Categorical Feature Selection

how to sel; ect relevant features if both the target and features are categorical

When both target and features are categorical, you're dealing with categorical-categorical feature selection, which is a bit different from the usual numeric target or numeric feature scenarios.

Here's a breakdown of **methods you can use to select relevant features** in this case:



🔧 1. Chi-Square Test (χ²)

Use when: Both the feature and target are categorical

What it does: Measures whether the distribution of the feature is independent of the target.

Use case:

```
python
from sklearn.feature_selection import chi2
from sklearn.feature_selection import SelectKBest
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OneHotEncoder
# Label encode your features and target
X_encoded = X.apply(LabelEncoder().fit_transform)
y_encoded = LabelEncoder().fit_transform(y)
# Apply Chi2
selector = SelectKBest(score_func=chi2, k='all')
chi2_scores = selector.fit(X_encoded, y_encoded)
# Get scores
feature_scores = pd.DataFrame({
    'Feature': X.columns,
    'Chi2 Score': chi2_scores.scores_
}).sort_values(by='Chi2 Score', ascending=False)
```

📊 2. Mutual Information (Information Gain)

Use when: You want to measure non-linear dependency between categorical variables.

Use case:

```
python

from sklearn.feature_selection import mutual_info_classif

X_encoded = X.apply(LabelEncoder().fit_transform)
y_encoded = LabelEncoder().fit_transform(y)

mi_scores = mutual_info_classif(X_encoded, y_encoded, discrete_features=True)

feature_scores = pd.DataFrame({
    'Feature': X.columns,
    'Mutual Info': mi_scores
}).sort_values(by='Mutual Info', ascending=False)
```

3. Cramér's V (Association Strength)

Use when: You want to understand the strength of association between two categorical variables. It's based on Chi-square.

Values range from 0 (no association) to 1 (strong association)

Use case:

```
import pandas as pd
import numpy as np
from scipy.stats import chi2_contingency

def cramers_v(x, y):
    confusion_matrix = pd.crosstab(x, y)
    chi2 = chi2_contingency(confusion_matrix)[0]
```

```
n = confusion_matrix.sum().sum()
phi2 = chi2/n
    r,k = confusion_matrix.shape
    phi2corr = max(0, phi2 - ((k-1)*(r-1))/(n-1))
    rcorr = r - ((r-1)**2)/(n-1)
    kcorr = k - ((k-1)**2)/(n-1)
    return np.sqrt(phi2corr / min((kcorr-1), (rcorr-1)))

# Example:
scores = {col: cramers_v(X[col], y) for col in X.columns}
pd.Series(scores).sort_values(ascending=False)
```

4. Feature Importance from Tree-Based Models (e.g., CatBoost)

Even for **categorical variables**, models like **CatBoost** can handle them natively without encoding and provide feature importances.

Use case:

```
python

from catboost import CatBoostClassifier

model = CatBoostClassifier(verbose=0)
model.fit(X, y, cat_features=list(X.columns))

# Get feature importance
importances = model.get_feature_importance()
pd.Series(importances, index=X.columns).sort_values(ascending=False)
```

▼ Final Recommendation Based on Use Case

Use Case	Method
Fast, simple test	Chi-Square Test

Use Case	Method
Non-linear relationships	Mutual Information
Association strength measure	Cramér's V
Embedded method (model-based)	CatBoost

If you're building a pipeline:

Start with Mutual Info + Chi2 \rightarrow then validate with Cramér's V \rightarrow use a model like CatBoost or LightGBM for feature ranking \rightarrow finalize top features.

Let me know if you want a **wrapper method**, e.g., **forward selection using mutual information or classification accuracy**, or if you'd like this in a **Jupyter-ready notebook**.

expalin Mutual Information (Information Gain)

Sure — let's dive deep into **Mutual Information (MI)**, also called **Information Gain**, especially for **categorical features and categorical target**.

Q What Is Mutual Information (MI)?

Mutual Information (MI) measures how much knowing one variable reduces the uncertainty of the other.

