# TabCaps: A Capsule Neural Network for Tabular Data Classification with BoW Routing

Jintai Chen , Kuanlun Liao, Yanwen Fang, Danny Z. Chen, Jian Wu*

# Background

## Tabular Data

columns = attributes for those observations

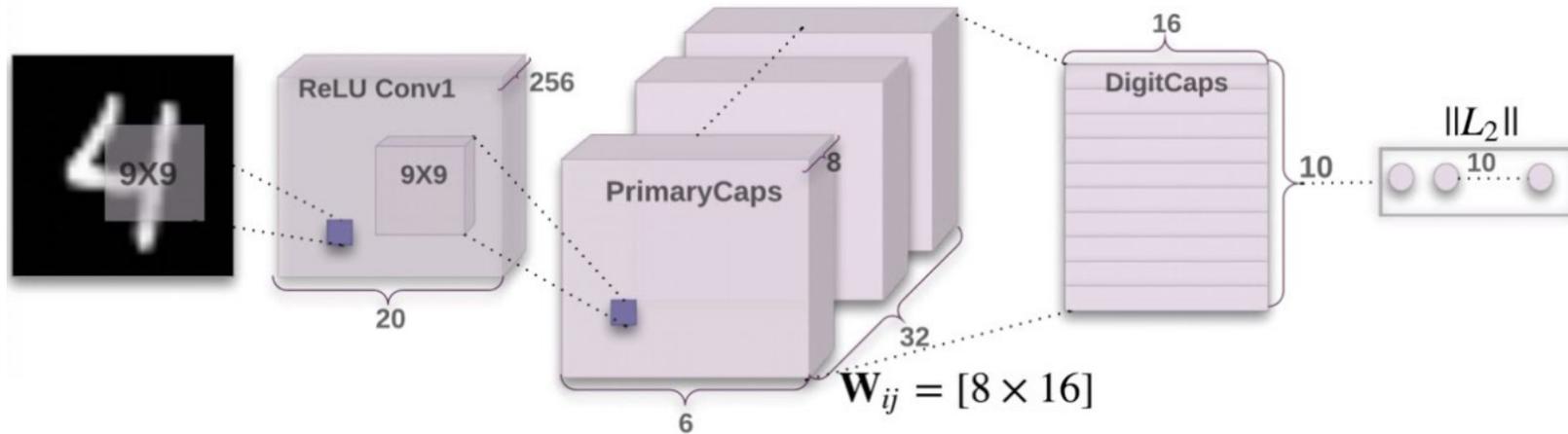| Player | Minutes | Points | Rebounds | Assists |
|--------|---------|--------|----------|---------|
| A | 41 | 20 | 6 | 5 |
| B | 30 | 29 | 7 | 6 |
| C | 22 | 7 | 7 | 2 |
| D | 26 | 3 | 3 | 9 |
| E | 20 | 19 | 8 | 0 |
| F | 9 | 6 | 14 | 14 |
| G | 14 | 22 | 8 | 3 |
| I | 22 | 36 | 0 | 9 |
| J | 34 | 8 | 1 | 3 |

Rows = observations

**Tabular data are represented by heterogeneous scalar features.**

**These features are aligned but their relations are unknown.**

**Mining interactions between heterogeneous features requires a higher sample complexity.**
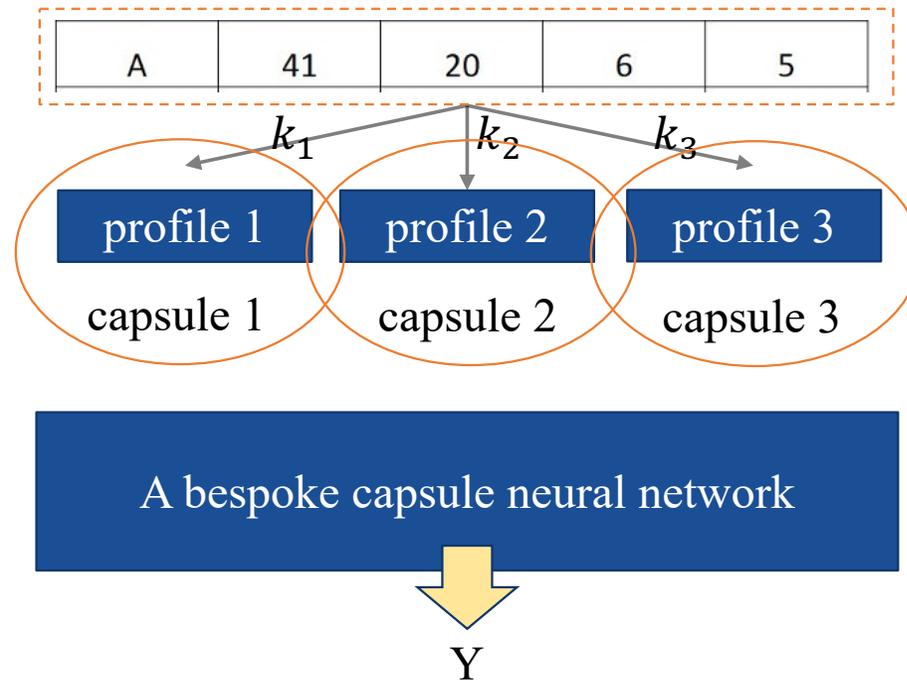
# Background



