

Check the 2022 data and incorporate it into gfiphc

Andrew Edwards

Last compiled on 20 April, 2023

First read the survey section of the latest IPHC report ([here](#)) for background if it's published (2022 report published 28th March 2023). In 2022 (as for 2020 and 2021) only the first 20 hooks were evaluated, so those data are not easily imported into GFBio. Going to incorporate into gfiphc here. Likely need this as a template for future years: resave this file with new year, and change all 2022's to the subsequent year, go through the code somewhat manually to check the output as you go along (in Emacs do Alt-query-replace to change years but read carefully as going along), and then finally render the full document to make the .pdf. This code includes some manual checks to make sure the data look okay. The planned stations for the 2022 survey are shown in this IPHC sampling manual ([click here](#)); page 4 again has Vancouver [Island] Outside, showing that not all stations were intended to be fished there. The 2022 annual report notes issues with the survey (such as ship availability due to higher Sablefish quota, crew availability) that may have reduced the number of stations, though this was referring to the full survey so need to do the maps here to understand. After doing these analyses there are 171 usable stations in 2022, though only 113 are standard; some being up in inlets, but for spatiotemporal analyses we can include the useful non-standard ones.

For comparison first look at 2013 data included in gfiphc:

```
load_all()
> i Loading gfiphc
setData2013
> # A tibble: 170 x 8
>   year station  lat  lon avgDepth effSkateIPHC E_it20 usable
>   <int> <chr>   <dbl> <dbl>   <int>      <dbl>   <dbl> <chr>
> 1  2013  2001    48.3 -126.     76        5.96    1.19 Y
> 2  2013  2002    48.3 -126.     93        5.90    1.19 Y
> 3  2013  2003    48.5 -125.     79        5.90    1.19 Y
> 4  2013  2004    48.5 -126.     56        5.96    1.20 Y
> 5  2013  2005    48.5 -126.     58        6.02    1.20 Y
> 6  2013  2006    48.5 -126.    110        5.78    1.16 Y
> 7  2013  2007    48.7 -125.     35        5.96    1.20 Y
> 8  2013  2008    48.7 -125.     35        5.90    1.20 Y
> 9  2013  2009    48.7 -126.     67        5.90    1.19 Y
> 10 2013  2010    48.7 -126.     41        5.96    1.20 Y
```

```

> # ... with 160 more rows
countData2013
> # A tibble: 1,304 x 4
>   year station spNameIPHC      specCount
>   <int> <chr>   <chr>          <int>
> 1  2013 2001    Spiny Dogfish      61
> 2  2013 2001    Empty Hook        57
> 3  2013 2001    Pacific Halibut     2
> 4  2013 2002    Spiny Dogfish      59
> 5  2013 2002    Empty Hook        56
> 6  2013 2002    Pacific Halibut     5
> 7  2013 2003    Sablefish (Blackcod) 1
> 8  2013 2003    Longnose Skate      4
> 9  2013 2003    Arrowtooth Flounder  7
> 10 2013 2003    Spiny Dogfish      13
> # ... with 1,294 more rows

```

We want to get the new data into the same format as those (columns with same names and classes, even though in retrospect some classes aren't ideally chosen, but also retaining retrieved and observed hooks for the set data). Two data sets are needed because later gfiphc code summarises catches of a particular species at the station level, and needs to create counts of zeros for the species of interest (and such zeros are not included in IPHC output).

Set-level information

For 2020, Maria was sent the file 2020 IPHCtoDFO_dataExtraction-Maria.xls for set details, but this is multiple sheets and more complex than needed. So I tried extracting directly from the IPHC website (which they want us to do in the future anyway), using the following instructions, which worked for 2020 and 2021:

Go to <https://www.iphc.int/data/fiss-data-query> and select the following options:

1. Year Range – 2022 to 2022.
2. Area 2B
3. Purpose Codes – All
4. IPHC Charter Regions – All
5. Maps – Nothing
6. Select non-Pacific halibut species – deselect All (yes, deselect).

Download tab on bottom right (see instructions above question 4), and select CrossTab. Select “Set and Pacific Halibut data” and .xlsx format (I tried .csv format but it didn't save with commas, strangely). Save in this folder as **set-and-halibut-data-2022.xlsx**. Open

in Excel and Export as .csv, `set-and-halibut-data-2022.csv`, and when trying to quit Excel say no to save changes (not sure if that matters).