In simple terms:

It tells you how much information a feature gives you about the target.

If a feature and the target are **independent**, MI = 0.

If they're **strongly related**, MI > 0.



Given two discrete random variables X (feature) and Y (target):

$$\operatorname{MI}(X;Y) = \sum_{x \in X} \sum_{y \in Y} P(x,y) \cdot \log \left(rac{P(x,y)}{P(x) \cdot P(y)}
ight)$$

- ullet P(x,y): joint probability of feature value x and target value y
- P(x), P(y): marginal probabilities of feature and target

This formula basically compares the **joint distribution** to what we'd expect under **independence**.

Why Use MI for Feature Selection?

- Handles non-linear relationships
- Doesn't assume normality or linearity
- Works with categorical variables
- Can capture combinatorial relationships between feature and label values

📊 Intuition with an Example

Let's say we have:

Color	Likes_Burger
Red	Yes
Red	Yes
Blue	No
Blue	No
Green	Yes

Color	Likes_Burger
Green	No

Now if "Color" is strongly associated with "Likes_Burger," the MI will be **high**, because knowing the color gives you **lots of info** about the burger preference.

If Likes_Burger is randomly distributed regardless of Color, MI = 0.

K How to Compute MI in Practice

You can compute mutual information using scikit-learn like this:

```
from sklearn.feature_selection import mutual_info_classif
from sklearn.preprocessing import LabelEncoder

# Assume X is your feature dataframe, y is the categorical target

# Label encode categorical columns
X_encoded = X.apply(LabelEncoder().fit_transform)
y_encoded = LabelEncoder().fit_transform(y)

# Compute MI
mi = mutual_info_classif(X_encoded, y_encoded, discrete_features=True)

# Attach scores to features
import pandas as pd
mi_scores = pd.Series(mi, index=X.columns).sort_values(ascending=False)
print(mi_scores)
```

You can then select features with the highest MI scores.

🔽 Mutual Information vs. Chi-Square

Metric	Handles Non-linear?	Output Type	Assumes Independence?
Mutual Information	✓ Yes	Info value (0–1+)	X No (measures dep.)
Chi-Square	X Limited	Test statistic	✓ Yes (tests H0)

Pro Tips

- Use **mutual_info_classif** for classification targets
- Use mutual_info_regression for regression targets
- You can combine it with **SelectKBest** to choose top k features

Would you like me to:

- Visualize how MI values differ across features?
- Show a wrapper method using MI for forward feature selection?
- Compare MI scores with model feature importances?

Let me know your direction.

how does chi square help in feature selection

Great question.

Let's break down how the Chi-Square test helps in feature selection when both features **and target are categorical** — step by step.



Objective of Chi-Square in Feature Selection

Chi-Square tests whether a categorical feature and the target variable are independent.

If a feature is **independent** of the target \rightarrow it's not useful \rightarrow low score \rightarrow **can be dropped** If a feature is **dependent** on the target \rightarrow it's useful \rightarrow high score \rightarrow **keep it**

What is the Chi-Square Test?

It compares the **observed frequencies** (what the data says) to the **expected frequencies** (what we'd expect if there was no relationship).

Notice Formula:

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

- O: Observed frequency from data
- ullet E: Expected frequency under independence

📊 Step-by-Step Intuition

Let's say you have:

1	
Gender	Purchased
Male	Yes
Male	No
Female	Yes
Female	No

Chi-square will:

- 1. Build a **contingency table** (cross-tab)
- 2. Compute expected frequencies assuming Gender and Purchased are independent
- 3. Compare observed vs expected using the formula

4. Output a **Chi-square score** (higher → stronger relationship)

How It Helps in Feature Selection

You calculate the Chi-square statistic for **each feature** with respect to the target.

