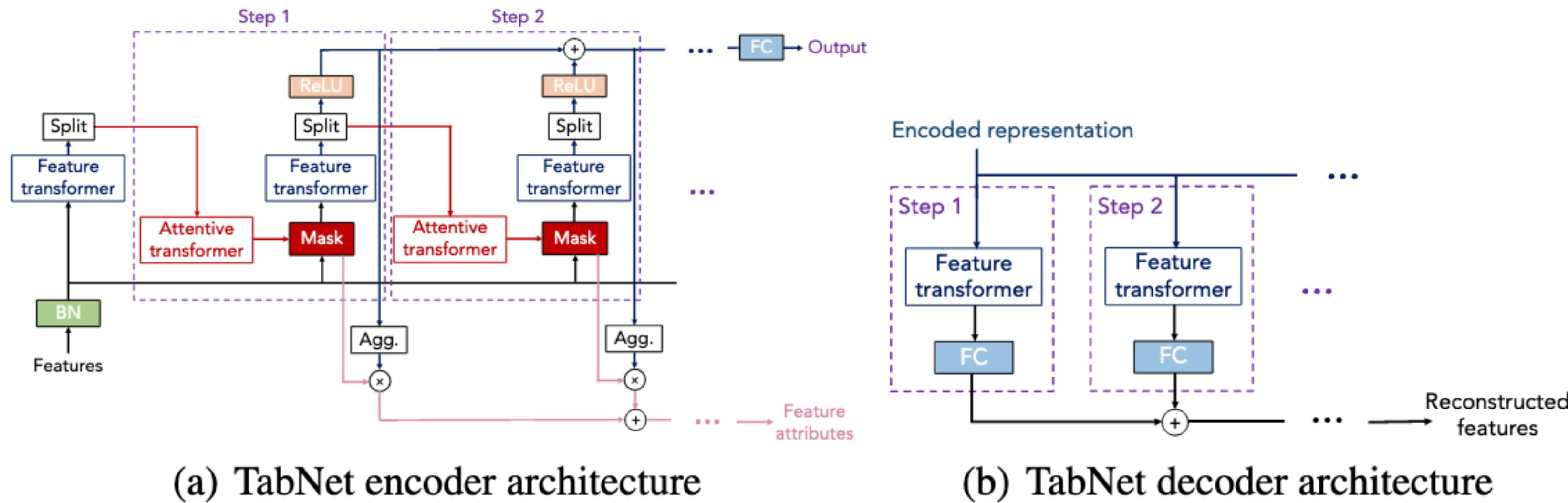
Conclusion: GBDT is still better

# Timeline



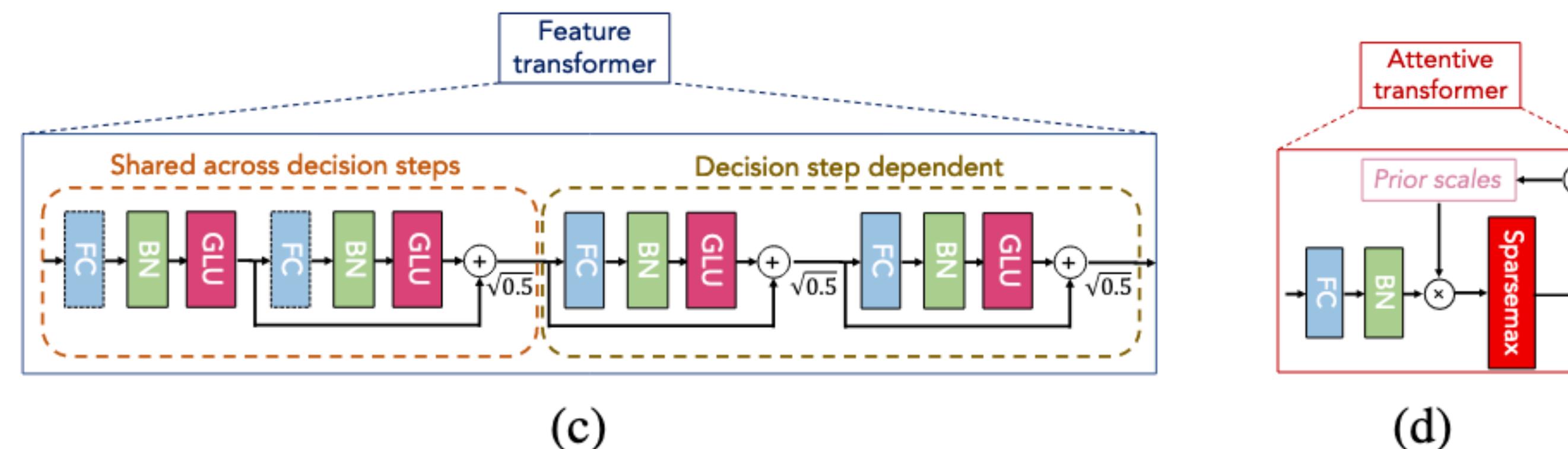
# Tabnet

## Architectures for tabular data



(a) TabNet encoder architecture

(b) TabNet decoder architecture



# Tabnet

## Architectures for tabular data

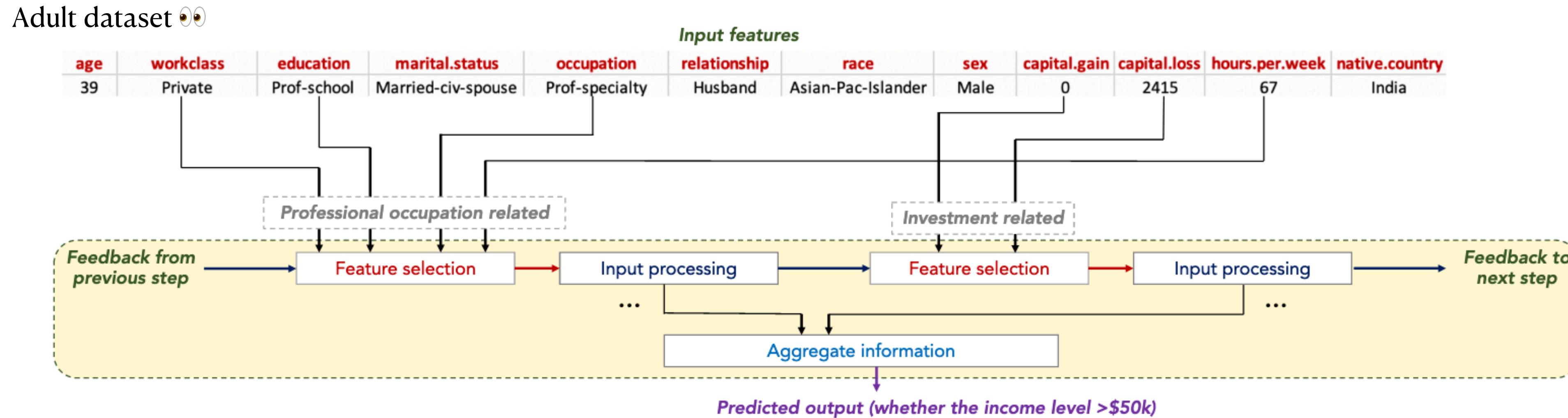


Figure 1: TabNet’s sparse feature selection exemplified for Adult Census Income prediction (Dua and Graff 2017). Sparse feature selection enables interpretability and better learning as the capacity is used for the most salient features. TabNet employs multiple decision blocks that focus on processing a subset of input features for reasoning. Two decision blocks shown as examples process features that are related to professional occupation and investments, respectively, in order to predict the income level.

# GrowNet

## Architectures for tabular data

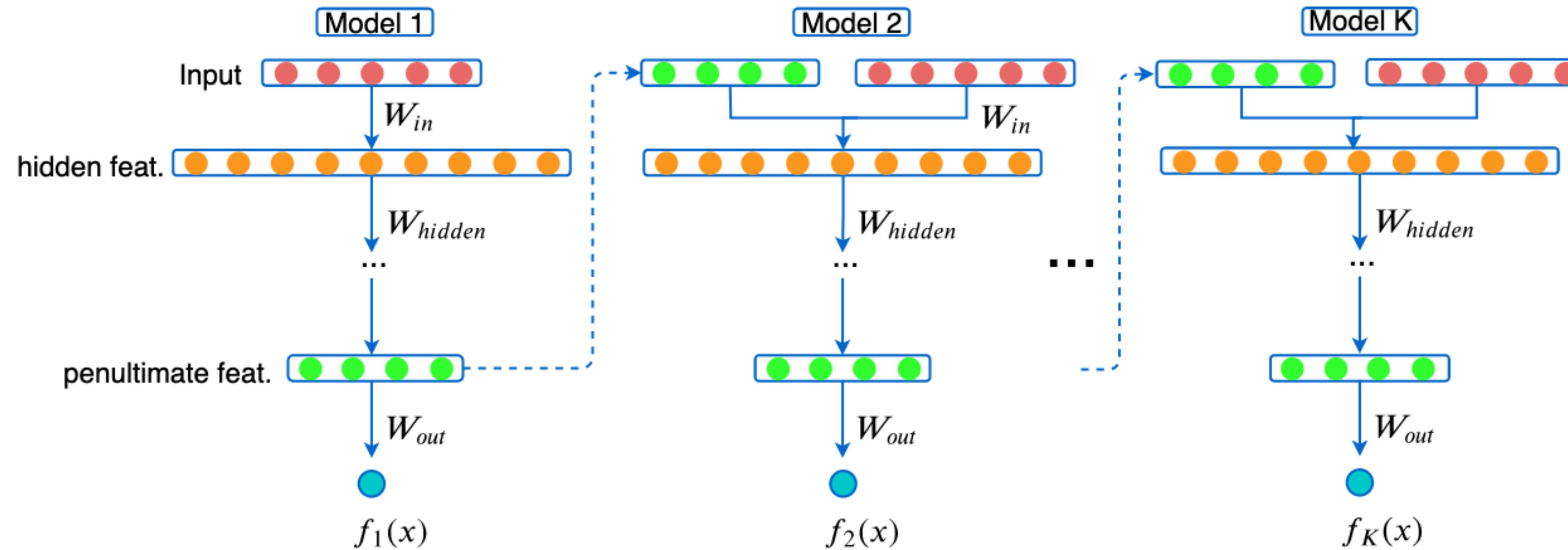


Figure 1: GrowNet architecture. After the first weak learner, each predictor is trained on combined features from original input and penultimate layer features from previous weak learner. The final output is the weighted sum of outputs from all predictors,  $\sum_{k=1}^{K} \alpha_k f_k(x)$ . Here Model K means weak learner K.