Repeat but with all non-halibut data (select All in number 6 and choose non-halibut in the CrossTab pop up)), and save as `non-halibut-data-2022.xlsx` and export as .csv in Excel, `non-halibut-data-2022.csv`. Importantly, this file (but not the first one) contains the numbers of observed hooks, needed in our calculations.

Load data for new year:

```
sets_raw <- readr::read_csv("set-and-halibut-data-2022.csv") %>%
  dplyr::mutate_if(is.character, factor)
> Rows: 174 Columns: 44
> -- Column specification -----
> Delimiter: ","
> chr (7): Vessel code, IPHC Reg Area, IPHC Charter Region, Purpose Code, Dat...
> dbl (32): Row number, Year, Stlkey, Station, Setno, IPHC Stat Area, BeginLat...
> lgl (3): Sigma-t, Oxygen_umol, Oxygen_sat
>
> i Use `spec()` to retrieve the full column specification for this data.
> i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

Now load the original 2020 data (do not change the 2020 here) to then test that the column names and types do not change in future years, and then check columns match `sets_raw`:

```
sets_raw_2020 <- readr::read_csv("set-and-halibut-data-2020.csv") %>%
  dplyr::mutate_if(is.character, factor)
> Rows: 198 Columns: 33
> -- Column specification -----
> Delimiter: ","
> chr (6): Vessel code, IPHC Reg Area, IPHC Charter Region, Purpose, Date, Eff
> dbl (24): Row number, Year, Stlkey, Station, Setno, IPHC Stat Area, BeginLat...
> lgl (1): Ineffcde
>
> i Use `spec()` to retrieve the full column specification for this data.
> i Specify the column types or set `show_col_types = FALSE` to quiet this message.

# For 2021 these were different - uncomment for future for first test
# testthat::expect_equal(names(sets_raw_2020),
#                           names(sets_raw))
# setdiff commands give the columns

# testthat::expect_equal(sapply(sets_raw_2020, typeof),
#                           sapply(sets_raw, typeof))

# Columns in 2020 not in new data:
setdiff(names(sets_raw_2020), names(sets_raw))
```

```

> [1] "Purpose"

# Columns in new data not in 2020:
setdiff(names(sets_raw), names(sets_raw_2020))
> [1] "Purpose Code"      "Profiler Lat"
> [3] "Profiler Lon"      "Profiler Bottom Depth (m)"
> [5] "Temp C"            "Max Pressure (db)"
> [7] "pH"                "Salinity PSU"
> [9] "Sigma-t"           "Oxygen_ml"
> [11] "Oxygen_umol"       "Oxygen_sat"

# For 2021 and 2022 looks like Purpose became Purpose Code, but are the same type:
summary(sets_raw_2020$Purpose)
>      Deep expansion Shallow expansion      Standard grid
>           3           30           165
summary(sets_raw$"Purpose Code")
> Standard Grid
>           174
testthat::expect_equal(typeof(sets_raw_2020$Purpose),
                        typeof(sets_raw$"Purpose Code"))

```

Those extra columns in 2021 and 2022 look related to oceanographic data, beyond the scope of gflphc, so can just ignore shortly.

Want to check the overlapping columns have the same type:

```

overlap_col_names <- intersect(names(sets_raw_2020),
                               names(sets_raw))

# testthat::expect_equal(sapply(dplyr::select(sets_raw_2020,
#                                           overlap_col_names),
#                               typeof),
#                         sapply(dplyr::select(sets_raw,
#                                           overlap_col_names),
#                               typeof))
# Error: sapply(dplyr::select(sets_raw_2020, overlap_col_names), typeof) not equal to
# 1/32 mismatches
# x[12]: "logical"
# y[12]: "integer"
# Get above error in 2021 and 2022, so for 2022 use the same wrangling as 2021
dplyr::select(sets_raw_2020,
              overlap_col_names[12])
> # A tibble: 198 x 1
>   Ineffcde
>   <lgl>
> 1 NA

```

```

> 2 NA
> 3 NA
> 4 NA
> 5 NA
> 6 NA
> 7 NA
> 8 NA
> 9 NA
> 10 NA
> # ... with 188 more rows
dplyr::select(sets_raw,
              overlap_col_names[12])
> # A tibble: 174 x 1
>   Ineffcde
>   <fct>
> 1 <NA>
> 2 <NA>
> 3 <NA>
> 4 <NA>
> 5 <NA>
> 6 <NA>
> 7 <NA>
> 8 <NA>
> 9 <NA>
> 10 <NA>
> # ... with 164 more rows

```

These are all NA's anyway (see below) and don't get saved, so no worries.