- A high score means the feature has a strong relationship with the target → select it
- A low score means the feature is likely noise or irrelevant → discard it

X Code Example (scikit-learn)

```
python
from sklearn.feature_selection import chi2
from sklearn.feature_selection import SelectKBest
from sklearn.preprocessing import LabelEncoder
import pandas as pd
# Example: X has categorical features, y is categorical target
# Label encode features & target
X_encoded = X.apply(LabelEncoder().fit_transform)
y_encoded = LabelEncoder().fit_transform(y)
# Run Chi-Square Test
chi2_scores, p_values = chi2(X_encoded, y_encoded)
# Output feature scores
feature_scores = pd.DataFrame({
    'Feature': X.columns,
    'Chi2 Score': chi2_scores,
    'p-value': p_values
}).sort_values(by='Chi2 Score', ascending=False)
print(feature_scores)
```

o When to Use Chi-Square for Feature Selection



Bonus Tip

You can combine it with **SelectKBest**:

```
python

from sklearn.feature_selection import SelectKBest, chi2
selector = SelectKBest(score_func=chi2, k=5) # Top 5 features
X_new = selector.fit_transform(X_encoded, y_encoded)
```

Chi-Square vs. Mutual Information

Metric	Handles Non-linear?	Gives p-values?	Measures strength?	Measures association?
Chi-Square	X No	✓ Yes	✓ Test statistic	Association
Mutual Info	✓ Yes	X No	✓ Info gain	Association

Want to see a **full example notebook** comparing chi-square vs mutual info with real data? Or want me to create a function to automate feature ranking using chi-square?

please do with an example

Chi-square will:

Build a contingency table (cross-tab)

Compute expected frequencies assuming Gender and Purchased are independent

Compare observed vs expected using the formula

Output a Chi-square score (higher → stronger relationship)

Perfect — let's walk through a **Chi-Square feature selection example from scratch** using Python and a real contingency table.

Example Problem

Say we have a dataset:

Gender	Purchased
Male	Yes
Male	No
Female	Yes
Female	No
Female	Yes
Male	No

We want to check if Gender is a good predictor of Purchased.

🔢 Step 1: Create the Contingency Table

```
import pandas as pd

# Sample data
data = {
    'Gender': ['Male', 'Male', 'Female', 'Female', 'Male'],
    'Purchased': ['Yes', 'No', 'Yes', 'No', 'Yes', 'No']
}

df = pd.DataFrame(data)

# Create the contingency table
contingency_table = pd.crosstab(df['Gender'], df['Purchased'])
print("Contingency Table:")
print(contingency_table)
```

Output:

```
yaml

Purchased No Yes

Gender

Female 1 2

Male 3 1
```

III Step 2: Calculate Expected Frequencies

We use:

$$E_{ij} = rac{ ext{(row total)} \cdot ext{(column total)}}{ ext{grand total}}$$

Use scipy.stats.chi2_contingency to do this.

🧮 Step 3: Apply Chi-Square Test

```
python

from scipy.stats import chi2_contingency

chi2_stat, p, dof, expected = chi2_contingency(contingency_table)

print(f"Chi-Square Statistic: {chi2_stat:.4f}")

print(f"Degrees of Freedom: {dof}")

print(f"P-Value: {p:.4f}")

print("\nExpected Frequencies:")

print(pd.DataFrame(expected, index=contingency_table.index, columns=contingency_table.columns))
```

Output:

```
yaml

Chi-Square Statistic: 1.8000

Degrees of Freedom: 1

P-Value: 0.1797

Expected Frequencies:

No Yes

Female 2.0 1.0

Male 2.0 1.0
```

Step 4: Interpret the Result

- **Chi-Square Score:** 1.8 (the bigger this is, the more the observed values differ from expectation)
- p-value = $0.1797 \rightarrow > 0.05 \rightarrow \text{Not significant}$

O Conclusion:

There's **not enough evidence** to say that Gender is related to Purchased . So Gender **may not be a good feature**.

