

Selecting ML/DL Models



what is the best way to check the dataset fits on what ml/dl model? like pycaret maybe



When determining the best machine learning (ML) or deep learning (DL) model for a given dataset, there are several methods, tools, and frameworks you can use. PyCaret, as you mentioned, is a low-code ML library that automates much of the process of selecting, training, and tuning models. However, there are also other strategies and tools to help you choose the right model for your dataset.

Here are some approaches and tools for checking dataset compatibility with ML/DL models:

1. PyCaret (Low-Code AutoML)

- **Description:** PyCaret is a low-code machine learning library that automates the entire workflow. You can quickly load your dataset, and PyCaret will train multiple models, perform cross-validation, and rank them based on various metrics (accuracy, AUC, RMSE, etc.).
- **Best for:** Beginners or when you want a fast comparison across multiple models without deep customization.
- **How to use:**

python

```
from pycaret.classification import *
clf = setup(data, target='target_column')
best_model = compare_models()
```

2. H2O.ai (AutoML for ML and DL)

- **Description:** H2O.ai's AutoML framework provides automated model selection, hyperparameter tuning, and evaluation across a variety of models, including deep learning models. It supports ML models (Random Forest, GBM) and DL models (deep neural networks).
- **Best for:** Users who want to incorporate deep learning into automated model selection.
- **How to use:**

python

```
import h2o
from h2o.automl import H2OAutoML
h2o.init()

# Load dataset into H2O frame
data = h2o.import_file("dataset.csv")

# Train AutoML
aml = H2OAutoML(max_models=20, seed=1)
aml.train(y="target", training_frame=data)
```

```
# View leaderboard
lb = aml.leaderboard
lb.head()
```

3. MLJAR (AutoML for Machine Learning)

- **Description:** MLJAR is a simple AutoML library that provides automated training and comparison of multiple machine learning models. It automatically handles preprocessing, feature engineering, model tuning, and provides detailed reports.
- **Best for:** Quick experimentation with machine learning models for tabular data.
- **How to use:**

python

```
from supervised.automl import AutoML
automl = AutoML()
automl.fit(X_train, y_train)
```

4. Scikit-Learn Model Comparison (Manual Approach)

- **Description:** If you prefer a manual approach, you can use scikit-learn to train and compare different machine learning models. You will need to code the process of model training, validation, and comparison, but it allows for greater customization.
- **Best for:** Those who want more control over the modeling process or want to better understand each model's performance on their dataset.
- **How to use:**

python

```
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.svm import SVC

models = {
    'Random Forest': RandomForestClassifier(),
    'SVM': SVC(),
}

for name, model in models.items():
    scores = cross_val_score(model, X_train, y_train, cv=5)
    print(f"{name}: {scores.mean()}")
```

5. TensorFlow/Keras for Deep Learning Models

- **Description:** If you have a dataset suited for deep learning (images, text, sequences), you can use TensorFlow or Keras to build, train, and evaluate various deep learning architectures (CNNs, RNNs, etc.).
- **Best for:** Datasets where deep learning models like convolutional neural networks (CNNs) or recurrent neural networks (RNNs) are needed (e.g., image, speech, or text data).
- **How to use:**

python

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense

model = Sequential([
    Dense(128, activation='relu', input_shape=(input_dim,)),
    Dense(64, activation='relu'),
    Dense(1, activation='sigmoid')
])
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model.fit(X_train, y_train, epochs=10, validation_split=0.2)
```

6. TPOT (Genetic Programming for Model Selection)

- **Description:** TPOT is an AutoML tool that uses genetic programming to optimize machine learning pipelines. It automates model selection, hyperparameter tuning, and even preprocessing.
- **Best for:** When you want to optimize both model and pipeline construction.
- **How to use:**

python

```
from tpot import TPOTClassifier
tpot = TPOTClassifier(verbosity=2, generations=5)
tpot.fit(X_train, y_train)
```

7. Hyperopt or Optuna for Hyperparameter Tuning

- **Description:** Hyperopt and Optuna are libraries for hyperparameter optimization. While they don't automate model selection per se, they can help tune parameters for various models, which is a critical step for getting the best performance.
- **Best for:** Users who want to focus on hyperparameter tuning across different models.
- **How to use (Optuna):**

python