```

sets_raw
> # A tibble: 174 x 44
>   Row numb~1 Year Stlkey Vesse~2 Station Setno IPHC ~3 IPHC ~4 IPHC ~5 Purpo~6
>   <dbl> <dbl> <dbl> <fct> <dbl> <dbl> <fct> <dbl> <fct> <fct>
> 1      1 2022 2.02e7 BDP      2324      1 2B      133 Charlo~ Standa~
> 2      2 2022 2.02e7 BDP      2150      2 2B      121 Charlo~ Standa~
> 3      3 2022 2.02e7 BDP      2318      3 2B      133 Charlo~ Standa~
> 4      4 2022 2.02e7 BDP      2147      4 2B      133 Charlo~ Standa~
> 5      5 2022 2.02e7 BDP      2148      5 2B      121 Charlo~ Standa~
> 6      6 2022 2.02e7 BDP      2145      6 2B      121 Charlo~ Standa~
> 7      7 2022 2.02e7 BDP      2142      7 2B      121 Charlo~ Standa~
> 8      8 2022 2.02e7 BDP      2139      8 2B      121 Charlo~ Standa~
> 9      9 2022 2.02e7 BDP      2297      9 2B      121 Charlo~ Standa~
> 10     10 2022 2.02e7 BDP      2295     10 2B      121 Charlo~ Standa~
> # ... with 164 more rows, 34 more variables: Date <fct>, Eff <fct>,
> #   Ineffcde <fct>, BeginLat <dbl>, BeginLon <dbl>, `BeginDepth (fm)` <dbl>,

```

```

> #   EndLat <dbl>, EndLon <dbl>, `EndDepth (fm)` <dbl>, `MidLat fished` <dbl>,
> #   `MidLon fished` <dbl>, `AvgDepth (fm)` <dbl>, `Lat - Grid target` <dbl>,
> #   `Lon - Grid target` <dbl>, `O32 Pacific halibut count` <dbl>,
> #   `U32 Pacific halibut count` <dbl>, `O32 Pacific halibut weight` <dbl>,
> #   `U32 Pacific halibut weight` <dbl>, `No. skates set` <dbl>, ...
summary(sets_raw)
>   Row number      Year      Stlkey      Vessel code      Station
>   Min.      : 1.00   Min.      :2022   Min.      :20220009   BDP:138   Min.      :2001
>   1st Qu.: 44.25   1st Qu.:2022   1st Qu.:20220263   PEN: 36   1st Qu.:2084
>   Median : 87.50   Median :2022   Median :20220356                      Median :2150
>   Mean    : 87.50   Mean    :2022   Mean    :20220402                      Mean    :2193
>   3rd Qu.:130.75   3rd Qu.:2022   3rd Qu.:20220683                      3rd Qu.:2282
>   Max.     :174.00   Max.     :2022   Max.     :20220742                      Max.     :3210
>
>   Setno      IPHC Reg Area IPHC Stat Area      IPHC Charter Region
>   Min.      : 1.00   2B:174      Min.      : 60.0   Charlotte      :75
>   1st Qu.: 44.25                      1st Qu.: 92.0   Goose Is.      :32
>   Median : 83.50                      Median :112.0   St. James      :36
>   Mean    : 75.29                      Mean    :107.4   Vancouver Outside:31
>   3rd Qu.:105.00                      3rd Qu.:122.0
>   Max.     :138.00                      Max.     :142.0
>
>   Purpose Code      Date      Eff      Ineffcde      BeginLat
>   Standard Grid:174   13-Jul-22: 6      N: 3      DO : 3      Min.      :48.36
>                      22-Jul-22: 6      Y:171     NA's:171     1st Qu.:51.34
>                      23-Jul-22: 6                      Median :52.54
>                      24-Jul-22: 6                      Mean    :52.33
>                      12-Jul-22: 5                      3rd Qu.:53.68
>                      14-Jul-22: 5                      Max.     :55.35
>                      (Other) :140
>   BeginLon      BeginDepth (fm)      EndLat      EndLon
>   Min.      : -133.7   Min.      : 6.00   Min.      :48.32   Min.      : -133.7
>   1st Qu.: -131.1   1st Qu.: 39.00   1st Qu.:51.36   1st Qu.: -131.1
>   Median : -130.1   Median : 69.50   Median :52.52   Median : -130.1
>   Mean    : -129.8   Mean    : 82.06   Mean    :52.32   Mean    : -129.8
>   3rd Qu.: -128.7   3rd Qu.:108.75   3rd Qu.:53.66   3rd Qu.: -128.7
>   Max.     : -125.1   Max.     :375.00   Max.     :55.32   Max.     : -125.1
>
>   EndDepth (fm)      MidLat fished      MidLon fished      AvgDepth (fm)
>   Min.      : 9.00   Min.      :48.34   Min.      : -133.7   Min.      : 7.00
>   1st Qu.: 43.00   1st Qu.:51.35   1st Qu.: -131.1   1st Qu.: 45.50
>   Median : 71.00   Median :52.52   Median : -130.1   Median : 66.50
>   Mean    : 83.64   Mean    :52.33   Mean    : -129.8   Mean    : 81.34
>   3rd Qu.:114.00   3rd Qu.:53.67   3rd Qu.: -128.7   3rd Qu.:113.50

```

