

Using Big Data to Get Us from Where We Think We Are to Where We Might Want to Go

Tidewater Big Data Enthusiasts
Chuck Cartledge
Developer

April 23, 2020 at 12:14 Noon

Contents

List of Figures	ii
List of Tables	ii
1 Introduction	1
2 Approach	1
2.1 Data sources	1
2.2 Database organization	4
2.3 Levenshtein distance algorithm	6
3 Results	11
4 Conclusion	19
A A PostgreSQL timeline	22
B Venn diagrams	23
C “The Dark Knight” actors	34
D Explaining SQL database concepts using Excel	38
E Misc. files	43

List of Figures

1	Barnes & Noble and Home Depot web pages.	2
2	Pareto chart of genres.	7
3	Movie poster for “Cupid Rides the Range”	20
4	Movie poster for “Dark Knight”	21
5	PostgreSQL Elephant.	22
6	Representative Venn diagrams.	24
7	Venn diagram of the Biography genre.	25
8	Venn diagram of the Biography and Crime genres (<i>Biography</i> \cup <i>Crime</i>).	26
9	Venn diagram of the Biography and Crime Western genres (<i>Biography</i> \cup <i>Crime</i> \cup <i>Western</i>).	27
10	Venn diagram of the Biography and Crime Western and Adventure genres (<i>Biography</i> \cup <i>Crime</i> \cup <i>Western</i> \cup <i>Adventure</i>).	28
11	Venn diagram of the Biography and Crime Western and Adventure and Documentary genres (<i>Biography</i> \cup <i>Crime</i> \cup <i>Western</i> \cup <i>Adventure</i> \cup <i>Documentary</i>). .	29
12	Colorful Venn diagram of Horror and Fantasy and Biography and Drama and Comedy.	30
13	An Excel workbook with PostgreSQL data (movies worksheet).	39
14	An Excel workbook with PostgreSQL data (actors worksheet).	40
15	An Excel workbook with PostgreSQL data (genre worksheet).	41
16	An Excel worksheet showing a the “Filter” interface.	42

List of Tables

1	Number of records downloaded from IMDb	4
2	PostgreSQL tables used for this demonstration.	4
3	Number of database entries by genre.	5
4	Levenshtein algorithm.	8
5	Levenshtein converting GUMBO to GAMBOL.	10
6	Normal output for “Copid rides the range”	11
7	Detailed output for “Copid rides the range”	13
8	Normal output for “The Dark Knight”	17
9	Choosing different genre combinations.	32
10	Pascal’s Triangle up to $n = 5$	33
11	Pascal’s Triangle up to $n = 5$ as binomials.	33
12	Pascal’s Triangle up to $n = 5$ reformatted coefficients.	33

13	Cross referencing Excel terms to PostgreSQL terms.	38
14	demo.R command line arguments.	43

1 Introduction

When we wander through the Internet, we leave behind digital traces of our passage. Some of our traces are embedded in our digital pictures[5], in our cell phone usage patterns[2], and from the cookies that we leave on our personal computers[3]. We are going to explore how some of these systems work as “recommenders” by using the Internet Movie Database(IMDb) to recommend movies that we might enjoy based on “misinformation” that we provide.

Our exploration will be to design a recommendation system that will understand and correct our misinformation, use the corrected information to identify specific relevant attributes, and using those attributes, recommend other new items we may be interested in. With data from the IMDb, we will:

1. Start with a misspelled movie title,
2. Have the database recommend a title that is close to our misspelled one,
3. Have the database identify which genre, or genres the movie belongs to,
4. Have the database combine the genres with a subset of the actors from the movie, and

to recommend other movies that we might be interested in viewing. These recommended movies will have some of the same actors from our original movie, and will belong to the same genre. We will be doing this using a venerable piece of structured query language (SQL) (see Section D on page 38) software called PostgreSQL (see Section A on page 22).

The idea being; if we liked the original movie, then we should like some of the recommended ones. This idea is behind many of the recommender systems that we use and observe everyday (see Figure 1 on the following page).

2 Approach

Now we will start delving into the details of how our movie recommendations system works, and conclude with some examples.

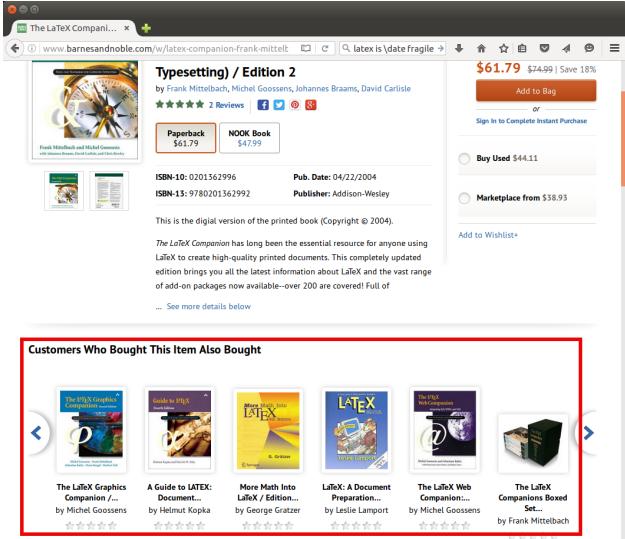
2.1 Data sources

Our raw movie database information is available for download from the Internet Movie Database (IMDb) website¹. Rather than download individual entries from the IMDb, we will download large files specific to data types of interest from the IMDb Alternative Interfaces website². The IMDb Alternative Interfaces in turn points to an FTP site where we can download data types of interest³. The three types of data we need are:

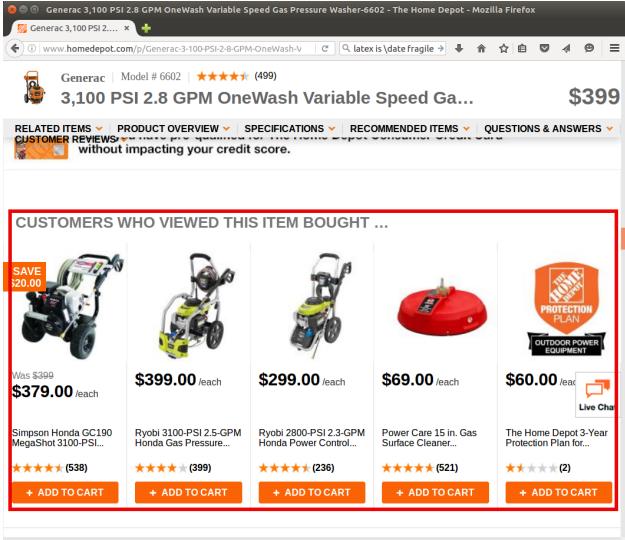
¹<http://www.imdb.com/>

²<http://www.imdb.com/interfaces>

³<ftp://ftp.fu-berlin.de/pub/misc/movies/database/>



(a) Barnes & Noble



(b) Home Depot

Figure 1: Barnes & Noble and Home Depot web pages. These websites recommend additional items based on the user's past activities and the activities of other users. The recommendation areas are outlined in red.

1. actors (a list of actors from each movie) in the file actors.list.gz (approximately 295 GB compressed)
2. genres (which assigns each movie to one or more genres) in the file genres.list.gz (approximately 17 GB compressed)
3. movies (in the IMDb general sense) in the file movies.list.gz (approximately 37 GB compressed)

These files will allow us to go from a movie title to actors in that movie to which genres that movie is part of, and then to find other movies in the same genres that had some of the same actors as the original movie. Conceptually; if we liked movie A then we should like movie B because it has some of the same actors, and is the same type as A.

Simplifying assumptions

This report is a “proof of concept” level effort, vice a production level effort. As such, there are a number of simplifying assumptions and design decisions that were made. These include:

1. Only deal with ASCII characters. The IMDb has movies (and other data) about productions made around the world. In many cases, the movie names contain characters that are outside the ASCII character set. Those characters are silently removed during the database creation.
2. Only deal with movies. The IMDb contains things other than traditional theater movies, including:
 - Shorts,
 - Videos,
 - TV shows,
 - Games,
 - Animations, and
 - other forms of electronic entertainment.
3. The number of common actors is fixed at 2. When the number of actors in a movie is small (< 10), it is likely that they will be together in another movie. When the number of actors is large (> 100), it is unlikely that they will be together in another movie. In a production system, the number of common actors could be a user entered value.
4. The Levenshtein distance (see Section 2.3 on page 6) is fixed at 4. In a production system, this could be a user entered value.

Table 1: Number of records downloaded from IMDb The gross number of records include a small number (< 300) lines of header and other information that is not loaded into the database.

File name	Gross number of records	Records in db.
actors.list.gz	17,855,640	6,117,750
actresses.list.gz	10,882,460	3,264,670
genres.list.gz	2,069,761	1,531,462
movies.list.gz	3,439,898	794,234

Table 2: PostgreSQL tables used for this demonstration. The name data type in the table is evident from the field name. Each table has another field (titleBase64) whose purpose may not be self evident. Because the IMDb is an international database, movie titles may contain characters that have special meaning to PostgreSQL, and would be difficult to load and process. Therefore the titles are all base 64 encoded so that special characters are preserved, but the encoding results in pure ASCII representation.

Table name	Column name	Characteristics
actors	actor	varchar(1000)
	titleBase64	varchar(1000)
genre	genre	varchar(1000)
	titleBase64	varchar(1000)
movies	title	varchar(1000)
	titleBase64	varchar(1000)

The simplifications resulted in fewer records being added to the recommendations database (see Table 1).

2.2 Database organization

The database consists of three tables, one for each type of data downloaded from the IMDb (see Table 2). A Pareto chart of the genres in the database was created (see Figure 2 on page 7). From the chart it is easy to see that almost 50% of the entries are classified as “short” or “drama” or both. An textual distribution is provided as well (see Table 3 on the next page).

Table 3: Number of database entries by genre.

Genre	Count
Action	44,836
Adult	7,781
Adventure	25,769
Animation	37,826
Biography	18,083
Comedy	177,313
Crime	32,598
Documentary	140,445
Drama	264,613
Family	30,753
Fantasy	25,100
Film-Noir	612
Game-Show	51
History	15,337
Horror	37,999
Musical	12,521
Music	18,441
Mystery	22,041
News	7,140
Reality-TV	228
Romance	51,450
Sci-Fi	20,959
Short	452,549
Sport	9,015
Talk-Show	116
Thriller	51,766
War	12,455
Western	13,665

With so many different genres, and so many movies categorized into each genre; a natural question is: how many movies share common categories? For instance, how many movies are categorized as Biography, Crime, and Documentary? While the question is easily stated, and easily computed there are nuances that contribute to making the question difficult to answer as the number of categories increases. For the three representative categories, there are 8 possible ways that movies could be categorized. These are

1. In Biography only.
2. In Crime only.
3. In Documentary only.
4. In Biography and Crime, but not in Documentary.
5. In Crime and Documentary, but not in Biography.
6. In Biography and Documentary, but not in Crime.
7. In Biography and Documentary and Crime.
8. Not in Biography, nor Documentary, nor Crime.

It turns out that the number of possible combinations is 2^n . Our proof of concept has 28 genres:

$$2^{28} = 268,435,456 \text{ possible combinations (see Section B on page 23)}$$

2.3 Levenshtein distance algorithm

“..., the Levenshtein distance is a string metric for measuring the difference between two sequences. Informally, the Levenshtein distance between two words is the minimum number of single-character edits (i.e. insertions, deletions or substitutions) required to change one word into the other.”

Wikipedia Staff [12]

Levenshtein distance (LD) may also be referred to as edit distance. The LD can be viewed as a metric for comparing two strings[6]. The greater the LD is, the more different the two strings are.