```
import optuna

def objective(trial):
    classifier = trial.suggest_categorical('classifier', ['SVM', 'RandomForest'])
    if classifier == 'SVM':
        model = SVC(C=trial.suggest_float('C', 0.1, 10))
    else:
        model = RandomForestClassifier(n_estimators=trial.suggest_int('n_estimators',
10, 100))

    score = cross_val_score(model, X_train, y_train, cv=5).mean()
    return score

study = optuna.create_study(direction='maximize')
study.optimize(objective, n_trials=100)
```

General Workflow for Choosing the Best Model:

1. **Dataset Analysis:** Start by analyzing your dataset to understand its nature (e.g., tabular, image, text). This determines the types of models you'll consider.
2. **Model Training/Comparison:**
 - For **tabular data**: Start with PyCaret, H2O.ai, or scikit-learn for traditional ML models (Random Forest, XGBoost, etc.).
 - For **image, text, or sequence data**: Consider deep learning frameworks like TensorFlow, Keras, or PyTorch.
3. **Model Performance Evaluation:** Use cross-validation to get reliable metrics (accuracy, AUC, F1 score, etc.) and compare model performance.
4. **Hyperparameter Tuning:** Use Hyperopt, Optuna, or built-in grid search (e.g., in PyCaret) to fine-tune model performance.

When to Use Which Approach:

- **PyCaret or H2O.ai:** If you want quick results and are working with tabular data.
- **Manual scikit-learn/TensorFlow:** If you need full control and customization.
- **AutoML (TPOT, MLJAR):** If you prefer an automated pipeline with some flexibility.
- **Hyperparameter tuning libraries (Optuna, Hyperopt):** If you need fine-tuning for a specific model after choosing it.

All these tools can make the process smoother, but your final choice of model should depend on the dataset's characteristics and your problem's complexity.



total merge dataset

1	50414234	ADI NUGRAHA YUNANDA	3.3	Software Engineer
2	50414279	ADITIA SEPTIANDI	3.1	Data Engineer
3	50414496	AGUS ANIN FALAHUDIN	3.25	Data Analyst
4	50414623	AHMAD ZULFIKAR	3.38	Software Engineer
5	50414625	ALZA ICHSAN	3.72	Mobile App
	NPM	Nama Lengkap	IPK	Klasifikasi Profes

heres my dataset



Your dataset contains 196 rows and 95 columns. It includes various features such as student names, grades in different subjects, and an important column labeled "IPK" (presumably representing a cumulative GPA or similar metric). The dataset seems to focus on academic performance, with many columns representing scores in various subjects or fields, which are numerical (integers).

Potential Steps to Fit the Dataset to an ML/DL Model:

1. Problem Definition:

- **Classification:** If you're trying to predict a discrete label (e.g., whether a student passes, classification of professions), it would be a classification problem.
- **Regression:** If you're trying to predict a continuous value (e.g., IPK or GPA), then it is a regression problem.

Based on the dataset structure, you could either:

- **Predict student performance (e.g., "IPK")** as a regression problem.
- **Predict a classification label** like "Klasifikasi Profesi" using student scores.

2. Feature Selection/Engineering:

- Remove irrelevant columns such as student names (e.g., `Nama Lengkap` and `NPM`), as these are not useful for prediction.
- Scale numeric columns if necessary (most algorithms like gradient boosting or neural networks may require scaling).

3. Choosing the Model:

- **For Regression** (predicting `IPK`): Models like Random Forest, XGBoost, or Neural Networks could work well.
- **For Classification** (predicting `Klasifikasi Profesi`): Logistic Regression, Random Forest, SVM, or Deep Learning models.

Would you like me to assist you with setting up an automated process to try different models on this dataset, using tools like PyCaret or Scikit-learn, or would you prefer to start by cleaning and preprocessing the dataset for model training?



that dataset is been cleaned before



Great! Since your dataset is already cleaned, we can directly proceed with exploring the appropriate models for your task. Based on the dataset, here are the next steps for model selection:

1. Decide on the Problem Type:

- **Regression:** If you are predicting a continuous value such as `IPK` (GPA), then you will need a regression model.
- **Classification:** If you are predicting a discrete category like `Klasifikasi Profesi`, you will need a classification model.

Let me know which type of problem you'd like to solve (regression or classification), and I can either:

- Use **PyCaret** to automatically compare different models, or

- Set up manual comparisons using **Scikit-learn** or other libraries.

Would you like to proceed with PyCaret for model comparison, or would you prefer another method?



my lecturer said "prediksi profesi pekerjaan mahasiswa menggunakan Wasserstein GAN"



It sounds like your lecturer has suggested using a **Wasserstein GAN (WGAN)** to predict the profession (``Klasifikasi Profesi``) of students. This is an interesting and more advanced approach, as GANs (Generative Adversarial Networks) are generally used for generative tasks but can also be adapted for structured prediction tasks.

Overview of Wasserstein GANs:

A **Wasserstein GAN (WGAN)** is a type of Generative Adversarial Network that aims to stabilize training and improve performance by using the Wasserstein distance (also known as Earth Mover's Distance) as a loss function. WGANs are typically used for generating high-quality samples from complex data distributions.

However, using WGANs for prediction tasks (like predicting professions in this case) might require some adaptation, as GANs are generally used for unsupervised learning and generating new data.

How WGAN Can Be Adapted to Your Task:

1. **Conditional GAN (cGAN) Extension:** One possible approach is to use a **Conditional GAN**, which is a variant of GAN that conditions both the generator and discriminator on additional input. In your case, you could condition the model on student performance metrics (the various subject scores) and generate predictions for the profession (``Klasifikasi Profesi``).
2. **Use WGAN for Data Augmentation:** Another possible application is to use the WGAN to generate synthetic student-profession pairs based on the dataset you have. You could then use the generated data to enhance the training of a traditional classification model (e.g., Random Forest, Neural Networks).

Steps to Implement WGAN for Predicting Professions:

1. **Set up the WGAN framework:**
 - **Generator:** The generator could be used to generate possible student scores/features that correspond to a specific profession.
 - **Discriminator:** The discriminator distinguishes between real student-profession pairs and generated ones.
 - **Conditioning:** You can condition the generator on student performance features to generate predicted professions.
2. **Training WGAN:**
 - You will train the WGAN to minimize the Wasserstein distance between the generated predictions and the real data.
3. **Using the Generator for Prediction:**

- Once trained, you can use the generator to predict the profession for new student data by sampling from the generator's output.

Python Implementation:

You will need deep learning frameworks like **TensorFlow** or **PyTorch** to implement a WGAN. Here's an outline of how you can proceed with PyTorch:

python