# NODE

## Architectures for tabular data

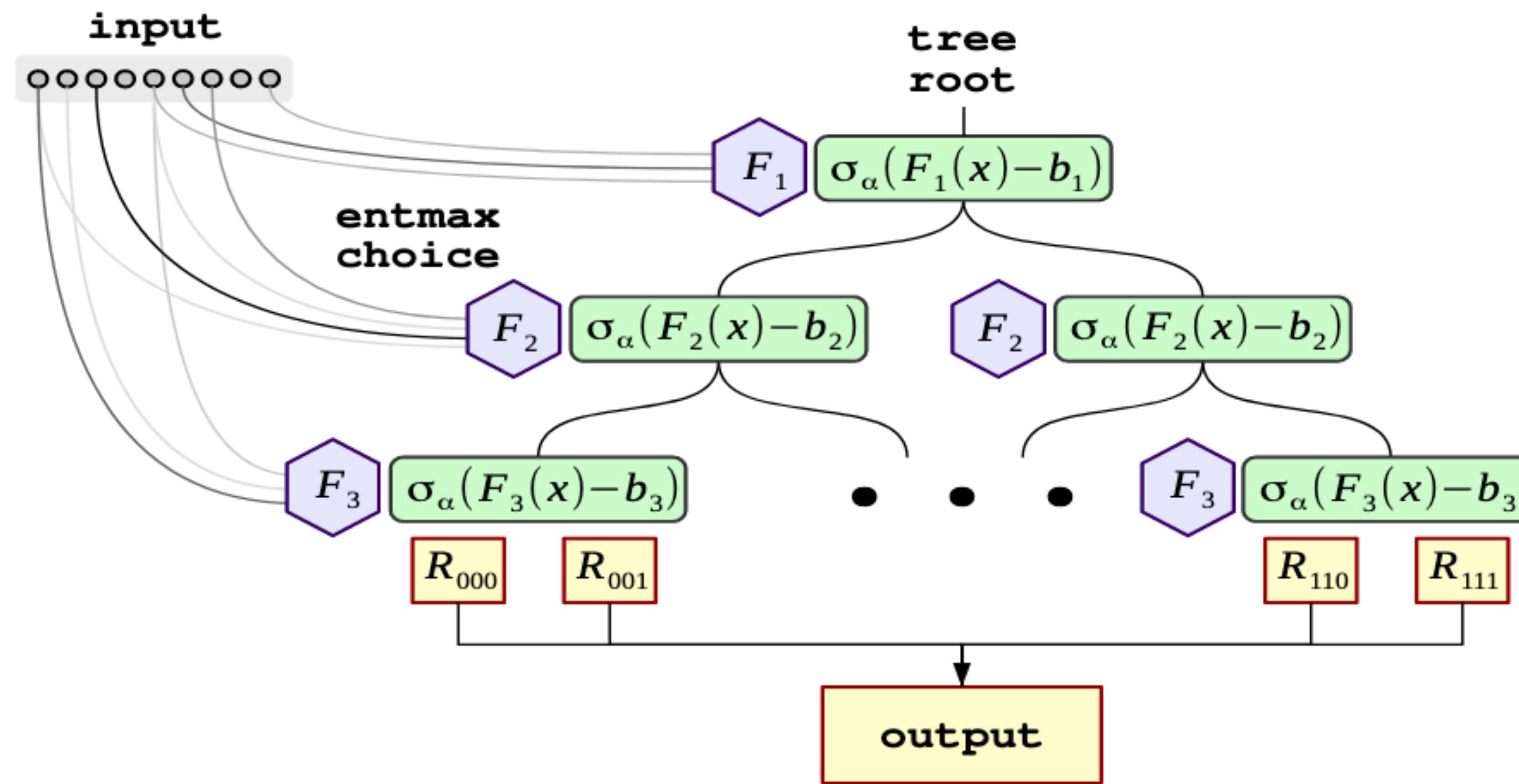


Figure 1: The single ODT inside the NODE layer. The splitting features and the splitting thresholds are shared across all the internal nodes of the same depth. The output is a sum of leaf responses scaled by the choice weights.

# NODE

## Architectures for tabular data

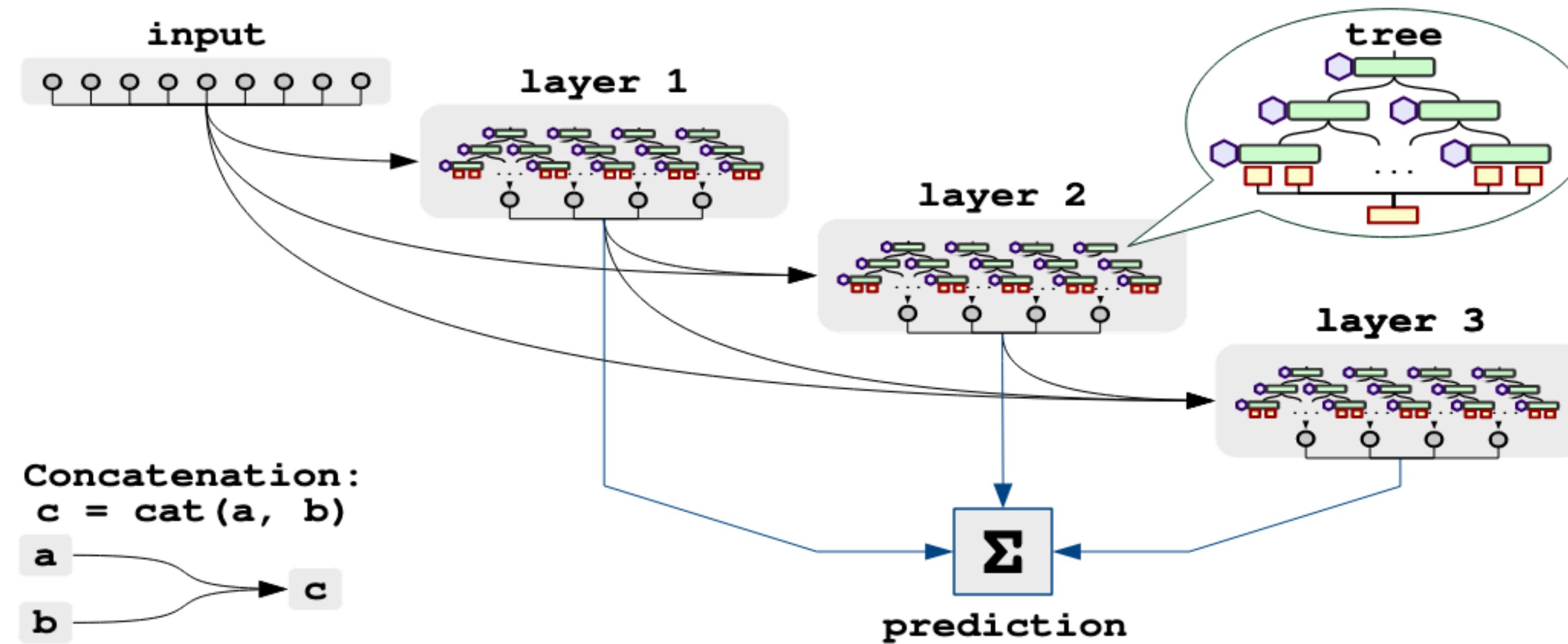
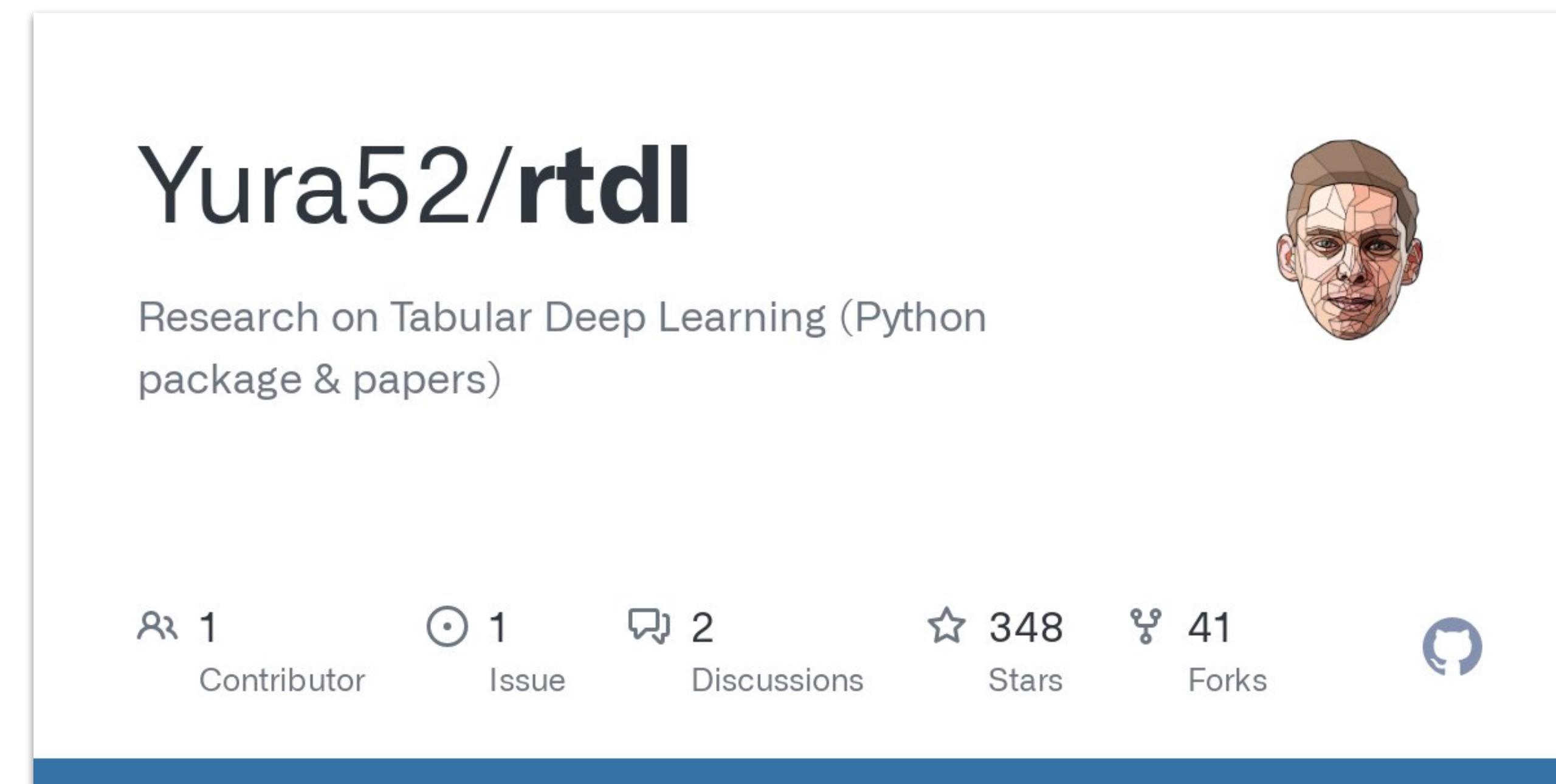


