

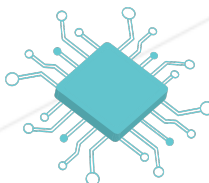


AP Computer Science Principles - Data Project

Data is Beautiful

Kai kai@42.us.org

Summary: Make your point by laying the facts out visually.



HACK
HIGH
SCHOOL



This work is licensed under a [Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-nc-sa/4.0/).

Contents

I	Introduction	2
I.1	Concepts	2
I.2	Hypothesis	2
II	Learn how to Handle Data	4
II.1	Read from a CSV	4
II.1.1	Graded Questions pt 1	6
II.2	Query an API	6
II.2.1	Register for API Key	7
II.2.2	Learn from the Documentation	7
II.2.3	Graded Questions pt 2	12
II.3	Put them together	13
II.3.1	Loop through CSV Again	13
II.3.2	Build the Query String	14
II.3.3	Use Abstraction on the Request function	15
II.4	Navigating the JSON Result	15
II.5	Store Results	17
II.5.1	Plan the Data Structure	17
II.5.2	Sorting by Month	18
II.5.3	Keep Track of Who's a Winner	20
III	Setting up a graph	23
III.1	Draw an X-axis and Title	23
III.1.1	An Opportunity for Abstraction	26
III.2	Draw bars with the example data	28
IV	Bonus: Multi-Year Slider	30

Chapter I

Introduction

I.1 Concepts

This assignment will proceed as a guided tutorial with discussion questions.

Consider three main topics in connection with this exercise:

- **Processing data:** Data is often stored in text files with simple formatting such as comma separated values (CSV) or [javascript object notation](#) (JSON). You can write code to loop through the lines of these files and extract the data.

A CSV is like a spreadsheet, but instead of being visually separated in boxes, each item is separated by a comma. You can export spreadsheets from Google Sheets or Microsoft Excel as a CSV.

- **Application Program Interface (API):** A way to send information back and forth from internet-connected databases.
- **Evaluating Hypotheses:** You must learn how to think scientifically about what conclusions you can be drawn from your data, and to understand when the data is insufficient to prove a hypothesis.
- **Displaying Data:** Have you ever stumbled upon an infographic on social media that made you pay more attention to an issue or change your mind about something you thought you knew? Artistic displays of data can be a more powerful way to explain your ideas than text or images alone. At the same time, if you are viewing someone else's data chart, you should be careful to fact check their sources and think logically about their conclusions.

I.2 Hypothesis


I will walk you through the process of collecting and visualizing data around a specific topic: [Can a Release Date Predict an Oscar Winner?](#)

In the article linked, author Zach Kram argues that which month a movie is released influences how likely it is to win an Oscar award.

Take some time to read the article. What do you think of his graphs and hypothesis? Do his visual data representations make the point of the article easy for you to understand?

In the tutorial below we will attempt to build a better visualization for his article. We will combine data from an API for [The Movie Database \(TMDB\)](#) and a CSV provided by [Datahub.io](#). The visualization will show a bar chart counting how many Oscar nominees were released in each month of the year. Oscar winners will be graphed on the same axis in a different color than the other nominees.

As the basic learning experience we will graph the release date distribution for just one year. If you like this project and decide to go further, we can add a slider that allows the viewer to interact by changing the year displayed.



APCSP/images/graph_2018.png

Chapter II

Learn how to Handle Data

Let's start by practicing two new coding techniques: Reading information from a CSV, and requesting information from an API.

We will combine the data from both sources together and store them in an object in our code.


II.1 Read from a CSV

Comma Separated Values (CSV): serves the same purpose as a spreadsheet, but it's displayed in plain text instead of a grid. Each "cell" is separated by a comma. Each row is on a new line.

Look in the p5js reference for the section on "IO", short for Input/Output. This will show you several ways to add data into your project (and, if you wanted to, how to export results).

We are going to use the [loadtable\(\)](#) function to read a CSV file. Let's first make sure you are comfortable adding a file to the project and reading its contents.