Capsule Neural Networks uses **"capsules"** to package scalar features as units.

The capsule features represent more concrete semantics.

Since mining feature relations is complex on tabular data, how about packaging them together and **conducting no interactions?**

# IDEA: Use Capsule and Conduct NO Feature Interactions



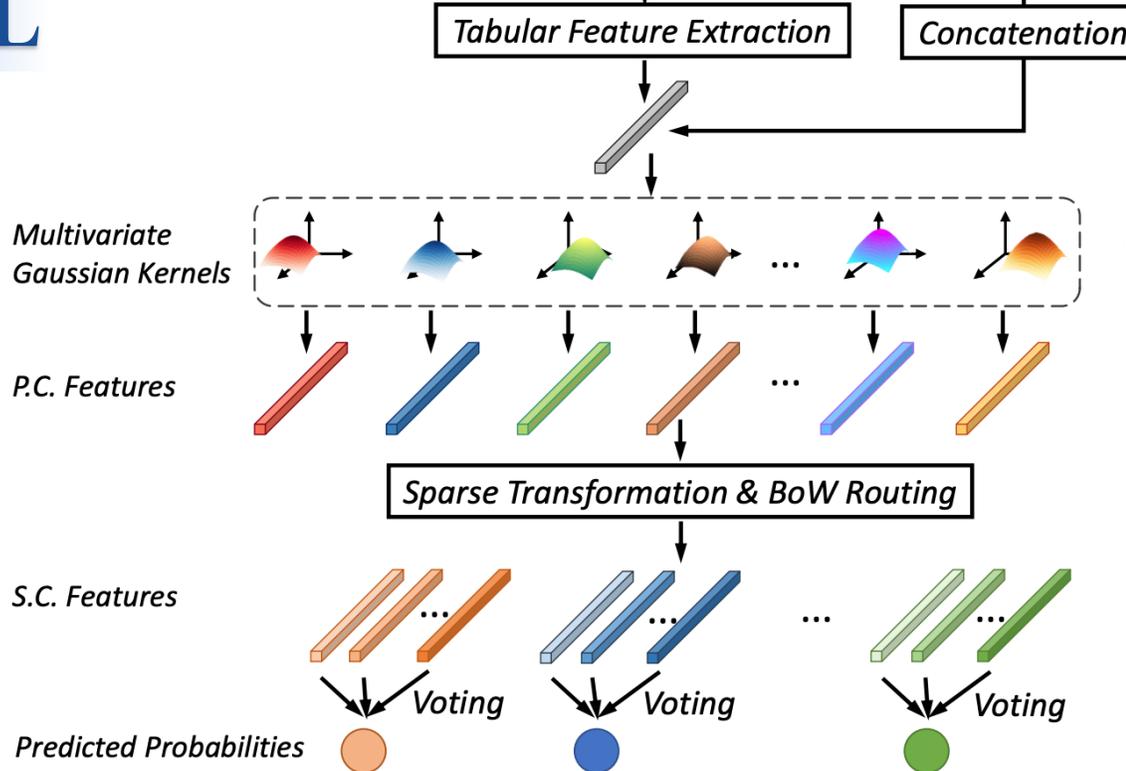Each function **k** **feature-wisely transforms a sample into a vector.**

Each capsule learns a **profile** (the vector) of the sample.

