

Scraping new papers for EGRET:

Getting Started: copied from osprey howtoscraper2019.tex

You'll need:

1. Excel or other program that makes .xls or .csv files
2. ImageJ download for free from here <https://imagej.net/Welcome>. You'll also need to add the Figure_Calibration.class, which will help for giving x and y calibrations to images.
 - (a) To add the the Figure_Calibration.class: In ImageJ go to plugins, select *add plug in*, then navigate to the Figure_Calibration.class file that you downloaded on your computer, click on it and follow through a few clicks to add the plugin.
 - (b) If you have some trouble getting the measurements to show up after calibrating, try switching to the pointer and clicking. You might need to set the preferences on your pointer tool to auto-measure.
 - i. A clarification on above from Tim Savas (original OPSREE lead data enterer): *After doing the figure calibration and selecting the yellow pointer tool, you start clicking inside the figure but no points appear. The reason for this is that the "rectangle" you drew for the figure calibration is still masking the figure, and until you click out of it, you can't draw points under it. It's hard to see! So to get rid of the invisible rectangle, just click the mouse once outside of its edge. Side note: After drawing all of your points onto a figure, you can press Command-M to bring up the resulting table of values. I do this in the video tutorial, and whenever I scrape, but didn't describe the key command!*

Now here's what to do:

1. Copy the excel file data/egret for git repo (egret/data) and make your own extension using your initials, for example, Jane Doe would write "egret_dmb". This will be the spread sheet you enter your data into and then in the future, someone will save the tab as a csv and merge all of our files together into the master data.
2. Familiarize yourself with each tab:
 - (a) **meta_general**: metadata for each sheet
 - (b) **source**: list of the paper we are working with. Bibliographic information and notes on usefulness for our purposes. Note the "assigned.to" column, which tells you which figure or table to focus on. You may find upon reviewing the text the paper does not suite out selection criteria, in which case you would change the 'A' (accept) in your accept_reject column to 'R' for reject. This tab should also be updated if a paper is not published in english.
 - (c) **data_detailed**: Detailed data for the experiment, with all relevant information filled out.
 - (d) **scratch**: For temporary formatting and manipulating data scraped from ImageJ.
 - (e) The two most important tabs to fill out are **source** and **data_detailed**.
3. Read your paper and update out the information in the "source" column if needed and fill the "data_detailed tab." Be sure your datasetID and study info agrees with the source tab, if not—figure out what is wrong and fix it.
4. Read the paper to decide if the study is eligible for inclusion in EGRET. These selection criteria include:
 - Seeds are germinated under experimental conditions
 - Not conducted under natural conditions in the field

- Species are not crops
 - Experimental treatments are not confounded
5. Take a screen shot of the figure and import into ImageJ, following detailed instructions from Tim's OSPREE data scraping video from [git/ospree/notes/howtoscrrape/Data Scraping Tutorial.mp4](https://git/ospree/notes/howtoscrrape/Data%20Scraping%20Tutorial.mp4) or the general instructions outlined in the egret wiki. Use the scratch tab to get data into the right format, and then copy into data_detailed.

A few more “how to’s”, trouble shooting etc

- Generally, we try to follow the language used in the paper, for example if a level of a treatment is referred to as chilling, but is actually warm, we still put it in the chilling column.
- **Our definition of a crop:** we defined crops as anything that is bred and produced at scale and for consumption. We still included species that were edible but foraged and species that are likely to grow in natural stands.
- **For dealing with lats and lons:** Tool to convert to decimal latitude and longitude: <https://www.fcc.gov/media/radio/dms-decimal> and remember to add NEGATIVE to your longitude if it's West. Also, you can check where things are by just typing in lat and long into Google maps.
- **On entering response times:** If the figure is a time-series or curve of the percent germination over time, the days to germination will be the values from the figures's x-axis. Otherwise the value will be the final day that germination was observed. Remember for time-series to also include the zero point, or to include days along the x-axis that have zero on the y-axis as this is still part of the curve.
- **On dealing with error:** If it's figures records error and it is *clear enough* to scrape, error can be recorded. Often times the SE bars are in the way of each other or not quite discernible, in which case we've decided to avoid them. But if the bars are clear, record them. Values can go in "resp_error" and just SE in "error type."
- **For dealing with cold storage conditions:** we defined storage as the period between when the seeds were obtained or collected and when they started applying treatments to them. See below for more detail on how we decided to deal with studies that has cold storage that could be physiologically similar to chilling.
- The paper pdf's were saved in a Google drive folder. Unfortunately, the folder went missing in the summer of 2022. We believe what happened was that the folder was made using a student account that was deleted once the owner graduated. Since we were not notified that the folder went missing until the end of the field season, it was too late to request it be temporarily restored and it was never backed up. We reached out to everyone and were able to get many of the pdf's back in a new folder (now owned by Lizzie). We have now re-downloaded the papers into a new folder and backed it up.