Go ahead and open the reference page for `loadTable()`. Look at the example code; in fact, let's copy and paste the entire example onto a new blank sketch to see how it works.



APCSP/images/loadTableExample.png


If you run the sketch at this point, it will do nothing except say "Loading..." forever. Why is that?

Notice the name of the CSV that it is trying to load on line 14. We need to create a CSV file in our project folder, and put the correct name of the CSV into that function.

Do you remember loading an image in APCSP part 05? Adding a CSV works just like that:

1. Find the arrow below the "play" button and to the left of "sketch.js". Click on it to expand the sidebar.
2. Click on the down arrow next to project-folder and choose "add file".
3. Call the file data.csv and press enter. Close the file explorer.

4. Copy the example data from lines 4-7 of the code and paste it into the data.csv file. Erase the comment marks.



APCSP/images/zoodata.png

Now, go back to the code in sketch.js and change "assets/mammals.csv" to "data.csv". Run the code again! Now you should get some console output from the results of the code in the setup() function. Read the code until you understand which line of code is producing which line(s) of output.

II.1.1 Graded Questions pt 1

Write down the answers to these questions:

1. Why is loadTable() called inside the preload() function? (Hint: Check the reference documentation for preload(!))
2. What code would you use if you wanted to print out only the 2nd item in the 2nd row of the CSV?

II.2 Query an API

Application Program Interface (API): A set of commands used to request information from a database server.

Server: A computer whose purpose is to stay on and connected to the network so you can request information from it.

Save or erase your sketch so far and start with a new one. We will practice making an API call before combining the parts.


II.2.1 Register for API Key

The first thing you'll want to do is register for your own API key. This is like a password that you use to identify yourself to the API.

1. Make an account at <https://www.themoviedb.org/>. You'll need to go into your email and confirm your account.
2. (HELP!) - If you don't have access to your email, you can partner with a classmate and try sharing their API key.
3. Log into your TMDb account and click on your first initial in the top right corner; select Settings.
4. Click on the "API" section. Then, click the link under "Request an API Key". Choose the Developer option and agree to the terms of use.
5. Fill out the form: Type of Use = Education. Application Name = APCSP Data Project. Application URL = codeforfun.org or paste a link to your p5js sketch. Application Summary = "Graphing release dates of Oscar nominees". Fill out your name and email. You can use the address of the school if you prefer: 6600 Dumbarton Circle, Fremont, CA 94555.
6. Now you have an API key! You have two of them, in fact. We will be using the Version 3 access key, so copy the shorter one (v3 auth).
7. Paste this API key as a comment at the top of your blank sketch. You'll need it later.

II.2.2 Learn from the Documentation

I'll take you on a brief tour of the steps I used to figure out how to use this API from p5js code. First, consult the [API Documentation](https://developers.themoviedb.org/) at developers.themoviedb.org.




APCSP/images/tmdb-apidocs.png

The links on the left are mostly **endpoints**, indicating the type of information that the API has to offer. We are going to use the endpoint "Search".

Click on "SEARCH" in the sidebar and then, under Search, click "Search Moves". On this page the documentation will show you what parameters you can use, and what the response will look like.

Switch to the "Try it out" tab and expand the "Code Generation" section. Choose Javascript. Nice of them - they already wrote some code for us to use!




APCSP/images/tmdb-searchmovie.png

NOW, let's fastforward and say that when I tried using their example code I got a bunch of errors. I tried learning more about the function it suggested that we use, XMLHttpRequest. What worked for me was to go to the official Javascript documentation (Mozilla) and use one of their [code examples for XMLHttpRequest\(\)](#). You may also want to skim the [main page](#) of documentation on XMLHttpRequest, but don't worry about trying to understand all of it.

Copy the code example for XMLHttpRequest. Inside your blank p5js sketch, create a new function called request() and paste that code inside of it. Make two edits:

1. Change xhr.responseText to 'json' instead of 'text'
2. Delete the line where it prints out xhr.responseText.

Then, add a call to `request()` from `setup()`.