Conduct no feature interactions and directly learning the semantics at data level, so we call it **Data-Level Learning.**

# MODEL

| | Pregnancies | Glucose | Blood Pressure | Skin Thickness | Insulin | BMI | Diabetes Pedigree | Age |
|---|---|---|---|---|---|---|---|---|
| Raw Features | Yes | 138 | 62 | 35 | 125 | 33.6 | 0.127 | 47 |

**Tabular Feature Extraction**   **Concatenation**

*Multivariate Gaussian Kernels*

*P.C. Features*

**Sparse Transformation & BoW Routing**

*S.C. Features*

Voting   Voting   Voting

*Predicted Probabilities*

① Add some features to obtain more features

② $u_i = k_i(x; \mu_i, \Sigma_i) = \dfrac{\exp(-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1}(x - \mu_i))}{(2\pi)^{m/2}|\Sigma_i|}$

③ $\hat{u}_{j|i} = u_i W'_{j|i}, \quad W'_{j|i} = entmax_\alpha(W_{j|i}),$

$u_j = Routing(\hat{u}_{j|i}).$

④ $l_k = \sum_{j \in G_k} ||v_j||_2 / ||G_k||$

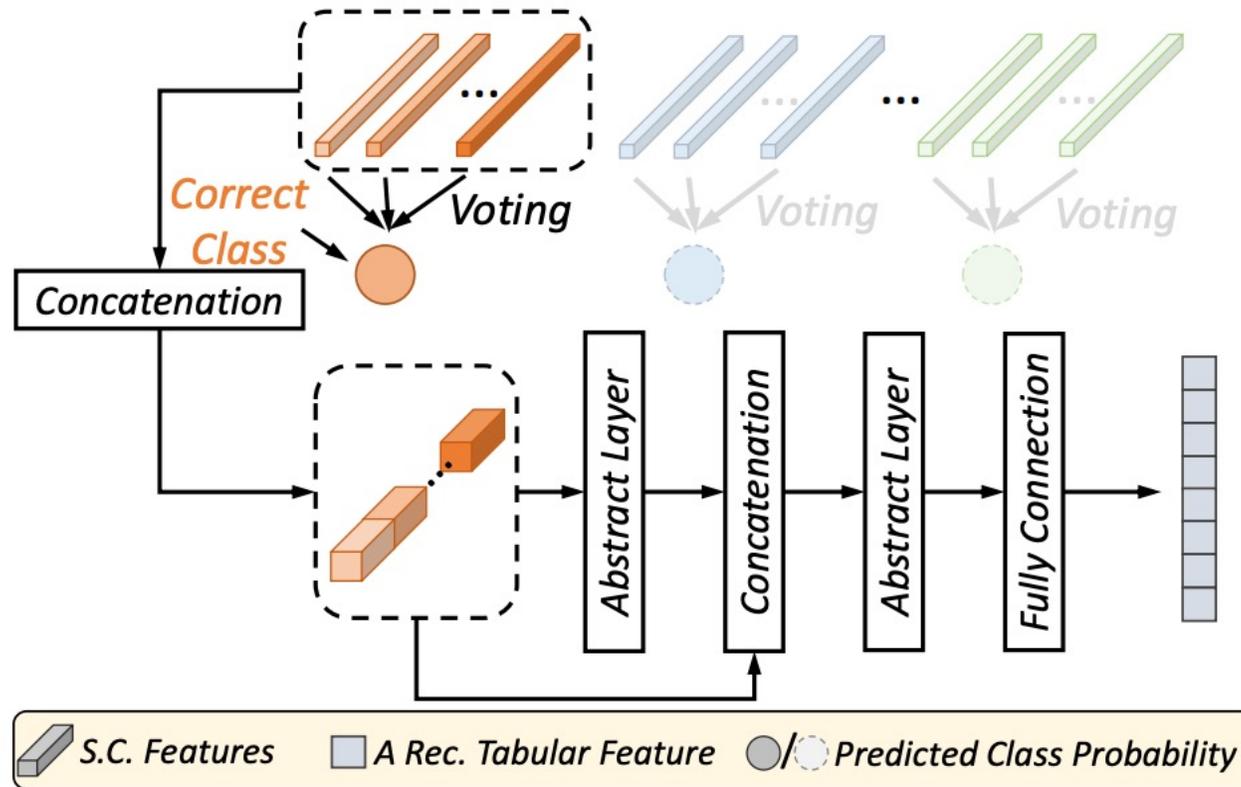20% of $v_j$ are dropped out in training

① **feature extension by Abstract Layer or MLP (automatic feature engineering)**

② **Gaussian kernels as function $k$**

③ **Transformation and Routing for selective capsule-feature-fusion**

④ **voting for prediction results (TIPS: dropout is helpful!)**

**The corresponding capsule features of correct class are stacked for reconstruction.**
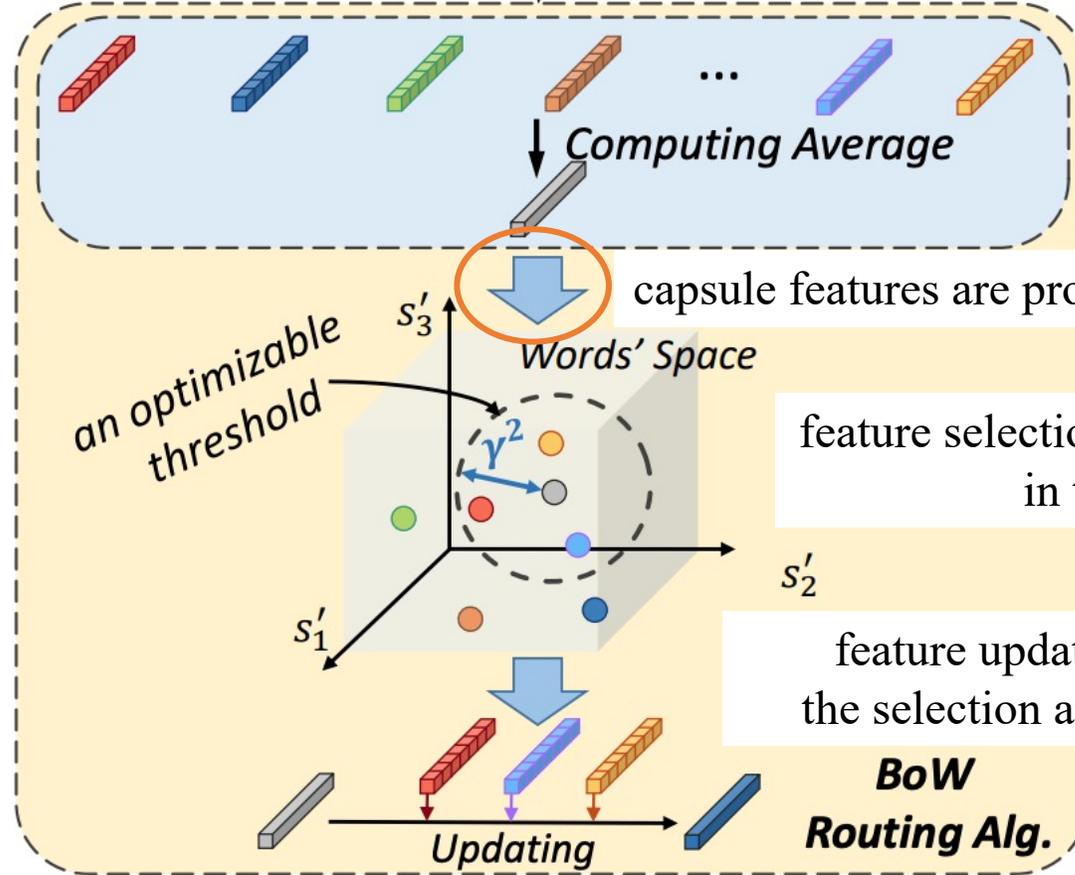
# BoW-Routing



P.C. Votes $\hat{u}_{j|\cdot}$

Initialized S.C. Features $v'_j$

Computing Average

capsule features are projected to a new feature space

an optimizable threshold

Words' Space

$s'_3$

$\gamma^2$

$s'_1$

$s'_2$

feature selection and similarity computing in the feature space

P.C. Votes Selecting

feature update according to the selection and the similarity

Features Updating of Final S.C.

Updating

BoW Routing Alg.