```
import torch
import torch.nn as nn
import torch.optim as optim

# Define the generator
class Generator(nn.Module):
    def __init__(self, input_size, hidden_size, output_size):
        super(Generator, self).__init__()
        self.main = nn.Sequential(
            nn.Linear(input_size, hidden_size),
            nn.ReLU(True),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(True),
            nn.Linear(hidden_size, output_size),
            nn.Softmax(dim=1) # Predicting probabilities for professions
        )

    def forward(self, x):
        return self.main(x)

# Define the discriminator
class Discriminator(nn.Module):
    def __init__(self, input_size, hidden_size, output_size=1):
        super(Discriminator, self).__init__()
        self.main = nn.Sequential(
            nn.Linear(input_size, hidden_size),
            nn.ReLU(True),
            nn.Linear(hidden_size, hidden_size),
            nn.ReLU(True),
            nn.Linear(hidden_size, output_size),
        )

    def forward(self, x):
        return self.main(x)

# Hyperparameters
input_size = 90 # Number of student features (you can adjust this)
hidden_size = 128
output_size = 1 # Predicting profession category

# Initialize models
G = Generator(input_size, hidden_size, output_size)
D = Discriminator(input_size, hidden_size)

# Optimizers
optimizerG = optim.RMSprop(G.parameters(), lr=0.00005)
optimizerD = optim.RMSprop(D.parameters(), lr=0.00005)

# Loss function for WGAN (Wasserstein distance)
def wasserstein_loss(output, target):
    return torch.mean(output * target)