`APCSP/images/xmlhttp.png`

We are going to edit the `"/server"` parameter and replace it with a long link describing the web address of the database, plus the endpoint that we want to pull from, the parameters we will use, and our API key.

Go back to the tab where you are logged into TMDb, or [click here](#). Make sure you are in the "Try it out" tab. Find your API key from earlier and paste it in the "Variables" section where it says "Your TMDb API key".

For the "query" parameter, type the name of a movie that you like.


Click "Send Request" and look at the results. You might have one or more results. In JSON format (which is basically a lot of arrays and objects packed together), the results tell you information such as: its popularity score, original language, plot summary, and... release date!

If you see multiple movies listed in the results, try filling in the "year" parameter for the actual year of the movie you were looking for. Run the query again. Did that help reduce the number of hits?

OK, now copy the long link next to the "SEND REQUEST" button, starting with

<https://api.themoviedb.org/3/search/movie?....>

Notice how this link is made up of different logical parts.



APCSP/images/apirequest.png

Back in your p5js editor, paste that link inside `xhr.open()` to replace `'/server'`. Make sure you wrap the link in quotation marks.

Place a call to `request()` inside of `setup()`. Then, run the code and see what it outputs!

APCSP/images/legosuccess.png

II.2.3 Graded Questions pt 2

1. Which variable in your API query code contains the main results of your query? What data type is that variable?
2. If you wanted to add another parameter to the query string without using the TMDB API docs, how would you do it? Say you have this query string and you want to also specify that "director" should be "Miyazaki".

```
https://api.themoviedb.org/3/search/multi?api_key=
6bb429redactedredactedcf95113e&genre=anime&language
=en-US&page=1
```

II.3 Put them together

II.3.1 Loop through CSV Again

OK, it's time to revisit your CSV reading abilities, use them on a CSV that contains movie info, and search for the release date of that movie in the API. Still with me? Let's go. Here is a CSV with the nominees for two main categories in 2018.

```
year,category,winner,entity
2018,BEST PICTURE,False,Black Panther
2018,BEST PICTURE,False,BlacKkKlansman
2018,BEST PICTURE,False,Bohemian Rhapsody
2018,BEST PICTURE,False,Vice
2018,BEST PICTURE,False,The Favourite
2018,BEST PICTURE,False,Green Book
2018,BEST PICTURE,False,Roma
2018,BEST PICTURE,False,A Star Is Born
2018,DIRECTING,False,Cold War
2018,DIRECTING,False,Roma
2018,DIRECTING,False,Vice
2018,DIRECTING,False,BlacKkKlansman
2018,DIRECTING,False,The Favourite
```

I have put it in the same format as the historical data, and assigned "False" to the "winner?" column for all. Go ahead and pick your favorite movie in each category and change that value to "True". Since we don't know the winners for this year we'll just make it up. Later, if you are interested enough to expand this project, we'll graph all the historical truth.

You can keep working on your sketch with the API code. Copy and paste this CSV data into a file on your project, like we did with the mammals example. Set up the `loadTable()` command using the `preLoad()` function that we discussed before. Then pull up the `loadTable()` reference page to remember which commands you can use to navigate the CSV data.

Write some code that does the following:

- Loops through each line of the CSV (yes, you might want to check that example again). (Hint, there's no need for a `*double*` for-loop.)
- On each line, it grabs the year and title and stores them as two variables.
- Ideally this would be in a new function called `read_csv()` that is called from the `setup()` function.
- Also, create a global variable for your `apikey` and store it there.


II.3.2 Build the Query String

I'll help you with this next part. The goal is to ask the TMDb database for the release date (not just the year, but the day and month) of each movie on the Oscar nominee list.

To do this, we need to generate a new query string for each nominee. We'll do that by adding (concatenating) the constant parts of the string together with the variable parts of the string.

You have used string concatenation before... in print statements, for debugging! This time we are going to save the result of "string addition" in a variable, instead of printing it out.

Oh, and as standard procedure you need to use the Javascript function `encodeURIComponent()` to replace all the spaces in the request with percent signs.