**Why the BoW Routing requires no iterations?**

**Previous CapsNets for images capture some unknown object parts in initializing capsule features and need routing-by-agreement. However, our data-level learning learns concrete profiles of the entire data and thus we believe that our routing does not need agreement.**
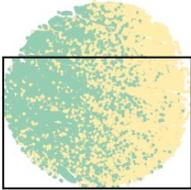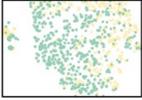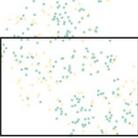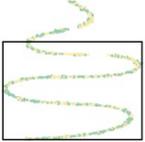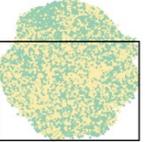
# Experiments

Table 1: **Classification Performances.** The best and second best performances of deep learning approaches are respectively marked in **bold** and underlined. Note that the reported log-loss values (the lower the better) are with a $100\times$ factor. The model size (# param.) and inference speed (fps) are on the *Diabetes* dataset. The performances are reported as "mean$\pm$std".

| Method | Click | Diabetes | EEG | Gas | Heart | Hill | Higgs | Epsilon | # param. | fps |
|---|---|---|---|---|---|---|---|---|---|---|
| XGboost | $62.253_{\pm0.02}$ | $14.338_{\pm0.03}$ | $14.117_{\pm0.02}$ | $2.087_{\pm0.06}$ | $32.371_{\pm0.04}$ | $69.049_{\pm1e\text{-}3}$ | $53.158_{\pm0.01}$ | $26.748_{\pm1e\text{-}3}$ | – | – |
| Catboost | $64.273_{\pm0.08}$ | $14.777_{\pm0.07}$ | $18.423_{\pm0.12}$ | $2.064_{\pm0.05}$ | $30.043_{\pm0.14}$ | $69.174_{\pm0.06}$ | $53.273_{\pm0.05}$ | $27.228_{\pm2e\text{-}3}$ | – | – |
| TabNet | $\underline{62.303}_{\pm0.03}$ | $17.964_{\pm0.04}$ | $45.340_{\pm0.04}$ | $4.647_{\pm0.04}$ | $44.967_{\pm0.01}$ | $87.804_{\pm0.07}$ | $54.668_{\pm0.03}$ | $26.743_{\pm0.02}$ | 3.4M | 73.1 |
| Net-DNF | $67.633_{\pm0.02}$ | $13.767_{\pm0.02}$ | $17.386_{\pm0.01}$ | $\mathbf{1.229}_{\pm0.04}$ | $55.371_{\pm0.03}$ | $\underline{15.787}_{\pm0.03}$ | $53.417_{\pm0.02}$ | $27.122_{\pm0.03}$ | 8.5M | 175.2 |
| NODE | $63.206_{\pm0.05}$ | $45.951_{\pm0.03}$ | $47.654_{\pm0.04}$ | $38.774_{\pm0.04}$ | $46.541_{\pm0.04}$ | $69.220_{\pm0.04}$ | $61.864_{\pm0.08}$ | $27.838_{\pm0.54}$ | 13.4M | 145.2 |
| FT-Transformer | $70.487_{\pm0.02}$ | $\underline{12.382}_{\pm0.03}$ | $\mathbf{7.446}_{\pm0.06}$ | $2.258_{\pm0.04}$ | $\mathbf{27.547}_{\pm0.05}$ | $20.084_{\pm0.03}$ | $\underline{53.310}_{\pm0.02}$ | $\underline{25.958}_{\pm0.85}$ | 9.3M | 284.7 |
| DANet-24 | $73.708_{\pm0.02}$ | $13.338_{\pm0.02}$ | $9.301_{\pm0.04}$ | $2.171_{\pm0.02}$ | $49.643_{\pm0.04}$ | $24.763_{\pm0.03}$ | $\mathbf{53.033}_{\pm0.01}$ | $26.431_{\pm0.01}$ | 5.5M | 54.9 |
| FCNN w/ mixup | $63.863_{\pm0.07}$ | $12.715_{\pm0.05}$ | $9.572_{\pm0.07}$ | $2.083_{\pm0.06}$ | $36.742_{\pm0.02}$ | $56.005_{\pm0.05}$ | $56.787_{\pm0.04}$ | $27.467_{\pm0.03}$ | 0.7M | 594.3 |
| FCNN w/ lasso | $87.005_{\pm0.17}$ | $41.071_{\pm0.75}$ | $31.852_{\pm0.05}$ | $4.141_{\pm0.06}$ | $44.881_{\pm0.06}$ | $69.302_{\pm0.01}$ | $132.102_{\pm0.07}$ | $32.282_{\pm0.02}$ | 0.7M | 568.8 |
| Vector CapsNet | $64.135_{\pm0.05}$ | $52.635_{\pm0.03}$ | $53.587_{\pm0.06}$ | $161.547_{\pm0.03}$ | $58.516_{\pm0.04}$ | $51.591_{\pm0.02}$ | $62.654_{\pm0.02}$ | $54.252_{\pm0.02}$ | 0.4M | 318.5 |
| TABCAPS (Ours) | $\mathbf{62.054}_{\pm0.04}$ | $\mathbf{12.043}_{\pm0.03}$ | $\underline{8.130}_{\pm0.05}$ | $\underline{2.013}_{\pm0.03}$ | $\underline{34.047}_{\pm0.02}$ | $\mathbf{14.301}_{\pm0.04}$ | $53.776_{\pm0.03}$ | $\mathbf{25.821}_{\pm0.02}$ | 0.2M | 501.1 |

