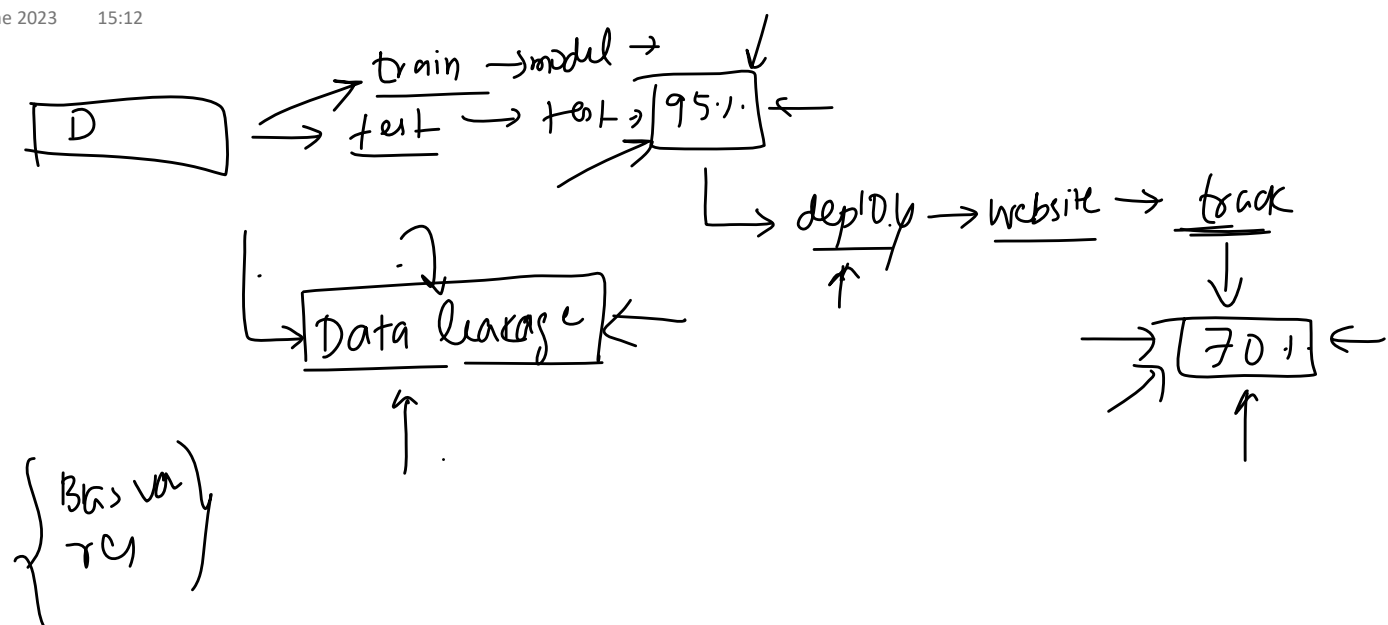


The Problem

12 June 2023 15:12

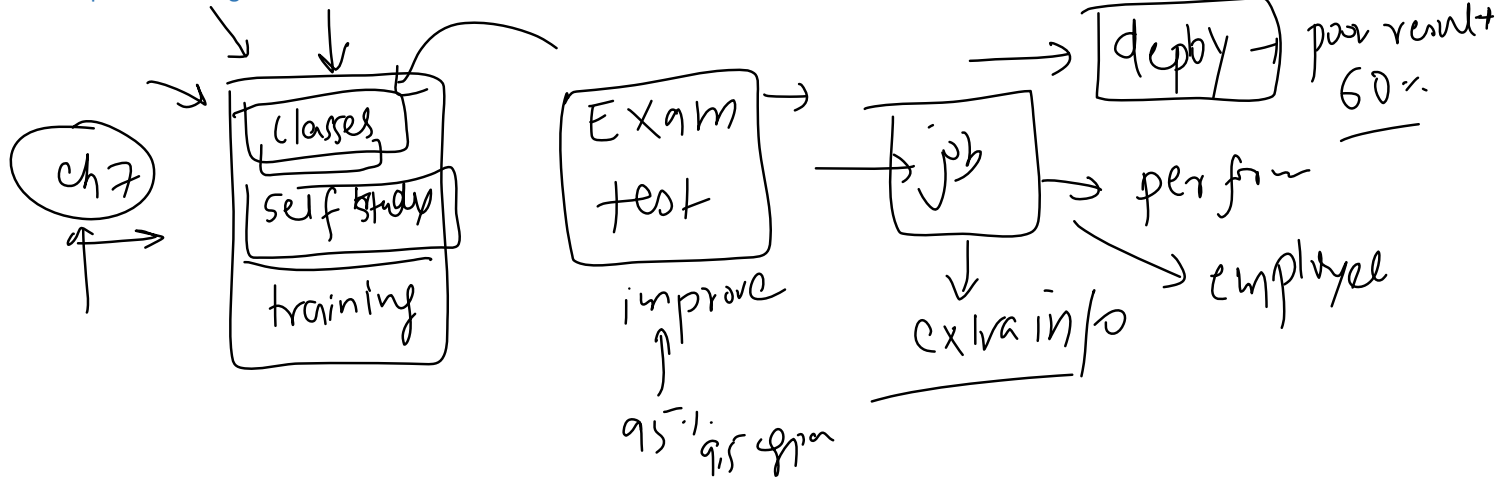


What is Data Leakage?