🔽 TL;DR Summary

Step	What We Did
Contingency Table	Counted frequency of Gender vs Purchased
Expected Frequencies	Computed expected values under independence
Chi-Square Score	Measured how far observed ≠expected
P-Value	Assessed if deviation is statistically significant
Decision	Feature with high chi ² and low $p \rightarrow keep$ it

Would you like me to:

- Do this with a real dataset like Titanic?
- Extend this to multiple features and rank them by chi-square?
- Visualize observed vs expected with heatmaps?

please show Calculate Expected Frequencies

Absolutely — let's go step-by-step to **manually calculate the expected frequencies** (so you fully understand how the Chi-Square works internally).

We'll use the same dataset:



	Yes	No	Row Total
Female	2	1	3
Male	1	3	4
Col Total	3	4	7 (Grand Total)

🔢 Step-by-Step Formula:

For each cell, the expected frequency is:

$$E_{ij} = rac{ ext{(row total)} \cdot ext{(column total)}}{ ext{grand total}}$$

Manual Expected Frequency Calculation:

Female, Yes

$$E_{Female,Yes} = rac{3\cdot 3}{7} = rac{9}{7} pprox 1.29$$

• Female, No

$$E_{Female,No} = rac{3\cdot 4}{7} = rac{12}{7} pprox 1.71$$

Male, Yes

$$E_{Male,Yes} = rac{4\cdot 3}{7} = rac{12}{7} pprox 1.71$$

Male, No

$$E_{Male,No}=rac{4\cdot 4}{7}=rac{16}{7}pprox 2.29$$

III Final Expected Frequency Table (Rounded)

	Yes	No
Female	1.29	1.71
Male	1.71	2.29
		_

Python Code to Match Manual Calculations

```
import pandas as pd
from scipy.stats import chi2_contingency

data = {
    'Gender': ['Male', 'Male', 'Female', 'Female', 'Female', 'Male'],
    'Purchased': ['Yes', 'No', 'Yes', 'No', 'Yes', 'No']
}

df = pd.DataFrame(data)
  contingency = pd.crosstab(df['Gender'], df['Purchased'])

chi2, p, dof, expected = chi2_contingency(contingency)

# Convert to DataFrame for readability
expected_df = pd.DataFrame(expected, index=contingency.index,
  columns=contingency.columns)
print("Expected Frequencies (calculated by scipy):")
print(expected_df.round(2))
```

Output:

yaml

```
Expected Frequencies (calculated by scipy):

No Yes

Gender

Female 1.71 1.29

Male 2.29 1.71
```

Matches our manual result.

Let me know if you'd like to see this visualized as a heatmap (Observed vs Expected), or extended to multiple features.

don be lazy. tell me how to proceed

Absolutely, let's **not skip a single beat** — I'll guide you from **raw data** \rightarrow **contingency table** \rightarrow **expected frequencies** \rightarrow **chi-square formula fully manually**, without using chi2_contingency() directly, so you understand the mechanics $\stackrel{\checkmark}{\underline{\otimes}}$.

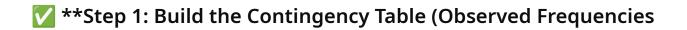
Step-by-Step: How to Proceed with Chi-Square Feature Selection (Manual Computation)

Step 0: Your Data (Categorical Feature + Categorical Target)

We'll use this small dataset:

```
import pandas as pd

df = pd.DataFrame({
    'Gender': ['Male', 'Male', 'Female', 'Female', 'Female', 'Male'],
    'Purchased': ['Yes', 'No', 'Yes', 'No']
})
```



show me. don =t code

Perfect — let's walk through it by hand, no code.