APCSP/images/make_query_string.png

II.3.3 Use Abstraction on the Request function

Modify the `request()` function so that it takes in the query string as a parameter and uses that inside of `xhr.open()`.

Modify `read_csv()` so that for each row, it generates a query string, and then calls `request(query)`.


I'll let you look at my code if you need to catch up.



APCSP/images/catch-up.png

II.4 Navigating the JSON Result


In the console, click the arrow next to one of the Object printouts to examine its contents.



APCSP/images/json-pretty.png

This is basically an **object** with many **variables** including "page" and "results". "results" is an array of objects. It could contain more than one movie, but it just has one in this example. Each movie in the result array is an object that has variables such as "vote_count", "id", and "video".

We want to pull out contents the "release_date" variable (at the bottom of the image). It's the first time we have worked with this type of data, so here's some code that will do it. Add the `find_release_month` function to your sketch and call it from inside the `request()` function, right after printing out the response.



APCSP/images/read-json.png

This code uses dot notation to reach inside of an object, and bracket notation to reach inside of an array. Look at what's happening:

- "response" is an object
- "response.results" is an array of objects
- "response.results[0]" is the first object in that array
- "response.results[0].release_date" is the variable labelled release_date on that object.

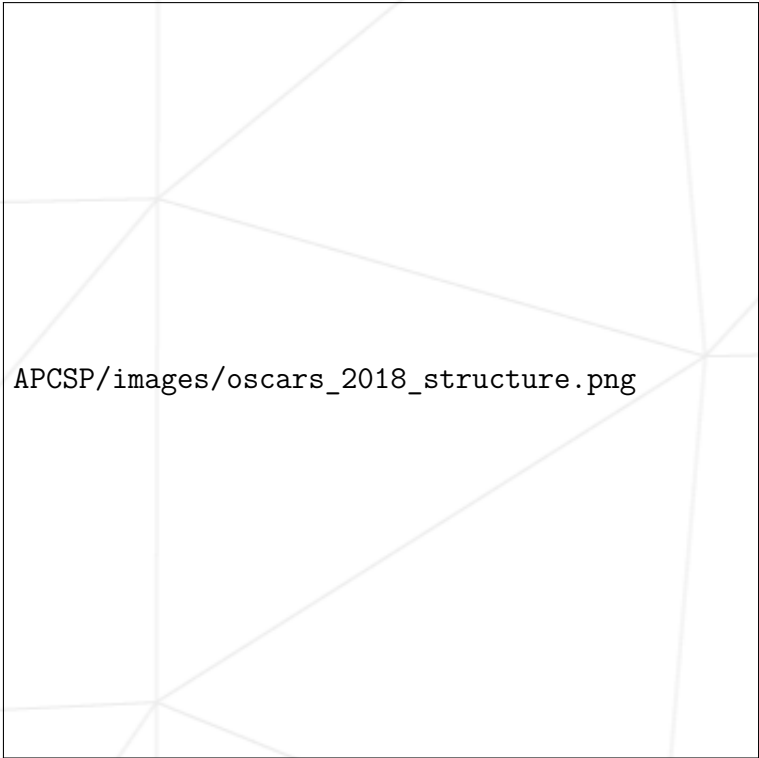
The result is in the format "YYYY-MM-DD", and we'll grab the month from it later.

II.5 Store Results

II.5.1 Plan the Data Structure

Similar to this [Visualization of the results rolling two dice](#), I'm planning to draw a bar chart simply consisting of vertical lines along an axis.

We need to store the release dates and movie info into a format that makes it easy to draw a graph. There are lots of valid ways to do this, but after thinking about it a lot, for this class I decided to use the following structure:



APCSP/images/oscars_2018_structure.png

The data for this year will be stored as one global variable, `oscars_2018`. That variable will be an object that is built from a constructor function. The constructor function is called `Year()`. Every `Year` object has two double arrays: one for all the nominees, and one for just the winners.

When a new `Year` object is created, the code loops through a counter for each month (1 to 12, or 0 to 11) and creates that many inner arrays in both the "winners" and "nominees" arrays.

