

- งานทั่วไป

- /noe งานบ้าน

DS : be able to perform ML pipeline

- data prep (impute missing, train/test split)
- traditional ML Decision tree, Linear, Logistic, GridSearch, Pipeline, k-means, SVM
- model evaluation / parameter tuning
- MLflow

Visual : Hint ข้อสอบส่วน อ.วีระ ฝากมาแจ้งนะครับ

ข้อสอบส่วนของ อ.วีระ

ส่วนแรก Spatial analysis จะให้นิสิตวิเคราะห์ข้อมูลเชิงพื้นที่ โดยจะมีโปรแกรม ให้บางส่วน นิสิตจะต้องเขียน โปรแกรมเพิ่มเติมเพื่อตอบคำถาม สิ่งที่ควรรู้คือ streamlit, pydeck, dbscan

ส่วนที่สอง network analysis จะให้นิสิตวิเคราะห์ข้อมูลเครือข่าย โดยใช้โปรแกรม gephi โดยให้ตอบคำถามคล้ายๆ กับที่ให้โจทย์ตัวอย่างไว้ใน github ของรายวิชา (ใน folder code)

MLFlow

Spark, RedDB, Kafka, Web

Gephi, Spatial analysis (pydeck, dbscan)

+ (1) One-hot vector (dummy codes)

■ Dummy coding = (n-1) dummy codes

Branch	BranchNum	D_BKK	D_Patum	D_Non
BKK	1	1	0	0
Patumtani	2	0	1	0
Nontaburi	3	0	0	1

X

Y/BK/BK

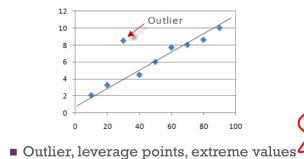
3 level = 2 dummy
ถ้าค่าที่ 3

reference

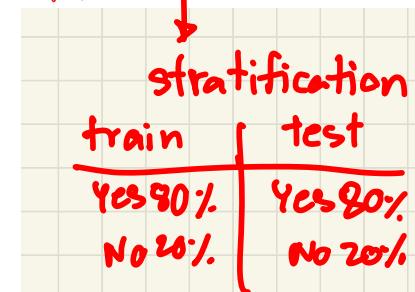
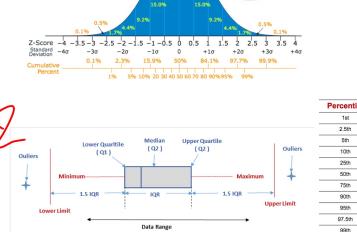


+ 3) Preparing features for ML (cont.):

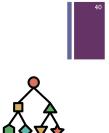
3.3 Truncate outliers



truncation

 $\mu \pm 3\sigma$ 

+ K-Fold Cross Validation
How to fix overfitting issue on test



iteration 1: Test Train Train Train Train
iteration 2: Train Test Train Train Train
iteration 3: Train Train Test Train Train
iteration 4: Train Train Train Test Train
iteration 5: Train Train Train Train Test

overall performance = mean(fold) = 80.8%

<https://towardsdatascience.com/cross-validation-explained-evaluating-estimator-performance-e51e5430f88>

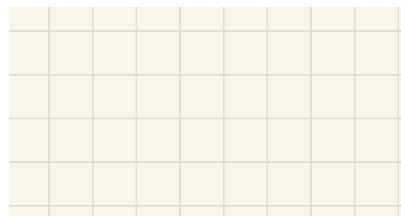
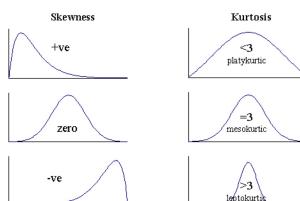
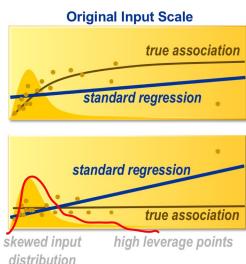
+ 3) Preparing features for ML (cont.):

3.4 Feature transformation

Skewness $\rightarrow \text{true} : \text{standard}$

Example: Salary, Balance in bank account

Solutions: Log, Binning



+ 3) Preparing features for ML (cont.):