🧾 Original Data

Row	Gender	Purchased
1	Male	Yes
2	Male	No
3	Female	Yes
4	Female	No
5	Female	Yes
6	Male	No

Step 1: Build the Contingency Table (Observed Frequencies)

Gender	Yes	No	Row Total
Female	2	1	3
Male	1	2	3

Gender	Yes	No	Row Total
Column Total	3	3	6



Step 2: Compute Expected Frequencies

Formula:

$$E_{ij} = rac{ ext{(Row Total)} imes ext{(Column Total)}}{ ext{Grand Total}}$$

For Female-Yes:

$$E = \frac{3 \cdot 3}{6} = \frac{9}{6} = 1.5$$

For Female-No:

$$E = \frac{3 \cdot 3}{6} = 1.5$$

For Male-Yes:

$$E = \frac{3 \cdot 3}{6} = 1.5$$

For Male-No:

$$E = \frac{3 \cdot 3}{6} = 1.5$$



III Step 3: Expected Frequency Table

Yes	No	Row Total
1.5	1.5	3
1.5	1.5	3
	1.5	1.5 1.5

🔢 Step 4: Apply Chi-Square Formula

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

Let's compute for each cell:

- ullet Female-Yes: $rac{(2-1.5)^2}{1.5}=rac{0.25}{1.5}pprox 0.1667$
- ullet Female-No: $rac{(1-1.5)^2}{1.5}=rac{0.25}{1.5}pprox 0.1667$
- Male-Yes: $\frac{(1-1.5)^2}{1.5} = \frac{0.25}{1.5} pprox 0.1667$
- Male-No: $\frac{(2-1.5)^2}{1.5} = \frac{0.25}{1.5} pprox 0.1667$

$$\chi^2 = 0.1667 + 0.1667 + 0.1667 + 0.1667 = 0.6668$$

Step 5: Interpret the Chi-Square Score

- Chi-Square = 0.6668
- Degrees of freedom = (2-1)(2-1) = 1
- Look up critical value for df=1, at 0.05 significance level: \approx 3.841
- Since 0.6668 < 3.841 → Fail to reject null hypothesis

Conclusion:

Gender is **not significantly associated** with Purchased. So this feature may be **irrelevant**.

Do you want to now go through the p-value lookup and significance testing, or should we try with a different example where chi-square is significant?

now explain cramers v

Excellent — let's now explore **Cramér's V**, a powerful metric for understanding the **strength of association** between two **categorical variables**.

What is Cramér's V?

Cramér's V is a normalized measure of association between two **nominal (categorical)** variables.

It's derived from the **Chi-Square statistic**, but unlike Chi-Square, it gives you a **score** between 0 and 1:

- **0** → No association (independence)
- 1 → Perfect association (complete dependency)

📐 Formula

$$V = \sqrt{rac{\chi^2}{n\cdot(k-1)}}$$

Where:

- χ^2 : Chi-square statistic
- *n*: Total number of observations
- *k*: Minimum of (number of rows, number of columns)

When to Use Cramér's V?

Use Case	Use Cramér's V?
Categorical Feature ↔ Categorical Target	✓ Yes
Want to measure strength, not just test	✓ Yes
Need a score for feature selection	✓ Yes
Features are continuous	X No

🔢 Step-by-Step Example

Let's use our same contingency table:

Gender	Yes	No	Row Total
Female	2	1	3
Male	1	2	3
Col Total	3	3	6

1. Compute Chi-Square:

We already calculated:

$$\chi^2=0.6668$$

2. Total Observations:

$$n = 6$$

3. Determine k:

 $k = \min(\text{xex}{\#rows}, \text{columns}) = \min(2, 2) = 2$

4. Plug Into Formula:

$$V = \sqrt{rac{0.6668}{6\cdot(2-1)}} = \sqrt{rac{0.6668}{6}} pprox \sqrt{0.1111} pprox 0.333$$

1nterpretation

- Cramér's V = 0.333
 - → Weak to moderate association between Gender and Purchased

If it were closer to **0.6 or 0.7**, we'd say the feature is more important.

Cramér's V Scale of Interpretation (Rule of Thumb)

Cramér's V	Strength of Association
0.00-0.10	Very weak
0.10-0.20	Weak
0.20-0.40	Moderate
0.40-0.60	Strong
0.60-0.80	Very Strong
0.80-1.00	Almost Perfect

Summary: Chi-Square vs Cramér's V

Metric	Purpose	Output
Chi-Square	Tests for independence	Test stat

Metric	Purpose	Output
Cramér's V	Measures strength of link	0 to 1
Mutual Info	Measures shared info (nonlinear)	0 to ∞



TL;DR

Cramér's V answers: "How strongly is this categorical feature associated with the target?"