Data cleaning and decisions

1. Cleaning code are located in the analysis/cleaning/source folder in the egret repo.
2. Early notes about data scraping can be found at <https://github.com/lizzieinvancouver/egret/issues/4>, with some highlights being:
 - We are trying to use camel case where possible for both script names and column names.
 - Adding a column for the rare cases where studies start measuring germination day relative to the day of first germination and not the day they start monitoring germination (i.e. day of germination temperatures, imbibing etc)

- In general we excluding papers not in english unless the person scraping the data could read the language fluently. But a few papers not were still included accidentally.
 - Decided to just add seed mass as a yes/no for now, we will have to go back and scrape it some day if we want it
 - When dealing with error, enter it as zero if bars are not distinguishable from the points, or as "indistinguishable" if multiple overlapping points have overlapping error bars that can not be scraped with certainty (see issue of a visual example of this).
3. We decided at our first retreat which columns we wanted to clean, the list of columns and who was assigned to clean them is outlined at <https://github.com/lizzieinvancouver/egret/issues/14>, key points from this issue are discussed below for each cleaning script they pertain to. This is also the issue were we list updated column names.
 4. JS and missingData.csv had 10,000+ empty rows in their csv data files for some reason, which was making the code run slowly. DL manually corrected this in the csv files.
 5. The following outlines some of the more complex cleaning decisions and issues:
 - **cleanChemical.R:** Chemical names should be cleaned and standardized, updating names of brand name products to clarify what they are (e.g. include explanations such as "—herbicide"). Names were updated to remove the treatment duration (which was moved to the trt.duration column where it belongs). Also see issue #65.
 - **cleanScarification.R:**
 - We initially thought to move the hot water scarification to the soaking column, but it was confirmed that this was in addition to soaking, so it was left as is.
 - There was also a question about "cold scarification", which had values duplicated in the chilling column (which seemed odd), this is not mentioned in the script, which I think means that is was explicitly called this in the paper.
 - Issue #14 has details about some of the add entries (e.g. partial scarification) that were double checked. Also see issue #41 and #62
 - **cleanCoordinates.R:**
 - For 123 studies there was no specific lat and longs given, but we still estimated one based on the location textually described the paper.
 - We are only interested in source population locations, not where studies were conducted.
 - To infer the coordinates for these studies, we used Google Earth to drop a pin in the centre of the named region (county, city, etc) and get a general value for these studies.
 - We added an additional column to denote when we needed to take these less specific site locations. I didn't change the columns but made detailed notes in the script.
 - To check that coordinates seem right we plotted them on a map and assess one by one if the locations made sens according to their source.population
 - See issue #21 for cleaning coordinates, #48 on number of studies with multiple provenances, #47 mapping code.
 - If available we kept a minimum of two decimals per lat/long, but in some cases, vague coordinates of a region were provided with no decimals and we kept those indicating a less specific site location.
 - If seeds were from three distinct locations, but they were mixed in the germination trial, we took the mean latitude and longitude of these locations.
 - **cleanGerminationTempDuration.R:**