LD has been used in:

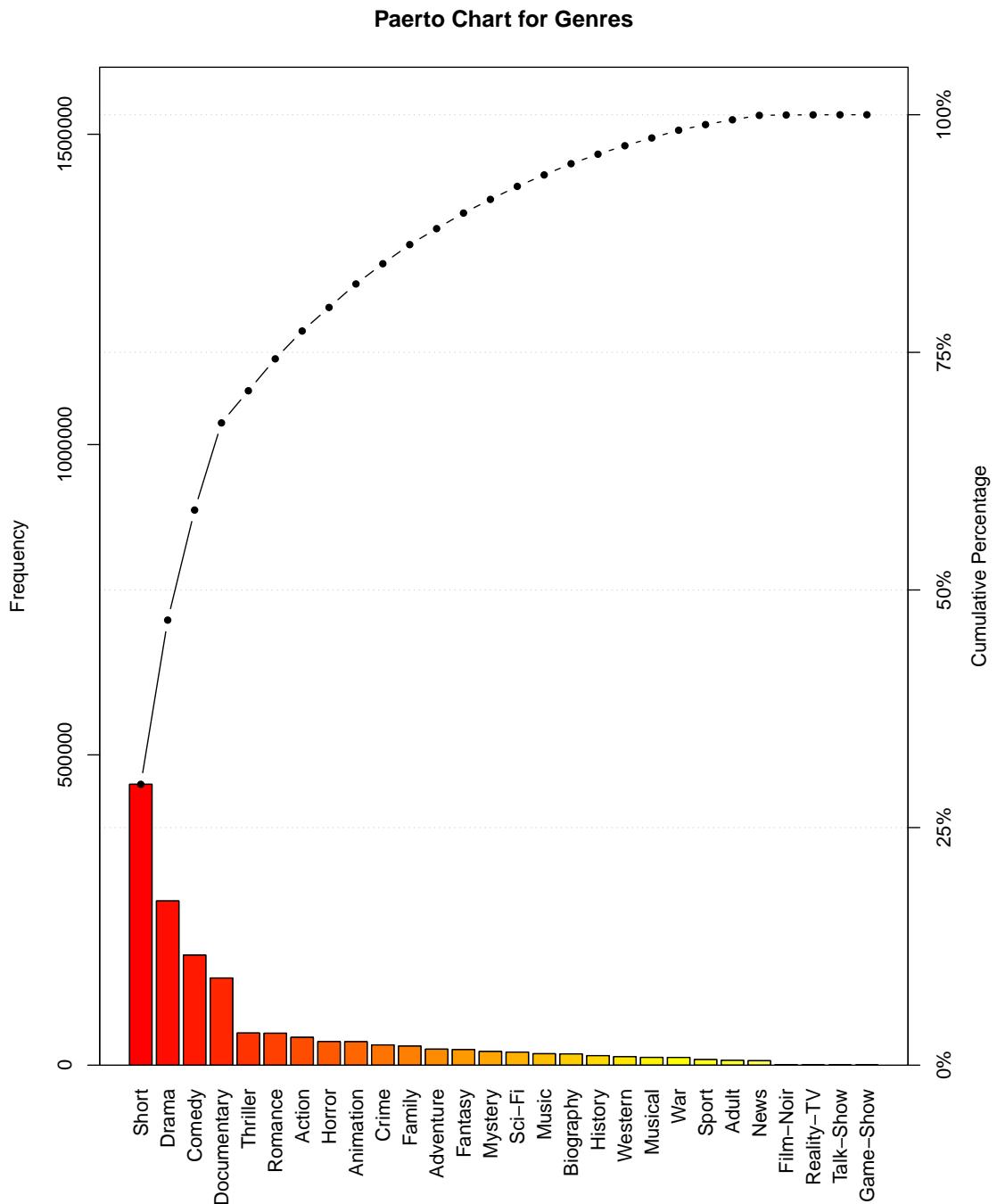


Figure 2: Pareto chart of genres. The three most common categorizations (Short, Drama, and Comedy) account for almost 60% of all categorizations.

- Spell checking
- Speech recognition
- DNA analysis
- Plagiarism detection

The following algorithm and example are taken from[4]. Assumptions:

1. The source string is called: s
2. The target string is called: t
3. The **cost** of all edit operations (insertions, deletions or substitutions) is 1. (Note: This may not be the real case in a real system.)

The goal is to compute the LD from s to t .

Mathematically[12]; LD between strings s, t (of length $| s |$ and $| t |$ respectively) is given by $lev_{s,t}(| s |, | t |)$. Where:

$$lev_{s,t}(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) = 0, \\ \min \begin{cases} lev_{s,t}(i - 1, j) + cost \\ lev_{s,t}(i, j - 1) + cost \\ lev_{s,t}(i - 1, j - 1) + cost_{(s_i \neq t_j)} \end{cases} & \text{otherwise.} \end{cases}$$

where $1_{(s_i \neq t_j)}$ is the indicator function equal to 0 when $s_i = t_j$ and equal to 1 otherwise. An algorithmic explanation can be found in (see Table 4).

Table 4: Levenshtein algorithm.

Step	Description
1	<ul style="list-style-type: none"> • Set n to be the length of s. • Set m to be the length of t. • If $n = 0$, return m and exit. • If $m = 0$, return n and exit. • Construct a matrix containing $0 \dots m$ rows and $0 \dots n$ columns.

(Continued on the next page.)

Table 4. (Continued from the previous page.)

Step	Description
2	<ul style="list-style-type: none"> • Initialize the first row to $0 \dots n$. • Initialize the first column to $0 \dots m$.
3	Examine each character of \mathbf{s} (i from 1 to n).
4	Examine each character of \mathbf{t} (j from 1 to m).
5	<ul style="list-style-type: none"> • If $s[i]$ equals $t[j]$, the cost is 0. • If $s[i]$ doesn't equal $t[j]$, the cost is 1.
6	<p>Set cell $lev(i, j)$ of the matrix equal to the minimum of:</p> <ul style="list-style-type: none"> • The cell immediately above plus 1: $lev(i - 1, j) + 1$. • The cell immediately to the left plus 1: $lev(i, j - 1) + 1$. • The cell diagonally above and to the left plus the cost: $lev(i - 1, j - 1) + \text{cost}$
7	After the iteration steps (3, 4, 5, 6) are complete, LD is found in cell $lev(n, m)$..

An example showing the lev matrix for changing the start string \mathbf{s} GUMBO to the target string \mathbf{t} GAMBOL provided (see Table 5 on the following page).

We will use LD to find the closest (lowest LD) between the a movie title entered by the user, and all movie titles in our database.

Table 5: Levenshtein converting GUMBO to GAMBOL. The LD (number of insertions, deletions or substitutions) is shown in the lower left hand cell of the last table as 2.

		Steps 1 and 2					Steps 3 to 6 (i = 1)					Steps 3 to 6 (i = 2)						
		G	U	M	B	O	G		U	M	B	O	G		U	M	B	O
		0	1	2	3	4	0	1	2	3	4	5	0	1	2	3	4	5
G	1						G	1	0				G	1	0	1		
A	2						A	2	1				A	2	1	1		
M	3						M	3	2				M	3	2	2		
B	4						B	4	3				B	4	3	3		
O	5						O	5	4				O	5	4	4		
L	6						L	6	5				L	6	5	5		
Steps 3 to 6 (i = 3)																		
		G	U	M	B	O	G		U	M	B	O	G		U	M	B	O
		0	1	2	3	4	0	1	2	3	4	5	0	1	2	3	4	5
G	1	0	1	2			G	1	0	1	2	3	G	1	0	1	2	3
A	2	1	1	2			A	2	1	1	2	3	A	2	1	1	2	3
M	3	2	2	1			M	3	2	2	1	2	M	3	2	2	1	2
B	4	3	3	2			B	4	3	3	2	1	B	4	3	3	2	1
O	5	4	4	3			O	5	4	4	3	2	O	5	4	4	3	2
L	6	5	5	4			L	6	5	5	4	3	L	6	5	5	4	3
Steps 3 to 6 (i = 4)																		
		G	U	M	B	O	G		U	M	B	O	G		U	M	B	O
		0	1	2	3	4	0	1	2	3	4	5	0	1	2	3	4	5
G	1	0	1	2	3		G	1	0	1	2	3	G	1	0	1	2	3
A	2	1	1	2	3		A	2	1	1	2	3	A	2	1	1	2	3
M	3	2	2	1	2		M	3	2	2	1	2	M	3	2	2	1	2
B	4	3	3	2	1		B	4	3	3	2	1	B	4	3	3	2	1
O	5	4	4	3	2		O	5	4	4	3	2	O	5	4	4	3	2
L	6	5	5	4	3		L	6	5	5	4	3	L	6	5	5	4	3
Steps 3 to 6 (i = 5)																		
		G	U	M	B	O	G		U	M	B	O	G		U	M	B	O
		0	1	2	3	4	0	1	2	3	4	5	0	1	2	3	4	5
G	1	0	1	2	3	4	G	1	0	1	2	3	G	1	0	1	2	3
A	2	1	1	2	3	4	A	2	1	1	2	3	A	2	1	1	2	3
M	3	2	2	1	2	3	M	3	2	2	1	2	M	3	2	2	1	2
B	4	3	3	2	1	2	B	4	3	3	2	1	B	4	3	3	2	1
O	5	4	4	3	2	1	O	5	4	4	3	2	O	5	4	4	3	2
L	6	5	5	4	3	2	L	6	5	5	4	3	L	6	5	5	4	3

3 Results

Interactions with the database are initiated by the user. Some misspelled movie names will result in recommendations (small number of actors, lots of movies being made at roughly the same time, small industry, etc.), while others won't (large cast, belonging to multiple conflicting genres, etc.). We will go through a couple of example iterations, showing the “normal” output, and the more detailed SQL queries and responses (see Section E on page 43).

First we will show a “normal” interaction where the user enters the misspelled title: Copid rides the range. Normal interactions are listed (see Table 6), and then detailed ones are listed (see Table 7 on page 13). Normal interactions are listed (see Table 8 on page 17) for “The Dark Knight” identified as the intended movie when “The Dirk Knight” (the default movie title) is entered.

Table 6: Normal output for “Copid rides the range”

Line	Output	Explanation
1	You entered 'Copid rides the range' and we chose 'Cupid Rides the Range' as our best match. Out of: 1. Cupid Rides the Range	A repeat of what the user entered to provide feedback.
2	'Cupid Rides the Range' belongs to 3 genres including: 1. Music 2. Western 3. Short	The genres that the movie was classified as part of.
3	These 10 actors appeared in 'Cupid Rides the Range': Phelps, Norman ; Strange, Glenn (I) ; McKenzie, Robert (I) Worden, Hank ; Ros, Elvira ; Phelps, Willie Whitley, Ray ; Card, Ken ; Phelps Brothers, The Phelps, Earl ;	A list of actors and actresses.
4	We will pick 2 of them.	Actors are selected at random and may vary between runs.

(Continued on the next page.)

Table 6. (Continued from the previous page.)

