# Recap

Linear Reg case study

Cross validation

Hyper parameter

model selection

# The Problem

Cross val $\longrightarrow$ why

supervised ml

$\longleftarrow X \longrightarrow$ | y

model train $\longrightarrow$ new data

evaluate $\longleftarrow$

old $\longrightarrow$ labeo

$X \longrightarrow y$ $\longrightarrow$ model

evaluate

deploy $\longrightarrow$ new data

# The Hold-out Approach

→ train_test_split

1000 customers



→ 0.75        0.25

training      test

train → model  → y_pred / y_true

1) shuffle
2) train_test

model evaluation

↓

Cross validation

# Problem with Hold-out Approach

*more data better models*

*High nl*

1. **Variability:** The performance of the model can be very sensitive to how the data is divided into training and testing sets. If the split is unfortunate, the training set may not be representative of the overall distribution of data, or the test set might contain unusually easy or difficult examples. This leads to high variance in the estimation of the model's performance.
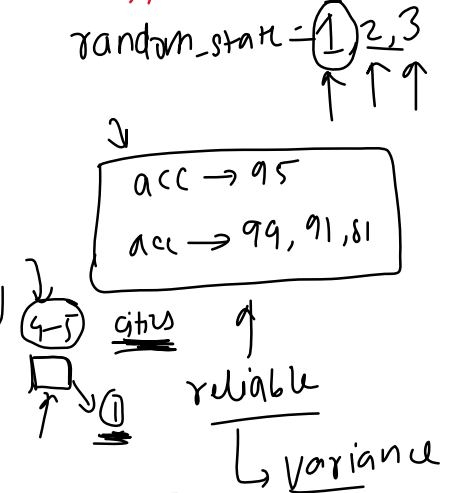
2. **Data inefficiency:** The holdout method only uses a portion of the data for training and a different portion for testing. This means that the model doesn't get to learn from all available data, which can be particularly problematic if the dataset is small.

3. **Bias in performance estimation:** If some classes or patterns are over- or under-represented in the training set or the test set due to the random split, it can lead to a biased performance estimation.

   *High_bias ←*

4. **Less reliable for hyperparameter tuning:** If the holdout method is used for hyperparameter tuning, there's a risk of overfitting to the test set because information might leak from the test set into the model. This means that the model's performance on the test set might be overly optimistic and not representative of its performance on unseen data.

*① shuffle*

*random_state = ① 2,3*

$$acc \rightarrow 95$$
$$acc \rightarrow 99, 91, 81$$

*→ data leakage*

*4-5 girus*

*reliable*
*↳ variance*

*100%*
*80%*

*80% training*
*20% percent*

*0.72, 0.76, 0.77 best → model deploy*

*80% → bias ↑↓*

*Bias variance trade off*
*↳ loss → bias + var + noise*

*→ bias*

*correct acc%*

# Why is hold-out approach used then?

52  62    8८    92        81.9  87.2  87.4

smaller

1. **Simplicity**: The holdout method is straightforward and easy to understand. You simply divide your data into two sets: a training set and a test set. This simplicity makes it appealing, especially for initial exploratory analysis or simple projects.

→ project
→ model

2. **Computational Efficiency**: The holdout method is computationally less intensive than methods like k-fold cross-validation. In k-fold cross-validation, you need to train and test your model k times, which can be computationally expensive, especially for large datasets or complex models. With the holdout method, you only train the model once.

→ cross-valid → multiple models

large

3. **Large Datasets**: For very large datasets, even a small proportion of the data may be sufficient to form a representative test set. In these cases, the holdout method can work quite well.

large dataset → high variance random

↳ diminish

1L row → random
↳ change

1 L → 10 sample

pop avg height

1700 sample2
↳ more data

# Cross Validation

given → f1
       → f2
       → f3 → test

The idea of cross-validation is to divide the data into several subsets or "folds" The model is then trained on some of these subsets and tested, on the remaining ones. This process is repeated multiple times, with different subsets used for training and validation each time. The results from each round are usually averaged to estimate the model's overall performance.

avg sal

final

Resampling → sampling → pop → Observa    avg

Cross val

Bootstrap

pop → Observa → samples → estimate

dataset (150)

iris → Sample → pop

100 → s1 mod1
    → s2 m2
    → s3 m3
    → s4 m4

avg
↓ crossscore

Leave One Out CV
( LOO CV)