3.4 Feature transformation - Log

from $y =$
GAVW $=$

Spending	Spending with outliers	LOG10(Spending)	LOG10(Spending with outliers)
2,500.00	2,500.00	3.40	3.40
2,900.00	2,900.00	3.46	3.46
3,200.00	3,200.00	3.51	3.51
4,000.00	4,000.00	3.60	3.60
4,500.00	4,500.00	3.65	3.65
6,200.00	6,200.00	3.79	3.79
10,000,000.00			7.00
mean	3,883.33	4,131,900.00	3.57
			4.06

+ 3) Preparing features for ML (cont.):

3.5 Feature engineering (cont.)

- Feature engineering
- Calculated variables
- Behavior from transactional data (RFM/RFA)



Data leaking from testing data in training data

યોગ્ય વિભાગીની ટેસ્ટ ડેટા ના ટ્રેન ડેટા-ટ૆સ્ટ

- Split by **subjects/videos** rather than individual images
- For **time series data**, split by period of time

ટ્રેન ડેટાનું કંતાર
અનેકોડો ટ્રેન ડેટાનું



Remark: Random Seed

- The experiment must be able to reconstruct (replicate).

- All randoms must be assigned a **random seed**.
 - `random.seed(12345)`
 - `random_state` option

જો seed હોય

Feature selection Helper

```
python Copy code
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import f_classif # or any other scoring function

# Create SelectKBest object with scoring function
selector = SelectKBest(score_func=f_classif, k=3) # Select top 3 features

# Fit the selector to your data
selector.fit(X, y)

# Transform your data to include only the selected features
X_selected = selector.transform(X)
```

parameters

```
python Copy code
from sklearn.feature_selection import SelectPercentile
from sklearn.feature_selection import f_classif # or any other scoring function

# Create SelectPercentile object with scoring function
selector = SelectPercentile(score_func=f_classif, percentile=20) # Keep top 20% of features

# Fit the selector to your data
selector.fit(X, y)

# Transform your data to include only the selected features
X_selected = selector.transform(X)
```

```
python Copy code
from sklearn.feature_selection import SelectKBest, f_classif, mutual_info_classif

# Create SelectKBest objects with different scoring functions
selector_f_classif = SelectKBest(score_func=f_classif, k=3)
selector_chi2 = SelectKBest(score_func=chi2, k=3)
selector_mutual_info = SelectKBest(score_func=mutual_info_classif, k=3)

# Fit and transform using each selector
X_selected_f_classif = selector_f_classif.fit_transform(X, y)
X_selected_chi2 = selector_chi2.fit_transform(X, y)
X_selected_mutual_info = selector_mutual_info.fit_transform(X, y)
```

```
python Copy code
from sklearn.feature_selection import SelectKBest, f_regression, mutual_info_regression

# Create SelectKBest objects with different scoring functions
selector_f_regression = SelectKBest(score_func=f_regression, k=3)
selector_mutual_info_regression = SelectKBest(score_func=mutual_info_regression, k=3)

# Fit and transform using each selector
X_selected_f_regression = selector_f_regression.fit_transform(X, y)
X_selected_mutual_info_regression = selector_mutual_info_regression.fit_transform(X, y)
```

Chi-square (χ^2) is a statistical test commonly used to determine whether there is a significant association between two categorical variables. It assesses whether the observed frequencies of categorical data differ from the expected frequencies under the assumption of independence between the variables.

Classification Problems:

For classification tasks, where the target variable is categorical, some common `score_func` options are:

→ `f_classif`:

Computes the ANOVA F-value for the provided samples.
Useful when the features are numeric.

Assumes the features are normally distributed and independent.

→ `chi2`:

Computes the chi-squared statistic between each non-negative feature and the target.
Suitable for categorical features.

Assumes the features are non-negative.

→ `mutual_info_classif`:

Estimates mutual information for a discrete target variable.
Works for both continuous and categorical features.

Captures any kind of statistical dependency, linear or non-linear.

Regression Problems:

For regression tasks, where the target variable is continuous, some common `score_func` options are:

→ `f_regression`:

Computes the F-value between each feature and the target.
Suitable for linear regression tasks.

Assumes the features are normally distributed and independent.

→ `mutual_info_regression`:

Estimates mutual information for a continuous target variable.
Useful for capturing non-linear dependencies.

Works for both continuous and categorical features.

