

TabLLM: Few-shot Classification of Tabular Data with Large Language Models

Hegselmann et al., AISTATS 2023

Presenter: *Haeyoung Lee*

December 02, 2025

Seoul National University

Motivation

- Real-world records are predominantly stored as *tabular* data across domains (e.g., healthcare, climate, finance), yet obtaining labels for supervised classification is difficult.
- Unlike in computer vision and NLP, deep learning has not consistently outperformed strong *gradient-boosted tree* ensembles on tabular data, which typically have weak locality, mixed feature types, and relatively few columns.
- LLMs pretrained on massive text corpora exhibit strong few-shot generalization across diverse domains and can achieve competitive performance with little or no labeled training data by leveraging the knowledge encoded in their parameters.
- **Key question:** Can we leverage an LLM for tabular prediction?
- **TabLLM** prompts an LLM with (i) a natural-language serialization of each table row and (ii) a short task description; multiple serialization strategies are explored, and zero-shot and few-shot classification is evaluated.

Notation

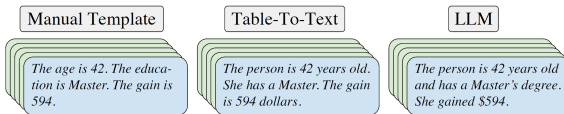
- Tabular dataset with n rows and d columns (features).
- x_i : a d -dimensional feature vector (the i -th row).
- $y_i \in C$: class label; C is the set of classes.
- $D = \{(x_i, y_i)\}_{i=1}^n$: dataset of row-label pairs.
- $F = \{f_1, \dots, f_d\}$: feature/column names (e.g., “age”, “education”).
- k : number of labeled training examples.
- D_k : a subset of size k sampled from D (with replacement) used for fine-tuning.
- $\text{serialize}(F, x)$: a function that takes the column names F and feature values x for a row as inputs and creates a textual representation.
- p : task-specific prompt (short description/question appended after serialization).
- LLM: an large language model with vocabulary V .
- $\text{LLM}((\text{serialize}(F, x), p))$: the prompted output of the LLM.
- Verbalizer: manually specified mapping from LLM output tokens (answer choices) to class labels in C .

Overview

1. Tabular data with k labeled rows

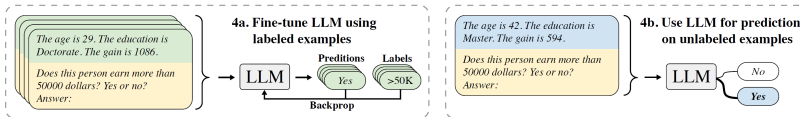
age	education	gain	income
39	Bachelor	2174	≤50K
36	HS-grad	0	>50K
64	12th	0	≤50K
29	Doctorate	1086	>50K
42	Master	594	

2. Serialize feature names and values into natural-language string with different methods



3. Add task-specific prompt

Does this person earn more than 50000 dollars? Yes or no? Answer:



• Steps:

1. Serialize a row (F, x) into a natural language string.
2. Append a task-specific prompt.
3. Obtain output probabilities from the LLM for verbalizer tokens (e.g., "Yes" / "No").
4. In few-shot setting, do parameter-efficient fine-tuning (T-Few) on k labeled examples.

From LLM Text to Class Label

Income Dataset:

```
answer_choices: 'No ||| Yes'
jinja: '{{serialization}}'

Does this person earn more than 50000
dollars per year? Yes or no?
Answer:
|||
{{ answer_choices[label] }}
```

Car Dataset:

```
answer_choices: 'Unacceptable |||
Acceptable ||| Good ||| Very good'
jinja: '{{serialization}}'

How would you rate the decision to buy
this car? Unacceptable, acceptable,
good or very good?
Answer:
|||
{{ answer_choices[label] }}
```

End Of Life Task:

```
answer_choices: 'No ||| Yes'
jinja: '{{serialization}}'

Does this patient die in the next nine
months? Yes or no?
Answer:
|||
{{ answer_choices[label] }}
```

- An LLM outputs a probability distribution over tokens, not class IDs.
- A **verbalizer** maps each class in C to a verbalizer token (e.g., $\leq 50K \mapsto \text{No}$, $> 50K \mapsto \text{Yes}$).
- **Inference:** compute the normalized likelihood of each verbalizer answer and choose the most probable class.