Figure 2: The NODE architecture, consisting of densely connected NODE layers. Each layer contains several trees whose outputs are concatenated and serve as input for the subsequent layer. The final prediction is obtained by averaging the outputs of all trees from all the layers.

# Revisiting deep learning for tabular data

## Our paper from NeurIPS 21

- But where are the results?
- Here: <https://github.com/Yura52/tabular-dl-revisiting-models>



# Revisiting deep learning for tabular data

Our paper from NeurIPS 21

$$\begin{aligned} \text{MLP}(x) &= \text{Linear}(\text{MLPBlock}(\dots(\text{MLPBlock}(x)))) \\ \text{MLPBlock}(x) &= \text{Dropout}(\text{ReLU}(\text{Linear}(x))) \end{aligned} \tag{1}$$

$$\begin{aligned} \text{ResNet}(x) &= \text{Prediction}(\text{ResNetBlock}(\dots(\text{ResNetBlock}(\text{Linear}(x))))) \\ \text{ResNetBlock}(x) &= x + \text{Dropout}(\text{Linear}(\text{Dropout}(\text{ReLU}(\text{Linear}(\text{BatchNorm}(x)))))) \\ \text{Prediction}(x) &= \text{Linear}(\text{ReLU}(\text{BatchNorm}(x))) \end{aligned} \tag{2}$$

# Revisiting deep learning for tabular data

## Our paper from NeurIPS 21

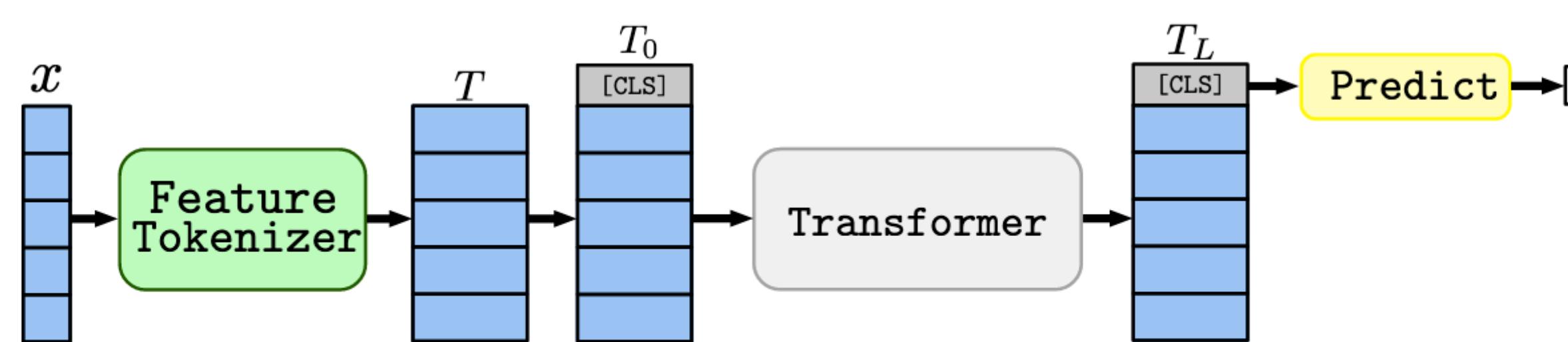
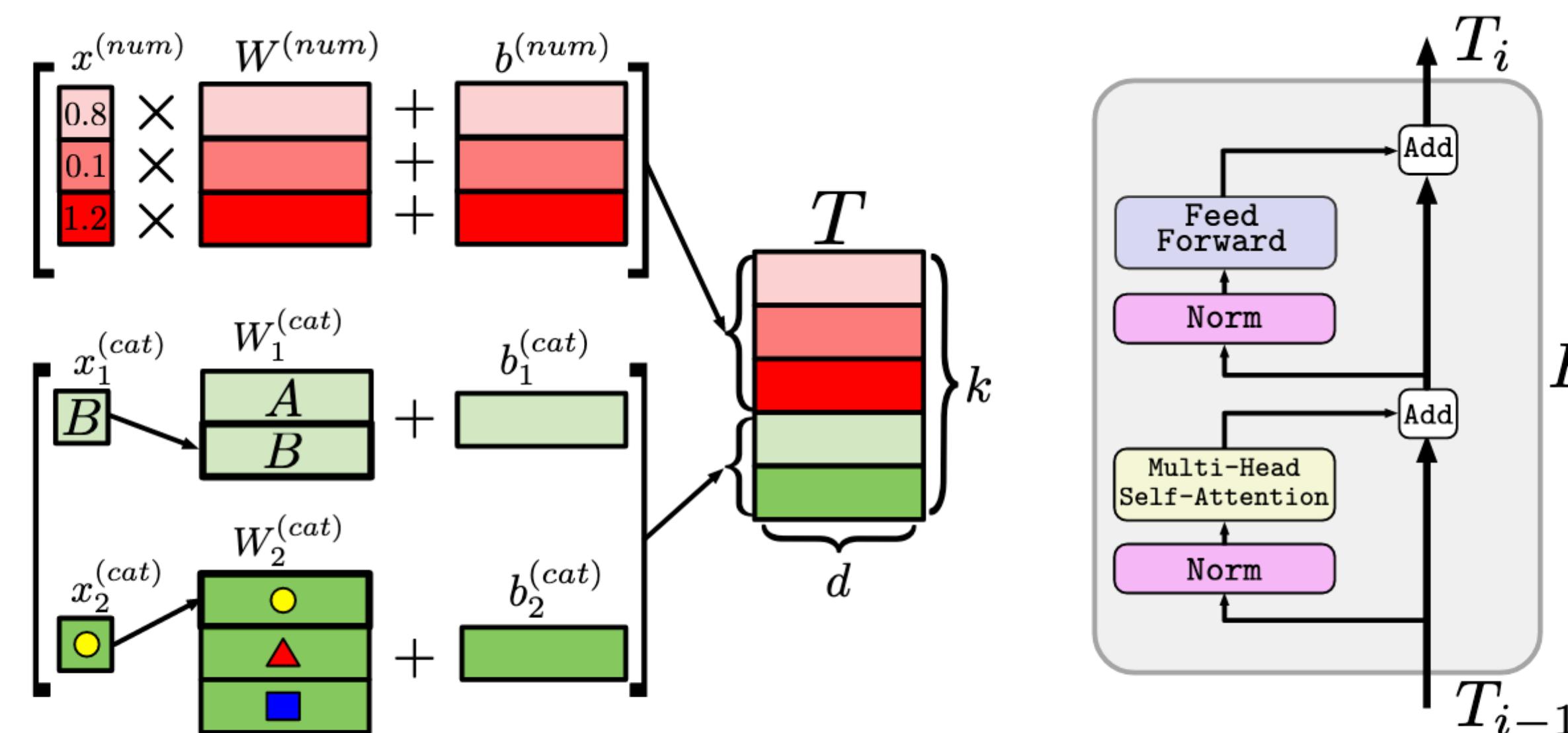


Figure 1: The FT-Transformer architecture. Firstly, Feature Tokenizer transforms features to embeddings. The embeddings are then processed by the Transformer module and the final representation of the [CLS] token is used for prediction.