Line	Output	Explanation
	We have selected these actors:	
	Phelps, Norman Phelps, Willie	
5	There are 21 recommendations.	Recommendations were found. Sometimes there won't be any.
	They are (in no particular order):	
	1. Border G-Man	
	2. Prairie Papas	
	3. Ranch House Romeo	
	4. A Buckaroo Broadcast	
	5. Bandits and Ballads	
	6. Trouble in Texas	
	7. Molly Cures a Cowboy	
	8. Rawhide	
	9. Sagebrush Serenade	
	10. Painted Desert	
	11. The Renegade Ranger	
	12. Trouble in Sundown	
	13. Gun Law	
	14. Rhythm Wranglers	
	15. Six-Gun Gold	
	16. A Western Welcome	
	17. Hittin' the Trail	
	18. Where the West Begins	
	19. Cactus Capers	
	20. Tex Rides with the Boy Scouts	
	21. The Utah Trail	

Table 7: Detailed output for “Copid rides the range”

Line	Output	Explanation
	\connect chuck;	
1	select distinct title, titlebase64, levenshtein(title, 'Copid rides the range') from movies where levenshtein (title, 'Copid rides the range') < 4 order by levenshtein limit 10 ;	Commands sent to the database.
	\q	
	You are now connected to database 'chuck' as user 'chuck'.	
2	title — titlebase64 — levenshtein +————— +————— Cupid Rides the Range — Q3VwaWQgUmlkZXMGdGhlIFJhbmdl — 3 (1 row)	Response from the database.
3	You entered 'Copid rides the range' and we chose 'Cupid Rides the Range' as our best match. Out of: 1. Cupid Rides the Range	Feedback to the user.
4	\connect chuck; select distinct genre from genre where titlebase64 = 'Q3VwaWQgUmlkZXMGdGhlIFJhbmdl';	Commands sent to the database.
	You are now connected to database 'chuck' as user 'chuck'.	
	genre — 5 Music Western Short (3 rows)	
		Response from the database.

(Continued on the next page.)

Table 7. (Continued from the previous page.)

Line	Output	Explanation
6	'Cupid Rides the Range' belongs to 3 genres including: 1. Music 2. Western 3. Short	Feedback to the user.
7	\connect chuck select distinct actor from actors where titlebase64 = 'Q3VwaWQgUmlkZXMGdGhlIFJhbmdl'; \q	Commands sent to the database.
	You are now connected to database 'chuck' as user 'chuck'. actor	
8	Phelps, Norman Strange, Glenn (I) McKenzie, Robert (I) Worden, Hank Ros, Elvira Phelps, Willie Whitley, Ray Card, Ken Phelps Brothers, The Phelps, Earl (10 rows)	Response from the database.
9	These 10 actors appeared in 'Cupid Rides the Range': Phelps, Norman Strange, Glenn (I) McKenzie, Robert (I) Worden, Hank Ros, Elvira Phelps, Willie Whitley, Ray Card, Ken Phelps Brothers, The	Feedback to the user.

(Continued on the next page.)

Table 7. (Continued from the previous page.)

Line	Output	Explanation
	Phelps, Earl	
	We will pick 2 of them.	
	We have selected these actors:	
	Phelps, Norman Phelps, Willie	
	\connect chuck	
10	select distinct titlebase64 from genre where titlebase64 in (select distinct a1.titlebase64 from actors as a1 join actors as a2 on a1.titlebase64 = a2.titlebase64 and a2.actor = 'Phelps, Willie' and a1.actor = 'Phelps, Norman');	Commands sent to the database.
	\q	
	You are now connected to database 'chuck' as user 'chuck'.	
	titlebase64	
	Qm9yZGVyIEctTWFu	
	Q3VwaWQgUmlkZXMGdGhlIFJhbmdl	
	UHJhaXJpZSBQYXHcw==	
	UmFuY2ggSG91c2UgUm9tZW8=	
	QSBCdWNrYXJvbyBCcm9hZGNhc3Q=	
	QmFuZGl0cyBhbmQgQmFsbGFkcw==	
	VHJvdWJsZSBpbIBUZXhhcw==	
	TW9sbHkgQ3VyZXMgYSBDb3dib3k=	
	UmF3aGlkZQ==	
	U2FnZWJydXNoIFNlcMVuYWRI	
11	UGFpbnRlZCBEZXNlcnQ=	Response from the database.
	VGhlIFJlbmVnYWRIIFJhbndlcg==	
	VHJvdWJsZSBpbIBTdW5kb3du	
	R3VuIExhdw==	
	Umh5dGhtIFdyYW5nbGVycw==	

(Continued on the next page.)

Table 7. (Continued from the previous page.)

Line	Output	Explanation
	U2l4LUd1biBHb2xk	
	QSBZXZN0ZXJuIFdlbGNvbWU=	
	SGl0dGluJyB0aGUgVHJhaWw=	
	V2hlcmUgdGhlIFdlc3QgQmVnaW5z	
	Q2FjdHVzIENhcGVycw==	
	VGV4IFJpZGVzIHdpdGggdGhlIEJveSBTY291dHM=	
	VGhlIFV0YWggVHJhaWw=	
	(22 rows)	
	There are 21 recommendations.	
	They are (in no particular order):	
	1. Border G-Man	
	2. Prairie Papas	
	3. Ranch House Romeo	
	4. A Buckaroo Broadcast	
	5. Bandits and Ballads	
	6. Trouble in Texas	
	7. Molly Cures a Cowboy	
	8. Rawhide	
	9. Sagebrush Serenade	
12	10. Painted Desert	Final data presented to the user.
	11. The Renegade Ranger	
	12. Trouble in Sundown	
	13. Gun Law	
	14. Rhythm Wranglers	
	15. Six-Gun Gold	
	16. A Western Welcome	
	17. Hittin' the Trail	
	18. Where the West Begins	
	19. Cactus Capers	
	20. Tex Rides with the Boy Scouts	

(Continued on the next page.)

Table 7. (Continued from the previous page.)

Line	Output	Explanation
21.	The Utah Trail	(Last page.)

Table 8: Normal output for “The Dark Knight”

Line	Output	Explanation
1	You entered 'The Dark Knight' and we chose 'The Dark Knight' as our best match. Out of: 1. The Dark Knight 2. The Dim Knight 3. The Dark Light	Feedback to the user.
2	'The Dark Knight' belongs to 5 genres includ- ing: 1. Action 2. Fantasy 3. Short 4. Drama 5. Crime	Feedback to the user.
3	These 261 actors appeared in 'The Dark Knight': Lesley, David White, Michael Jai List of actors redacted (see Section C on page 34) Ritz, Sara Campbell, Tommy Pierce, Ernest Malkin, Simon Webb, Bronson Hoving, John McAllister, Lisa (I) Chin Han List of actors redacted (see Section C on page 34)	Feedback to the user.

(Continued on the next page.)

Table 8. (Continued from the previous page.)

Line	Output	Explanation
	Kalesperis, Adam Shah, Vivek (II)	
	We will pick 2 of them.	
4	We have selected these actors: Campbell, Tommy Chin Han	Feedback to the user.
5	Sorry no recommendations possible.	Feedback to the user.

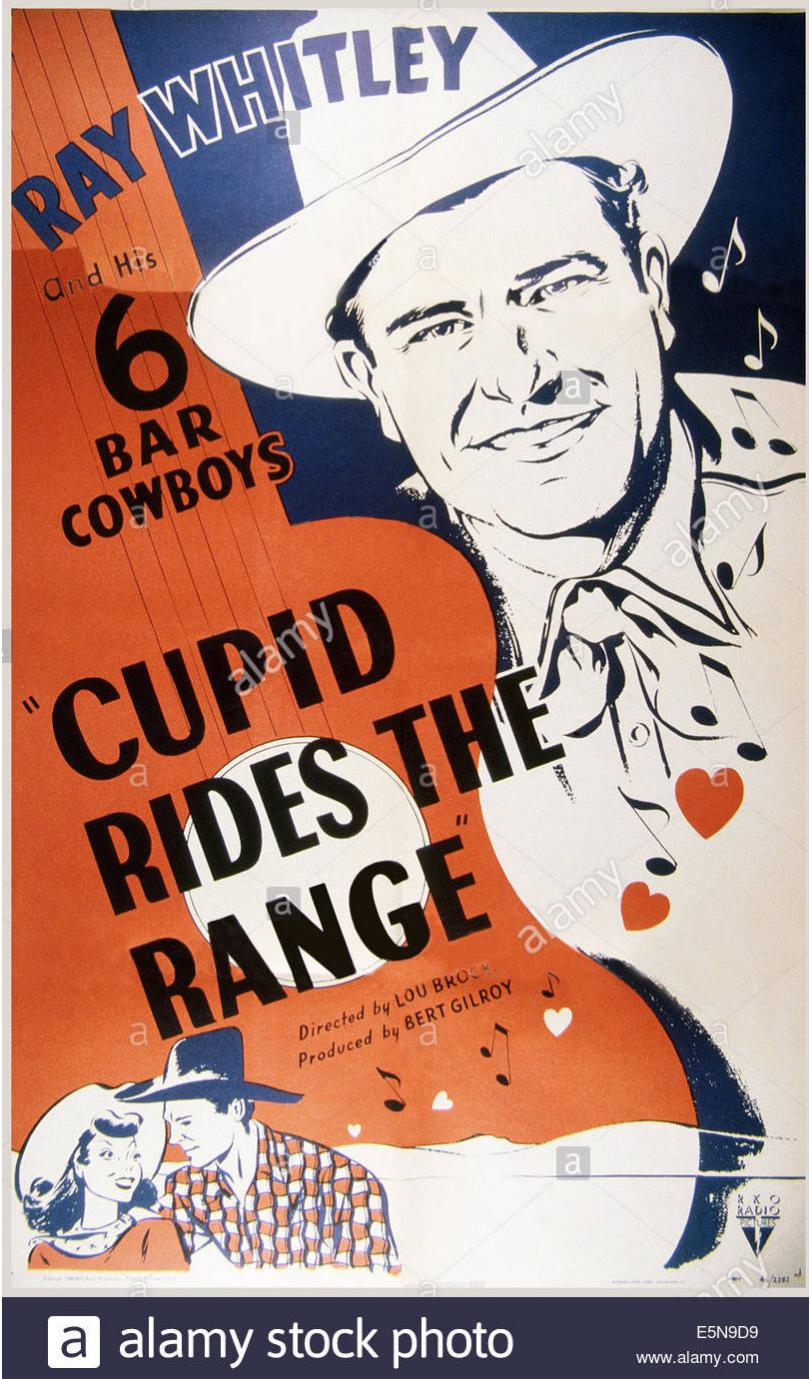
4 Conclusion

We created a movie recommendation system that mimics recommendation systems in that we encounter all over the Internet. We demonstrated how the system would respond when given misspelled movie titles. Recommendations were based on having at least two actors from the original movie in another movie that was in same genres as the original. In some cases, recommendations were found, and in others none were found.

One misspelled movie title, “Copid rides the range” (see Figure 3 on the next page) resulted in 21 recommendations. Another misspelled movie title “Dirk Knight” (see Figure 4 on page 21) resulted in 0 recommendations.

The proof of concept recommendation system worked. It is not a production class system, and could be improved by:

1. Selecting the common actors based on their placement in the credits, vice random selection.
2. Supporting the full Unicode character set, vice only ASCII.
3. Incorporation the user’s past genre history when limiting searches, vice using the complete set of genre assignments from the original movie.



a alamy stock photo

E5N9D9
www.alamy.com

Figure 3: Movie poster for “Cupid Rides the Range” Movie staring Ray Whitley (1939).Image from [7]



Figure 4: Movie poster for “Dark Knight” Movie starring Christian Bale, and Heath Ledger (2008). Image from [8].

A A PostgreSQL timeline

PostgreSQL is old[9].