9 Serialization Strategies

Income Dataset (List Template):

```
- Age: 30
- Race: Asian-Pac-Islander
- Sex: Female
- Marital status: never married
- Rel. to head of the household: Own
- Native country: Taiwan
- Occupation: execution and management
- Work class: private sector employee
- Capital gain last year: 0
- Capital loss last year: 0
- Education: bachelor's degree
- Work hours per week: 52
```

Income Dataset (Text Template):

The Age is 30. The Race is Asian-Pac-Islander. The Sex is Female. The Marital status is never married. The Relation to head of the household is Own-child. The Native country is Taiwan. The Occupation is execution and management. The Work class is private sector employee. The Capital gain last year is 0. The Capital loss last year is 0. The Education is bachelor's degree. The Work hours per week is 52.

Income Dataset (Table-To-Text):

The age was 30 years, and was the youngest ever to enter the age. The race was held in the Asian-Pac-Islander, and was won by the race. The sex of the village was Female. The first female to be married is Marital status never reported. the family has the head of the household. The Chinese: native region of Taiwan. He was the executioners of the execution and management of the city of New York City. the private sector employee is a private sector employee. The capital was Capital of the State of India. The capital loss of the state was 0.5%. The bachelor's degree in Education was bachelor's degree. the week 52 was the 52-hour week.

Income Dataset (Text T0):

Kim is a 30-year-old Asian-Pacific Islander. She is never married and has never had children. The man is the owner of the house and he is the only child. A woman is executing a contract as a private sector employee. The company had a capital loss of \$ 0 last year. The man has a bachelor's degree and works 52 hours a week.

Income Dataset (Text GPT-3):

The person is 30 years old, Asian-Pac-Islander, female, never married, and an own-child relation to the head of the household. The person is from Taiwan and is an execution and management occupation in the private sector employee work class. The person has 0 dollars in capital gain and 0 dollars in capital loss from the previous year. The person has a bachelor's degree and works 52 hours per week.

- **Templates (manual):**

- ▶ List Template: "feature: value, ..."
- ▶ Text Template: "The feature is value., ..."

- **LLM-based generation:**

- ▶ Table-to-Text: an LLM fine-tuned on a table-to-text generation task. Feed each (column, value) tuple separately and concatenate the generated outputs.
- ▶ Text T0: use T0 with 11B params, split a row into *pairs* of two (column, value) tuples, prompt each pair with "Write this information as a sentence:", then combine the resulting sentences.
- ▶ Text GPT-3: use GPT-3 via API; provide the *full list* of features at once with "Rewrite all list items in the input as a natural text."

9 Serialization Strategies

- **Ablations (based on the List Template):**
 - ▶ **List Only Values:** *remove column names* (values only) to test whether feature/column names contribute to performance.
 - ▶ **List Permuted Names:** *permute the column names* so that each value is paired with an incorrect name. (e.g., sex:30, age:male)
Goal: test how important the correct name–value association is.
 - ▶ **List Permuted Values:** *permute values* within each column using a fixed, column-specific mapping applied to all examples.
Categorical ex.: Bachelor→PhD, Master→12th, ...
Continuous ex.: bin into 10 uniform ranges, then permute bins (e.g., age 20–30→70–80).
Goal: test whether the LLM relies on fine-grained value information beyond column names.
 - ▶ **List Short:** keep *at most 10 features*; used for the healthcare dataset to satisfy the LLM input-length limit and test the effect of reduced information.

Experimental Setup: LLM and Fine-tuning

- **LLM: T0 (11B params)**; a T5-style encoder–decoder model that is further **instruction-tuned on a large variety of task-specific prompts**, which yields strong **zero-shot generalization** across diverse tasks.
 - ▶ Input budget: **1024 tokens** (roughly 400 words).
- **Few-shot fine-tuning**: Fine-tune on the k -shot set using the parameter-efficient **T-Few** recipe.
 - ▶ Public tabular datasets: **30 epochs** for all few-shot runs.
 - ▶ Healthcare dataset: fewer epochs for larger k (e.g., 10 epochs up to 256 shots; 3 epochs for 1,024+ shots) to reduce runtime and overfitting.
- **Zero-shot baseline**: **GPT-3** with no fine-tuning.

Experimental Setup: Datasets

- **9 Public tabular benchmarks :**

- ▶ Requirements: $\leq 50,000$ rows, ≤ 30 columns for cost and token budget of T0.
- ▶ Require *textual* feature names; exclude datasets with *derived* features (e.g., mean pixel values).
- ▶ Included datasets (rows, features):
Bank (45,211, 16), Blood (748, 4), California (20,640, 8), Car (1,728, 8), Credit-g (1,000, 20), Income (48,842, 14), Jungle (44,819, 6), Diabetes (768, 8), Heart (918, 11).