1

Leave p out CV

cross_val_score
↓
api ←

2 → K-fold CV
 ├ K-fold
 ├ Stratified
3├ nested k-fold
 └ repeated k-fold

time based
→ Time series

msc

test data → iron
          ↓
          oLswr

# Leave One Out Cross Validation (LOOCV)

dataset $\rightarrow$ 1000 rows $\rightarrow$ $\underline{\text{models}}$ $\rightarrow$ model
$(n)$            $(1000)$        evaluation

[diagram of dataset split into 1000 segments, labeled 1000]

[segment with n-1 rows, 999] $\rightarrow$ 999 rows $\rightarrow$ training $\rightarrow$ $\boxed{m\,1}$
                                1 row $\rightarrow$ test

[segment] 999 rows $\rightarrow$ m2
          1 row

⋮

[segment]   $\underline{m\,1000}$         $\overline{LOOCV}$
            ⤵                    ⤷
accuracy

$m_1 + m_2 \cdots \cdots + m_{1000} \rightarrow \underline{avg} \rightarrow \dfrac{\text{final accuracy}}{\text{score}}$

## Advantages:

1. **Use of Data**: LOOCV uses almost all of the data for training, which can be beneficial in situations where the dataset is small and every data point is valuable.
2. **Less Bias**: Since each iteration of validation is performed on just one data point, LOOCV is less biased than other methods, such as k-fold cross-validation. The validation process is less dependent on the random partitioning of data.
3. No Randomness: There's no randomness in the train/test split, so the evaluation is stable, without variation in the results due to different random splits. (no shuffling)

$\rightarrow$ LOOCV      $\underline{80\% \text{ data}} \rightarrow$ high bias
     ⤷
  $n-1$ row      $\boxed{\dfrac{n-1 \text{ n}}{n}} \rightarrow$ high %
                      ⤵
                 reduce bias

## Disadvantages:

1. Computational Expense: LOOCV requires fitting the model N times, which can be computationally expensive and time-consuming for large datasets.

2. High Variance: LOOCV can lead to higher variance in the model performance since the training sets in all iterations are very similar to each other.

3. Inappropriate Performance Metric: Performance metrics like R^2 are not appropriate to be used with LOOCV as they are not defined when the validation set only has one sample.

4. Not Ideal for Imbalanced Data: In classification problems, if you have imbalanced classes, LOOCV may not provide a reliable estimate of model performance because the single validation sample in each iteration may not

[diagram with bias arrows]   bias
                              ⤵
                             reduce

5000 row data
     ⤵
5 or models
         ↑
accuracy & lot of time

4. Not Ideal for <u>Imbalanced Data</u>: In classification problems, if you have imbalanced classes, LOOCV may not provide a reliable estimate of model performance because the single validation sample in each iteration may not be representative of the overall class distribution.

When to use:

1. <u>Small datasets</u>: LOOCV is most beneficial when you have a limited amount of data. With small datasets, you want to use as much data as possible for training to get a reliable model, which is exactly what LOOCV offers by using all but one data point for training.

2. <u>Balanced datasets</u>: LOOCV might not perform well on imbalanced datasets, especially in classification problems, because the training set might end up missing some classes. Thus, it's more appropriate to use LOOCV when you have a balanced dataset.

3. <u>Need for less biased performance estimate</u>: <u>Since LOOCV uses nearly all the data for training</u>, <u>it gives a less biased estimate</u> of model performance compared to other methods like k-fold cross-validation.

accuracy { lot of time, e ot of space }

1 data point

accura → variance

o ver

1000 → 990  10
            Y    N

→ 990 test sets
       └→ N component
540 - 460
  Y      N

model
Var  → bias
low    high

LOOCV

K-Fold Cross Validation $\rightarrow$ most used cv technique  general   pure wrcups

10 June 2023   12:34

variant                                Linear Reg $\rightarrow$ case study

Data leaks

All Data                               $\downarrow$

100                                    $\rightarrow$ Cross validation $\rightarrow$ $\frac{1}{2}$

$\downarrow$         5 equal folds      $\rightarrow$ Hyperparameter
20 rows                                    ③

Training data         Test data                              ②

K = 5, 10                              Data Leakage

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

① 

Split 1  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5   m1 $\rightarrow$ global results
Split 2  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5   m2 $\rightarrow$                model
Split 3  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5   m3 $\rightarrow$  Finding Parameters
Split 4  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5   m4        avg $\rightarrow$ avg
Split 5  Fold 1  Fold 2  Fold 3  Fold 4  Fold 5   m5 $\rightarrow$              acer