# Revisiting deep learning for tabular data

Our paper from NeurIPS 21

Adult dataset 

	CA ↓	AD ↑	HE ↑	JA ↑	HI ↑	AL ↑	EP ↑	YE ↓	CO ↑	YA ↓	MI ↓	rank (std)
TabNet	0.510	0.850	0.378	0.723	0.719	0.954	0.8896	8.909	0.957	0.823	0.751	7.5 (2.0)
SNN	0.493	0.854	0.373	0.719	0.722	0.954	0.8975	8.895	0.961	0.761	0.751	6.4 (1.4)
AutoInt	0.474	<b>0.859</b>	0.372	0.721	0.725	0.945	0.8949	8.882	0.934	0.768	0.750	5.7 (2.3)
GrowNet	0.487	<b>0.857</b>	–	–	0.722	–	0.8970	8.827	–	0.765	0.751	5.7 (2.2)
MLP	0.499	0.852	0.383	0.719	0.723	0.954	0.8977	8.853	0.962	0.757	0.747	4.8 (1.9)
DCN2	0.484	0.853	0.385	0.716	0.723	0.955	0.8977	8.890	0.965	0.757	0.749	4.7 (2.0)
NODE	0.464	<b>0.858</b>	0.359	0.727	0.726	0.918	0.8958	<b>8.784</b>	0.958	<b>0.753</b>	<b>0.745</b>	3.9 (2.8)
ResNet	0.486	0.854	<b>0.396</b>	0.728	0.727	<b>0.963</b>	0.8969	8.846	0.964	0.757	0.748	3.3 (1.8)
FT-T	<b>0.459</b>	<b>0.859</b>	0.391	<b>0.732</b>	<b>0.729</b>	0.960	<b>0.8982</b>	8.855	<b>0.970</b>	0.756	0.746	1.8 (1.2)

# Revisiting deep learning for tabular data

Our paper from NeurIPS 21

Adult dataset 

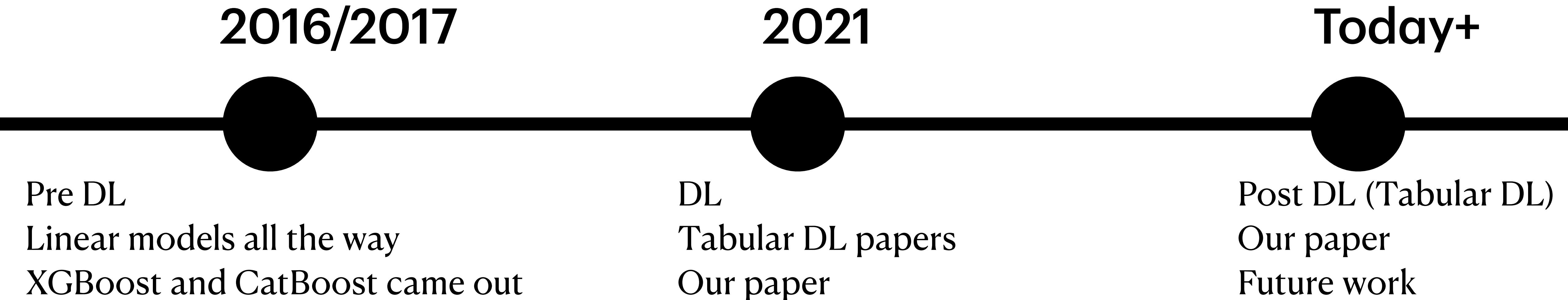
	CA ↓	AD ↑	HE ↑	JA ↑	HI ↑	AL ↑	EP ↑	YE ↓	CO ↑	YA ↓	MI ↓
Default hyperparameters											
XGBoost	0.462	<b>0.874</b>	0.348	0.711	0.717	0.924	0.8799	9.192	0.964	0.761	0.751
CatBoost	<b>0.428</b>	0.873	0.386	0.724	0.728	0.948	0.8893	8.885	0.910	0.749	0.744
FT-Transformer	0.454	0.860	<b>0.395</b>	<b>0.734</b>	<b>0.731</b>	<b>0.966</b>	<b>0.8969</b>	<b>8.727</b>	<b>0.973</b>	<b>0.747</b>	<b>0.742</b>
Tuned hyperparameters											
XGBoost	0.431	0.872	0.377	0.724	0.728	–	0.8861	8.819	0.969	<b>0.732</b>	0.742
CatBoost	<b>0.423</b>	<b>0.874</b>	0.388	0.727	0.729	–	0.8898	8.837	0.968	0.740	<b>0.741</b>
ResNet	0.478	0.857	0.398	0.734	0.731	0.966	0.8976	8.770	0.967	0.751	0.745
FT-Transformer	0.448	0.860	<b>0.398</b>	<b>0.739</b>	<b>0.731</b>	<b>0.967</b>	<b>0.8984</b>	<b>8.751</b>	<b>0.973</b>	0.747	0.743

# Concurrent work

**Tabular DL go brrr (not really)**

- Well tuned MLPs outperform GBDTs <https://arxiv.org/abs/2106.11189>
- SAINT: transformers for tabular data problems <https://arxiv.org/abs/2106.01342>

# Timeline



# Embeddings for numerical features

## Our work

$$z_i = f_i \left( x_i^{(num)} \right) \in \mathbb{R}^{d_i}$$

$\text{MLP}(z_1, \dots, z_k) = \text{MLP}(\text{concat}[z_1, \dots, z_k])$

$\text{concat}[z_1, \dots, z_k] \in \mathbb{R}^{d_1 + \dots + d_k}$

	GE ↑	CH ↑	EY ↑	CA ↓	HO ↓	AD ↑	OT ↑	HI ↑	FB ↓	SA ↑	CO ↑	MI ↓
MLP	0.632	0.856	0.615	0.495	3.204	0.854	0.818	0.720	5.686	0.912	<b>0.964</b>	0.747
MLP-L	<b>0.639</b>	<b>0.861</b>	0.635	0.475	3.123	0.856	<b>0.820</b>	0.723	5.684	0.916	0.963	0.748
MLP-LR	<b>0.642</b>	<b>0.860</b>	<b>0.660</b>	<b>0.471</b>	<b>3.084</b>	0.857	<b>0.819</b>	<b>0.726</b>	<b>5.625</b>	0.923	0.963	<b>0.746</b>
MLP-AutoDis	<b>0.649</b>	0.857	0.634	0.474	3.165	<b>0.859</b>	0.807	<b>0.725</b>	<b>5.670</b>	<b>0.924</b>	<b>0.963</b>	–

# Embeddings for numerical features

## Our work

$$z_i = f_i \left( x_i^{(num)} \right) \in \mathbb{R}^{d_i}$$

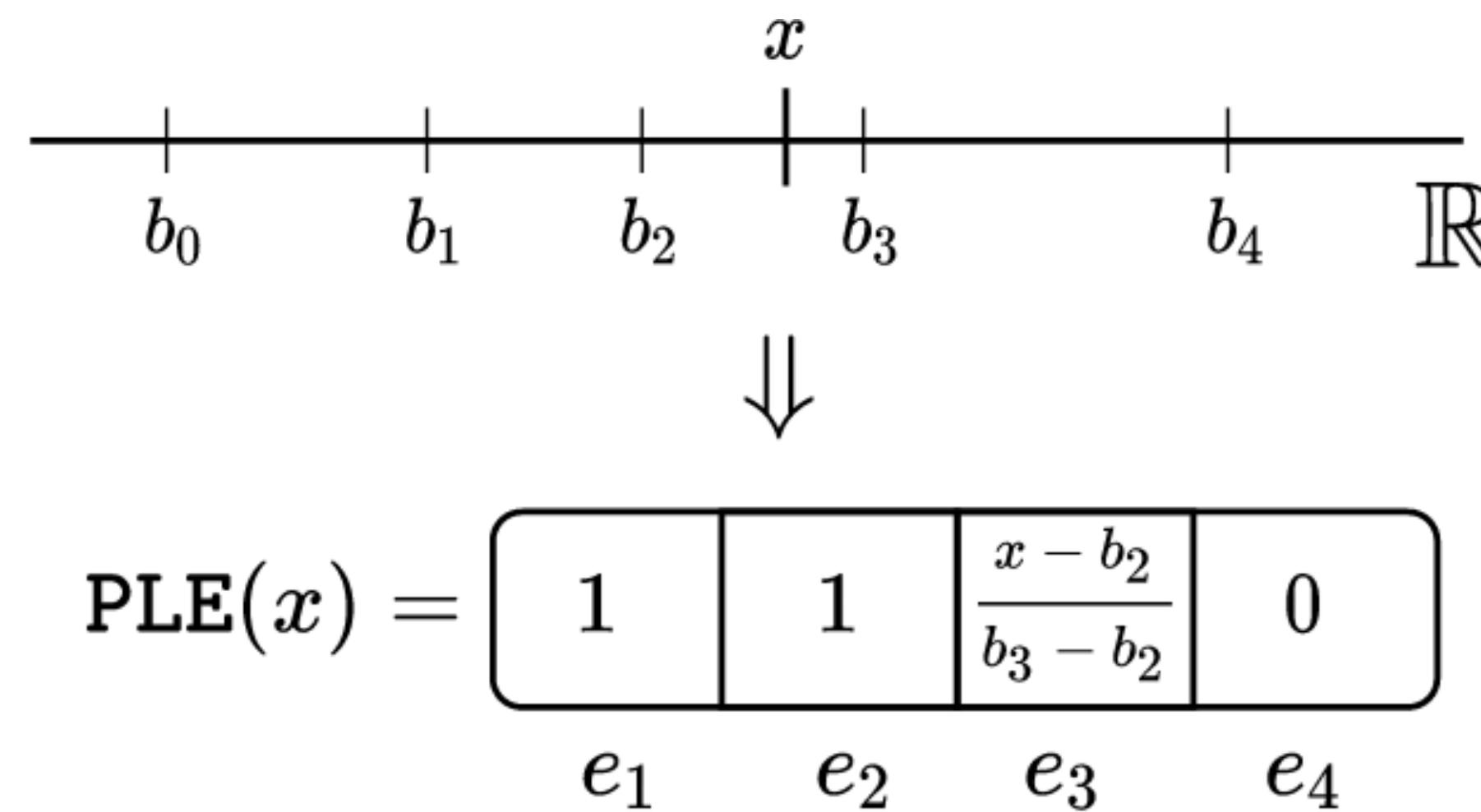
$\text{MLP}(z_1, \dots, z_k) = \text{MLP}(\text{concat}[z_1, \dots, z_k])$

