Crazy EDA.

Steps before model training

```python
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

df = pd.DataFrame(scaler.fit_transform(df), columns=df.columns)
```

Train

All Numerical Columns.

① Cat → Num
② Feature Scaling

\# Sample - $m$

dimension/ \# features - $d$

Predictors - $(m, d)$ matrix $X$

$I$th sample - $X^{[i]}$ — one row

$j$th feature - $X_j$ — one column

True output/target - $y^i$

Predicted output/target $\hat{y}^i$

# Goal of generalization in ML

m samples - $\{x^i, y^i\}^m$

Historical Data. $\longrightarrow$ Train an ML Model

$$x^i \xrightarrow[\text{Model}]{\text{ML}} \hat{y}^i \quad \longleftarrow \text{Predicted value}$$

Ideally, $y^i \approx \hat{y}^i \longrightarrow$ Predicted Output

$\quad\quad\quad \hookrightarrow$ True Output

$\Longrightarrow$ For all <u>m samples</u> $\boxed{y^i \approx \hat{y}^i} \xrightarrow{\hspace{1cm}}$ GOOD TRAINING
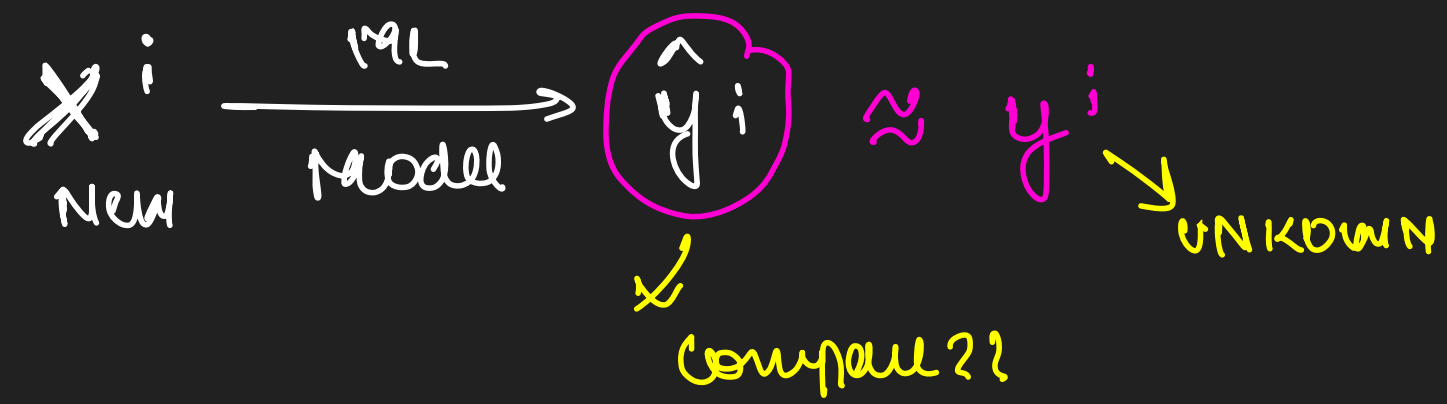
DO WE WANT TO PERFORM PREDICTION ON TRAINING

for training data y is already present?

Why to even do prediction.

Goal of ML: GENERALISATION.

You model should perform well
both on training and (NEW) data.

True output
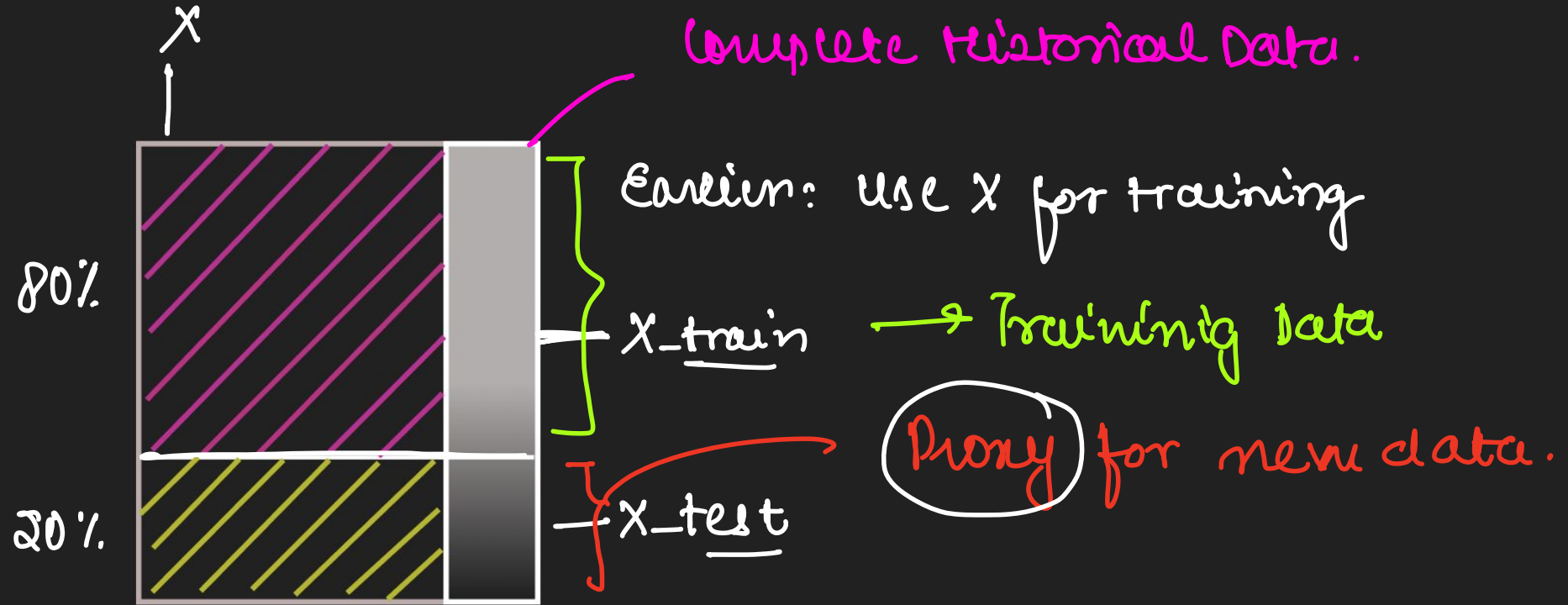
$$X^i \xrightarrow[\text{Model}]{ML} \hat{y}^i \approx y^i$$

New

compare??

UNKOWN

check if model is GENERALISING??

DATA??

**PROBLEM: No labels for new data.**

**Because we don't know the ground truth (y) for New Sample**

X

Complete Historical Data.

Earlier: Use X for training

80%

X_train → Training Data

Proxy for new data.

20%

X_test

Now: Use only X_train for training

# Phases for Model Development

① **Training** - use trainin split for model training

$$y_{train} \approx \hat{y}_{train} \quad - \quad \text{Model is learning}$$

② **Testing / Evaluation** - test split

$$y_{test} \approx \hat{y}_{test} \quad - \quad \text{Model is generalising}$$

```python
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.3,
                                                    random_state=100
```

0.2

random split

80%    20%

200
200
seed.