- Started at UC Berkley as POSTGRES in 1986.
- DARPA and Army Research Office (ARO) project
- Demoed in 1987, presented in 1988, released in 1989
- Evolved into Postgres95
- Name changed in 1985 to PostgreSQL to reflect origin and new SQL capability
- Open source software available 1 Aug. 1996 [10]
- Online presence at PostgreSQL.org since 22 Oct. 1996
- Under active development, major releases approx. yearly
- Official name PostgreSQL, nicknamed Postgres
- Postgres is Post(In)gres (yet another SQL database)[11]

An old and on going project. PostgreSQL is very widely available.

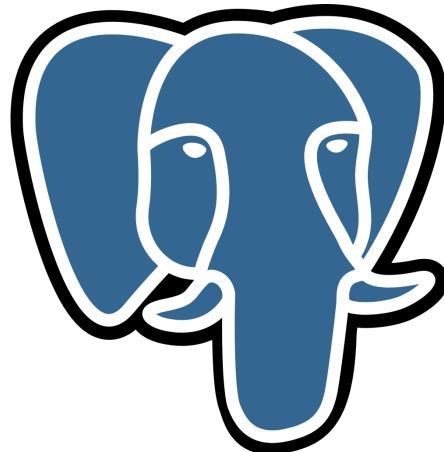


Figure 5: PostgreSQL Elephant.

B Venn diagrams

“A Venn diagram (also known as a set diagram or logic diagram) is a diagram that shows all possible logical relations between a finite collection of different sets.”

Wikipedia Staff [13]

Venn diagrams are used to visualize how different sets of data intersect (or don't intersect). We'll use three different sets of data A , B , and U for our general discussion before delving into questions about our proof of concept genre database.

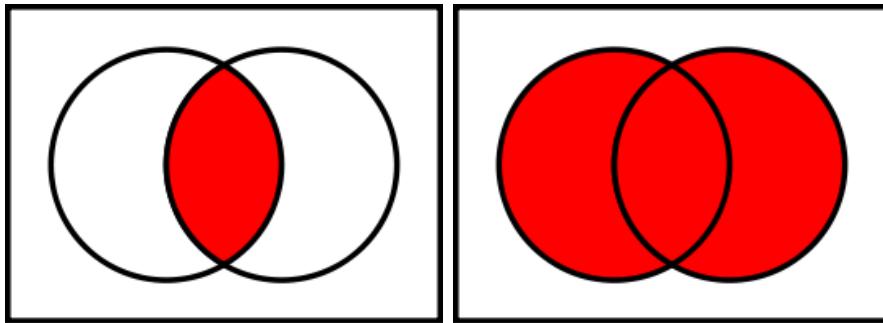
A and B represent sets of data, that have some members in their respective sets in common. These sets could be animals, vegetables, minerals, numbers, books, or anything that we want. U is the universe of all things, including all of A and all of B . Using these three sets (see Figure 6 on the next page), we can say somethings about them.

A series of Venn diagrams were created to visualize the question of how many movies were in the data sets Biography and Documentary and Crime. The Western and Adventure genres were added to the data sets of interest in order to show how complex Venn diagrams can get with only a small number of sets.

Monochromatic Venn diagrams can convey the overlap of the data sets (genres) (see Figures 7 on page 25 through 11 on page 29), but there are times when a polychromatic presentation helps (see Figure 12 on page 30).

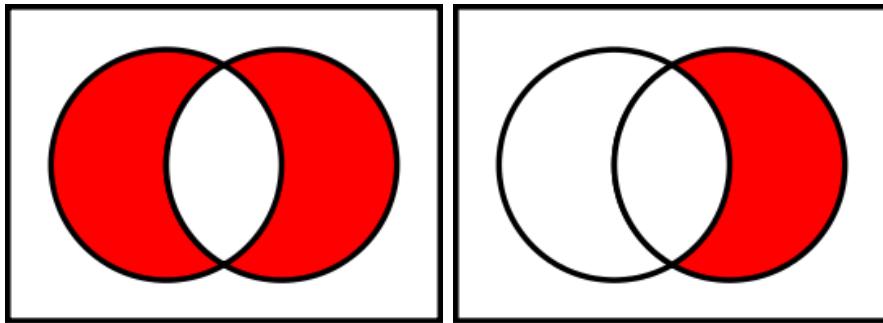
The number of ways of picking unordered subsets from a larger set is known as the binomial coefficient. If the larger set contains the elements *Biography*, *Crime*, *Documentary*, *Adventure* and we were interested in counting the number of subsets that had exactly 2 elements, then the subsets would be:

1. *Biography, Crime*
2. *Biography, Documentary*
3. *Biography, Adventure*
4. *Crime, Documentary*
5. *Crime, Adventure*
6. *Documentary, Adventure*



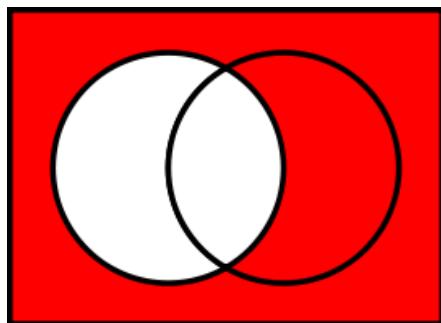
(a) Intersection of two sets: $A \cap B$

(b) Union of two sets: $A \cup B$



(c) Symmetric difference of two sets:
 $A \Delta B$

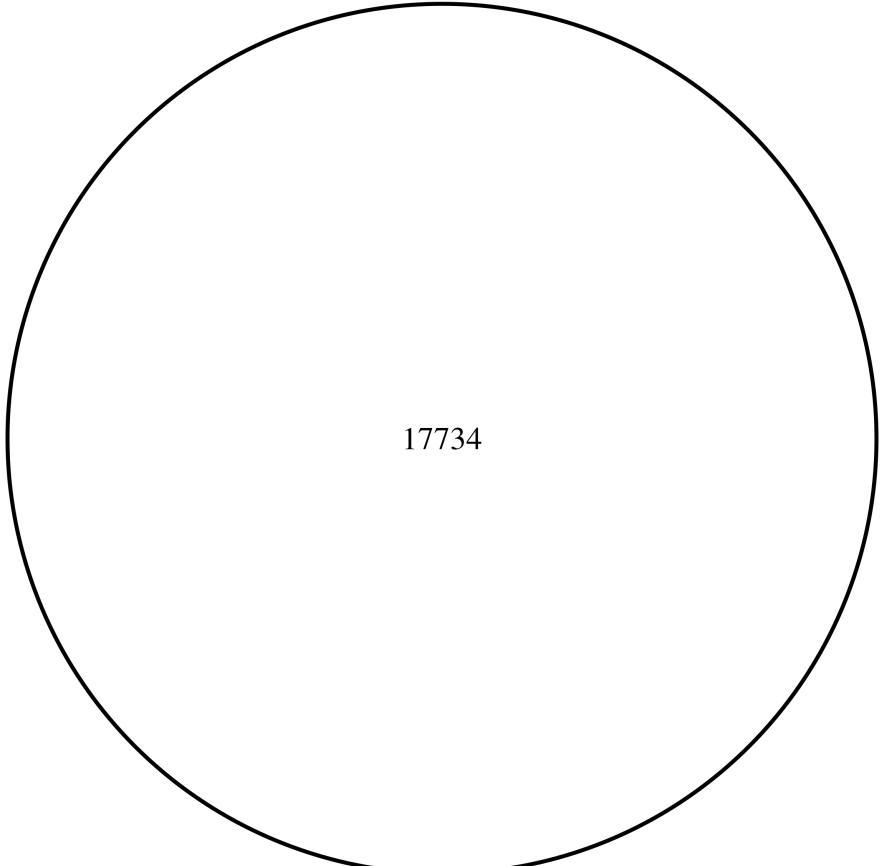
(d) Relative complement of A (left) in
 B (right): $A^c \cap B = B \setminus A$



(e) Absolute complement of A in U :
 $A^c = U \setminus A$

Figure 6: Representative Venn diagrams. The corresponding mathematical notation is in the caption of each subfigure.

Biography



17734

Figure 7: Venn diagram of the Biography genre.

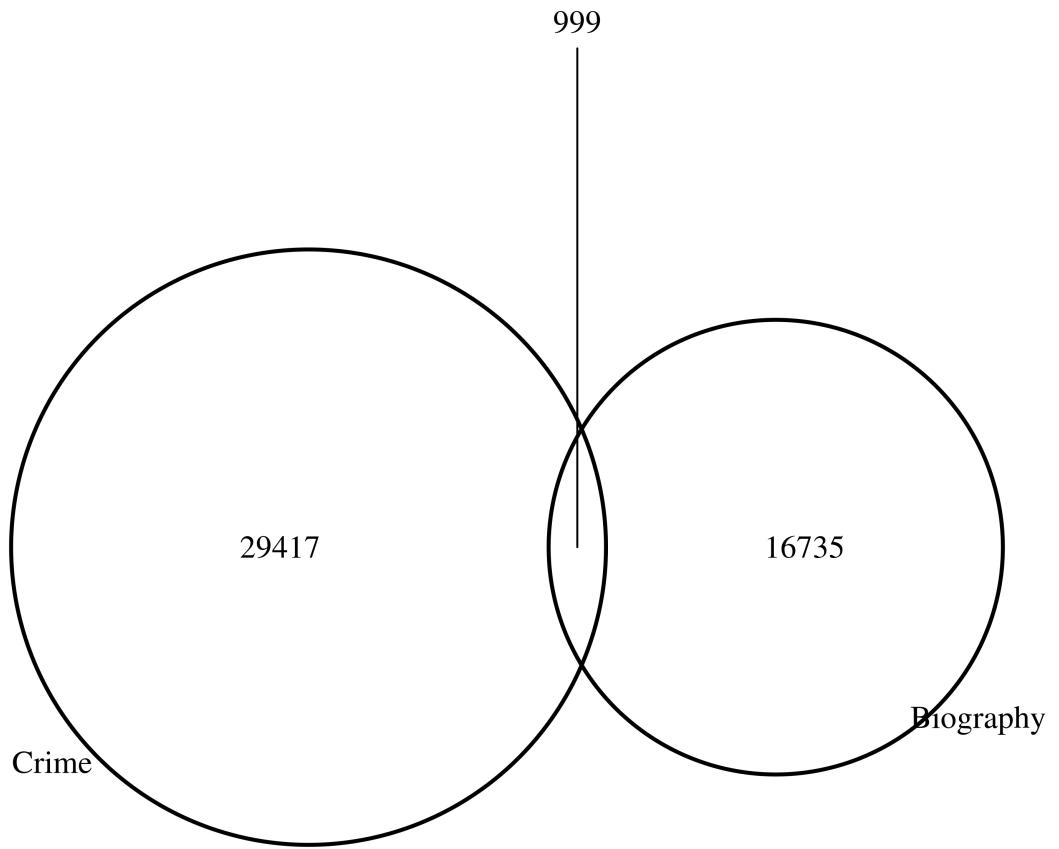


Figure 8: Venn diagram of the Biography and Crime genres ($\text{Biography} \cup \text{Crime}$). Showing there are 999 common movies in both genres.