# Training loop (simplified)
num_epochs = 10000
for epoch in range(num_epochs):
    # Train Discriminator
    for _ in range(5):
        real_data = get_real_student_data() # This should be your dataset
        fake_data = G(torch.randn(real_data.size(0), input_size))
```

```
D_real = D(real_data)
D_fake = D(fake_data)

D_loss = -torch.mean(D_real) + torch.mean(D_fake)
optimizerD.zero_grad()
D_loss.backward()
optimizerD.step()

# Clamp weights of discriminator
for p in D.parameters():
    p.data.clamp_(-0.01, 0.01)

# Train Generator
fake_data = G(torch.randn(batch_size, input_size))
D_fake = D(fake_data)
G_loss = -torch.mean(D_fake)

optimizerG.zero_grad()
G_loss.backward()
optimizerG.step()

# Print the loss
if epoch % 100 == 0:
    print(f"Epoch {epoch}, D Loss: {D_loss.item()}, G Loss: {G_loss.item()}")
```

Key Considerations:

- **Data Preparation:** Since WGAN is a generative model, you will need to make sure that the features of your dataset (the student scores) are properly scaled.
- **Evaluation:** After training, you can use the generator to generate predictions of professions for new student data. This could be done by sampling the generator multiple times and taking the profession with the highest probability.

Next Steps:

1. **Setup:** Ensure that you have the necessary libraries (`torch`, `torchvision`) installed for WGAN implementation.
2. **Training:** Train the model on your dataset.
3. **Evaluation:** Evaluate the generator's performance in predicting professions.

Let me know if you'd like to proceed with implementing this or need any help adjusting the model architecture for your dataset.



but first can i test the dataset first? like maybe theres model better then Wasserstein GAN?, maybe with H2O.ai? for dl? Wasserstein GAN is deep learning right?



Yes, you can absolutely test your dataset first with traditional machine learning models or other deep learning approaches before diving into more complex models like Wasserstein GAN (which is indeed a deep learning model).

Using **H2O.ai** for testing different models, including deep learning models, is a great way to benchmark and compare performance quickly. H2O.ai's **AutoML** functionality can automatically train and evaluate multiple models, including deep learning models, without much manual intervention.

Here's how you can proceed with **H2O.ai** to test your dataset:

Steps to Use H2O.ai for AutoML

1. **Install H2O.ai** if you haven't already:

```
bash

pip install h2o
```

2. **Prepare the Dataset:** Load your dataset and remove any non-predictive columns like student IDs or names.
3. **Run H2O AutoML:** You can set it up to test multiple models (including traditional ML and deep learning).

Here's a full example of using H2O.ai to predict the profession (`Klasifikasi Profesi`) based on student scores:

Example Code for H2O.ai AutoML (Classification):

```
python

import h2o
from h2o.automl import H2OAutoML
import pandas as pd

# Initialize H2O cluster
h2o.init()

# Load your dataset (assuming your dataset is a pandas dataframe)
file_path = '/mnt/data/total_merge_dataset.csv'
data = pd.read_csv(file_path)

# Convert pandas DataFrame to H2O Frame
h2o_data = h2o.H2OFrame(data)

# Set target column and features
target = 'Klasifikasi Profesi' # Predicting profession
features = data.columns.tolist()
features.remove(target)
features.remove('Nama Lengkap') # Remove non-predictive columns like names or IDs
features.remove('NPM')

# Define training and test sets
train, test = h2o_data.split_frame(ratios=[.8])

# Run H2O AutoML
aml = H2OAutoML(max_runtime_secs=300, # Adjust this for more time if needed
                seed=1,
                include_algos=["DeepLearning", "XGBoost", "GBM", "StackedEnsemble"])
aml.train(x=features, y=target, training_frame=train)