$\text{concat}[z_1, \dots, z_k] \in \mathbb{R}^{d_1 + \dots + d_k}$

	GE ↑	CH ↑	EY ↑	CA ↓	HO ↓	AD ↑	OT ↑	HI ↑	FB ↓	SA ↑	CO ↑	MI ↓
MLP	0.632	0.856	0.615	0.495	3.204	0.854	0.818	0.720	5.686	0.912	<b>0.964</b>	0.747
MLP-L	<b>0.639</b>	<b>0.861</b>	0.635	0.475	3.123	0.856	<b>0.820</b>	0.723	5.684	0.916	0.963	0.748
MLP-LR	<b>0.642</b>	<b>0.860</b>	<b>0.660</b>	<b>0.471</b>	<b>3.084</b>	0.857	<b>0.819</b>	<b>0.726</b>	<b>5.625</b>	0.923	0.963	<b>0.746</b>
MLP-AutoDis	<b>0.649</b>	0.857	0.634	0.474	3.165	<b>0.859</b>	0.807	<b>0.725</b>	<b>5.670</b>	<b>0.924</b>	<b>0.963</b>	–

# Embeddings for numerical features

Our work



$$f_i(x) = \text{Periodic}(x) = \text{concat}[\sin(v), \cos(v)],$$

$$v = [2\pi c_1 x, \dots, 2\pi c_k x]$$

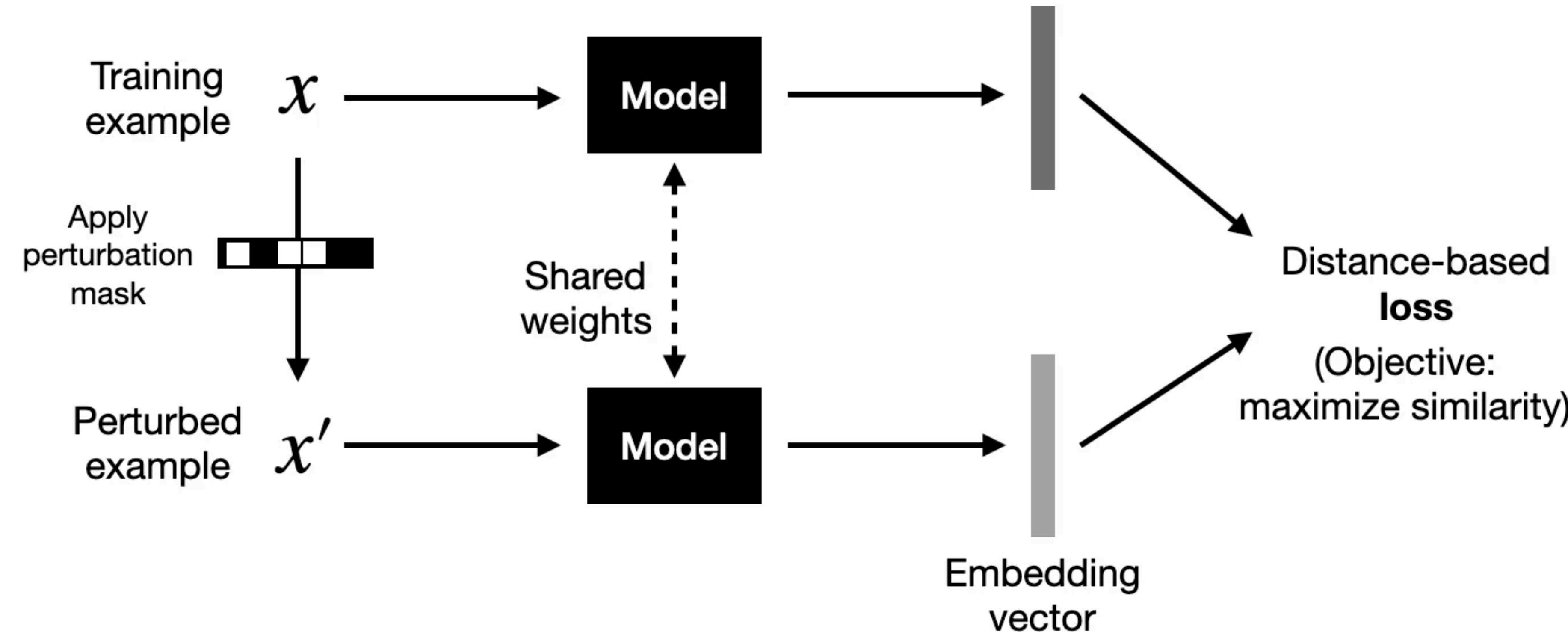
# Embeddings for numerical features

**Very sorry :(**

	GE ↑	CH ↑	EY ↑	CA ↓	HO ↓	AD ↑	OT ↑	HI ↑	FB ↓	SA ↑	CO ↑	MI ↓	Avg. Rank
CatBoost	0.692	0.861	0.757	0.430	3.093	0.873	0.825	0.727	5.226	0.924	0.967	<b>0.741</b>	6.8 ± 4.9
XGBoost	0.683	0.859	0.738	0.434	3.152	<b>0.875</b>	0.827	0.726	5.338	0.919	0.969	0.742	9.0 ± 5.7
MLP	0.665	0.856	0.637	0.486	3.109	0.856	0.822	0.727	5.616	0.913	0.968	0.746	15.6 ± 2.4
MLP-LR	0.679	0.861	0.694	0.463	3.012	0.859	0.826	0.731	5.477	0.924	0.972	0.744	10.2 ± 4.4
MLP-Q-LR	0.682	0.859	0.732	0.433	3.080	0.867	0.818	0.724	<b>5.144</b>	0.924	0.974	0.745	10.7 ± 4.6
MLP-T-LR	0.673	0.861	0.729	0.435	3.099	0.870	0.821	0.727	5.409	0.924	0.973	0.746	10.3 ± 3.8
MLP-PLR	<b>0.700</b>	0.858	0.968	0.453	<b>2.975</b>	0.874	<b>0.830</b>	<b>0.734</b>	5.388	<b>0.924</b>	0.975	0.743	4.9 ± 4.8
ResNet	0.690	0.861	0.667	0.483	3.081	0.856	0.821	0.734	5.482	0.918	0.968	0.745	12.1 ± 4.7
ResNet-LR	0.672	0.862	0.735	0.450	2.992	0.859	0.822	0.733	5.415	0.923	0.971	0.743	9.8 ± 4.3
ResNet-Q-LR	0.674	0.859	0.794	0.427	3.066	0.868	0.815	0.729	5.309	0.923	0.976	0.746	9.2 ± 4.8
ResNet-T-LR	0.683	0.862	0.817	<b>0.425</b>	3.030	0.872	0.822	0.731	5.471	0.923	0.975	0.744	7.8 ± 3.6
ResNet-PLR	0.691	0.861	0.925	0.443	3.040	<b>0.874</b>	0.825	0.734	5.400	0.924	0.975	0.743	5.2 ± 2.3
Transformer-L	0.668	0.861	0.769	0.455	3.188	0.860	0.824	0.727	5.434	0.924	0.973	0.743	10.6 ± 3.3
Transformer-LR	0.666	0.861	0.776	0.446	3.193	0.861	0.824	0.733	5.430	0.924	0.973	0.743	9.4 ± 4.1
Transformer-Q-LR	0.690	0.857	0.842	<b>0.425</b>	3.143	0.868	0.818	0.726	5.471	<b>0.924</b>	0.975	0.744	8.5 ± 5.5
Transformer-T-LR	0.686	0.862	0.833	<b>0.423</b>	3.149	0.871	0.823	0.733	5.515	0.924	<b>0.976</b>	0.744	7.2 ± 4.6
Transformer-PLR	0.686	<b>0.864</b>	<b>0.977</b>	0.449	3.091	0.873	0.823	0.734	5.581	<b>0.924</b>	0.975	0.743	6.0 ± 4.5