- Germination temperatures in 108 studies varied, either as a range or between day and night. These studies with multiple temperatures were manually checked and additional data columns entered for day and night temperatures.
 - For studies that just report temperature ranges, we took the mean of the two.
 - For studies that define conditions as "open air" we consider this to be equivalent to "ambient"
 - If it was not specified what temps were day and night, we assume the warmer temp is day.
 - We did not explicitly collect data on thermoperiodicity, but apparently all papers do say the warm treatment coincides with photoperiod, which allows us to calculate the means and GDD accordingly—but some papers have photoperiods longer than a 24h cycle—here we did a similar weighting as we would if it was 48h or 72h and assume that the daylight lengths reported are correct to get the number of night hours
 - Supra-optimal germination temps: The issue was raised whether greater mgt with temp is an artifact of high seed numbers rapidly increasing germ percents. How should we deal with this in the analysis?
 - See issue #18 for germination temp and duration cleaning and notes about unusual values and #40 regarding supra-optimal germination temps
- **cleanChillTempDuration.R:**
 - Some studies have alternating day and night temps, so we created unique columns for alternating temps and daylength.
 - If duration of chilling wasn't provided, we input NA for duration, but we kept the temperature.
 - Some studies include warm temperatures as a chilling level. Sometimes, when it was a chilling treatment following a warm stratification, we entered it in chillTemp and its corresponding chillDuration. E.g. "61 days at 20 C then 61 days at 5C".
 - There are 30 studies that are missing chilling temperature and durations. These pdf's were manually checked and notes made to check whether this is true or improperly entered data. It was discovered that in some of these studies, there are no chilling treatments, but there is moderately cold storage temperatures.
 - We discussed whether this storage could function as chilling, but decided that it depended on 1. whether the storage is wet, which would allow the seeds to be imbibed, 2. storage is between -20 and 10 degree C.
 - Option 1: Create new columns that include cold and wet storage in our chilling calculations for all studies, not just those that have no reported chilling treatments.
 - Option 2: Combine the chilling and storage columns: <https://github.com/lizzieinvancouver/egret/tree/main/analyses/analyseSeedCues/cleaningSeedCues>
 - Code to calculate chilling units is in analysesSeedCues
 - See issue #23, or #44 to check papers without chilling for which storage temps may be cold. See issue #39 for details on storage
 - chillTemp: denotes the temperature of the chilling treatment. If it was alternating, we denoted it as X then Y
 - chillDuration: same as for chillTemp, but for the number of days
 - chillTempUnc: uncertainty of temperature when given in the article. E.g 5±1Celsius. 5 would be the chillTemp and 1 would be the chillTempUnc
 - chillTempCycle: - period of each alternating cycle for the entire chillDuration of that chilling condition Example: a chillTemp of "10 and 20" with a chillTempCycle of "16 and 8" and a chillDuration of "14" means 16 hours at 10°C then 8 hours at 20°C and so on for 14 days; or, a chillTempCycle of "168 and 168" with a corresponding chillTemp of "4 and 22" and chillDuration of "112" means 168 hours (or 7 days) at 4°C then 168 hours at 22°C for 112 days

- chillLightCycle : (hours/day) - daily light. Example: "12 then 12 then 0" means 12 hours of light for the first and second chilling conditions, and the third chilling condition is implemented in the dark
- **cleanPhotoperiod.R:**
 - 326 studies have light/dark treatments
 - created new "lightDark" column
 - Photoperiod data was in both the photoperiod column and in the other.treatment column
 - A total of 113 papers have no photoperiods given. If studies done in complete darkness, photoperiod should be zero, but may have been entered as "NA". Also see issue #42
- **cleanStorage.R:** Here we created a general column which grouped similar treatments and one that included more diverse treatments. The cleaning of storage is now linked to chilling and whether studies without chilling have moist storage. Also see issue #39.
- **cleanYearGermination.R:** Over 100 studies do not specify what year the study was conducted. It was discussed and our meeting Aug 2 as to whether it was worth our manually checking the paper. It was decided that it was not, although if we do change our mind, it would be possible for a number of studies to infer it based on the language used in some papers (e.g. seeds were stored for x months after collection"). Also see issue #63, closed for now.
- **cleanResponseVar.R:** code cleaning the response variables.
 - Decided for now we are most interested in percent germination and mean germination time. See issues #25 and #33.
 - There are some issues with the response variables units, particularly for germination rate or speed, this is discussed in issue #26.
 - For germination curves, we extracted the maximum percent.germination reported for each study.
- **cleanMisc.R:** In this code we are removing redundant data that was reported in both tables and figures of studies, and studies for which we can not use the data (treatments are not properly described). Here we also clean minor columns, such as seed.mass.
- **cleanSpecies.R:** We cleaned the species name using WorldFlora package.
- **Manual removal of crops**
 - While we we tried to avoid crops in our datascaping, we double checked manually.
 - Dan inspected the list of species, looking up the latin binomials of any unfamiliar species to see if they were crops.
 - The *criteria* for being removed as a crop were species that a) have a long history of domestication and therefore likely to have had artificial selection on germination or b) a cultivar.
 - Species used in agriculture (i.e., cover crops) were not dropped, as we do not expect there is likely to be extensive artificial selection on germination behavior.