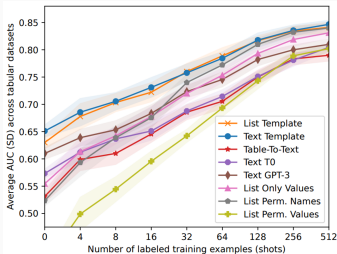
- **Healthcare claims:**

- ▶ 3 binary tasks: End-of-Life (EoL), Surgery, Likelihood of Hospitalization (LoH).

Baselines

- **Classical baselines:**
 - ▶ Logistic Regression (LR)
 - ▶ Gradient-boosted trees: XGBoost, LightGBM
- **Deep tabular baselines:**
 - ▶ TabNet (attention over columns), SAINT (attention over rows & columns), NODE (differentiable tree ensemble)
 - ▶ TabPFN (pretrained Bayesian NN on synthetic tabular data)
- **Details:**
 - ▶ Hyperparameter tuning for all baselines *except* TabPFN.
 - ▶ Few-shot setting: no separate validation set \Rightarrow **4-fold CV** on the k shots for model selection.
 - ▶ Categorical features: **one-hot encoding** (ordinal encoding tested but worse).
- **Healthcare claims:** only LR + LightGBM (runtime constraints).
- **Metric:** Area under the curve (AUC).

Results: Effects of Serialization



Car Dataset (Text Template):

The Buying price is low. The Doors is three. The Maintenance costs is low. The Persons is more than four. The Safety score is medium. The Trunk size is medium.

Car Dataset (Text GPT-3):

This car a good choice for those who are looking for a low-priced vehicle with low maintenance costs. It is also a good choice for families or groups of friends who need a car with a bit more space than a smaller car. The safety score is medium, so it is not the best choice for those who are looking for a car with the highest safety rating.

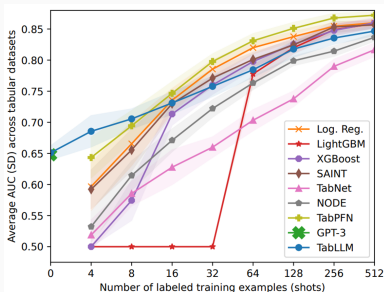
- **Text Template** performs best overall.
- In the **zero-shot** setting, **Text Template** outperforms the List Template,
 - ▶ likely because its natural-sentence format is closer to T0's training distribution.
 - ▶ This gap largely vanishes by **~8 shots**, suggesting that sophisticated serializations might be unnecessary when some training data exists.
- **LLM-based serialization is less reliable:**
 - ▶ generated text may **hallucinate** additional content or **omit** features, which can introduce misleading cues and lead to lower performance than fixed templates.
- **Ablations converge with more data:**
 - ▶ Only Values and Permuted Names perform poorly at **zero / very few shots**, but **match the main templates** once enough training examples are provided.

Results: Effects of Serialization

Method	Number of Shots							
	0	16	64	256	1,024	4,096	16,384	all
End of Life (EoL)								
TabLLM (T0 + List Template)	0.70	0.74	0.78	0.78	0.79	0.81	0.81	—
TabLLM (T0 + Text Template)	0.63	0.71	0.74	0.76	0.78	0.79	0.80	—
TabLLM (T0 + List Short)	0.68	0.71	0.76	0.79	0.80	0.81	0.82	—
TabLLM (T0 + List Perm. Names)	0.62	0.66	0.70	0.74	0.75	0.77	0.79	—
Method	EoL			Surgery	LoH			
Age, sex, and race	0.59			0.57	0.65			
Least frequent conditions	0.57			0.64	0.67			
Least frequent procedures	0.59			0.59	0.65			
Least frequent concepts (cond. + proc.)	0.55			0.55	0.66			
Most frequent conditions	0.67			0.66	0.69			
Most frequent procedures	0.59			0.58	0.65			
Most frequent concepts (cond. + proc.)	0.62			0.61	0.65			
Oldest conditions	0.65			0.66	0.69			
Oldest procedures	0.59			0.58	0.65			
Oldest concepts (cond. + proc.)	0.60			0.60	0.67			
Most recent conditions	0.65			0.66	0.69			
Most recent procedures	0.55			0.59	0.65			
Most recent concepts (cond. + proc.)	0.59			0.60	0.66			

- **Healthcare claims: List Template** slightly outperforms Text Template. **List Short** is only slightly worse, suggesting robustness to incomplete feature sets.
- **Claims-specific finding:** selecting the **most frequent conditions per patient** works best.