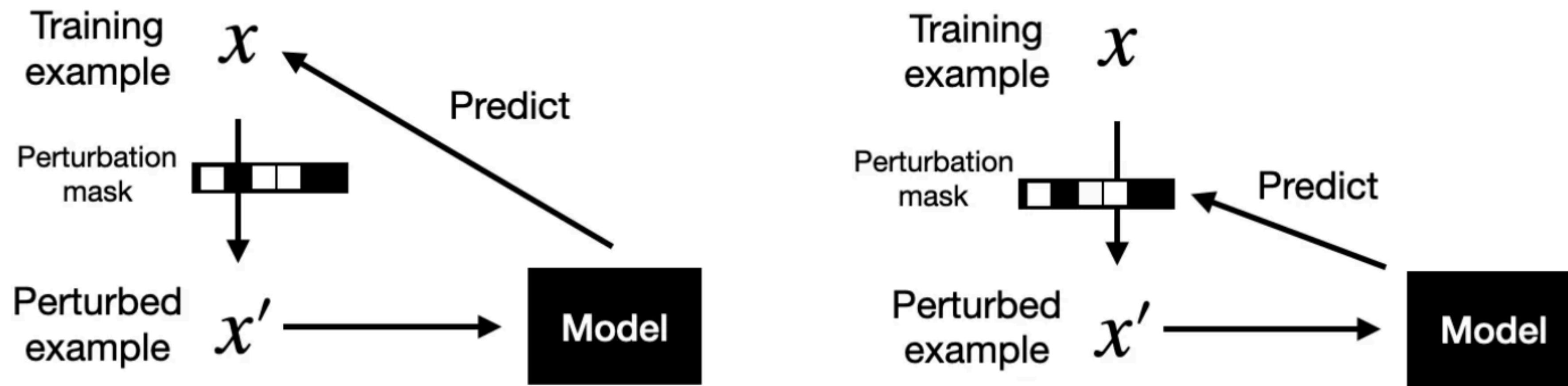
# Pretraining in tabular deep-learning

<https://magazine.sebastianraschka.com/p/ahead-of-ai-3-ainnouncements>



# Pretraining in tabular deep-learning

<https://magazine.sebastianraschka.com/p/ahead-of-ai-3-ainnouncements>



A simplified depiction of the two self-prediction processes. Left: the model predicts the original training example from the perturbed example. Right: the model predicts that perturbation mask given the perturbed example.

# Pretraining in tabular deep-learning

<https://magazine.sebastianraschka.com/p/ahead-of-ai-3-ainnouncements>

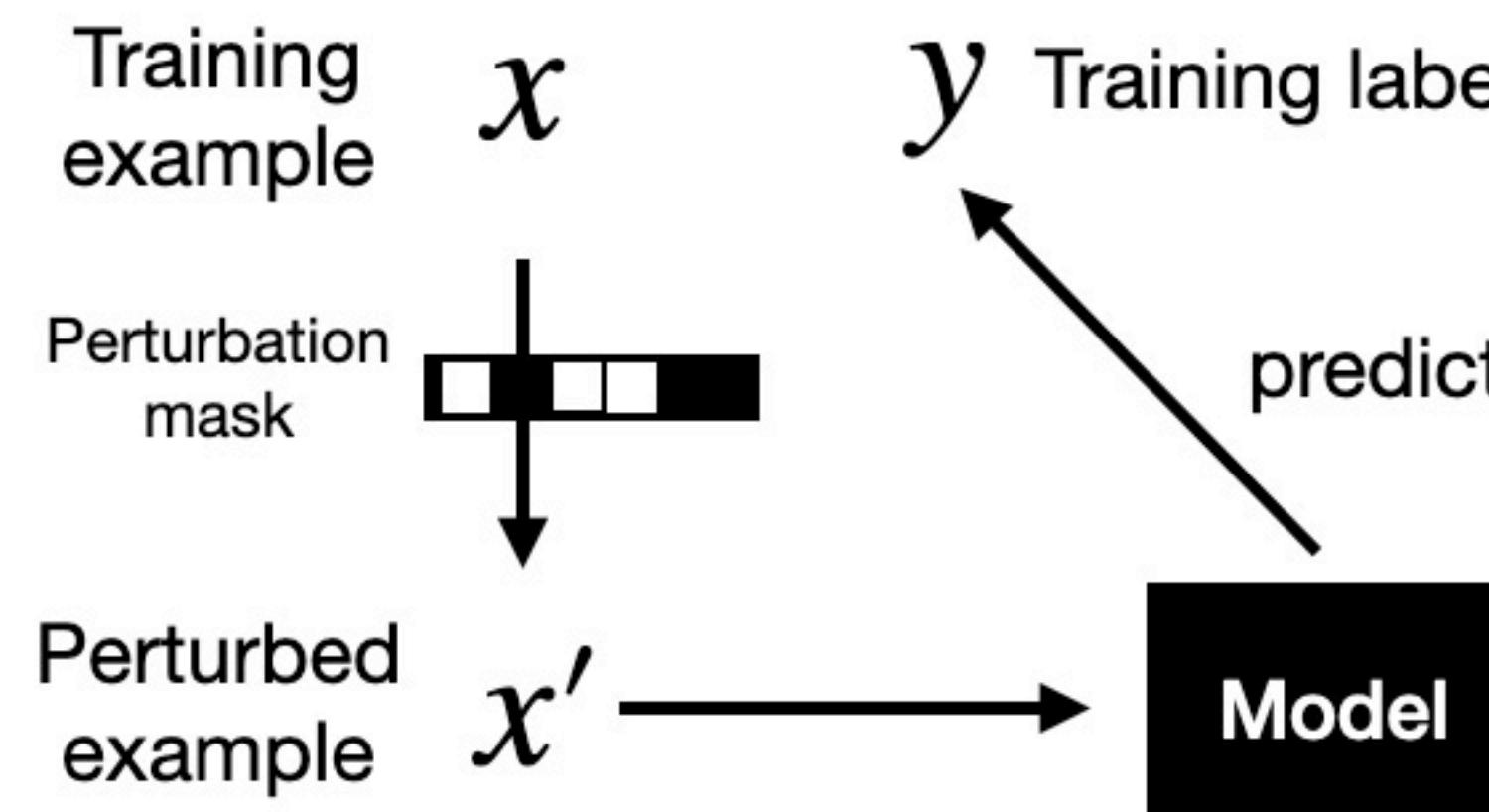
Table 2: Results for pretraining deep models with different objectives. We report metrics averaged over 15 seeds, bold entries correspond to results that are statistically significantly better (we use Tukey HSD test). The comparisons are separate for different models.  $\uparrow$  corresponds to accuracy and ROC-AUC metrics,  $\downarrow$  corresponds to RMSE and log-loss for OT. "no pretraining" stands for the supervised baseline, initialized with random weights

	GE $\uparrow$	CH $\uparrow$	CA $\downarrow$	HO $\downarrow$	OT $\downarrow$	HI $\uparrow$	FB $\downarrow$	AD $\uparrow$	WE $\downarrow$	CO $\uparrow$	MI $\downarrow$	
MLP												
no pretraining	0.635	0.849	0.506	3.156	0.479	0.801	5.737	0.908	1.909	0.963	0.749	
contrastive	0.672	<b>0.855</b>	0.455	<b>3.056</b>	0.469	<b>0.813</b>	5.697	0.910	1.881	0.960	0.748	
rec	0.662	0.853	<b>0.445</b>	<b>3.044</b>	<b>0.466</b>	0.805	<b>5.641</b>	<b>0.910</b>	<b>1.875</b>	<b>0.965</b>	<b>0.746</b>	
mask	<b>0.691</b>	<b>0.857</b>	0.454	3.113	0.472	<b>0.814</b>	<b>5.681</b>	<b>0.912</b>	1.883	0.964	0.748	
MLP-PLR												
no pretraining	0.668	<b>0.858</b>	0.469	<b>3.008</b>	0.483	0.809	<b>5.608</b>	0.926	1.890	0.969	0.746	
rec	0.667	0.852	<b>0.439</b>	<b>3.031</b>	<b>0.472</b>	0.808	<b>5.571</b>	<b>0.926</b>	<b>1.877</b>	<b>0.971</b>	<b>0.745</b>	
mask	<b>0.685</b>	<b>0.863</b>	<b>0.434</b>	<b>3.007</b>	0.477	<b>0.818</b>	<b>5.586</b>	<b>0.927</b>	1.911	<b>0.970</b>	0.748	
MLP-T-LR												
no pretraining	0.634	<b>0.866</b>	0.444	3.113	0.482	0.805	5.520	0.925	1.897	0.968	0.749	
rec		<b>0.652</b>	0.857	<b>0.424</b>	3.109	<b>0.472</b>	0.808	<b>5.363</b>	0.924	<b>1.861</b>	<b>0.969</b>	<b>0.746</b>
mask		<b>0.654</b>	<b>0.868</b>	<b>0.424</b>	<b>3.045</b>	<b>0.472</b>	<b>0.818</b>	5.544	<b>0.926</b>	1.916	<b>0.969</b>	0.748

