# Database Programming

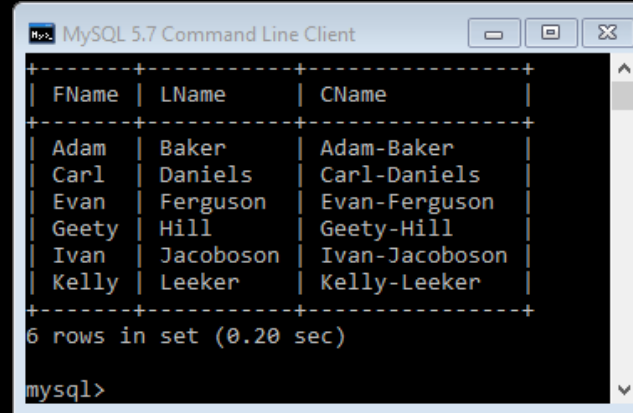## LECTURE 8: REGEX AND TEXT SEARCHING

Professor HG Locklear
hlocklear@pace.edu

# Pattern Matching

▶ SQL **pattern matching** enables you to use _ to match any single character and **%** to match an arbitrary number of characters (including zero characters).

▶ In MySQL, SQL patterns are case-insensitive by default.

▶ **LIKE 'b%'** anything that begins with **b**

▶ **LIKE '%fy'** anything that ends with **fy**

▶ **LIKE '%w%'** anything that contains **w**

▶ **LIKE '_____'** anything that has five characters

```
MySQL 5.7 Command Line Client                    [_][□][✕]
+--------+-----------+-----------------+
| FName  | LName     | CName           |
+--------+-----------+-----------------+
| Adam   | Baker     | Adam-Baker      |
| Carl   | Daniels   | Carl-Daniels    |
| Evan   | Ferguson  | Evan-Ferguson   |
| Geety  | Hill      | Geety-Hill      |
| Ivan   | Jacoboson | Ivan-Jacoboson  |
| Kelly  | Leeker    | Kelly-Leeker    |
+--------+-----------+-----------------+
6 rows in set (0.20 sec)

mysql>
```

```
MySQL 5.7 Command Line Client                    [_][□][✕]
Database changed
+--------+-----------+-----------------+
| FName  | LName     | CName           |
+--------+-----------+-----------------+
| Adam   | Baker     | Adam-Baker      |
+--------+-----------+-----------------+
1 row in set (0.04 sec)

+--------+-----------+-----------------+
| FName  | LName     | CName           |
+--------+-----------+-----------------+
| Adam   | Baker     | Adam-Baker      |
+--------+-----------+-----------------+
1 row in set (0.00 sec)

+--------+-----------+-----------------+
| FName  | LName     | CName           |
+--------+-----------+-----------------+
| Adam   | Baker     | Adam-Baker      |
| Carl   | Daniels   | Carl-Daniels    |
| Evan   | Ferguson  | Evan-Ferguson   |
| Ivan   | Jacoboson | Ivan-Jacoboson  |
+--------+-----------+-----------------+
4 rows in set (0.00 sec)

+--------+-----------+-----------------+
| FName  | LName     | CName           |
+--------+-----------+-----------------+
| Evan   | Ferguson  | Evan-Ferguson   |
| Ivan   | Jacoboson | Ivan-Jacoboson  |
+--------+-----------+-----------------+
2 rows in set (0.00 sec)

+--------+-----------+-----------------+
| FName  | LName     | CName           |
+--------+-----------+-----------------+
| Adam   | Baker     | Adam-Baker      |
+--------+-----------+-----------------+
```

```
30  SELECT * FROM codeNames WHERE FName LIKE 'A%' ;
31  SELECT * FROM codeNames WHERE FName LIKE 'a%' ;
32  SELECT * FROM codeNames WHERE FName LIKE '%a%';
33  SELECT * FROM codeNames WHERE FName LIKE '%n' ;
34  SELECT * FROM codeNames WHERE LName LIKE '_____';
```

# REGEX Operators

▶ A **Regular Expression** is a powerful way of specifying a pattern for a complex search.

▶ Regular Expression Operators

   ▶ **NOT REGEXP** (Negation of REGEXP)

   ▶ **REGEXP** (True if a String matches the regular expression...false otherwise)

   ▶ **RLIKE** (True if a String matches the regular expression...false otherwise)

▶ A regular expression pattern match succeeds **if the pattern matches anywhere** in the value being tested.