As we go through the data I want to identify the release month and whether the movie was a winner or not. Then, I will copy the name of the movie into the corresponding spot in these double arrays.



Check Your Understanding: What information will be stored in `oscars_2018.nominees[3]`?

II.5.2 Sorting by Month

In order to fill up the data structure, let's add some logic to the `find_release_month()` function.

I want you to write this code yourself. Here it is in pseudocode:

- Split the release date string by the "-" delimiter. (HINT: research how to use Javascript's `string.split()` function!)

- The second item in the resulting array represents the release month. Use the Javascript function `parseInt()` to transform it from a string to a number.
- This month is represented in the way humans count, from 1 to 12. Subtract one from the month number so that we can use it with our array of twelve arrays, which are numbered from 0 to 11.


Can you take that psuedocode and turn it into real code? Use your brain. Type it in to the `find_release_month()` function.

Next I am going to create a new function, called `store_result(month, title, winner)`. Here is the code for it:



```
APCSP/images/store_result.png
```

Copy this function into your sketch and we'll take a minute to understand what it does. One thing that will help is to add the command `print(oscars_2018);` into your `setup()` function. This will print out a representation of the object, while all of its arrays are still empty:



APCSP/images/oscars_object_empty.png

The `store_result` function takes in three parameters: `month` (as a number), `title` (as a string), and `winner` (it should be a boolean). All movies will be added to the `nominees` array. `oscars2018.nominees[n]` takes us to the 'nth' array, which represents the 'nth' month - January = 0, February = 1, March = 2 and so on.

`push()` is a basic Javascript function which allows us to add an item into the chosen array.

Then, the `"if(winner)"` statement checks if the **winner** boolean is **true** and also adds this movie to the winner array if it is.

II.5.3 Keep Track of Who's a Winner

It's time to go back and fix a detail that we ignored so far: we need to extract the "winner" variable from the CSV, and pass it from function to function until it is used to sort films in our data object.

Trace the flow from one function to another in the code so far. What order do the functions run in?

1. `preload()` runs first (and variables are initiated about the same time, so `Year()` runs somewhere around here).
2. `setup()` calls `read_csv()`.
3. `read_csv()` makes calls to `make_query_string()` and `request()`.
4. `make_query_string()` runs and returns a value.
5. `request()` runs next, and makes a call to `find_release_month()`.
6. `store_data()` hasn't been called by any code yet.

I suggest that we make the call to `store_data()` from `find_release_month()`. But, `store_data()` needs a parameter for the **winner**! `find_release_month` has no such information available at the moment.

It does NOT make sense to store `winner` as a global variable, because it is different for every movie! We need to keep track of it alongside the other movie-specific information.

Please make the following changes to the code:

- Inside `read_csv`, add the line `"var winner = (table.getString(r, 2) == 'True');"`. The `"=="` statement checks if the string from the table is equal to the string `True`. If that evaluates to `True`, then the variable `"winner"` will receive the `True` boolean. Otherwise, it is a `False` boolean.
- Modify the `request()` function definition so it takes two parameters now, `query_string` and `winner`. Send `winner` to `request()` when you call `request()`.
- Modify the `find_release_month()` function so that it takes two parameters, `response` and `winner`. Send `winner` to `find_release_month()` when you call that function.
- Also, add a line to `find_release_month()` which extracts the title from the response object and saves it as a variable called `"title"`. This one line of code will be very similar to the line of code that extracts the release date.
- Add a call to `store_result(month, title, winner)` at the end of `find_release_month()`.
- Run the code and see if it works! If you are stuck, try debugging with print statements. Verify that it works properly by printing out the `oscars_2018` variable at the end of `store_result()`.



There is an advanced topic hidden in the code here: **asynchronous Javascript**. Learn about it more when you have time. What it means for us today is that we cannot return a value from the `request()` function and jump back into the `read_csv()` algorithm. `request()` takes a relatively long time to run each time, and the other code in the program will not wait for it before moving on. This could result in the "response" variable being undefined. That's one reason I am telling you to write the functions in a specific way.