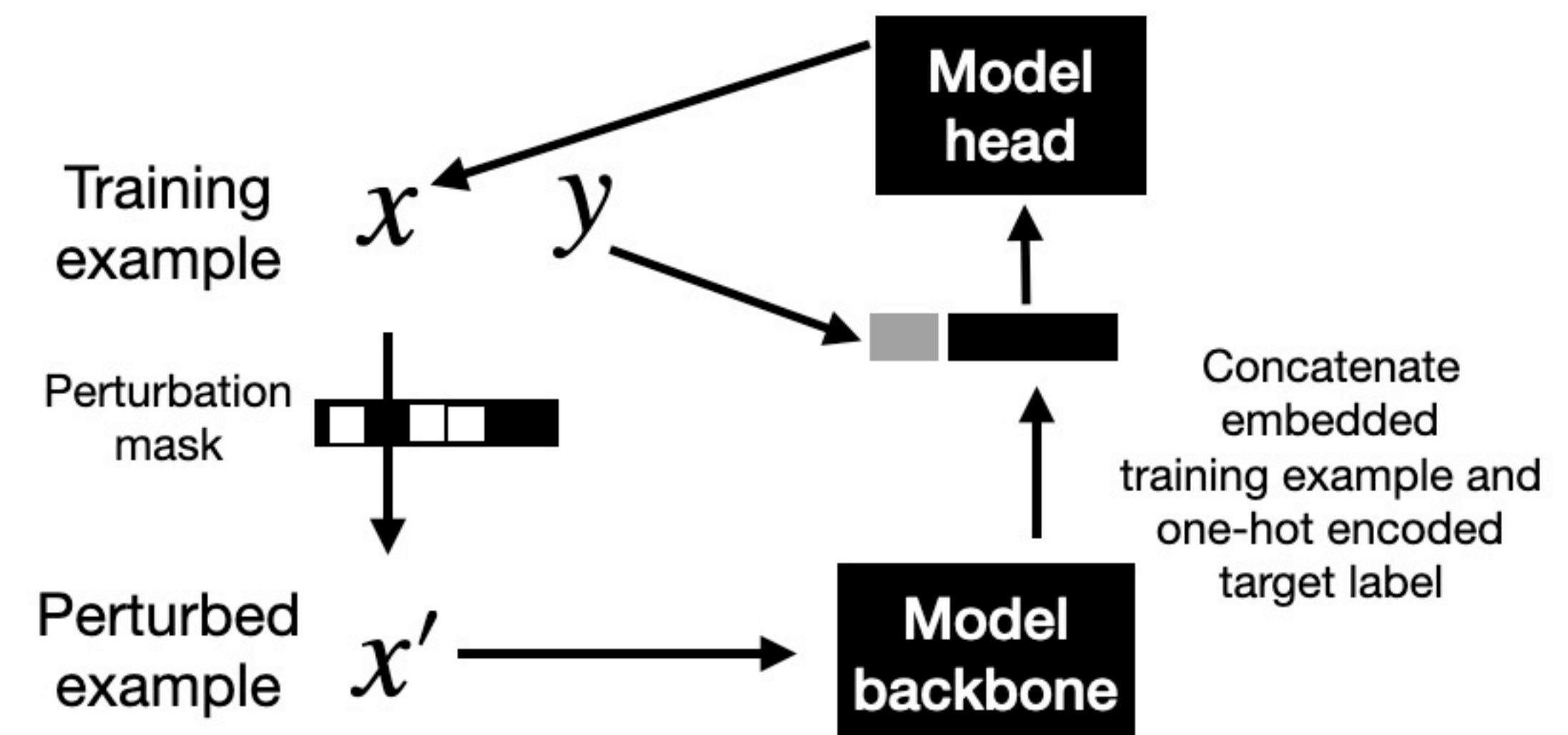
# Pretraining in tabular deep-learning

<https://magazine.sebastianraschka.com/p/ahead-of-ai-3-ainnouncements>

“sup”



“rec + target”



# Pretraining in tabular deep-learning

<https://magazine.sebastianraschka.com/p/ahead-of-ai-3-ainnouncements>

2 datasets where gradient boosting  
outperforms the deep learning-  
based methods

Table 4: Comparison of pretrained models to GBDT. Notation follows Table 2. Results represent ensembles of models. Bold entries correspond to the overall statistically significant best entries.

	GE ↑	CH ↑	CA ↓	HO ↓	OT ↓	HI ↑	FB ↓	AD ↑	WE ↓	CO ↑	MI ↓	Avg. Rank
CatBoost	0.692	<b>0.864</b>	0.430	3.093	0.450	0.807	5.226	0.928	1.801	0.967	<b>0.741</b>	2.6 ± 1.4
XGBoost	0.683	0.860	0.434	3.152	0.454	0.805	5.338	0.927	<b>1.782</b>	0.969	<b>0.742</b>	3.0 ± 1.5
MLP												
no pretraining	0.656	0.852	0.482	3.055	0.467	0.805	5.666	0.910	1.850	0.968	0.747	4.8 ± 1.1
mask + target	0.709	0.860	0.414	2.949	0.457	<b>0.828</b>	5.551	0.916	1.809	0.969	0.746	2.8 ± 1.2
rec + sup	0.709	0.859	0.419	2.951	<b>0.442</b>	0.817	5.531	0.913	1.801	0.973	0.745	2.5 ± 1.2
MLP-P-LR												
no pretraining	0.695	<b>0.864</b>	0.454	2.953	0.470	0.814	5.324	0.928	1.835	<b>0.974</b>	0.744	2.6 ± 1.2
mask + target	<b>0.719</b>	<b>0.866</b>	<b>0.407</b>	2.952	0.458	<b>0.828</b>	5.373	<b>0.930</b>	1.849	0.973	0.745	2.1 ± 1.2
rec + sup	<b>0.737</b>	0.862	0.424	2.964	0.449	0.811	<b>5.124</b>	<b>0.929</b>	1.813	<b>0.974</b>	0.744	2.0 ± 1.0
MLP-T-LR												
no pretraining	0.662	<b>0.868</b>	0.437	3.028	0.472	0.808	5.424	0.927	1.850	0.972	0.747	3.7 ± 1.1
mask + target	0.673	<b>0.868</b>	<b>0.410</b>	<b>2.894</b>	0.460	<b>0.827</b>	5.458	<b>0.930</b>	1.849	0.972	0.746	2.4 ± 1.4
rec + sup	0.705	<b>0.866</b>	0.425	3.057	<b>0.444</b>	0.814	5.422	0.927	1.811	<b>0.974</b>	0.746	2.5 ± 1.1

# Diffusion models 😱 and tabular data generation

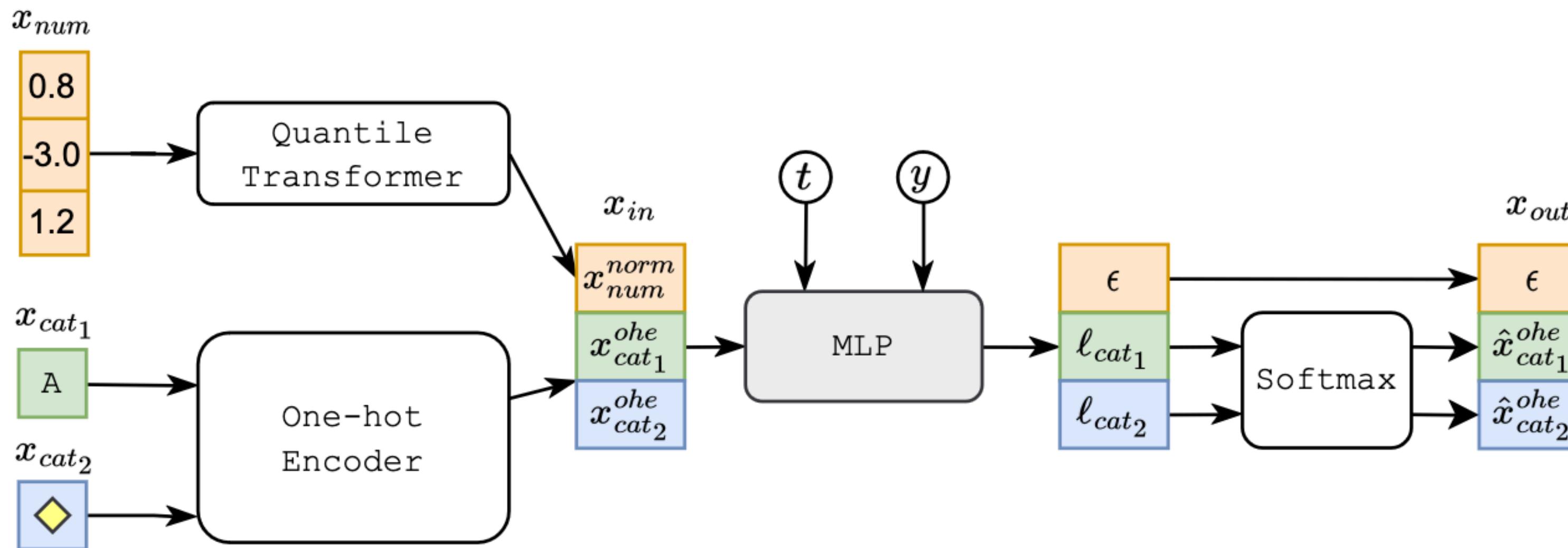


Figure 1: TabDDPM scheme for classification problems;  $t$ ,  $y$  and  $\ell$  denote a diffusion timestep, a class label, and logits, respectively.

