# Beginner-Friendly Guide to TabPFN(Prepared by Siddharth Yadav)

## 1. What is TabPFN?

TabPFN stands for Tabular Prior-data Fine-tuned Network. It is a transformer-based model developed by Microsoft Research for tabular data classification. Unlike traditional machine learning models, TabPFN does not require training on your dataset. Instead, it performs zero-shot learning using knowledge gained from millions of synthetic tasks it was pretrained on.

## 2. How TabPFN Works

TabPFN treats classification as a Bayesian inference problem. During its pretraining phase, it was exposed to millions of synthetic tasks. When provided with a new dataset, TabPFN uses its internal 'prior' to make probabilistic predictions for unlabeled data points based on the given small labeled subset. It performs this via a forward pass without any model training.

## 3. When to Use TabPFN

TabPFN is especially useful when:
- You have a small tabular dataset (≤128 training samples)
- You need fast predictions without training
- You want to test hypotheses quickly
- You are working with classification (not regression)
- You want state-of-the-art performance on small tabular tasks

## 4. Installing TabPFN

You can install TabPFN using pip:
pip install tabpfn

Or, clone the official repository for more control:
git clone https://github.com/automl/TabPFN
cd TabPFN
pip install -e .

## 5. Example Code to Use TabPFN

```
import torch
from tabpfn import TabPFNClassifier
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score
```

```
# Load dataset
X, y = load_iris(return_X_y=True)
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=128, test_size=10, stratify=y)

# Preprocess data
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Load TabPFN model
model = TabPFNClassifier(device='cuda' if torch.cuda.is_available() else 'cpu')

# Fit and predict
model.fit(X_train, y_train)
preds = model.predict(X_test)
print('Accuracy:', accuracy_score(y_test, preds))
```

## 6. Important Notes and Limitations

- TabPFN supports only classification problems (not regression).
- It works best for datasets with up to 128 training samples and 100 features.
- Categorical variables must be encoded before input.
- Running on CPU with more than 1000 samples is blocked by default for performance reasons.

## 7. Overcoming CPU Limitations

To use TabPFN with larger datasets on CPU, set the following environment variable before running your code:
export TABPFN_ALLOW_CPU_LARGE_DATASET=1

Or within Python:
import os
os.environ['TABPFN_ALLOW_CPU_LARGE_DATASET'] = '1'

## 8. Resources

- GitHub: https://github.com/automl/TabPFN
- Paper: https://arxiv.org/abs/2207.01848
- TabPFN Client API (for cloud inference): https://github.com/PriorLabs/tabpfn-client