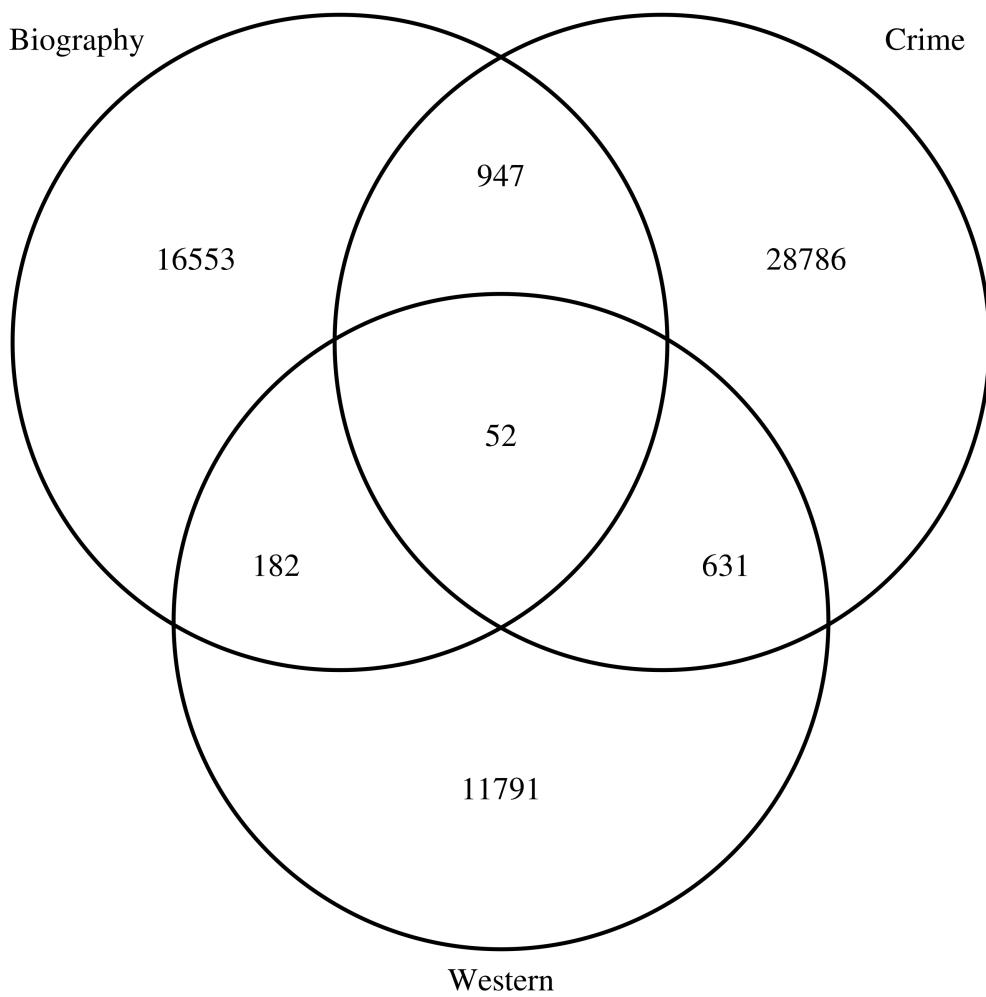


Figure 9: Venn diagram of the Biography and Crime Western genres ($\text{Biography} \cup \text{Crime} \cup \text{Western}$). Showing there are 52 common movies in all genres.

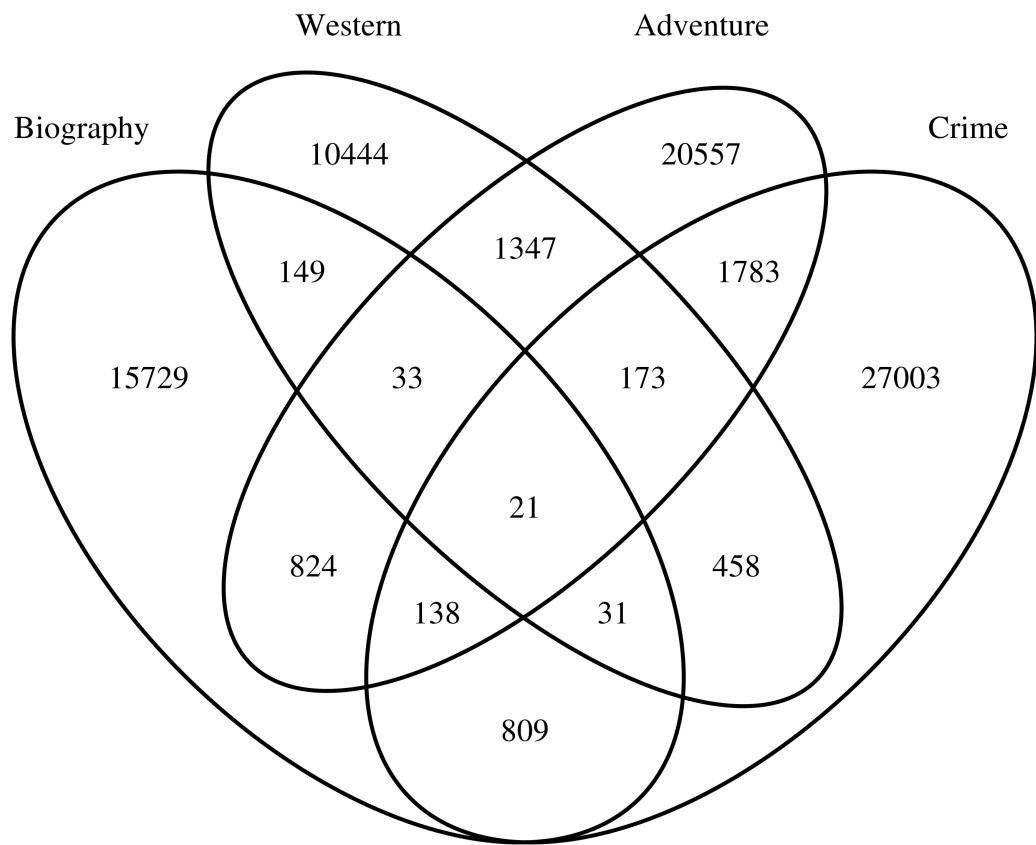


Figure 10: Venn diagram of the Biography and Crime Western and Adventure genres ($\text{Biography} \cup \text{Crime} \cup \text{Western} \cup \text{Adventure}$). Showing there are 21 common movies in all genres.

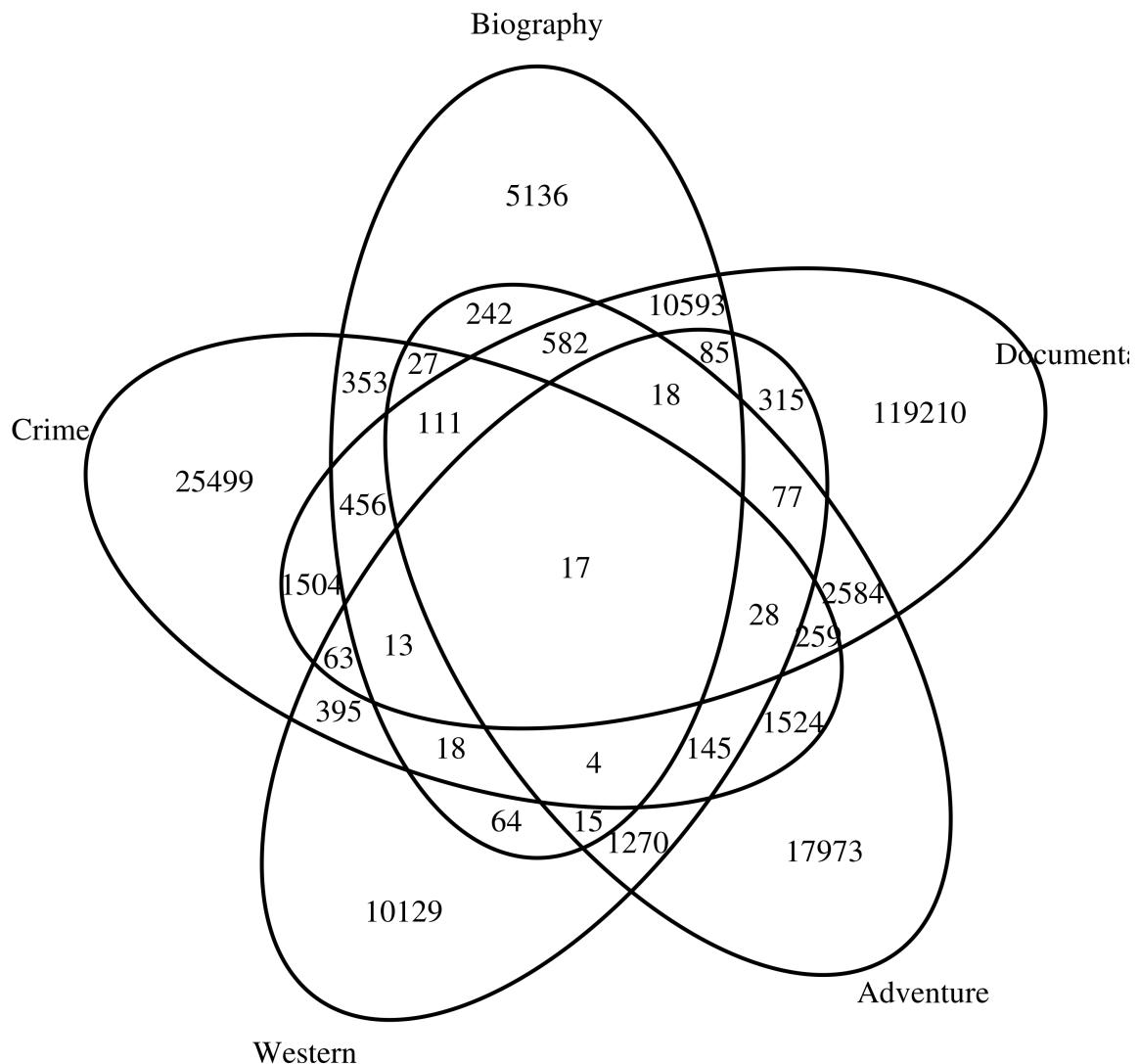


Figure 11: Venn diagram of the Biography and Crime Western and Adventure and Documentary genres ($\text{Biography} \cup \text{Crime} \cup \text{Western} \cup \text{Adventure} \cup \text{Documentary}$). Showing there are 17 common movies in all genres.