# Diffusion models 🧨 and tabular data generation

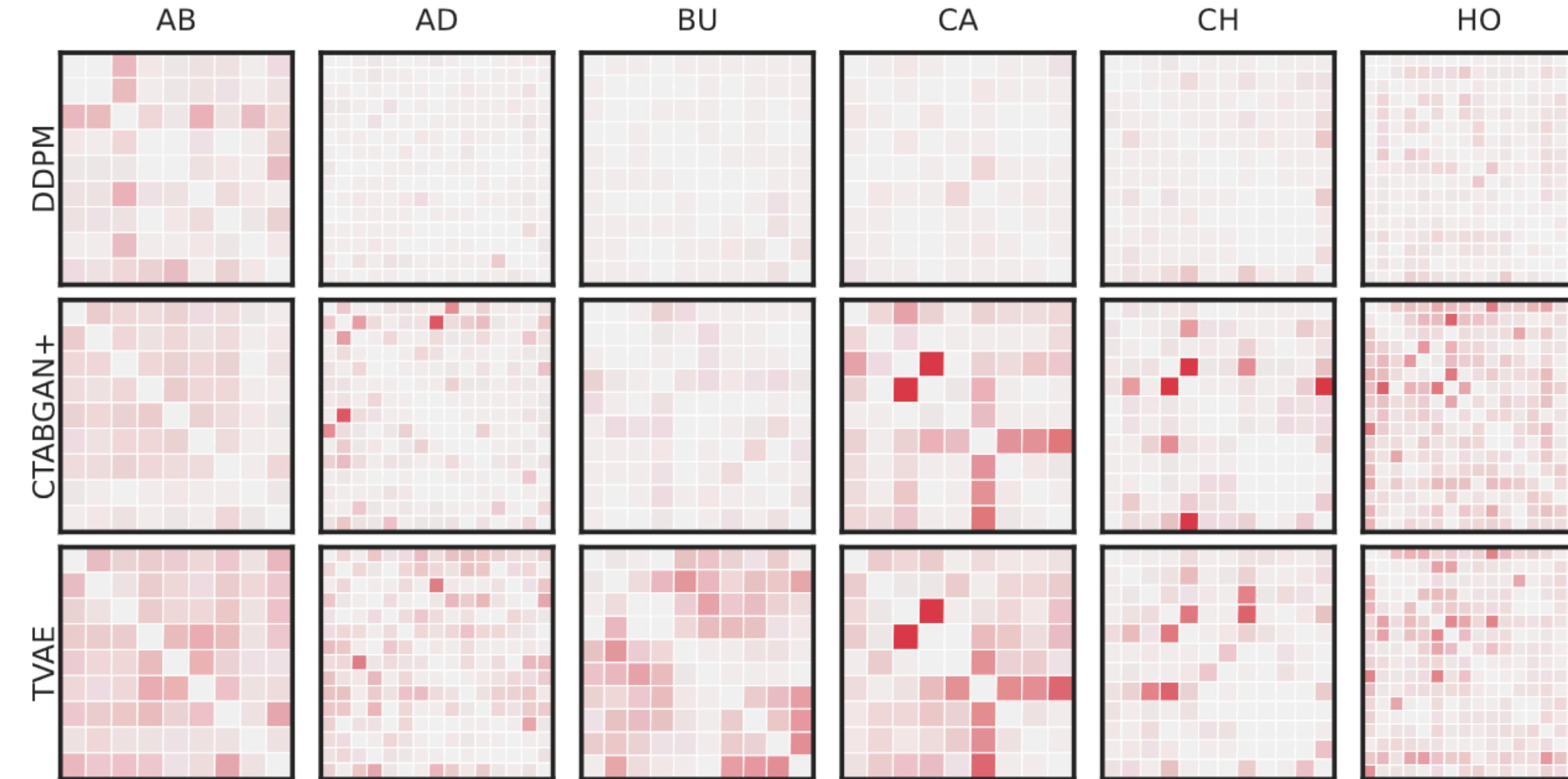


Figure 3: The absolute difference between correlation matrices computed on real and synthetic datasets. A more intensive red colour indicates a higher difference between the real and synthetic correlation values.

# Diffusion models 🥵 and tabular data generation

	AB ( <i>R</i> 2)	AD ( <i>F</i> 1)	BU ( <i>F</i> 1)	CA ( <i>R</i> 2)	CAR ( <i>F</i> 1)	CH ( <i>F</i> 1)	DE ( <i>F</i> 1)	DI ( <i>F</i> 1)
TVAE	0.238±.012	0.742±.001	0.779±.004	-13.0±1.51	0.693±.002	0.684±.003	0.643±.003	0.712±.010
CTABGAN	-	0.737±.007	0.786±.008	-	0.684±.003	0.636±.010	0.614±.007	0.655±.015
CTABGAN+	0.316±.024	0.730±.007	0.837±.006	-7.59±.645	<b>0.708±.002</b>	0.650±.008	0.648±.008	<b>0.727±.023</b>
SMOTE	<b>0.400±.009</b>	0.750±.004	0.842±.003	0.667±.006	0.693±.001	0.690±.003	0.649±.003	0.677±.013
TabDDPM	<b>0.392±.009</b>	<b>0.758±.005</b>	<b>0.851±.003</b>	<b>0.695±.002</b>	0.696±.001	<b>0.693±.003</b>	<b>0.659±.003</b>	0.675±.011
Real	0.423±.009	0.750±.006	0.845±.004	0.663±.002	0.683±.002	0.679±.003	0.648±.003	0.699±.012

	FB ( <i>R</i> 2)	GE ( <i>F</i> 1)	HI ( <i>F</i> 1)	HO ( <i>R</i> 2)	IN ( <i>R</i> 2)	KI ( <i>R</i> 2)	MI ( <i>F</i> 1)	WI ( <i>F</i> 1)
TVAE	<< 0	0.372±.006	0.590±.004	0.174±.012	0.470±.024	0.666±.006	<b>0.880±.002</b>	0.497±.001
CTABGAN	-	0.339±.009	0.539±.006	-	-	-	0.856±.003	0.656±.011
CTABGAN+	<< 0	0.373±.009	0.598±.004	0.222±.042	0.669±.018	0.197±.051	0.867±.002	0.653±.027
SMOTE	<b>0.651±.002</b>	<b>0.478±.005</b>	0.664±.003	0.394±.006	0.709±.008	<b>0.751±.005</b>	0.860±.001	<b>0.793±.004</b>
TabDDPM	0.527±.005	0.462±.005	<b>0.670±.002</b>	<b>0.426±.007</b>	<b>0.734±.007</b>	0.611±.013	0.850±.004	<b>0.792±.004</b>
Real	0.645±.005	0.431±.005	0.663±.002	0.415±.007	0.708±.007	0.768±.013	0.850±.004	0.684±.004

Table 3: The values of machine learning efficiency computed with regards to five weak classification/regression models. Negative scores denote negative R2, which means that performance is worse than an optimal constant prediction.

# Diffusion models 🥵 and tabular data generation

	AB (R2)	AD (F1)	BU (F1)	CA (R2)	CAR (F1)	CH (F1)	DE (F1)	DI (F1)
TVAE	0.433±.008	0.781±.002	0.864±.005	0.752±.001	0.717±.001	0.732±.006	0.656±.007	<b>0.714±.039</b>
CTABGAN	–	0.783±.002	0.855±.005	–	0.717±.001	0.688±.006	0.644±.011	<b>0.731±.022</b>
CTABGAN+	0.467±.004	0.772±.003	0.884±.005	0.525±.004	0.733±.001	0.702±.012	0.686±.004	<b>0.734±.020</b>
SMOTE	<b>0.549±.005</b>	0.791±.002	0.891±.003	<b>0.840±.001</b>	0.732±.001	0.743±.005	<b>0.693±.003</b>	0.683±.037
TabDDPM	<b>0.550±.010</b>	<b>0.795±.001</b>	<b>0.906±.003</b>	0.836±.002	<b>0.737±.001</b>	<b>0.755±.006</b>	<b>0.691±.004</b>	<b>0.740±.020</b>
Real	0.556±.004	0.815±.002	0.906±.002	0.857±.001	0.738±.001	0.740±.009	0.688±.003	0.785±.013

	FB (R2)	GE (F1)	HI (F1)	HO (R2)	IN (R2)	KI (R2)	MI (F1)	WI (F1)
TVAE	0.685±.003	0.434±.006	0.638±.003	0.493±.006	0.784±.010	0.824±.003	0.912±.001	0.501±.012
CTABGAN	–	0.392±.006	0.575±.004	–	–	–	0.889±.002	<b>0.906±.019</b>
CTABGAN+	0.509±.011	0.406±.009	0.664±.002	0.504±.005	0.797±.005	0.444±.014	0.892±.002	0.798±.021
SMOTE	<b>0.803±.002</b>	<b>0.658±.007</b>	<b>0.722±.001</b>	0.662±.004	<b>0.812±.002</b>	<b>0.842±.004</b>	0.932±.001	<b>0.913±.007</b>
TabDDPM	0.713±.002	0.597±.006	<b>0.722±.001</b>	<b>0.677±.010</b>	0.809±.002	<b>0.833±.014</b>	<b>0.936±.001</b>	<b>0.904±.009</b>
Real	0.837±.001	0.636±.007	0.724±.001	0.662±.003	0.814±.001	0.907±.002	0.934±.000	0.898±.006

Table 4: The values of machine learning efficiency computed with regards to the state-of-the-art tuned CatBoost model.

# The end

Let's talk

(If you think you might like doing research, apply for ML residency here:

<https://yandex.com/jobs/vacancies/yandex-research-ml-residency-1740>)

my contacts: @puhsuuu everywhere

# Papers mentioned

- <https://arxiv.org/abs/2106.11959> - Revisiting models
- <https://arxiv.org/abs/2203.05556> - Numerical feature embeddings
- <https://arxiv.org/abs/2207.03208> - Pre-training
- <https://arxiv.org/abs/2209.15421> - Generating tabular data with TabDDPM