

Topic 1

outlier removal methods

What is an outlier?

- A data point that is significantly different from other observations
- It is a point that doesn't belong with the main cluster of data

Why care?

They can seriously skew your results.

An outlier can pull the mean (average) & warp the line a model tries to fit, leading to incorrect predictions

Method 1: 3 sigma (3 σ) Rule

- Based on the properties of a normal distribution (the BELL CURVE)

Key Assumption: Your data should be roughly bell shaped for this to work well

The Rule: Any data point that falls outside of 3 std. deviation from the mean is considered an outlier

— Working —

calculate the mean (μ)
(Find the average of your data)

↓

calculate std. deviation (σ)

↓

Define boundaries

$x = \text{lower} \rightarrow \mu - 3\sigma$ & $\text{upper} \rightarrow \mu + 3\sigma$

↓

Any point $x < \text{lower bound}$ or $x \geq \text{upper bound}$
is an outlier

pro → simple & fast, cons → not robust

Method 2: IQR method

A robust method not influenced by extreme values. It looks at the middle 50% of the data

Key Assumption: Does not require the data to be normally distributed

Rule: Any pt. outside the range defined by $1.5 \times$ the IQR below Q_1 & above Q_3 is an outlier

Sort \uparrow → from smallest to largest

↓

Find quartile

Q_1 : The 25th percentile

Q_3 : The 75th percentile

↓

calculate IQR ($IQR = Q_3 - Q_1$)

↓

Define boundaries

$\text{Lower} = Q_1 - 1.5 \times IQR$ & $\text{upper} = Q_3 + 1.5 \times IQR$

↓

Identify

Pro \rightarrow Very Robust to outliers
cons \rightarrow slightly more steps to calculate

Topic - 2 : overfitting, underfitting & Generalization

The goal: A Generalized model

The model learns the depth underlying pattern from the training data while ignoring the random noise

A generalized model performs well not only on the data it was trained on, but more importantly on new, unseen data.

Analogy: A student who studies to truly understand the concepts (generalization) will ace both the practice test & the final exam

underfitting: (too simple model)

The model is not complex enough to capture underlying true trend in the data

This is a problem of high bias. The model makes overly simplistic assumptions about the data. The model has high error on both Training & Test data.

overfitting - Too complex model

The model learns the training data too perfectly. It starts to memorize the random noise & outliers as if they were real patterns.

This is a problem of high variance. The model is overly sensitive to the specific data it was trained on.

The model has very low error on the training data, high error on new data.

The Bias Variance trade off

Bias \rightarrow error from being too simple

Variance \rightarrow error from being too complex

The trade off \rightarrow As you \downarrow bias

As you \downarrow bias \rightarrow may cause \uparrow in variance
Goal is to find the perfect balance

TOPIC - 3 SVM Alg. used for classification

Maximum margin classifier

main job is to find the best possible boundary
that separates 2 classes

Key idea: It doesn't just find any
line that separates the data. It finds
the one that creates the widest possible
empty space b/w 2 classes. This is
called maximizing the margin.

Key terminology

Hyper plane: The decision boundary itself
 \hookrightarrow In 2D \rightarrow It is a line
3D \rightarrow it is a plane

Support vectors;

data points closest to hyper plane
They lie on the edge of margin

Margin: This is the empty space

\hookrightarrow goal of SVM to make empty space
as wide as possible



working

Plot the data points



Find a separating hypel plane



Identify support vectors & margin



Maximize the margin



classify new data

Margin

Hard margin

• very strict, data perfectly
strict linearly
separable

Soft margin

Much more flexible

The trade off

• Mistakes allowed = 0

• extremely sensitive
to outliers

Using the 'C'

Hyper parameter

High 'C' → Acts more
like strict
hard margin

Low 'C' - more
tolerant

Topic - 4

Hyper parameters tuning

What is hyper parameters tuning?

A setting or configuration for a model that is chosen before the training process begins.

The model does not learn hyper parameters. The data scientist sets it.

Eg:-

The C value in SVM

No. of trees in Random Forest

Learning rate for a neural network

Why tuning is necessary?

The performance of a model is extremely sensitive to its hyper parameters settings.

To find the specific combination of hyper parameters \rightarrow to find the best model.

Common tuning techniques

Grid Search

Create a grid of possible values you

want to test for each hyper parameter.

Grid search exhaustively trains & evaluates a model for every single possible combination pros \rightarrow you are guaranteed to find best combination cons \rightarrow slow & expensive

Randomized Search

It tests a fixed no. of random combinations from the range of values you provide

Pros → faster & more efficient

Cons → not guaranteed to find the absolute best combination

Grid Search Step

→ Define the hyperparameter grid

→ Create all combinations

→ Evaluate each combination

→ Select the best combination

Randomized Search

→ Define hyperparameter space

→ Set the no. of iterations

→ Randomly sample & evaluate

→ Select the best combination

Topic - 5 cross validation

Get a reliable & stable estimate of a model performance

It prevents data leakage

The k-fold cross validation process

The initial split

The very first step is to partition your entire dataset into a training set & a test set

The test set is locked away. All the following steps happen only on the training set.

2 Partition into ' k ' folds

The training set is divided into ' k ' equal sized, non overlapping subsets called 'folds' (k is usually 5 or 10)

3) The iteration loop

The process iterates ' k ' times in each run

- A single unique fold chosen as the validation set

- The remaining $k-1$ folds as the training data
- The model is trained on the training data & then scored on the validation fold

4) Aggregate the Scores

After ' k ' rounds you will have ' k ' performance scores

The final cross validation score is the average of these k scores