12 June 2023 15:13

Data leakage, in the context of machine learning and data science, refers to a problem where information from outside the training dataset is used to create the model. This additional information can come in various forms, but the common characteristic is that it is information that the model wouldn't have access to when it's used for prediction in a real-world scenario.

This can lead to overly optimistic performance estimates during training and validation, as the model has access to extra information. However, when the model is deployed in a production environment, that additional information is no longer available, and the performance of the model can drop significantly. This discrepancy is typically a result of mistakes in the experiment design.



Ways in which Data Leakage can occur

12 June 2023 15:13

input → fi → future → available

1. Target Leakage:

Target leakage occurs when your predictors include data that will not be available at the time you make predictions.

2. Multicollinearity with target col

3. Duplicated Data

4. Preprocessing Leakage → Train test contamination & Improper Cross Validation

5. Hyperparameter Tuning

credit card fraud detect

value | date | website | . . . | reversed_transaction | fraud |

→ 500 — — — — —

normal → 400

reversed_transaction → 1 (normal) / 0 (fraud)

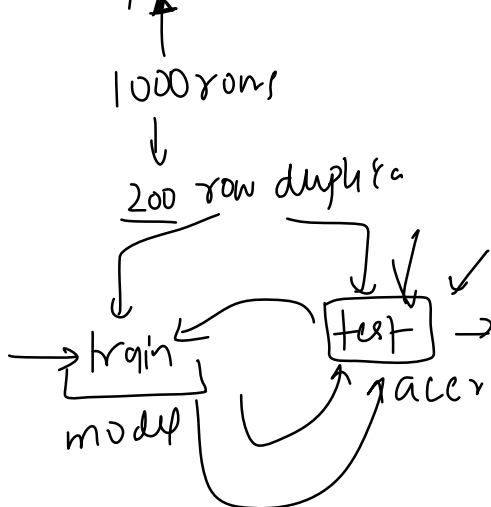
prediction → 80.1%

smartphone price predict

brand | ram | . . . | ratings | price

input → accuracy

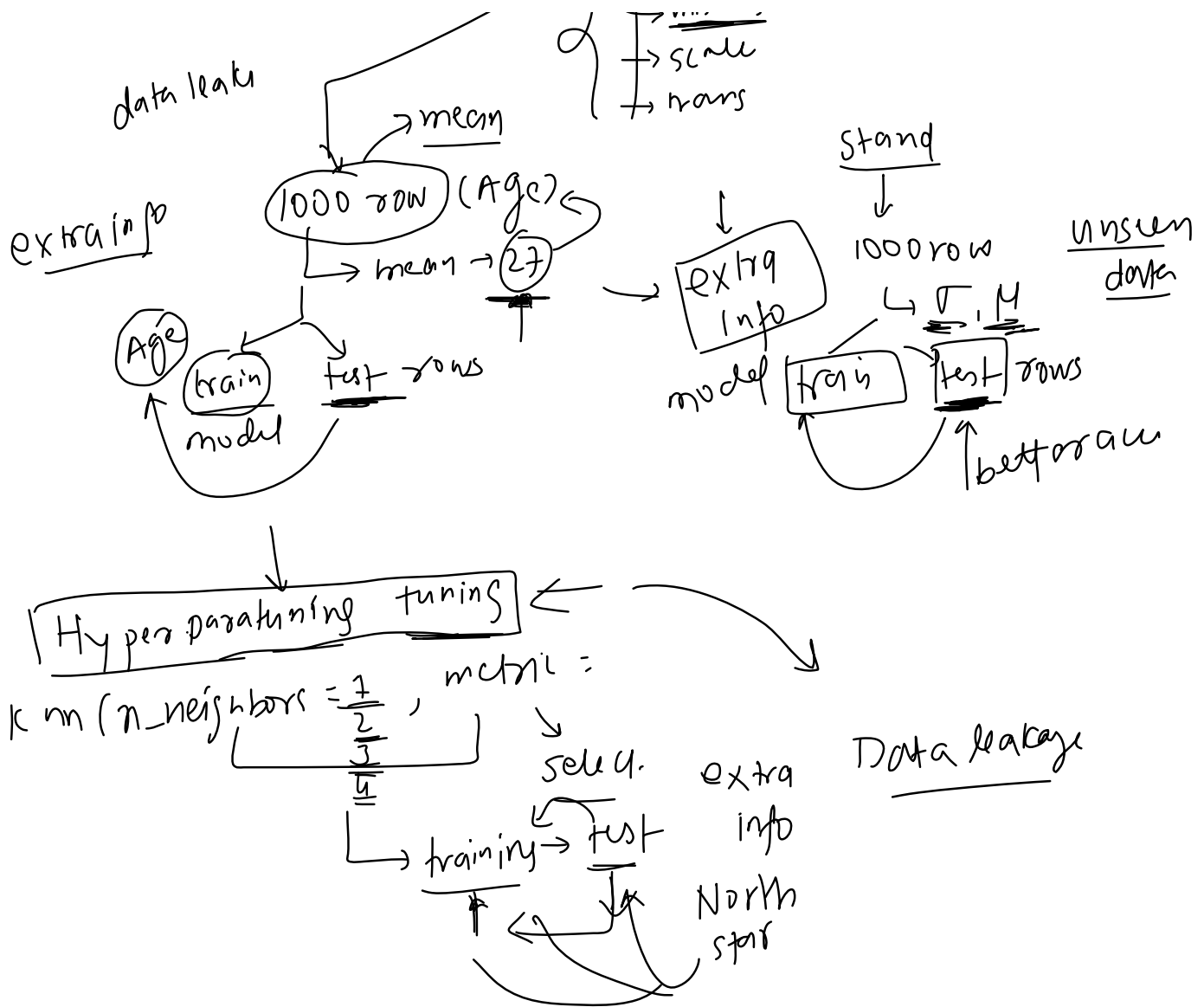
duplicates → sort of duplicate



extra info

preprocessing → model

missing
scale



How to detect?

12 June 2023 15:13

$f_1 f_2 \dots f_3$

1. Review Your Features: Carefully review all the features being used to train your model. Do they include any data that wouldn't be available at the time of prediction, or any data that directly or indirectly reveals the target? Such features are common sources of data leakage.



2. Unexpectedly High Performance: If your model's performance on the validation or test set is surprisingly good, this could be a sign of data leakage. Most predictive modelling tasks are challenging, and exceptionally high performance could mean that your model has access to information it shouldn't.



↓ producing