ML

Linear Regression

```

import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score

# Generate random data
np.random.seed(42)
X = np.random.rand(100, 1) * 10
y = 3 * X.squeeze() + np.random.randn(100) * 2 # True relationship: y = 3*X + noise

# Split data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Create and train the model
model = LinearRegression().fit(X_train, y_train)

# Make predictions
y_pred = model.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print("Mean Squared Error (MSE):", mse)
print("R-squared (R²) Score:", r2)

```

Confusion Matrix

		Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60	
Actual: YES	FN = 5	TP = 100	105	
	55	110		

■ True Positive (TP)
 ■ True Negative (TN)
 ■ False Positive (FP)
 ■ False Negative (FN)

 ■ Accuracy = $(TP + TN) / \text{total}$
 ■ Misclassification = $(FP + FN) / \text{total}$



Precision, Recall, F1

		Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60	
Actual: YES	FN = 5	TP = 100	105	
	55	110		

■ Precision = correctly predict = $TP / (TP + FP)$
 ■ Recall = coverage = $TP / (TP + FN)$
 ■ F1 = $(2 * \text{precision} * \text{recall}) / (\text{precision} + \text{recall})$



Precision, Recall, F1: Average

		Predicted: NO	Predicted: YES	
Actual: NO	TN = 50	FP = 10	60	
Actual: YES	FN = 5	TP = 100	105	
	55	110		

■ What is a positive class?
 ■ 1) Direct target marketing? Yes
 ■ 2) Intelligent diagnosis to predict "Corona" or "Flu" - both are important.

 ■ If there is no positive (all classes are important)
 ■ Macro Average (Micro Average) by the amount of data in that class

 ■ Precision = correctly predict = $TP / (TP + FP)$
 ■ Recall = coverage = $TP / (TP + FN)$
 ■ F1 = $(2 * \text{precision} * \text{recall}) / (\text{precision} + \text{recall})$



Sequential Feature Selection (SFS) is a technique used for selecting a subset of features from a larger set of features based on their contribution to the performance of a machine learning model. It can be used forward, backward, or bidirectional.

Steps in Sequential Feature Selection:

- Initialization: Start with an empty set of features.
- Selection Step: Add one feature to the set based on a criterion (forward, backward, or bidirectional).
- Forward Selection: Start with an empty set and add the feature that provides the best improvement in model performance.
- Backward Selection: Start with all features and remove the feature that provides the worst improvement in model performance.
- Step Criterion: Sequential Feature Selection uses a step criterion to determine if a feature is added or removed.
- Stop Criterion: Sequential Feature Selection uses a stop criterion to determine when to stop the search process.
- Improves Model Performance: Sequential Feature Selection helps in identifying the most relevant features, potentially improving the model's performance.
- Excludes Overfitting: Helps in reducing overfitting by selecting only the most informative features.

Let's say we have a dataset with multiple features and a target variable. We can use forward selection as follows:

```

# Import scikit-learn
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.feature_selection import SequentialFeatureSelector as SFS

```

```

# Load Iris dataset
iris = load_iris()
X, y = iris.data, iris.target

```

```

# Split data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

```

```

# Create a simple KNN classifier model
model = KNeighborsClassifier()

```

```

# Fit the model
model.fit(X_train, y_train)

```

```

# Get the selected features indices and names
selected_features_idx = SFS(model, k_features=3).fit(X_train, y_train).get_support(indices=True)
selected_features_name = [feature_names[i] for i in selected_features_idx]

```

```

print('Selected Features Indices:', selected_features_idx)
print('Selected Features Names:', selected_features_name)

```

In this case, the selected features are Sepal Length, Sepal Width, and Petal Length.

We use the fit dataset and test it after training and testing.

We create a logistic regression model.

We use the selected features from model to perform forward selection.

The k_features parameter specifies the number of features to select (in this case, 3).

Important Notes:

Computational cost: Sequential feature selection can be computationally expensive, especially for large datasets or when using complex models.

Training Parameters: You might need to tune parameters such as the number of features to select, the scoring metric, and cross-validation settings.

Model Selection: Sequential Feature Selection is a powerful technique for feature selection, especially when dealing with datasets with many features.

Performance: Sequential Feature Selection can help in improving model performance, reducing overfitting, and getting insights into the most relevant features for the task at hand.

learn Install User Guide API Examples More »

scikit-learn 0.22.1 Other version

Please cite us if you use the software.

sklearn.metrics.confusion_matrix

sklearn.metrics.confusion_matrix

Examples using

sklearn.metrics.confusion_matrix

sklearn.metrics.confusion_matrix(x, y_true, y_pred, labels=None, sample_weight=None, normalize=None)

[source]

Compute confusion matrix to evaluate the accuracy of a classification.