   ▶ **This differs from a LIKE pattern match, which succeeds only if the pattern matches the entire value.**

```
48   SELECT * FROM codeNames WHERE FName REGEXP 'l';
49   SELECT * FROM codeNames WHERE FName NOT REGEXP 'l';
50   SELECT * FROM codeNames WHERE FName RLIKE 'l';
51   SELECT * FROM codeNames WHERE FName REGEXP 'L'; --
```

MySQL 5.7 Command Line Client

```
+-------+----------+--------------+
| FName | LName    | CName        |
+-------+----------+--------------+
| Carl  | Daniels  | Carl-Daniels |
| Kelly | Leeker   | Kelly-Leeker |
+-------+----------+--------------+
2 rows in set (0.01 sec)

+-------+----------+----------------+
| FName | LName    | CName          |
+-------+----------+----------------+
| Adam  | Baker    | Adam-Baker     |
| Evan  | Ferguson | Evan-Ferguson  |
| Geety | Hill     | Geety-Hill     |
| Ivan  | Jacoboson| Ivan-Jacoboson |
+-------+----------+----------------+
4 rows in set (0.01 sec)

+-------+----------+--------------+
| FName | LName    | CName        |
+-------+----------+--------------+
| Carl  | Daniels  | Carl-Daniels |
| Kelly | Leeker   | Kelly-Leeker |
+-------+----------+--------------+
2 rows in set (0.00 sec)

+-------+----------+--------------+
| FName | LName    | CName        |
+-------+----------+--------------+
| Carl  | Daniels  | Carl-Daniels |
| Kelly | Leeker   | Kelly-Leeker |
+-------+----------+--------------+
2 rows in set (0.00 sec)

mysql>
```

# Regular Expressions

- A **Regular Expression** describes a set of strings.

- The simplest regular expression is one that has no special characters in it.

- For example, the regular expression **hello** matches **hello** and nothing else.

- Nontrivial regular expressions use certain special constructs so that they can match more than one string.

- For example, the regular expression **hello|world** contains the **|** alternation operator and matches either the **hello** or **world**.

```
63  SELECT * FROM codeNames WHERE FName REGEXP 'l|m';
```

```
MySQL 5.7 Command Line Client

+--------+----------+--------------+
| FName  | LName    | CName        |
+--------+----------+--------------+
| Adam   | Baker    | Adam-Baker   |
| Carl   | Daniels  | Carl-Daniels |
| Kelly  | Leeker   | Kelly-Leeker |
+--------+----------+--------------+
3 rows in set (0.00 sec)
```

# Regular Expressions

▶ As a more complex example, the regular expression **B[an]\*s** matches any of the strings Bananas, Baaaaas, Bs, and any other string starting with a B, ending with an s, and containing any number of a or n characters in between.

▶ **^** Match the beginning of a String

▶ **$** Match the end of a String

▶ **.** Match any character (including newline)

▶ **\*** Match any sequence of 0 or more characters

▶ **+** Match any sequence of one or more characters

▶ **?** Match either 0 or 1 character

▶ **|** Match either of the sequences

▶ **()\*** Match 0 or more instances of the sequence in the parentheses

**Alternative Notation**
**{n},{m,n}**

**a\*** could be written as **a{0,}**
**a+** could be written as **a{1,}**
**a?** could be written as **a{0,1}**

# Regular Expressions

```
80   SELECT * FROM codeNames WHERE FName REGEXP '^A';
81   SELECT * FROM codeNames WHERE FName REGEXP 'n$';
82   SELECT * FROM codeNames WHERE FName REGEXP '.';
83   SELECT * FROM codeNames WHERE FName REGEXP 'e*';
```