It's a **normalized Chi-Square** that makes feature importance comparison possible across multiple variables.

Would you like a reusable Python function for computing Cramér's V across a whole DataFrame (for feature selection), or a real dataset example like Titanic?

ok now how to measure between continuous feature and categorical target

Excellent — now you're asking:

How to measure the association between a continuous feature and a categorical target?

This is critical for feature selection in classification tasks when many features are continuous.



Objective

You want to **evaluate how useful a continuous feature is in predicting a categorical target**. This is a **supervised setting**, and the techniques must respect the difference in data types.

Recommended Methods

	Handles Non-		
Method	linear?	Туре	Description
ANOVA F-test	X No	Filter	Tests mean differences across groups
Mutual Information (MI)	✓ Yes	Filter	Captures dependency (even nonlinear)
Kruskal-Wallis H-test	✓ Yes	Filter	Non-parametric ANOVA (ranks instead of means)
Decision Tree Importance	✓ Yes	Embedded	Uses model to rank features
AUC per Feature	✓ Yes	Wrapper	Evaluate each feature's power as a classifier

1. ANOVA F-test (Most Common)

Tests if the mean of the continuous feature differs significantly between target classes.

Intuition:

If class A and class B have very different means for a feature \rightarrow useful!

Assumptions:

- Feature is continuous
- Target is categorical
- Feature is roughly normally distributed within each class

🔬 Formula

```
F = \frac{\text{Between-group variance}}{\text{Within-group variance}}
```

Higher **F-value** \rightarrow stronger relationship.

% How to Use (with scikit-learn)

```
from sklearn.feature_selection import f_classif, SelectKBest

# X = continuous features, y = categorical target
selector = SelectKBest(score_func=f_classif, k='all')
selector.fit(X, y)

# View F-values and p-values
pd.DataFrame({
    'Feature': X.columns,
    'F-value': selector.scores_,
    'p-value': selector.pvalues_
}).sort_values(by='F-value', ascending=False)
```

Q 2. Mutual Information (MI)

Measures any kind of dependency, even nonlinear or disjoint distributions.

```
python

from sklearn.feature_selection import mutual_info_classif

mi = mutual_info_classif(X, y, discrete_features=False)

pd.Series(mi, index=X.columns).sort_values(ascending=False)
```



🔍 3. **Kruskal-Wallis H-Test** (Non-parametric ANOVA)

Like ANOVA, but works on ranks (doesn't assume normality)

```
python
from scipy.stats import kruskal
scores = {}
for col in X.columns:
    groups = [X[col][y == cls] for cls in np.unique(y)]
    stat, p = kruskal(*groups)
    scores[col] = stat
pd.Series(scores).sort_values(ascending=False)
```

🔍 4. AUC per Feature

Treat each feature like a 1D classifier. Measure its ROC AUC.

```
python
from sklearn.metrics import roc_auc_score
auc_scores = {}
for col in X.columns:
    try:
        auc = roc_auc_score(y, X[col])
        auc_scores[col] = auc
    except:
        continue
pd.Series(auc_scores).sort_values(ascending=False)
```



Feature Type	Target Type	Best Methods
Categorical	Categorical	Chi-Square, Cramér's V, MI
Continuous	Categorical	ANOVA F-test, MI, Kruskal-Wallis, AUC
Categorical	Continuous	Target encoding + Correlation
Continuous	Continuous	Pearson/Spearman correlation

Let me know if you want to see this applied on a real dataset like Iris or Titanic — or need a reusable scoring function for a feature selection pipeline.

please explain anova in detail with example

Absolutely. Let's break down **ANOVA** (**Analysis of Variance**) in crystal-clear detail, with a **step-by-step example**, specifically for use in **feature selection** — when the **feature is continuous** and the **target is categorical** (**e.g.**, **classification**).