# View the AutoML Leaderboard to see the best models
lb = aml.leaderboard
lb.head()

# Make predictions on the test set
predictions = aml.predict(test)
print(predictions)
```

Key Points:

- **Include Deep Learning:** H2O's AutoML will include deep learning models (like multi-layer perceptrons) in the model comparison if you specify ``include_algos=["DeepLearning", ...]``.
- **Leaderboard:** After running, you can view the leaderboard to see which models performed best (e.g., deep learning, XGBoost, GBM, etc.).
- **Testing Time:** You can limit the amount of time AutoML spends training (via ``max_runtime_secs``) to get results quickly.

Results Interpretation:

- Once AutoML completes, the leaderboard will display the top-performing models. If deep learning models perform well, you might want to explore more advanced deep learning techniques like Wasserstein GAN. Otherwise, a simpler model like **XGBoost** or **Random Forest** may be sufficient.
- **Evaluation Metrics:** The leaderboard will show metrics like accuracy, AUC (for binary classification), and log loss, allowing you to easily compare different models' performance.

Why Start with H2O.ai?

- **Ease of Use:** It allows you to quickly benchmark a variety of models (including deep learning) without the need to manually tune or configure each one.
- **Model Comparison:** It helps identify whether simpler machine learning models or traditional deep learning models might work better than the Wasserstein GAN.

Once you've tested your dataset and explored the performance of simpler models using **H2O.ai**, you can decide whether to pursue more complex approaches like the **Wasserstein GAN**.

Would you like to proceed with setting up this H2O.ai process, or would you prefer assistance with another step?



Let's try H2O.ai AutoML now.



python

```
import h2o
from h2o.automl import H2OAutoML
import pandas as pd

# Initialize H2O cluster
h2o.init()

# Load the dataset
file_path = '/mnt/data/total_merge_dataset.csv'
data = pd.read_csv(file_path)

# Convert pandas DataFrame to H2O Frame
h2o_data = h2o.H2OFrame(data)

# Set target column and features
target = 'Klasifikasi Profesi' # Predicting profession
features = data.columns.tolist()
```

```
features.remove(target)
features.remove('Nama Lengkap') # Remove non-predictive columns like names or IDs
features.remove('NPM')

# Define training and test sets
train, test = h2o_data.split_frame(ratios=[.8])

# Run H2O AutoML with a time limit of 5 minutes
aml = H2OAutoML(max_runtime_secs=300, # Adjust for time (300 seconds = 5 minutes)
                seed=1,
                include_algos=["DeepLearning", "XGBoost", "GBM", "StackedEnsemble"])
aml.train(x=features, y=target, training_frame=train)

# Get the AutoML leaderboard
lb = aml.leaderboard
lb.head()
```

```
-----
ModuleNotFoundError                                Traceback (most recent call last)
Cell In[2], line 1
----> 1 import h2o
      2 from h2o.automl import H2OAutoML
      3 import pandas as pd

ModuleNotFoundError: No module named 'h2o'
```

It seems that the H2O library is not installed in this environment. To proceed, you can follow these steps in your local environment:

1. Install H2O.ai:

```
bash
```

```
pip install h2o
```

2. **Run the Provided Code** in your local Python environment (with the modifications I provided earlier).

Alternatively, if you'd like, I can assist you with generating a script or guide you through any part of the process locally. Would you like to proceed that way or have any other preferences?



What was the best-performing model?



Since the environment here does not support running H2O.ai, I cannot execute the AutoML process directly. However, when you run the H2O AutoML process in your local environment, the best-performing model will be listed at the top of the **leaderboard** that is generated.

The **leaderboard** contains a ranking of models based on performance metrics such as accuracy, AUC (for classification tasks), or RMSE (for regression tasks). You can easily identify the best-performing model from the leaderboard.

Here is how you can retrieve the best model:

After Running H2O AutoML:

```
python
```

```
# View the leaderboard to check the best model
lb = aml.leaderboard
print(lb)

# Get the best model (leader)
best_model = aml.leader
print(best_model)
```

The **leaderboard** will rank models like:

- **GBM (Gradient Boosting Machine)**
- **XGBoost**
- **Deep Learning Models**
- **Stacked Ensembles**

The model at the top of the leaderboard is typically the best-performing one.

If you run this in your local setup, you'll be able to see which model (e.g., deep learning or a more traditional model) performed the best for predicting student professions. Would you like help interpreting the results once you've run it, or any additional guidance on the process?

ChatGPT can make mistakes. Check important info.