Results: TabLLM vs. Baselines



Dataset	Method	Number of Shots									
		0	4	8	16	32	64	128	256	512	all
Bank	XGBoost	—	0.50 ₀₀	0.56 ₀₀	0.68 ₀₄	0.76 ₀₁	0.83₀₂	0.85 ₀₀	0.88 ₀₀	0.90 ₀₁	0.94₀₀
	TabPFN	—	0.51 ₀₁	0.66₀₀	0.69₀₁	0.76 ₀₁	0.82 ₀₁	0.86₀₂	0.89₀₀	0.90 ₀₀	0.93 ₀₀
	TabLLM	0.63₀₁	0.59 ₀₃	0.64 ₀₀	0.65 ₀₀	0.64 ₀₀	0.69 ₀₀	0.82 ₀₀	0.87 ₀₀	0.88 ₀₁	0.92 ₀₁
	TabLLM	—	0.50 ₀₀	0.58 ₀₀	0.66 ₀₁	0.67 ₀₀	0.68 ₀₀	0.71 ₀₀	0.70 ₀₀	0.67 ₀₀	0.71 ₀₀
Blood	XGBoost	—	0.52 ₀₀	0.64 ₀₀	0.67₀₀	0.70₀₁	0.73₀₄	0.75₀₄	0.76₀₄	0.76₀₁	0.74₀₁
	TabPFN	—	0.58₀₀	0.66₀₀	0.66 ₀₀	0.68 ₀₁	0.68 ₀₀	0.68 ₀₀	0.70 ₀₀	0.68 ₀₁	0.70 ₀₀
	TabLLM	—	0.50 ₀₀	0.62 ₀₁	0.74 ₀₁	0.79 ₀₁	0.82 ₀₄	0.87 ₀₄	0.90 ₀₀	0.92 ₀₁	0.97₀₀
	TabLLM	—	0.63 ₀₁	0.63₀₁	0.80₀₀	0.83₀₁	0.89₀₁	0.91₀₄	0.92₀₀	0.93₀₀	0.94 ₀₀
Calhousing	XGBoost	—	0.63 ₀₀	0.67 ₀₀	0.70 ₀₀	0.77 ₀₄	0.81 ₀₂	0.83 ₀₀	0.86 ₀₂	0.86 ₀₂	0.93 ₀₀
	TabPFN	—	0.50 ₀₀	0.59 ₀₀	0.70 ₀₀	0.82 ₀₁	0.91 ₀₂	0.95 ₀₄	0.96 ₀₀	0.99 ₀₁	1.00 ₀₀
	TabLLM	0.61₀₁	0.63 ₀₀	0.67 ₀₀	0.70 ₀₀	0.77 ₀₄	0.81 ₀₂	0.83 ₀₀	0.86 ₀₂	0.86 ₀₂	0.93 ₀₀
	TabLLM	—	0.50 ₀₀	0.59 ₀₀	0.70 ₀₀	0.82 ₀₁	0.91 ₀₂	0.95 ₀₄	0.96 ₀₀	0.99 ₀₁	1.00 ₀₀
Car	XGBoost	—	0.64 ₀₀	0.75 ₀₀	0.87₀₀	0.92₀₂	0.97₀₀	0.99₀₁	1.00₀₀	1.00 ₀₀	1.00 ₀₀
	TabPFN	—	0.64 ₀₀	0.75 ₀₀	0.87₀₀	0.92₀₂	0.97₀₀	0.99₀₁	1.00₀₀	1.00 ₀₀	1.00 ₀₀
	TabLLM	0.82₀₂	0.83₀₀	0.85₀₀	0.86 ₀₁	0.91 ₀₂	0.96 ₀₂	0.98 ₀₀	0.99 ₀₀	1.00 ₀₀	1.00 ₀₀
	TabLLM	—	0.50 ₀₀	0.51 ₀₀	0.59 ₀₀	0.66 ₀₁	0.67 ₀₀	0.68 ₀₀	0.73 ₀₀	0.75 ₀₁	0.78₀₄
Credit-g	XGBoost	—	0.58 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.72₀₀	0.75₀₄	0.75 ₀₁	0.75 ₀₁
	TabPFN	—	0.58 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.72₀₀	0.75₀₄	0.75 ₀₁	0.75 ₀₁
	TabLLM	0.53₀₁	0.69₀₄	0.66₀₀	0.66₀₀	0.72₀₀	0.70 ₀₀	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
	TabLLM	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
Diabetes	XGBoost	—	0.50 ₀₀	0.59 ₀₀	0.72₀₀	0.69 ₀₀	0.73 ₀₀	0.78 ₀₀	0.80 ₀₁	0.80 ₀₁	0.84₀₄
	TabPFN	—	0.61 ₀₁	0.67₀₁	0.71 ₀₀	0.77₀₁	0.82₀₀	0.83₀₀	0.83₀₀	0.81₀₂	0.81₀₁
	TabLLM	0.68₀₀	0.61 ₀₀	0.63 ₀₀	0.69 ₀₀	0.68 ₀₁	0.73 ₀₀	0.79 ₀₀	0.76 ₀₂	0.78 ₀₄	0.80 ₀₀
	TabLLM	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
Heart	XGBoost	—	0.50 ₀₀	0.55 ₀₀	0.84 ₀₀	0.88 ₀₄	0.91 ₀₁	0.91 ₀₄	0.90 ₀₀	0.92 ₀₁	0.94 ₀₁
	TabPFN	—	0.84₀₀	0.88₀₂	0.87 ₀₀	0.91₀₂	0.92₀₂	0.92₀₂	0.92 ₀₀	0.92 ₀₂	0.92 ₀₂
	TabLLM	0.54₀₁	0.76 ₀₁	0.83 ₀₀	0.87 ₀₀	0.91 ₀₁	0.91 ₀₄	0.90 ₀₄	0.92 ₀₀	0.92 ₀₁	0.94 ₀₁
	TabLLM	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
Income	XGBoost	—	0.73 ₀₀	0.71 ₀₀	0.76 ₀₀	0.80 ₀₄	0.82 ₀₄	0.84 ₀₄	0.86 ₀₀	0.87 ₀₁	0.89 ₀₀
	TabPFN	—	0.73 ₀₀	0.71 ₀₀	0.76 ₀₀	0.80 ₀₄	0.82 ₀₄	0.84 ₀₄	0.86 ₀₀	0.87 ₀₁	0.89 ₀₀
	TabLLM	0.84₀₀	0.84₀₁	0.84₀₂	0.84₀₄	0.84₀₁	0.84₀₂	0.86₀₀	0.87 ₀₀	0.89₀₁	0.92 ₀₀
	TabLLM	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
Jungle	XGBoost	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
	TabPFN	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁
	TabLLM	0.69₀₀	0.64 ₀₁	0.64 ₀₂	0.65 ₀₁	0.71 ₀₁	0.78 ₀₂	0.81 ₀₂	0.84 ₀₀	0.89 ₀₁	1.00₀₁
	TabLLM	—	0.50 ₀₀	0.59 ₀₀	0.64 ₀₀	0.69 ₀₁	0.70 ₀₁	0.71 ₀₀	0.72 ₀₁	0.72 ₀₁	0.75 ₀₁