**The performances are competitive to or even better than other approach that conducts complex feature interactions.**
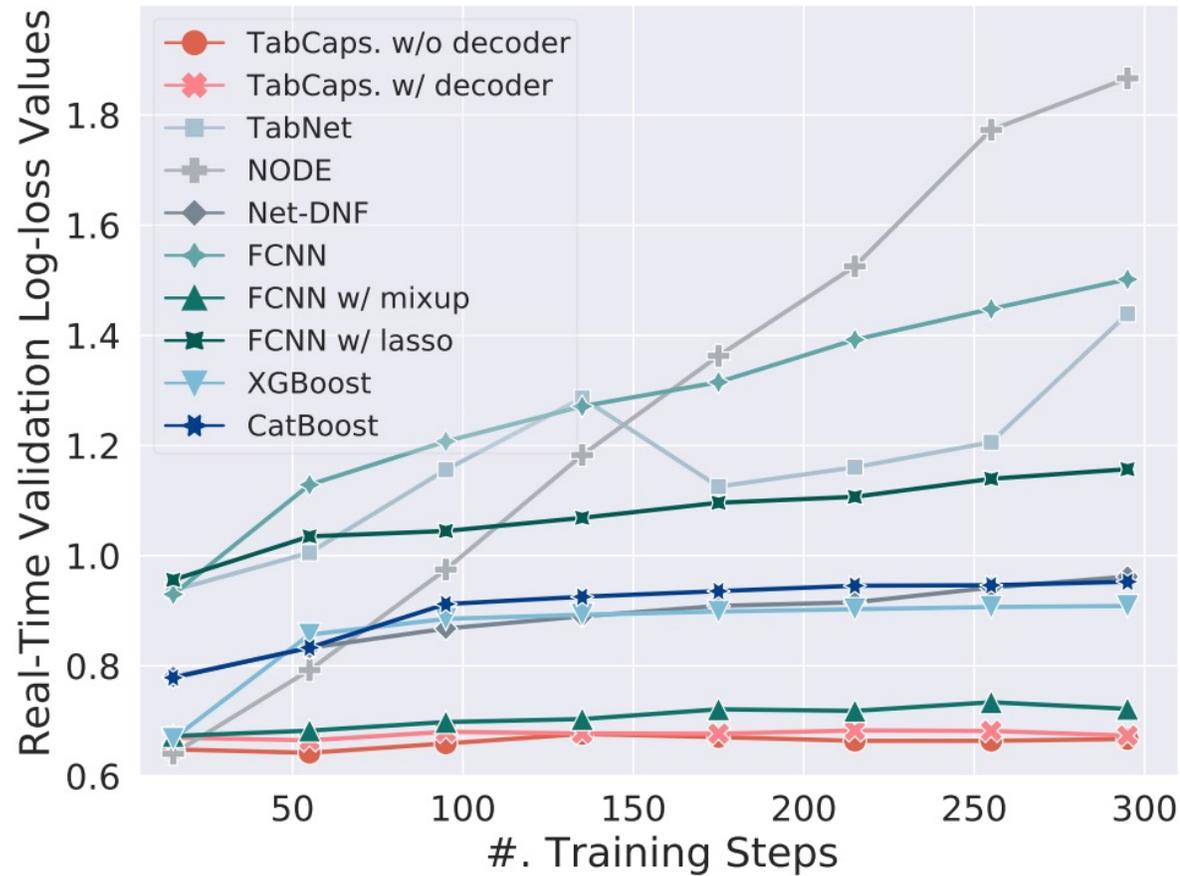
# Experiments