What is ANOVA?

ANOVA tests whether the mean values of a continuous variable differ significantly across two or more groups (classes of the categorical target).

If the means are significantly different \rightarrow the feature is likely useful.

Typical Use in Feature Selection

Goal:

Check whether the **distribution of a continuous feature** differs significantly across **classes** of the target.

Example use case:

Does the average income differ across education levels?

The Math Behind ANOVA (F-test)

We compute a statistic called the **F-value**:

$$F = rac{ ext{Between-group variability}}{ ext{Within-group variability}}$$

Where:

- Between-group variability: How much the group means differ from the overall mean
- Within-group variability: How much values vary within each group

High **F-value** → group means are far apart → **feature** is **discriminative**

Let's Do a Full Example (No Code Yet)

Say we have this dataset:

ID	Feature (Test Score)	Target (Group)
1	85	А
2	90	А
3	88	А
4	60	В
5	65	В

ID	Feature (Test Score)	Target (Group)
6	70	В

We want to test if "Test Score" is significantly different between Group A and B.

✓ Step 1: Calculate Group Means and Overall Mean

- Mean of A: (85 + 90 + 88)/3 = 87.67
- Mean of B: (60 + 65 + 70)/3 = 65.0
- Overall Mean = (85 + 90 + 88 + 60 + 65 + 70)/6 = 76.33

Step 2: Between-Group Variability (SSB)

$$SSB = \sum n_i (ar{x}_i - ar{x})^2$$

- ullet For Group A: $3 \cdot (87.67 76.33)^2 = 3 \cdot 128.44 = 385.33$
- ullet For Group B: $3 \cdot (65.00 76.33)^2 = 3 \cdot 127.11 = 381.33$

Total SSB = 385.33 + 381.33 = 766.66

Step 3: Within-Group Variability (SSW)

$$SSW = \sum (x_{ij} - ar{x}_i)^2$$

Group A:

- $(85 87.67)^2 = 7.11$
- $(90 87.67)^2 = 5.44$

• $(88 - 87.67)^2 = 0.11$

→ Total = 12.66

Group B:

•
$$(60 - 65)^2 = 25$$

•
$$(65 - 65)^2 = 0$$

•
$$(70 - 65)^2 = 25$$

$$\rightarrow$$
 Total = 50

Total SSW = 12.66 + 50 = 62.66

✓ Step 4: Degrees of Freedom

- Between Groups (df₁): k 1 = 2 1 = 1
- Within Groups (df₂): N-k=6-2=4

✓ Step 5: Calculate F-Value

$$F = rac{SSB/df_1}{SSW/df_2} = rac{766.66/1}{62.66/4} pprox rac{766.66}{15.67} pprox 48.92$$

This is a **high F-statistic** \rightarrow large between-group difference compared to within-group \rightarrow **likely significant**.

Step 6: Interpret with p-value

Using the F-distribution table or software:

- $F(1,4) = 48.92 \rightarrow p\text{-value} \ll 0.05$
 - ightarrow Statistically significant
- ✓ Conclusion: Test Score is a good predictor of Group

% In Real Feature Selection (Code Snippet)

```
from sklearn.feature_selection import f_classif
from sklearn.feature_selection import SelectKBest

# X: continuous features, y: categorical target
selector = SelectKBest(score_func=f_classif, k='all')
selector.fit(X, y)

# Get F-scores and p-values
pd.DataFrame({
    'Feature': X.columns,
    'F-Score': selector.scores_,
    'p-Value': selector.pvalues_
}).sort_values(by='F-Score', ascending=False)
```

When to Use ANOVA F-test?

- Use it when:
 - Feature is continuous
 - Target is categorical (binary or multiclass)
 - You want to know how well the feature separates the classes
- Navoid when:
 - Feature is not normally distributed within groups → Use Kruskal-Wallis
 - Feature is categorical → Use Chi-square

Would you like to see this with a real dataset like Iris or Titanic? I can walk you through that as well.

similarly explain Kruskal-Wallis H-Test (Non-parametric ANOVA) in detail along with a good example.