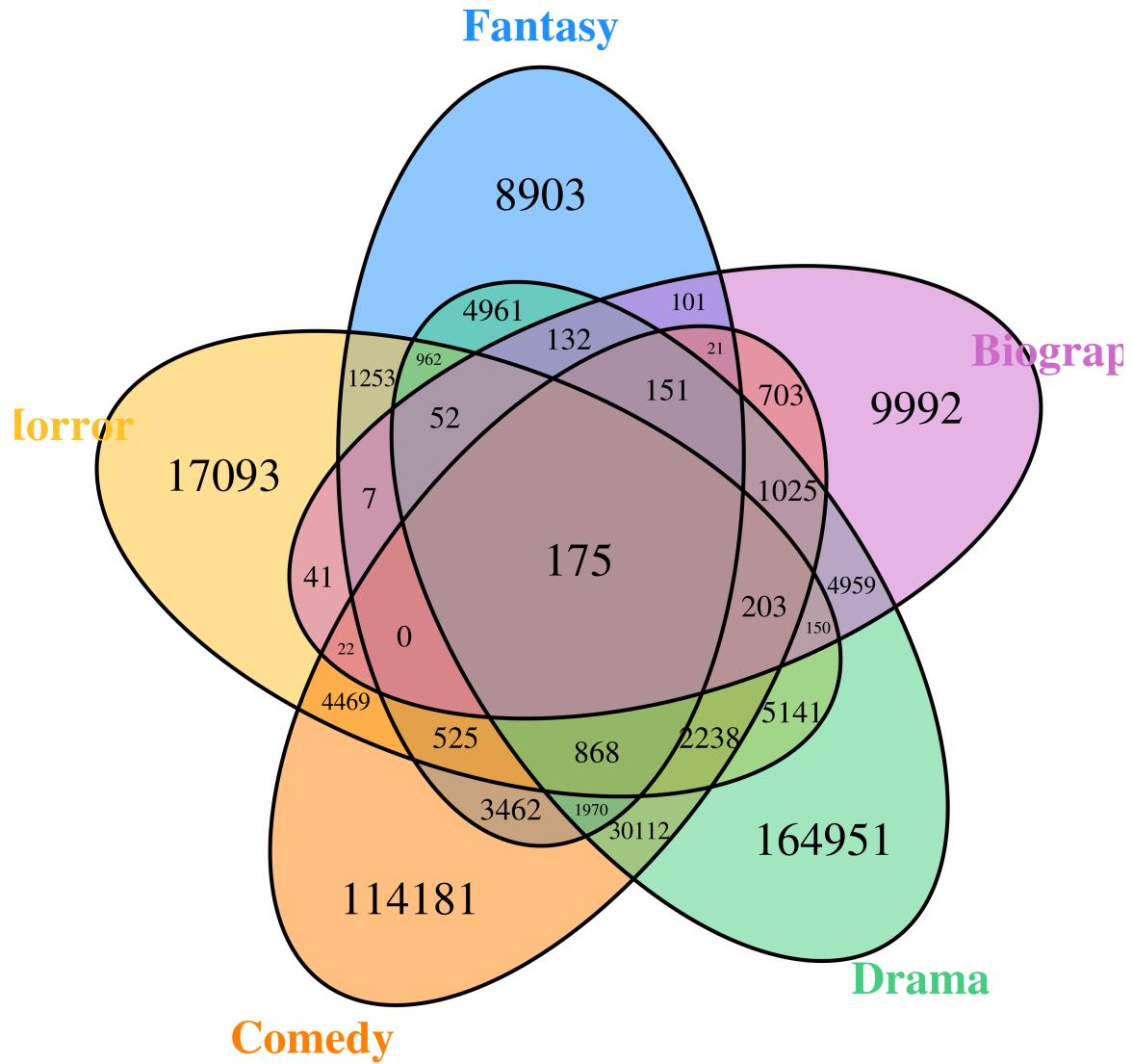


Figure 12: Colorful Venn diagram of Horror and Fantasy and Biography and Drama and Comedy. It is interesting to see that there are 0 movies common to all genres. This diagram was made possible by modifying the code from <http://www.inside-r.org/packages/cran/VennDiagram/docs/venn.diagram>. The modified code is attached to this report.

If we vary the number of elements in the subsets from 0 to 4 (because our example has 4 elements in the complete set), we end up with an interesting pattern of number of subsets based on the desired size of the subset (see Table 9 on the next page).

If we say that n is the number of elements in the larger set, and k is the number of elements in the subset, then we can say $n\text{choose}k$ which is often written $\binom{n}{k}$ (the binomial coefficient) which expands to $\frac{n!}{k!(n-k)!}$.

By looking at the pattern of numbers in the **Count** column (see Table 9 on the following page), and varying k (all the while keeping n fixed and larger than k), we start to see a pattern. This pattern is known as Pascal's Triangle. Pascal's Triangle

“[T]he triangular array of integers, with 1 at the apex, in which each number is the sum of the two numbers above it in the preceding row . . .

Borowski and Borwein [1]

A Pascal Triangle for $n = 5$ has been created (see Table 10 on page 33). While the coefficients are patently obvious afterwards, it is helpful to see they represent (see Table 11 on page 33). A slightly different presentation of Pascal's Triangle makes it easy to understand and compute the number of coefficients⁴ (see Table 12 on page 33).

Looking at enough of these cases, and counting the number of subsets for each $n\text{choose}k$ case where we vary k from 0 to n , we can arrive at the final formula:

$$\begin{aligned} \text{subsets} &= \sum_{k=0}^n \binom{n}{k} \\ &= \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{n} \\ &= \frac{n!}{0! * (n-0)!} + \frac{n!}{1! * (n-1)!} + \dots + \frac{n!}{n! * (n-n)!} \\ &= 2^n \end{aligned}$$

⁴www.cs.columbia.edu/~cs4205/files/CM4.pdf

Table 9: Choosing different genre combinations.

Choose	Chosen	Subsets
0	1. (None.)	1
1	1. <i>Biography</i> 2. <i>Crime</i> 3. <i>Documentary</i> 4. <i>Adventure</i>	4
2	1. <i>Biography, Crime</i> 2. <i>Biography, Documentary</i> 3. <i>Biography, Adventure</i> 4. <i>Crime, Documentary</i> 5. <i>Crime, Adventure</i> 6. <i>Documentary, Adventure</i>	6
3	1. <i>Biography, Crime, Documentary</i> 2. <i>Biography, Crime, Adventure</i> 3. <i>Biography, Documentary, Adventure</i> 4. <i>Crime, Documentary, Adventure</i>	4
4	1. <i>Biography, Crime, Documentary, Adventure</i>	1
Total		16

Table 10: Pascal's Triangle up to $n = 5$. Each coefficient is the sum of the two above to the left and to the right (assuming that blanks are treated as 0s). Python source code to create table from <http://www.bedroomlan.org/coding/pascals-triangle-latex>. Python code is attached to this report.

$n = 0:$	1					
$n = 1:$	1 1					
$n = 2:$	1 2 1					
$n = 3:$	1 3 3 1					
$n = 4:$	1 4 6 4 1					
$n = 5:$	1 5 10 10 5 1					

Table 11: Pascal's Triangle up to $n = 5$ as binomials.

$n = 0:$	$\binom{0}{0}$					
$n = 1:$	$\binom{1}{0} \quad \binom{2}{0}$					
$n = 2:$	$\binom{2}{0} \quad \binom{2}{1} \quad \binom{2}{2}$					
$n = 3:$	$\binom{3}{0} \quad \binom{3}{1} \quad \binom{3}{2} \quad \binom{3}{3}$					
$n = 4:$	$\binom{4}{0} \quad \binom{4}{1} \quad \binom{4}{2} \quad \binom{4}{3} \quad \binom{4}{4}$					
$n = 5:$	$\binom{5}{0} \quad \binom{5}{1} \quad \binom{5}{2} \quad \binom{5}{3} \quad \binom{5}{4} \quad \binom{5}{5}$					

Table 12: Pascal's Triangle up to $n = 5$ reformatted coefficients.

n	$\binom{n}{0}$	$\binom{n}{1}$	$\binom{n}{2}$	$\binom{n}{3}$	$\binom{n}{4}$	$\binom{n}{5}$	Σ
0	1						1
1	1	1					2
2	1	2	1				3
3	1	3	3	1			8
4	1	4	6	4	1		16
5	1	5	10	10	5	1	32

C “The Dark Knight” actors

Actors and actresses can have the same name as other actors and actresses. The IMDb resolves these by appending a Roman numeral in parentheses after the name.

These 261 actors (in alphabetical order) appeared in 'The Dark Knight':

- | | | |
|--------------------------------|----------------------------|------------------------------|
| 1. Ajala, David | 21. Birchard, Paul (I) | 41. Chin Han |
| 2. Albertson, Jeff | 22. Borchardt, Amanda (II) | 42. Cho, Matt (I) |
| 3. Aldaz, Tracy L. | 23. Brindowski, Adam | 43. Chu, Henry Milton |
| 4. Allen, Matthew W. | 24. Brody, Jon Lee | 44. Clear, Patrick |
| 5. Anderson, Christopher James | 25. Brooke, Peter (I) | 45. Clevenger, Kelli |
| 6. Armourae, Stephen | 26. Bucholz, Adam | 46. Cosey, David |
| 7. Armstrong, William (I) | 27. Bulcock, Philip | 47. Coster, Ritchie |
| 8. Azpeitia, Alisa | 28. Burns, Debbi | 48. Crane, Nancy |
| 9. Bach, Mike (I) | 29. Caballero, Joseph Luis | 49. Crocker, Martin |
| 10. Bakare, Ariyon | 30. Cabrera, Maritza | 50. Curnen, Monique Gabriela |
| 11. Baker, Wayne (I) | 31. Caine, Michael (I) | 51. Dastmalchian, David |
| 12. Bale, Christian | 32. Caiola, Shirin | 52. Daugherty, Rachel |
| 13. Ballantyne, Martin (I) | 33. Calmels, Fabrice | 53. Dawson, Bruce Allen |
| 14. Ballard, Doug | 34. Campbell, Tommy | 54. Day, Danielle (II) |
| 15. Barcelos, Alberto | 35. Carbonell, Nestor | 55. Dean, Ron (I) |
| 16. Bartlett, Tommy (I) | 36. Carrington, Nigel | 56. DeFaria, Peter |
| 17. Beam, Greg | 37. Chadwick, David (II) | 57. Derence, Sam |
| 18. Bennett, Blayne | 38. Chapman, Josh (I) | 58. Dillane, Richard |
| 19. Berard, Ryan | 39. Chen, Edison (I) | 59. di Makena, Matteo |
| 20. Bicknell, Andrew | 40. Chernicky, Laura | 60. Divizio, Richard |

- | | | |
|-------------------------------------|---------------------------------|--------------------------|
| 61. Domino, Tony | 85. Gaitsch, Thomas | 108. Hathaway, Alexander |
| 62. Doyle, Jessica (V) | 86. Gamble, Nathan | 109. Healy, Brendan P. |
| 63. Dunn, Sarah Jayne | 87. Ganyo, Scott | 110. Heaney, Craig |
| 64. Eckhart, Aaron | 88. Gayle, Lorna | 111. Hellman, Erik |
| 65. Edwards, Grahame (I) | 89. Geer, Daniel Richard | 112. Hepple, Victoria |
| 66. Edwards, Laine | 90. Gerard, Bill (I) | 113. Hinshelwood, Sophia |
| 67. Egan, R. Michael | 91. Giraud, Marisol | 114. Hodges, Jordon |
| 68. Elliott, David William
James | 92. Glanfield, Tim | 115. Hoving, John |
| 69. Ellis, Winston | 93. Goldring, Danny | 116. Huebsch, Matthew |
| 70. Farb, Aaron | 94. Gordon, Tom (XII) | 117. Ibrahim, Bill |
| 71. Farruggio, James | 95. Gorman, Michael An-
drew | 118. Jamroz, Gerard |
| 72. Feore, Aidan | 96. Gossen, Dan | 119. Jay, Erron |
| 73. Fichtner, William | 97. Griffin, Gordon (I) | 120. Jefferson, Daniel |
| 74. Fierro, James (I) | 98. Grover, Sharlene | 121. Jimenez, Ramses |
| 75. Fojtik, Gene | 99. Gull, R. Michael (I) | 122. Jones, Will (VI) |
| 76. Foster, Michael Corey | 100. Gunn, Hannah | 123. Kalesperis, Adam |
| 77. Foster, Reese | 101. Gyllenhaal, Maggie | 124. Kaliebe, Bob |
| 78. Fowler, Jim (IV) | 102. Hallam, Natalie | 125. Katt, Nicky |
| 79. Frederick, Jason (VI) | 103. Hall, Anthony Michael | 126. Keiser, Mark (I) |
| 80. Freeman, K. Todd | 104. Harper, Terrell | 127. Kierscht, Charlie |
| 81. Freeman, Morgan (I) | 105. Hartley, David (II) | 128. Knox, Jennifer |
| 82. Fuller, Jason | 106. Hartmann, Thomas
(VII) | 129. Kolasa, Keryaki |
| 83. Fulsher, Darren Elliot | 107. Harto, Joshua | 130. Kosik, Thomas |
| 84. Fultz, David | | 131. Kress, Don |
| | | 132. Kross, Ryan |
| | | 133. Krueger, Tim (I) |