3. Inconsistent Performance Between Training and Unseen Data: If your model performs significantly better on the training and validation data compared to new, unseen data, this might indicate that there's data leakage.

4. Model Interpretability: Interpretable models, or techniques like feature importance, can help understand what the model is learning. If the model places too much importance on a feature that doesn't seem directly related to the output, it could be a sign of leakage.



↪ feature imp

$f_1 \rightarrow 0.13$
 $f_2 \rightarrow 0.81$

How to remove Data Leakage

12 June 2023 15:14

train test cross val

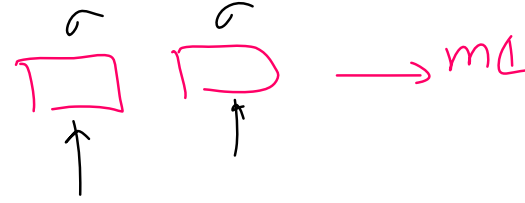
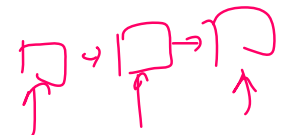
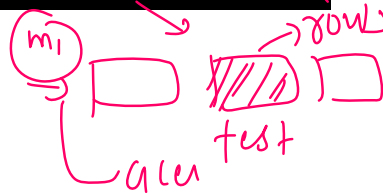
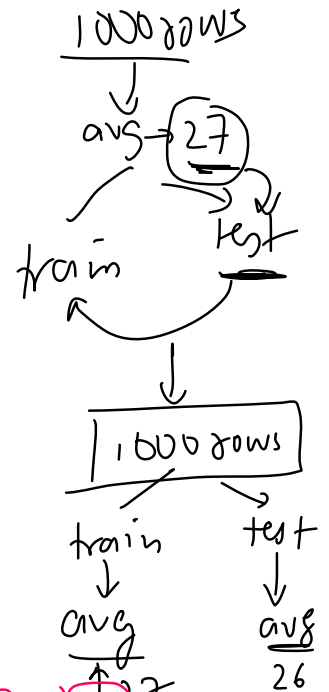
1. Understand the Data and the Task: Before starting with any kind of data processing or modelling, it's important to understand the problem, the data, and how the data was collected. You should understand what each feature in your data represents, and whether it would be available at the time of prediction.
2. Careful Feature Selection: Review all the features used in your model. If any feature includes information that wouldn't be available at the time of prediction, or that directly or indirectly gives away the target variable, it should be removed or modified.
3. Proper Data Splitting: Always split your data into training, validation, and testing sets at an early stage of your pipeline, before doing any pre-processing or feature extraction.
4. Pre-processing Inside the Cross-Validation Loop: If you're using techniques like cross-validation, make sure to do any pre-processing inside the cross-validation loop. This ensures that the pre-processing is done separately on each fold of the data, which prevents information from the validation set leaking into the training set.

```
# Incorrect way
X_normalized = normalize(X) # normalize the whole dataset
cross_val_score(model, X_normalized, y, cv=5) # perform cross-validation
```

```
# Correct way
pipeline = make_pipeline(normalizer, model)
cross_val_score(pipeline, X, y, cv=5) # per
```

6. Avoid Overlapping Data: If the same individuals, or the same time periods, appear in both your training and test sets, this can cause data leakage. It's important to ensure that the training and test sets represent separate, non-overlapping instances.

duplicate
remove
data leakage



Validation Set

12 June 2023 22:02

Deep lar