MySQL 5.7 Command Line Client

```
+-------+----------+-------------+
| FName | LName    | CName       |
+-------+----------+-------------+
| Adam  | Baker    | Adam-Baker  |
+-------+----------+-------------+
1 row in set (0.00 sec)

+-------+----------+----------------+
| FName | LName    | CName          |
+-------+----------+----------------+
| Evan  | Ferguson | Evan-Ferguson  |
| Ivan  | Jacoboson| Ivan-Jacoboson |
+-------+----------+----------------+
2 rows in set (0.00 sec)

+-------+----------+----------------+
| FName | LName    | CName          |
+-------+----------+----------------+
| Adam  | Baker    | Adam-Baker     |
| Carl  | Daniels  | Carl-Daniels   |
| Evan  | Ferguson | Evan-Ferguson  |
| Geety | Hill     | Geety-Hill     |
| Ivan  | Jacoboson| Ivan-Jacoboson |
| Kelly | Leeker   | Kelly-Leeker   |
+-------+----------+----------------+
6 rows in set (0.00 sec)

+-------+----------+----------------+
| FName | LName    | CName          |
+-------+----------+----------------+
| Adam  | Baker    | Adam-Baker     |
| Carl  | Daniels  | Carl-Daniels   |
| Evan  | Ferguson | Evan-Ferguson  |
| Geety | Hill     | Geety-Hill     |
| Ivan  | Jacoboson| Ivan-Jacoboson |
| Kelly | Leeker   | Kelly-Leeker   |
+-------+----------+----------------+
```
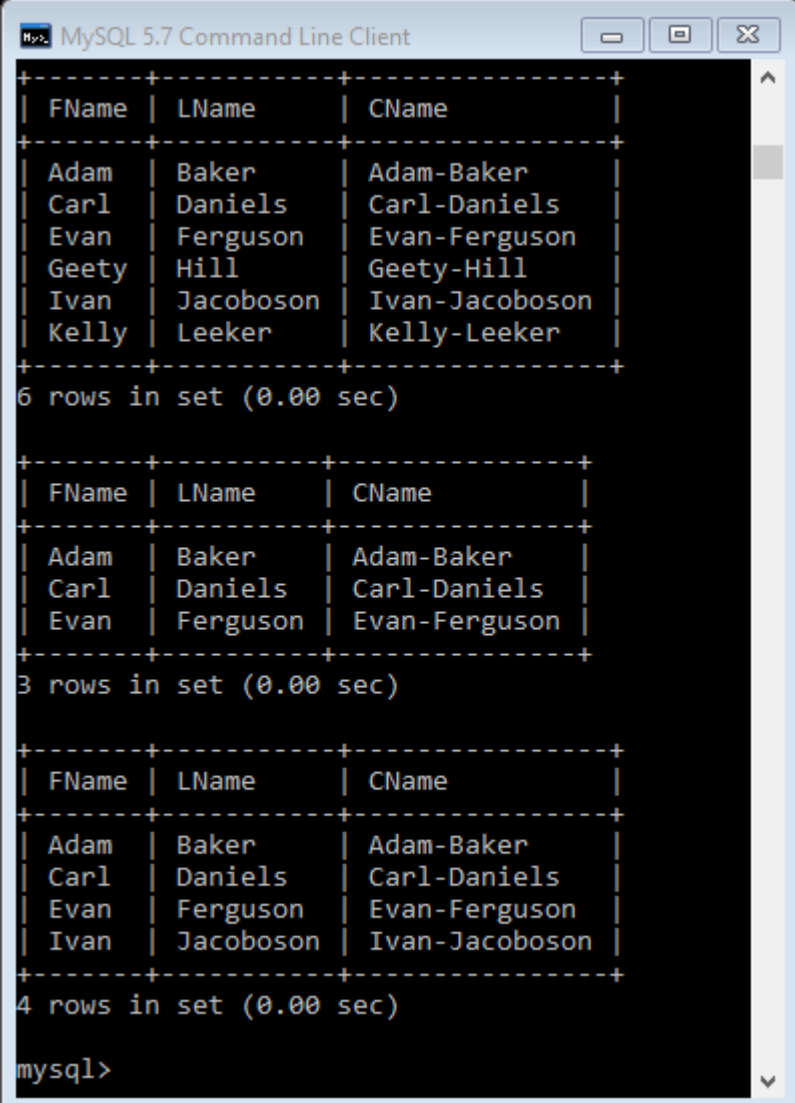
# Regular Expressions

```
84   SELECT * FROM codeNames WHERE FName REGEXP 'e+';  --
85   SELECT * FROM codeNames WHERE FName REGEXP 'e?';  -- z
86   SELECT * FROM codeNames WHERE FName REGEXP 'Ad|rl';
87   SELECT * FROM codeNames WHERE CName REGEXP '(rl-)*';
88   SELECT * FROM codeNames WHERE CName REGEXP '(rl-)+';
```