- | | | |
|---------------------------|-----------------------------|---------------------------|
| 134. Kupferer, Keith | 158. Marchesi, Al | 182. Ozbay, Sal |
| 135. Kuster, Michael | 159. Marshall, Peter (V) | 183. Patti, Chelsey |
| 136. Lambdin, Brandon | 160. Marston, Rob | 184. Pearson, Adam (II) |
| 137. Latham, Dan (II) | 161. Martino, J.R. | 185. Pedersen, Libby (II) |
| 138. Lazicki, Joseph | 162. Mawell, Xander | 186. Petschler, Chris |
| 139. Leahy, Patrick (I) | 163. Mazurk, Joseph | 187. Pierce, Ernest |
| 140. Ledger, Heath | 164. McAllister, Lisa (I) | 188. Pirie, Ian |
| 141. Lee, Carl (V) | 165. McComas, Tom (I) | 189. Plante, Rory |
| 142. Leitch, Matthew | 166. McElroy, Tom (I) | 190. Poe, Matthew (II) |
| 143. Lesley, David | 167. McEnany, Krista | 191. Query, Charles |
| 144. Lewis, Walter (I) | 168. McFarlane, Colin (I) | 192. Radz, Marc |
| 145. Lineham, Emily | 169. McGonagle, Ryan (I) | 193. Rapkin, Brian |
| 146. Lister, Tommy 'Tiny' | 170. McGraw, Melinda | 194. Rasque, Mitch |
| 147. Lopez, Debra (II) | 171. Mellor, James (II) | 195. Reeves, Buster |
| 148. Lovegren, Ivan | 172. Monk, Roger (II) | 196. Riotta, Vincent |
| 149. Lovejoy, Lateef | 173. Mosley, Gertrude | 197. Rippy, Matt |
| 150. Luther, Andy (II) | 174. Murphy, Cillian | 198. Ritz, Sara |
| 151. Luther, Scott | 175. Nadolski, David J. | 199. Rivera, Dale |
| 152. Lutz, Lanny | 176. Nelson, Joseph (XII) | 200. Rnic, Peter |
| 153. Lynn, Deborah (II) | 177. Nicoli, Vincenzo | 201. Roberts, Eric (I) |
| 154. Lynn, Noelle (I) | 178. Olawumi, Olumiji | 202. Rollins, Joshua (I) |
| 155. Macchi, Jonathan | 179. Oldman, Gary | 203. Rosen, Beatrice |
| 156. Makepeace, Richard | 180. Oliveira, Joseph (III) | 204. Rowell, Harry |
| 157. Malkin, Simon | 181. O'Neill, Matthew (I) | 205. Ryder, Gary (II) |

- | | | |
|----------------------------|---------------------------------|---------------------------|
| 206. Ryland, Jonathan | 225. Stern, January | 243. Tyrell, Michael (II) |
| 207. Saindon, Kelly | 226. Stern, Robert Patrick | 244. Ukena, Matt |
| 208. Satcher, Daryl | 227. Stevens, Willie | 245. Venn, Charles |
| 209. Scales, James | 228. Stewart, Douglas (VII) | 246. Venus, Alisa |
| 210. Schleef, Elisa | 229. Stone, Jordan (IV) | 247. Vieau, Michael |
| 211. Schweiner, Greg | 230. Stone, Robert (IX) | 248. Walsh, Kyle (II) |
| 212. Seybold, Jan | 231. Stoyanov, Michael (I) | 249. Warfield, James |
| 213. Shah, Amit (III) | 232. Strobel, Richard (I) | 250. Warman, John |
| 214. Shah, Vivek (II) | 233. Summers, Ronan | 251. Webb, Bronson |
| 215. Shalhoub, Amanda | 234. Szarabajka, Keith | 252. White, Michael Jai |
| 216. Shallenberger, Matt | 235. Tait, Sheila | 253. Whyte, Mike |
| 217. Sharp, Colin | 236. Tait, Tristan (I) | 254. Williams, Erik A. |
| 218. Shields, Michelle (I) | 237. Terracina, Nydia Rodriguez | 255. Wilson, Chris (XXIX) |
| 219. Smillie, William | 238. Thomas, Chris D. | 256. Wilson, Helene |
| 220. Smirnova, Sofiya | 239. Thurner, John | 257. Wolf, Lisa (IV) |
| 221. Snowden, John (I) | 240. Townsend, Tom (II) | 258. Wong, Wai (III) |
| 222. Solabarrieta, Lorea | 241. Tsou, Chuen | 259. Zahir, Essa |
| 223. Sora, Dwight | 242. Turk, John (I) | 260. Zahrn, Will |
| 224. Spielbauer, Bruce | | 261. Zaideman, Kevin |

D Explaining SQL database concepts using Excel

The world of structured query language (SQL) can be and sound very confusing, but you've probably been using most of SQL's concepts and ideas for a while, without realizing it. I will attempt to connect SQL ideas of tables, columns, rows, and selects using Excel. I'll use a table and a few figures to help make the connection (see Table 13).

Table 13: Cross referencing Excel terms to PostgreSQL terms.

Excel term	PostgreSQL
Workbook	Database
Worksheet	Table
Row	Row
Column	Column
	Variable type. All entries in a column must be of the same type (characters, numbers, undefined, etc.).
Filter	Select. SQL is built around the idea of a SELECT statement. Excel's filter command is not nearly as powerful.
1,048,576 rows	Unlimited
16,384 cols	250-1600 depending on column types

Our PostgreSQL database has three tables (see Figures 13 on the next page through 15 on page 41). Excel provides a way to filter data on a worksheet (see Figure 16 on page 42). The Excel filter command works on one worksheet, not across multiple worksheets. SQL supports SELECT command on one, two, or as many tables as are available.

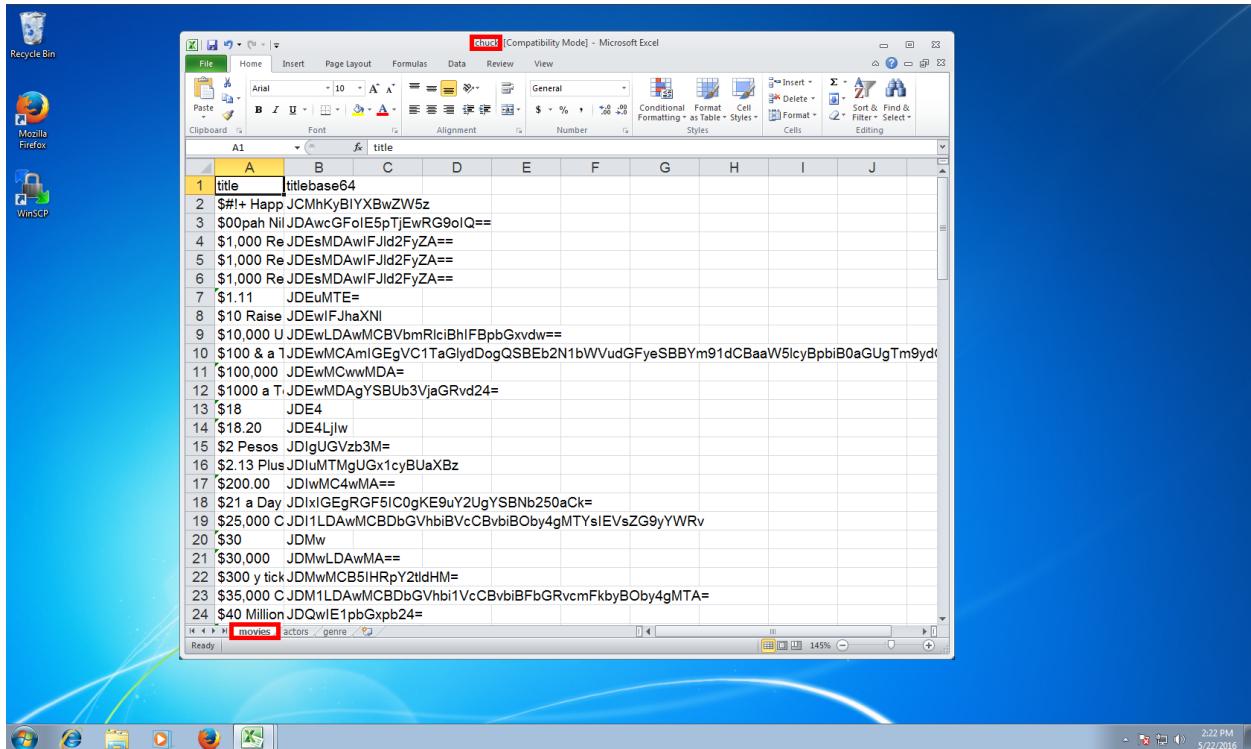
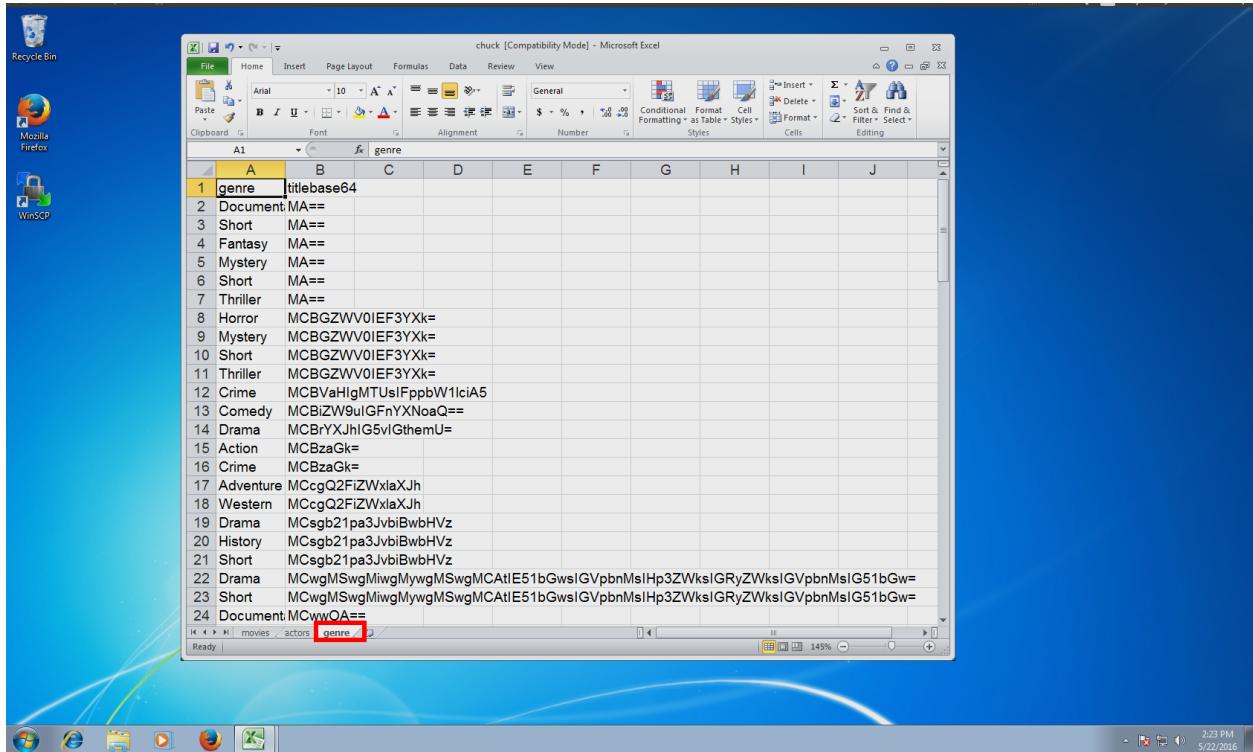


Figure 13: An Excel workbook with PostgreSQL data (movies worksheet). The workbook/database name and worksheet/table are outlined in red.