By definition a confusion matrix C is such that C_{ij} is equal to the number of observations known to be in group i and predicted to be in group j.

Thus in binary classification, the count of true negatives is $C_{0,0}$, false negatives is $C_{0,1}$, true positives is $C_{1,1}$ and false positives is $C_{1,0}$.

Read more in the User Guide.

Parameters: x : array-like of shape (n_{samples} , n_{classes})

Ground truth (correct) target values.

y_{true} : array-like, optional

Labels to consider. By default all labels.

y_{pred} : array-like, optional

Prediction results. By default all labels.

$sample_weight$: array-like of shape (n_{samples} ,), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

normals : {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix. This may be used to reorder or select a subset of labels. If None is given, those once in x , y_{true} or y_{pred} are used in sorted order.

y -like of shape (n_{samples}), default=None

normals : {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix relative to the (row) samples, predicted (columns) conditions or all the population. If None, confusion matrix will not be normalized.

learn Install User Guide API Examples More »

scikit-learn 0.22.1 Other version

Please cite us if you use the software.

sklearn.metrics.classification_report

sklearn.metrics.classification_report

Examples using

sklearn.metrics.classification_report

sklearn.metrics.classification_report(y_true, y_pred, labels=None, target_names=None, sample_weight=None, output_dict=False, zero_division='warn')

[source]

Build a text report showing the main classification metrics

Read more in the User Guide.

Parameters: y_{true} : array-like, or label names

Ground truth (correct) target values.

y_{pred} : array-like, or label names

Estimated targets as returned by a classifier.

$labels$: array, shape = [n_{labels}]

Optional list of label indices to consider.

$target_names$: list of strings

Optional display names matching the labels.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

$normalize$: {“raw”, “per”, “off”, default=None}

Normalizes confusion matrix.

$sample_weight$: array-like of shape (n_{samples}), default=None

Sample weights.

$digits$: int

Number of digits for formatting the returned values will be $\text{np.set_printoptions}(\text{precision}=2)$.

$output_dict$: bool

If True, returns a dictionary.

$zero_division$: {“warn”, “0”, “1”}

Divide by zero error value in division by zero condition.

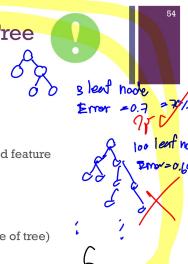
$normalize$: {“raw”, “per”, “off”, default=None}

Important Parameters in Decision Tree

Hyper Tuning

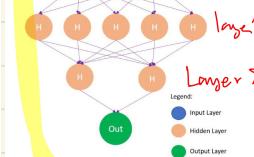
- Splitting measure (criterion): gini / entropy regularization
- Maximum depth: ~5-10 (depend on number of feature)
- Maximum leaf nodes: depend on number of class (target) and feature
- Minimum sample split: 5 - 20% (depend on number of data)
- Minimum impurity decrease: (default 0)
- Pruning adjustment (ccp_alpha): 0.001 - 0.01 (depend on size of tree)

Balance between: model performance & complexity



Example

tree_size = 1000



batch_size: int, default='auto'

Size of minibatches for stochastic optimizers. If the solver is 'lbfgs', the classifier will not use minibatch. When set to "auto", batch_size=min(200, n_samples).

learning_rate : {constant, invscaling, adaptive}, default='constant'

Learning rate schedule for weight updates.

```
mlp_clf = MLPClassifier(hidden_layer_sizes=(5,2),
                        max_iter = 300, activation = 'relu',
                        solver = 'adam')
```

```
param_grid = {
    'hidden_layer_sizes': [(150,100,50), (120,80,40), (100,50)],
    'max_iter': [50, 100, 150],
    'activation': ['tanh', 'relu'],
    'solver': ['sgd', 'adam'],
    'alpha': [0.0001, 0.05],
    'learning_rate': ['constant','adaptive'],}
```

```
grid = GridSearchCV(mlp_clf, param_grid, n_jobs=-1, cv=5)
grid.fit(trainX_scaled, trainY)
```

```
print(grid.best_params_)
```

Sure, let's break down CAP theorem and ACID in a beginner-friendly way:

CAP Theorem:

CAP theorem, also known as Brewer's theorem, deals with distributed data stores, like databases, in computer science. It states that it's impossible for a distributed system to simultaneously provide all three of the following guarantees:

1. **Consistency (C)**: Every read receives the most recent write or an error. In other words, all nodes in the system have the same data at the same time.