# Regular Expressions

```
96   SELECT * FROM codeNames WHERE FName REGEXP 'e{2,}';
97   SELECT * FROM codeNames WHERE FName REGEXP 'y{1,}';
98   SELECT * FROM codeNames WHERE FName REGEXP 'e{0,1}';
99   SELECT * FROM codeNames WHERE FName REGEXP 'e{0,2}';
100  SELECT * FROM codeNames WHERE FName REGEXP 'G{1,2}';
```

# Regular Expressions

```
107    SELECT * FROM codeNames WHERE FName REGEXP '^[A-K]';
108    SELECT * FROM codeNames WHERE FName REGEXP '^[a-e]';
109    SELECT * FROM codeNames WHERE FName REGEXP '[l-n]$';
```

# Full Text Searching

▶ In MySQL we have the concept of **full-text** searching.

▶ Full-Text searching is utilizing the MATCH function and the AGAINST operator to takes a comma-separated list that specifies the attributes to be searched and a string to search for and an optional modifier that indicates what type of search to perform.

▶ There are three types of full-text searches:

  ▶ **Natural Language Search**

    ▶ Interprets the search string as a phrase in natural human language.

    ▶ There are no special operators, with the exception of double quote (") characters.

  ▶ **Boolean Search**

    ▶ A boolean search interprets the search string using the rules of a special query language.

    ▶ It can also contain operators that specify requirements for the words being searched for.

  ▶ *Query Expansion Search (we will not explore this in our course)*

    ▶ Is a modification of a natural language search.

    ▶ The words from the most relevant rows returned by the search are added to the search string and the search is done again.

    ▶ The query returns the rows from the second search.

```
4
5  CREATE TABLE dna_SEQ(
6  ID INT AUTO_INCREMENT,
7  Sequence TEXT,
8  FULLTEXT (Sequence),
9  CONSTRAINT pk_dna PRIMARY KEY(ID)
10 );
11
12 INSERT INTO dna_SEQ (Sequence) VALUES('ATGU-GUAT-UUUU-TTAA');
13 INSERT INTO dna_SEQ (Sequence) VALUES('ATGU-GUAT-UTUU-TAAA');
14 INSERT INTO dna_SEQ (Sequence) VALUES('ATGU-GUAT-UUUU-corrupted');
15
16 SELECT *
17 FROM dna_SEQ
18 WHERE MATCH (Sequence) AGAINST ('corrupted' IN NATURAL LANGUAGE MODE);
19
20 SELECT *
21 FROM dna_SEQ
22 WHERE MATCH (Sequence) AGAINST ('TTAA' IN NATURAL LANGUAGE MODE);
```

**Must specify what full text is based on**

MySQL 5.7 Command Line Client

```
+----+--------------------------+
| ID | Sequence                 |
+----+--------------------------+
|  3 | ATGU-GUAT-UUUU-corrupted |
+----+--------------------------+
1 row in set (0.03 sec)

+----+--------------------------+
| ID | Sequence                 |
+----+--------------------------+
|  1 | ATGU-GUAT-UUUU-TTAA      |
+----+--------------------------+
1 row in set (0.00 sec)
```

```
   C:\Users\msgth\Desktop\AA_DBProgrammingFall2020\XML\searching.sql - Sublime Text (UNREGISTERED)
File  Edit  Selection  Find  View  Goto  Tools  Project  Preferences  Help

  xmltesting.sql           searching.sql

 4
 5  CREATE TABLE dna_SEQ(
 6  ID INT AUTO_INCREMENT,
 7  Sequence TEXT,
 8  Fragment TEXT,
 9  FULLTEXT (Sequence,Fragment),
10  CONSTRAINT pk_dna PRIMARY KEY(ID)
11  );
12
13  INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-UTY','ATGU-GUAT-UUUU-TTAA');
14  INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-YUT','ATGU-GUAT-UTUU-TAAA');
15  INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-XLN','ATGU-GUAT-UUUU-corrupted');
16
17  SELECT *
18  FROM dna_SEQ
19  WHERE MATCH (Sequence,Fragment) AGAINST ('XLN corrupted' IN NATURAL LANGUAGE MODE);
20
21  SELECT *
22  FROM dna_SEQ
23  WHERE MATCH (Sequence,Fragment) AGAINST ('VWS -TAAA' IN NATURAL LANGUAGE MODE);

Line 19, Column 46
```