- TabLLM performs best in **zero-shot / very-few-shot**.
 - ▶ **Zero-shot:** TabLLM achieves strong performance on most tasks and is **on par with GPT-3** despite T0 being much smaller (11B vs. 175B).
 - ▶ **Few-shot:** Performance improves with more shots; in the **very-few-shot** regime, TabLLM is often **substantially better** than most baselines.
- **TabPFN** is best overall; a small number of shots is often sufficient to catch up to TabLLM.
- LR is frequently the second-best baseline due to extensive parameter tuning.

- **Zero-shot capability & prior knowledge:** Simple template serializations (List/Text) yield **meaningful zero-shot** performance, indicating effective use of the LLM's prior knowledge.
- Sample efficiency varies by task: stronger on **semantically rich** datasets (e.g., INCOME) and weaker on **mostly numeric** datasets (e.g., BLOOD).
- TabLLM relies on the LLM's **pretrained semantics**; gains may be limited when column names/values are **out-of-domain** (e.g., gene identifiers).
 - ▶ T0's instruction-tuning contains no medical tasks, which may reduce gains on medical datasets.
- **Ethics & responsible use:** Since LLMs can inherit biases and stereotypes from historical training data, applying TabLLM to sensitive tasks (e.g., income or health) requires **careful interpretation**.

Thank you!