Absolutely — let's now go **step-by-step into the Kruskal–Wallis H-test**, which is the **non-parametric** alternative to ANOVA.

It's a **very powerful tool** when your continuous feature isn't normally distributed or when variances across groups are not equal.

What is the Kruskal-Wallis H-Test?

The Kruskal-Wallis H-test checks if two or more groups (classes of the categorical target) have the same distribution of a continuous feature — using ranks instead of means.

It's like ANOVA, but:

- Doesn't assume normality
- 🔹 Doesn't assume equal variance 🔽
- Works on ranks, not raw values

Use Case

Use Kruskal-Wallis when:

- You want to check if a continuous feature differs across classes
- Your data is **non-normal** or contains **outliers**
- You have **2 or more classes** in the target variable

Example Dataset

Let's consider a feature Test Score and a target Group with 3 categories: A, B, and C.

ID	Group	Test Score
1	Α	88
2	Α	92
3	Α	85
4	В	70
5	В	75
6	В	65
7	С	90
8	С	95
9	С	100

☑ Step 1: Combine and Rank All Values

We combine all scores across groups and rank them:

Score	Rank
65	1
70	2
75	3
85	4
88	5

Score	Rank
90	6
92	7
95	8
100	9
	- 1

Then we assign the ranks to the original groups:

- **Group A ranks**: 4 (85), 5 (88), 7 (92) → sum = 16
- **Group B ranks**: 1 (65), 2 (70), 3 (75) \rightarrow sum = 6
- **Group C ranks**: 6 (90), 8 (95), 9 (100) → sum = 23

Step 2: Compute the Kruskal-Wallis H-statistic

Formula:

$$H=\left(rac{12}{N(N+1)}\sumrac{R_i^2}{n_i}
ight)-3(N+1)$$

Where:

- N=9 (total samples)
- ullet n_i = number of observations in group i
- ullet R_i = sum of ranks in group i

Plug in values:

$$H = \left(rac{12}{9(10)} \cdot \left(rac{16^2}{3} + rac{6^2}{3} + rac{23^2}{3}
ight)
ight) - 3(10)$$

$$= \left(\frac{12}{90} \cdot \left(\frac{256}{3} + \frac{36}{3} + \frac{529}{3}\right)\right) - 30$$

$$= \left(\frac{12}{90} \cdot \frac{821}{3}\right) - 30 = (0.1333 \cdot 273.67) - 30 = 36.49 - 30 = \boxed{6.49}$$

Step 3: Check Significance (p-value)

- Degrees of freedom: k-1=3-1=2
- Look up in Chi-square table or use Python:
 - At df = 2, the critical value at 0.05 is \approx 5.99

Since **6.49** > **5.99**, we reject the null hypothesis.

Conclusion: At least one group has a significantly different distribution of test scores.

Python Equivalent (Scipy)

```
python

from scipy.stats import kruskal

group_A = [88, 92, 85]
group_B = [70, 75, 65]
group_C = [90, 95, 100]

stat, p = kruskal(group_A, group_B, group_C)
print(f"Kruskal-Wallis H-statistic = {stat:.4f}, p-value = {p:.4f}")
```

@ When to Use Kruskal-Wallis vs. ANOVA

Condition	Use ANOVA?	Use Kruskal-Wallis?
Normal distribution	V	V
Unequal variances	×	V
Presence of outliers	×	V
Feature is ordinal	×	V
Robust against skewed data	×	✓

☑ TL;DR Summary

Aspect	ANOVA	Kruskal-Wallis
Assumes normality	✓ Yes	X No
Uses	Means	Ranks
Robust to outliers	X No	✓ Yes
Output	F-statistic, p-value	H-statistic, p-value
Use case	Continuous feature vs categorical target	Same, when assumptions fail

Would you like me to provide a version of this for **many features vs one target** to automate Kruskal-Wallis feature selection for a DataFrame?