**Must specify what full text is based on**

```
MySQL 5.7 Command Line Client

+----+----------+-------------------------+
| ID | Sequence | Fragment                |
+----+----------+-------------------------+
|  3 | VWS-XLN  | ATGU-GUAT-UUUU-corrupted |
+----+----------+-------------------------+
1 row in set (0.04 sec)


+----+----------+-------------------------+
| ID | Sequence | Fragment                |
+----+----------+-------------------------+
|  2 | VWS-YUT  | ATGU-GUAT-UTUU-TAAA     |
|  1 | VWS-UTY  | ATGU-GUAT-UUUU-TTAA     |
|  3 | VWS-XLN  | ATGU-GUAT-UUUU-corrupted |
+----+----------+-------------------------+
3 rows in set (0.00 sec)
```

# Boolean

```
4
5  CREATE TABLE dna_SEQ(
6  ID INT AUTO_INCREMENT,
7  Sequence TEXT,
8  Fragment TEXT,
9  FULLTEXT (Sequence,Fragment),
10 CONSTRAINT pk_dna PRIMARY KEY(ID)
11 );
12
13 INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-UTY','ATGU-GUAT-UUUU-TTAA');
14 INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-XLN','ATGU-GUAT-UTUU-TAAA');
15 INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-XLN','UTGU-GUAT-UUUU-corrupted');
16
17 SELECT *
18 FROM dna_SEQ
19 WHERE MATCH (Sequence,Fragment) AGAINST ('+XLN -corrupted' IN BOOLEAN MODE);
20
21 SELECT *
22 FROM dna_SEQ
23 WHERE MATCH (Sequence,Fragment) AGAINST ('+VWS +UTGU' IN BOOLEAN MODE);
```

```
+ include
- exclude
```

MySQL 5.7 Command Line Client

```
+----+----------+------------------------+
| ID | Sequence | Fragment               |
+----+----------+------------------------+
|  2 | VWS-XLN  | ATGU-GUAT-UTUU-TAAA    |
+----+----------+------------------------+
1 row in set (0.05 sec)

+----+----------+------------------------+
| ID | Sequence | Fragment               |
+----+----------+------------------------+
|  3 | VWS-XLN  | UTGU-GUAT-UUUU-corrupted |
+----+----------+------------------------+
```

# Scoring

```
score

0.062016263604164124
0.031008131802082062
                    0
```

```
 4
 5  CREATE TABLE dna_SEQ(
 6  ID INT AUTO_INCREMENT,
 7  Sequence TEXT,
 8  Fragment TEXT,
 9  FULLTEXT (Sequence,Fragment),
10  CONSTRAINT pk_dna PRIMARY KEY(ID)
11  ) ENGINE = INNODB;
12
13  INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-UTY','ATGU-GUAT-UUUU-TTAA');
14  INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-XLN','ATGU-GUAT-UTUU-TAAA');
15  INSERT INTO dna_SEQ (Sequence,Fragment) VALUES('VWS-XLN','UUUU-GUAT-UUUU-corrupted');
16
17  SELECT ID, Sequence, Fragment, MATCH (Sequence,Fragment) AGAINST ('UUUU' IN BOOLEAN MODE) AS score
18  FROM dna_SEQ
19  ORDER BY score DESC;
```

Line 17, Column 1

**Must specify search engine**

**Ranking System is based on TF-IDF**
*https://en.wikipedia.org/wiki/Tf%E2%80%93idf*
**Specifies how relevant this result is based on the word you were searching for**

MySQL 5.7 Command Line Client

```
+----+----------+--------------------------+----------------------+
| ID | Sequence | Fragment                 | score                |
+----+----------+--------------------------+----------------------+
|  3 | VWS-XLN  | UUUU-GUAT-UUUU-corrupted | 0.062016263604164124 |
|  1 | VWS-UTY  | ATGU-GUAT-UUUU-TTAA      | 0.031008131802082062 |
|  2 | VWS-XLN  | ATGU-GUAT-UTUU-TAAA      |                    0 |
+----+----------+--------------------------+----------------------+
3 rows in set (0.00 sec)
```