$\rightarrow$ Evaluate $\rightarrow$ train
$\rightarrow$ Hold out set $\searrow$ test

5 models          Final evaluation    Test data          $\rightarrow$ problems
1 one each   n rows $\rightarrow$ n models $\rightarrow$ n-1        unseen   $\rightarrow$ use
fold                                        row training   1 testing   data   $\rightarrow$ Crossvalidation      costly
                                                                                    $\rightarrow$ LOOCV
             $\rightarrow$ avg                                                          what
                                                                                     draw
                                                                                     advant

$\rightarrow$ K fold  K-1 fold train       n-1 rows training $\rightarrow$ bias low
                high bias

**Advantages of K-Fold Cross Validation:** $\leftarrow$ K-fold $\rightarrow$ Variant

1. Reduction of Variance: By averaging over k different partitions, the variance
   of the performance estimate is reduced. This is beneficial because it means
   that the performance estimate is less sensitive to the particular random
   partitioning of the data.                                     no. of rows 10000

                                                    $\rightarrow$ LOOCV

2. Computationally Inexpensive: Take less time and space in comparison to
   LOOCV                                            $\rightarrow$ K=5,10 $\rightarrow$ 5 models

**Disadvantages of K-Fold Cross Validation:**

1. Potential for High Bias: If k is too small, there could be high bias if the test set
   is not representative of the overall population.   1 fold     r2 $\rightarrow$ 1 row
                                                                        $\rightarrow$ vary
2. May not work well with Imbalanced Classes: If the data has imbalanced
   classes, there's a risk that in the partitioning, some of the folds might not    multiple
   contain any samples of the minority class, which can lead to misleading
   performance metrics.                              1000 $\rightarrow$ 5 fold, 200 rows/fold

2. May not work well with imbalanced classes: If the data has imbalanced classes, there's a risk that in the partitioning, some of the folds might not contain any samples of the minority class, which can lead to misleading performance metrics.
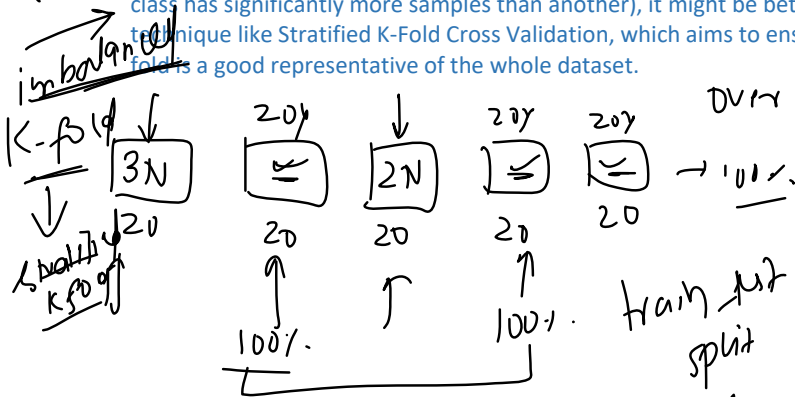
Vanilla K fold