Table 2: **Extreme generalization performances.** The best and second best performances of deep learning approaches are respectively marked in **bold** and underlined. Note that the reported log-loss values (the lower the better) are with a $100\times$ factor. The *Epsilon* dataset is not included due to its extremely high computation complexity in conducting t-SNE projection.

| Method | Click | Diabetes | EEG | Gas | Heart | Hill | Higgs |
|---|---|---|---|---|---|---|---|
| Training-Test Split | | | | | | | |
| XGboost | $66.070_{\pm0.03}$ | $65.886_{\pm0.09}$ | $70.654_{\pm0.02}$ | $31.504_{\pm0.04}$ | $35.650_{\pm0.01}$ | $69.657_{\pm0.09}$ | $54.557_{\pm0.04}$ |
| Catboost | $63.925_{\pm0.04}$ | $68.819_{\pm0.06}$ | $68.799_{\pm0.04}$ | $18.864_{\pm0.04}$ | $35.207_{\pm0.08}$ | $69.162_{\pm0.03}$ | $54.632_{\pm0.07}$ |
| TabNet | $115.907_{\pm0.11}$ | $225.22_{\pm0.08}$ | $79.666_{\pm0.07}$ | $158.618_{\pm0.03}$ | $44.967_{\pm0.06}$ | $89.114_{\pm0.08}$ | $55.763_{\pm0.11}$ |
| Net-DNF | $67.625_{\pm0.02}$ | $\underline{58.792}_{\pm0.05}$ | $68.261_{\pm0.04}$ | $15.124_{\pm0.03}$ | $55.371_{\pm0.07}$ | $48.301_{\pm0.04}$ | $55.738_{\pm0.06}$ |
| NODE | $\underline{63.839}_{\pm0.04}$ | $67.021_{\pm0.04}$ | $68.357_{\pm0.04}$ | $57.698_{\pm0.06}$ | $46.541_{\pm0.03}$ | $69.771_{\pm0.10}$ | $61.870_{\pm0.03}$ |
| FT-Transformer | $78.431_{\pm0.11}$ | $59.283_{\pm0.04}$ | $68.278_{\pm0.07}$ | $\mathbf{6.416}_{\pm0.06}$ | $\mathbf{26.132}_{\pm0.05}$ | $66.972_{\pm0.05}$ | $\mathbf{53.970}_{\pm0.10}$ |
| DANet-24 | $74.401_{\pm0.02}$ | $59.736_{\pm0.06}$ | $69.021_{\pm0.03}$ | $10.395_{\pm0.01}$ | $49.643_{\pm0.02}$ | $\underline{37.976}_{\pm0.04}$ | $\underline{54.182}_{\pm0.01}$ |
| FCNN mixup | $66.052_{\pm0.05}$ | $60.262_{\pm0.04}$ | $68.850_{\pm0.08}$ | $25.102_{\pm0.03}$ | $35.674_{\pm0.17}$ | $67.126_{\pm1e\text{-}3}$ | $55.847_{\pm0.01}$ |
| FCNN lasso | $106.123_{\pm3e\text{-}3}$ | $67.082_{\pm0.04}$ | $93.170_{\pm0.04}$ | $61.310_{\pm0.02}$ | $76.854_{\pm0.03}$ | $75.853_{\pm0.02}$ | $106.580_{\pm0.06}$ |
| Vector CapsNet | $64.724_{\pm0.05}$ | $66.009_{\pm0.02}$ | $\underline{67.845}_{\pm0.04}$ | $163.193_{\pm0.04}$ | $60.848_{\pm0.04}$ | $64.743_{\pm0.09}$ | $62.791_{\pm0.02}$ |
| TABCAPS (Ours) | $\mathbf{63.355}_{\pm0.04}$ | $\mathbf{58.409}_{\pm0.02}$ | $\mathbf{67.471}_{\pm0.01}$ | $\underline{8.750}_{\pm0.06}$ | $\underline{34.503}_{\pm0.05}$ | $\mathbf{17.887}_{\pm0.04}$ | $54.707_{\pm0.07}$ |

**We biasedly split train and test sets to inspect the generalization capability.**

**TabCaps performs well!**

# Experiments



We observe that overfitting often occurs on the Click data.

We demonstrate the model's ability to resist overfitting through comparison.

THANKS

Thank you for listening!