Scraping the USDA Seed Manual:

Following our first egret retreat in November 2023, we realized that the egret data had very little data for woody tree species. But we were very interested in merging both the egret and osprey datasets. We tried to better understand why this was the case and whether there were other sources for that specific dataset. In the end, one of the best sources we found was the USDA Seed manual: https://github.com/lizzieinancouver/egret/blob/main/analyses/scrapeUSDAseedmanual/input/USDA_woody_plant_seed_manual_2008.pdf.

Getting Started: copied from egret/analyses/scrapeUSDAseedmanual.md

- Used Amazon Textract—for early notes on this see: https://github.com/lizzieinvancouver/egret/blob/main/analyses/scrapeUSDAseedmanual/scraping/data_scraping_log.txt
- Produced files for each scraped data table—but each table is uniquely formatted as each section of the book has different authors
- Numbered folders are the page ranges from which tables were taken
- Tables were scraped if they contained data on 1. phenology, 2. germination, 3. traits —scraped data was sorted into each of these categories using `renameRelevantDataTablesScript.R`
- Sometime stratification temperatures are in the text or the notes of the tables which could not be scraped by Amazon Textract. We checked the column stratification duration and then went back to the original book to see if any important information is missing and added the data manually
- Went through the original book checking for tables did not get scraped by Amazon Textract and scraped those tables manually
- For efficiency we began cleaning the germination data—see #15

Data cleaning and decisions

- Our goal is to format the data in such a way it can be merged with the egret data, following outlines some of the more complex cleaning decisions:
- **cleanSpeciesUsda.R:** We cleaned the species name using WorldFlora package.
- **cleanMeanUsda.R:** We converted entries with a range of numeric values into three columns (min, max and avg), these columns include: pretreatment, coldStratDur, photoperiod, tempDay, tempNight, germDuration, samples, chillDuration and responseValue.
- **cleanResponseUsda.R:** for response variable we converted "percent.germ.total", "percent.germ", "germ.capacity", "mean.germ.capacity" and "percent.germ.15degC.incubated" to "perc.standard"
- **cleanChillUsda.R:**
 - We created a new column called chilling with 'Y/N' data, this includes both original pretreatment chill and cold stratification.
 - We combined chilling and cold stratification entries into a new class of columns called: chill.dur.Min and chill.dur.Max, and chill.dur.Avg
 - Entries with only one original value for chilling end up in chill.dur.Avg.
- **cleanScarificationUsda.R:**
 - We selected treatments that we believe are scarification in the pretreatment column (including acid, soaking, mechanical, etc.) and made two columns: one (scarifTypeGen) for the general type of scarification, and the other (scarifTypeSpe) for the specific scarification details.
 - For scarifTypeGen, we have four types: soaking, chemical, mechanical and unknowScarification
 - For scarifTypeSpe, we first include the general scarification type, followed by the specific scarification method if that information is available.
- **cleanPreTrtUsda.R**
 - In the pretreatment column, we subsetted all treatment might be related to chilling into a new column.

- We select all rows with non-NA values in the column pregermination treatment hot water soak temperature and assigned 'Hot water' treatment to these rows.
- **cleanEgretColUsda.R:** We updated the names of all columns to better fit with egret