↳ classification

→ 95% 5% → 100 row 5 fold
   Y    N            cv
       95y 5N

multiple
1000 → 5 fold  200 rows/fold
       variance    fo rows
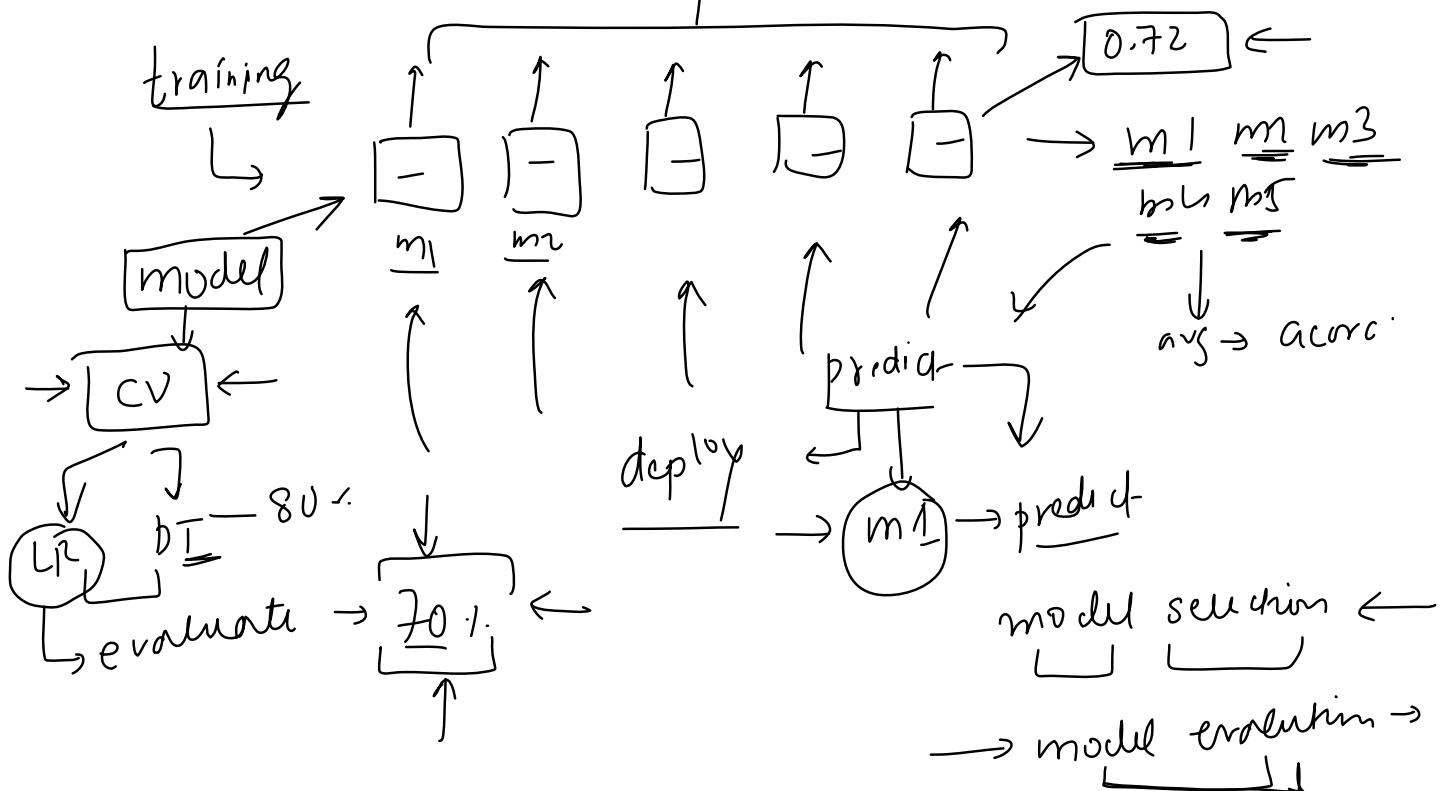hold-out variance
        ⬇
      K-fold

**When to use:**

1. When you have a sufficiently large dataset: K-Fold Cross Validation requires the model to be trained K times, so it can be computationally intensive. However, if your dataset is large enough, this increased computational cost can be justified by the more reliable estimate of model performance.

2. When your data is evenly distributed: K-Fold Cross Validation works best when the data is evenly distributed. If your dataset is imbalanced (i.e., one class has significantly more samples than another), it might be better to use a technique like Stratified K-Fold Cross Validation, which aims to ensure each fold is a good representative of the whole dataset.

imbalanced

K-fold
   ⬇
 stratified
 K fold

3N    20y    2N    20y  20y  over avg 9k%
20    20     20    20   20    → 100%
     100%        100%  train test
                       split

mpg) horse power
        ⬇
poly degree = 1,2,3...6

Hold out approach
        ⬇
(random state =
        ⬇
      0.72  ←

m1  m  m3
   b4 m5
     ⬇
   avg → acorc

training
   ↳

model → [ — ] [ — ] [ — ] [ — ] [ — ]
         m1   m2

CV
 ↓   ↓
LR   DI — 80%

evaluate → 70% ←

deploy →   predict
        → m1 → predict

model selection ←

→ model evaluation →

95% Y    5% N     95% Y    5% N

Imbalanced
↓
classification

| | class 1 Y 2 | class 2 N 4 | ... 3 | class n M |
|---|---|---|---|---|
| round 1 | | | | |
| round 2 | | | | |
| round 3 | | | | |
| round 4 | | | | |
| ⋮ | | | | |
| round K | k-1 : 1 | k-1 : 1 | | k-1 : 1 |

■ training data    ■ validation data

Y
95%

N
5%

K-fold

100%   avg →

1, 2, 3
↓  ↓  ↓
2   4   3

Class is not

realistic → avg → realistic
value

m1 → 70%   5 num

unseen data → test data

m L → 72%

unseen data – test