The screenshot shows a Microsoft Excel window titled "chuck [Compatibility Mode] - Microsoft Excel". The worksheet is named "actor". The data consists of 24 rows, each containing two columns: "actor" and "titlebase64". The "actor" column contains names like "Homo", "Steve", and "Short, Too". The "titlebase64" column contains long strings of base64 encoded data. The entire table is highlighted with a red border.

actor	titlebase64
1 actor	tNlreXRhaXRIZW4gbXVzZW8=
2 \$, Homo	U3V1cmkgaWxsdKNpb25p
3 \$, Homo	RS5SLIBTbHV0cw==
4 \$, Steve	MjAxMjBBV4kgQXdhcmRzlFNob3c=
5 Short, Too	MjAxMjBBV4gUmVkiENhcnBldCBTaG93
6 Short, Too	NXRoiEFubnVhbCBWSDegSGlwIehvcCBlb25vcnM=
7 Short, Too	QW1cmjlyW4gUGltca==
8 Short, Too	QmVhdHmsIFJoeW1icyAmlExpZmU6IFRoZSBUcmF2ZWxzlG9mlGEgVHJpYmUgQ2FsbGVkIFF1ZX
9 Short, Too	QnJlVtbpmcgT3VOiBUaGUgQWxjXRyXcgQ29uY2vdA==
10 Short, Too	R2FuZ3N0YSBSYX46IFRoZSBHbG9ja3vZW50YX5
11 Short, Too	R2V0IEi0IFdaZXJlIFvdSBGaxQgaW4gMQ==
12 Short, Too	R2V0IEi0IFdaZXJlIFvdSBGaxQgaW4gMQ==
13 Short, Too	R2hldHRvIFBoeXnpY3M=
14 Short, Too	R2hvc3RyaWRliHRoZSBXaGlw
15 Short, Too	SGLwlEhvcCBVbmNbnNvcmVkiFZvbC4gNDogTWlhWkgVmijZQ==
16 Short, Too	TWFjERyZTogTGvhZW5kIG9mlHRoZSBCYXk=
17 Short, Too	TWVuYWNlIEIJFNvY2ldHk=
18 Short, Too	T3pvbmUgV2zdCaZ
19 Short, Too	UGltcGFsYKrbp246IFJldHVbiBvZiB0aGUgVHJpbGw=
20 Short, Too	UGxhbmV0IFJvY2s6IFRoZSBTdG9yeSBvZiBlaXTSG9wlGFuZCB0aGUgQ3JhY2sgR2VuZXJhdGlvb
21 Short, Too	UG9ybmRvZ3M6IFRoZSBHZlbnR1cmVzIg9mlFNhZGll
22 Short, Too	Umh5bWUgJiBSZWfzb24=
23 Short, Too	U2NhcmZhY2U6IEdyZWFOZKN0IhpHMgb24gRFZE
24 Short, Too	U3RvcCBQZXBwZXlgUGFsbWVy

Figure 14: An Excel workbook with PostgreSQL data (actors worksheet). The worksheet/table are outlined in red.



The screenshot shows an Excel spreadsheet titled "chuck [Compatibility Mode] - Microsoft Excel". The active sheet is named "genre". The data is organized into two columns: "genre" and "titlebase64". The first column contains genre names, and the second column contains their base64 encoded binary representations. The first row is highlighted in yellow.

genre	titlebase64
Document	MA==
Short	MA==
Fantasy	MA==
Mystery	MA==
Short	MA==
Thriller	MA==
Horror	MCBGZWV0IEF3YXk=
Mystery	MCBGZWV0IEF3YXk=
Short	MCBGZWV0IEF3YXk=
Thriller	MCBGZWV0IEF3YXk=
Crime	MCBVaHlgMTUsIPpbW1ciA5
Comedy	MCBiZW9uIGFnYXnoaQ==
Drama	MCBrYXjhIG5vIGthemU=
Action	MCBzaGk=
Crime	MCBzaGk=
Adventure	MCcgQ2FzWxlaXjh
Western	MCcgQ2FzWxlaXjh
Drama	MCsgb21pa3JvbIBwbHVz
History	MCsgb21pa3JvbIBwbHVz
Short	MCsgb21pa3JvbIBwbHVz
Drama	MCwgMSwgMwgMywgMSwgMCAtE51bGwsIGVpbnMsIHp3ZWksIGRyZWksIGVpbnMsIG51bGw=
Short	MCwgMSwgMwgMywgMSwgMCAtE51bGwsIGVpbnMsIHp3ZWksIGRyZWksIGVpbnMsIG51bGw=
Document	MCwwOA==

Figure 15: An Excel workbook with PostgreSQL data (genre worksheet). The worksheet/table are outlined in red.

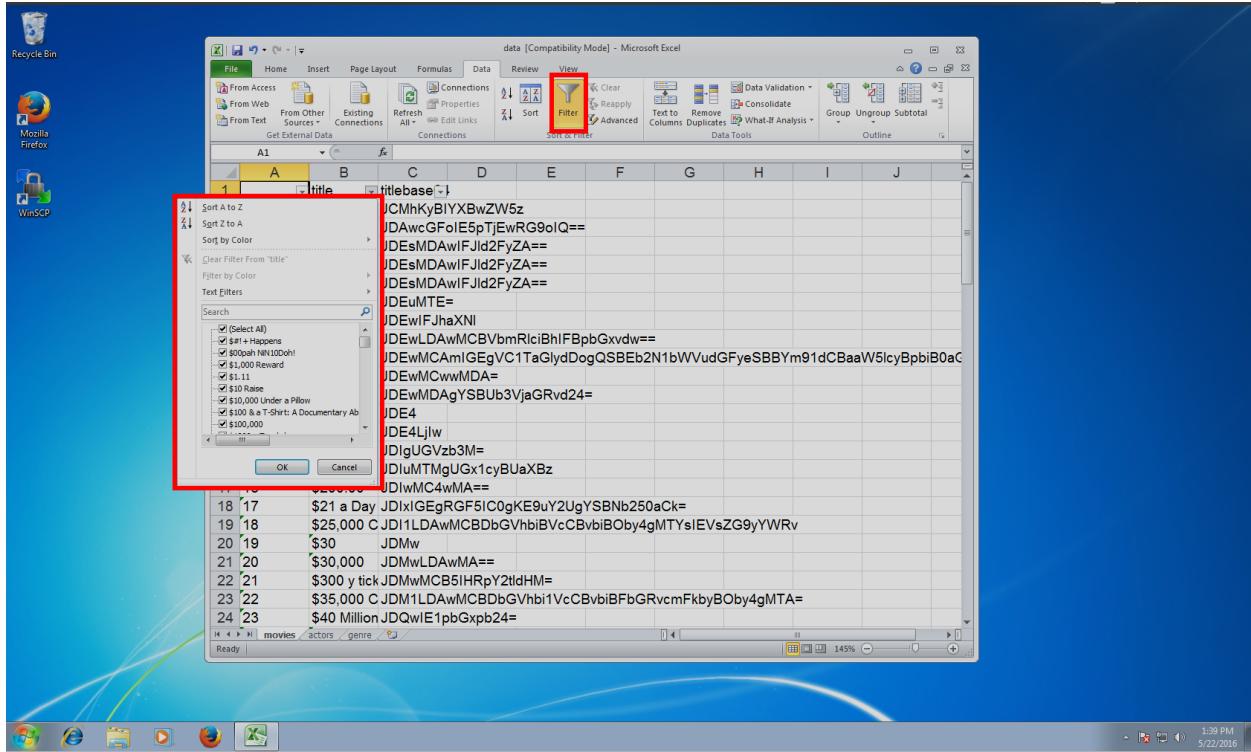


Figure 16: An Excel worksheet showing the “Filter” interface.

E Misc. files

The files used to create all these figures are attached to this report. They are:

1. chuck.xls - an Excel spreadsheet with 100 entries from the postgress movies, actors, and genre tables.
2. demo.R - an R program used to discover and make recommendations using the Internet Movie Database⁵. The PostgreSQL database used by demo.R can be created and populated by setting the appropriate command line arguments. 
3. pascalTriangle.py - a Python program to generate Pascal's Triangle suitable for use in a LATEXfile.
4. vennDiagrams.R - an R program used to create Venn diagrams based on the genres. 
5. writeExcel.R - an R program used to extract data from the database and make it available for Excel.

The demo.R supports a small set of command line arguments (see Table 14).

Table 14: demo.R command line arguments. Arguments can be entered in any order, but can only be used once per program execution.

Argument	Usage	Default value
m	The misspelled movie title. If the title is more than one word long then it should be enclosed in quotes (e.g., -m ‘‘The Terminator’’).	“The Dirk Knight”
p	Print SQL commands and results.	FALSE
r	Reset and restart the PostgreSQL database. This will drop and create all tables, and reload all data from source gzipped files in “well known” locations.	FALSE

(Continued on the next page.)

⁵<http://www.imdb.com/>

Table 14. (Continued from the previous page.)

Argument	Usage	Default value
t	The tex file for collecting data used in automated reports.	(NULL)
u	Select the user for the PostgreSQL database. Combining this option with the r option would create a unique database for a user.	chuck

F References

- [1] E.J. Borowski and J.M. Borwein, *Web-Linked Dictionary, Mathematics*, Harper Collins, 2005.
- [2] Francesco Calabrese, Giusy Di Lorenzo, Liang Liu, and Carlo Ratti, *Estimating origin-destination flows using mobile phone location data*, IEEE Pervasive Computing **10** (2011), no. 4, 0036–44.
- [3] David C Edelman, *Four ways to get more value from digital marketing*, McKinsey Quarterly **6** (2010).
- [4] Michael Gilleland, *Levenshtein distance, in three flavors*, <http://people.cs.pitt.edu/~kirk/cs1501/Pruhs/Fall2006/Assignments/editdistance/Levenshtein%20Distance.htm>, 2006.
- [5] Fabien Girardin, Francesco Calabrese, Filippo Dal Fiore, Carlo Ratti, and Josep Blat, *Digital footprinting: Uncovering tourists with user-generated content*, Pervasive Computing, IEEE **7** (2008), no. 4, 36–43.
- [6] Vladimir I Levenshtein, *Binary codes capable of correcting deletions, insertions, and reversals*, Soviet Physics Doklady, vol. 10, 1966, pp. 707–710.
- [7] Alamy Staff, *Stock photo - cupid rides the range*, <http://www.alamy.com/stock-photo-cupid-rides-the-range-ray-whitley-1939-72383173.html>, 2016.
- [8] IMDb Staff, *The dark knight*, http://www.imdb.com/title/tt0468569/?ref_=ttmd_ph_tt1, 2016.
- [9] PostgreSQL Staff, *A Brief History of PostgreSQL*, <http://www.postgresql.org/docs/9.0/static/history.html>, 2015.
- [10] Wikipedia Staff, *PostgreSQL*, <https://en.wikipedia.org/wiki/PostgreSQL>, 2015.
- [11] _____, *Ingres (database)*, [https://en.wikipedia.org/wiki/Ingres_\(database\)](https://en.wikipedia.org/wiki/Ingres_(database)), 2016.
- [12] _____, *Levenshtein distance*, https://en.wikipedia.org/wiki/Levenshtein_distance, 2016.
- [13] _____, *Venn diagram*, https://en.wikipedia.org/wiki/Venn_diagram, 2016.