2. **Availability (A)**: Every request to the system gets a response, even if one of the nodes in the system fails.

3. **Partition tolerance (P)**: The system continues to operate despite network partitions (communication failures) between nodes.

Explanation:

- **Consistency (C)** means that all nodes in a distributed system have the same data at the same time. So if you write data to one node and then immediately read from another, you should get the same data, not something outdated.

- **Availability (A)** means that every request to the system gets a response, even if some nodes are down or unreachable. Users can always read and write to the system.

- **Partition Tolerance (P)** means that the system continues to operate even if there are network issues causing some nodes to be unreachable by others. It ensures that the system can handle network failures without data loss.

Implications:

- If you choose Consistency and Partition Tolerance (CP), then during a network partition, the system will stop taking requests to maintain consistency.

- If you choose Availability and Partition Tolerance (AP), then during a network partition, the system will still be available but may return data that is not fully up-to-date (as it prioritizes availability over consistency).

- If you choose Consistency and Availability (CA), then during a network partition, the system will sacrifice availability to ensure that all nodes have consistent data.

ACID:

ACID is a set of properties that ensure that database transactions are processed reliably. It stands for:

1. **Atomicity**: Ensures that either all operations within a transaction succeed, or if any operation fails, the entire transaction fails, and the database state is left unchanged.

2. **Consistency**: Guarantees that only valid data is written to the database. If a transaction starts with the database in a consistent state, it will end with the database in a consistent state.

3. **Isolation**: Transactions are executed in isolation from each other, so the result of a transaction is not visible to other transactions until it is committed. This prevents interference between transactions.

4. **Durability**: Once a transaction is committed, its changes are permanent and will not be lost, even in the event of a system crash.

Explanation:

- **Atomicity**: Think of this as an "all or nothing" rule. If you're transferring money from one account to another, either the full amount is moved successfully, or nothing happens at all. No partial transfers.

- **Consistency**: This ensures that your data remains in a valid state at all times. For example, if you have a rule that a person's age can't be negative, any transaction that tries to set a person's age to -5 would be rejected.

- **Isolation**: Transactions should not interfere with each other. If two transactions are happening at the same time, they should not impact each other until they are both completed. This prevents scenarios like one transaction reading data from another.

- **Durability**: Once a transaction is confirmed (like clicking "submit" on an order), the changes are saved permanently, even if the system crashes immediately after. So, you won't lose your order confirmation even if the power goes out.

Conclusion:

- **ACID** properties are mainly concerned with single-node (or single-database) transactions, ensuring reliability and integrity.

- **CAP theorem** deals with distributed systems, emphasizing the trade-offs between consistency, availability, and partition tolerance.

In practical terms, designers and engineers often need to balance these principles based on the needs of their application. Some applications may prioritize data consistency and be willing to sacrifice availability during network issues at times.

network regularity

- Degree
- Avg path
- Clustering coefficient

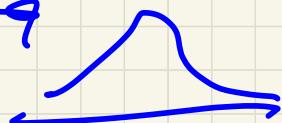
$$\text{clust} = \frac{\text{# triangles}}{\text{# triplets}}$$

$$0 < \text{clust} < 1$$

→ random network

Dist

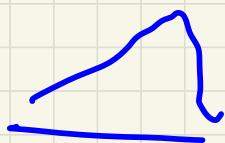
- (diff) random



- (more) regular

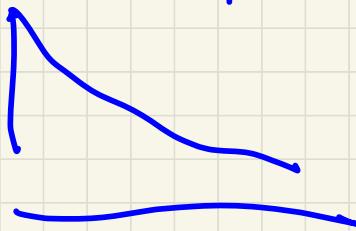
→ small-world network

Dist



→ Scale-Free

Common node has $\approx 1/\log n$
distribution



Central

Scale-Free degree central random

Small closeness center random

Small Between central
(and Nearest neighbor)
- Eigenvector (containing)

CSV = import Spreadsheet

Statistics → avg path length →
(v07) → avg Design
Dist

Modularity future
(v07)