APCSP/images/oscars_object_full.png

Above is what the Oscars 2018 object will look like if you are successful so far!

Chapter III

Setting up a graph


Similar to this [Visualization of the results rolling two dice](#), we can draw a simple bar chart using just lines and the coordinate system. If you want to improve this graph to use rectangles for the bars, do it as a bonus at the end.

III.1 Draw an X-axis and Title

Plan out the size and shape of the graph: We will place 12 vertical bars, one for each month. Let's say they will be spaced 40 pixels apart with 50 pixels padding on each side of the graph. That means the x axis should start at $x = ?$ and end at $x = ?$. Can you fill in the question marks?

Vertically, how tall will the bars be? Right now we have data for Best Picture and Best Director. Let's say we will show them together without separating them in the chart. Most categories always have 5 nominees, and the Best Picture can have 10 nominees. So we will be graphing 15 movies for each year if we use just these 2 categories. We should make space for a bar that is 15 units tall just in case, by some coincidence, all nominees were released in the same month. Let's say that each vertical unit is `<y_scale>` pixels tall and there is padding of `<top_margin>` pixels on top for the title.

Create global variables that store your choices for the following values: spacing between bars, how many vertical pixels are used per movie counted, margin on each side of the graph, top margin pixels for the title, and bottom margin pixels for the labels.



`APCSP/images/scaling_variables.png`

Then you can use the following code to create a canvas of the corresponding size and an x axis for labels. Notice that `setup()` is responsible both for triggering `read_csv` and for setting up the graph.

APCSP/images/chart_setup.png

III.1.1 An Opportunity for Abstraction

Phew, that was a lot of code to type or copy over! Did you notice some parts that are repetitive? As I was writing those equations I found myself calculating the **y coordinate** of the **x axis** multiple times. You can clean up the code by storing that value in another global variable.

1. Add another global variable, `var x_axis_y_intercept`, but don't assign a value to it.
2. At the beginning of `setup()`, set that variable equal to the `top_margin` plus `y_scale` times `max_bar_height`. We are counting down from the top and that is the maximum bar height, in pixels, that could occur.
3. Replace every place you see `top_margin + (y_scale * max_bar_height)` with `x_axis_y_intercept`. It's used in creating the canvas and drawing the axis line.

APCSP/images/x_axis_y_intercept.png

III.2 Draw bars with the example data


Now we'll create a function to draw the bars of the bar chart, and test it with the example data.

In pseudocode here is what I am going to do:

- Loop 12 times so I follow the same algorithm for each month:
- Check if there were any nominees released in the current month, and if there were, draw a blue line from 0 to the number of nominees.
- Check if there were any winners released in the current month, and if there were, draw a yellow line from 0 to the number of winners, on top of the previous blue line. This works because we will include winners in the count of nominees, storing the titles of winning films in both arrays.

You can type the following code to complete this part. Notice that the hard part is figuring out the Y coordinates of the line. Since 0 is at the top, we subtract the length of the bar from the Y coordinate of the X axis.

Another hard part is pausing the draw() loop and calling it at a specific time when you want the graph to be redrawn. Notice that I include noLoop() in setup (this will prevent the graph from drawing before the data is available), and I call redraw() every time that store_result updates the data. I moved the background, axis and title drawing steps from setup() to draw().



`APCSP/images/draw_bars.png`

Good Luck! Debug by thinking critically if you had any errors while typing it over! Feel free to ask your neighbor for help. Remember, the functions can be defined in the code in any order and it will work the same, so feel free cut and paste them into the positions that help you understand your code.

Chapter IV

Bonus: Multi-Year Slider

The bonus chapter of this assignment is a challenge:

Take the full CSV from datahub.io that I linked at the beginning of the PDF.

Creae a whole array of Year objects so that you can store multiple years worth of data.

Loop through the whole CSV and pull the information you are interested in.

Add a slider using the p5.dom library. In the draw function, check the value of `slider.value()` in order to decide which year's data to draw.