```

> Max.      :297.00    Max.      :55.33    Max.      :-125.1    Max.      :257.00
>
> Lat - Grid target Lon - Grid target 032 Pacific halibut count
> Min.      :48.33     Min.      :-133.7    Min.      : 0.00
> 1st Qu.:51.33       1st Qu.: -131.1    1st Qu.: 10.00
> Median :52.52       Median : -130.1    Median : 22.50
> Mean     :52.33      Mean     :-129.8    Mean     : 26.07
> 3rd Qu.:53.67       3rd Qu.: -128.7    3rd Qu.: 36.75
> Max.      :55.33     Max.      :-125.1    Max.      :130.00
>
> U32 Pacific halibut count 032 Pacific halibut weight
> Min.      : 0.00           Min.      : 0.0
> 1st Qu.: 2.25             1st Qu.: 220.8
> Median : 13.00           Median : 500.0
> Mean     : 27.62          Mean     : 605.1
> 3rd Qu.: 39.00           3rd Qu.: 811.5
> Max.      :183.00         Max.      :4027.0
>
> U32 Pacific halibut weight No. skates set No. skates hauled Avg no. hook/skate
> Min.      : 0.00           Min.      :8        Min.      :2.000    Min.      : 98.00
> 1st Qu.: 19.25           1st Qu.:8        1st Qu.:8.000    1st Qu.: 99.00
> Median : 105.50          Median :8        Median :8.000    Median : 99.00
> Mean     : 208.84         Mean     :8        Mean     :7.948    Mean     : 99.32
> 3rd Qu.: 297.75          3rd Qu.:8        3rd Qu.:8.000    3rd Qu.:100.00
> Max.      :1412.00        Max.      :8        Max.      :8.000    Max.      :101.00
>
> Effective skates hauled Soak time (min.) Profiler Lat Profiler Lon
> Min.      :2.110          Min.      :381.0    Min.      :48.36    Min.      :-131.7
> 1st Qu.:7.950            1st Qu.:466.2    1st Qu.:50.23    1st Qu.: -130.0
> Median :7.950            Median :548.0    Median :51.67    Median : -129.2
> Mean     :7.925          Mean     :565.5    Mean     :51.46    Mean     : -128.8
> 3rd Qu.:8.030            3rd Qu.:642.8    3rd Qu.:52.67    3rd Qu.: -127.7
> Max.      :8.110          Max.      :944.0    Max.      :53.99    Max.      : -125.2
>
>                                     NA's :96        NA's :96
> Profiler Bottom Depth (m) Temp C Max Pressure (db) pH
> Min.      : 16.00          Min.      : 5.203    Min.      : 7.00    Min.      :7.487
> 1st Qu.: 79.25            1st Qu.: 6.389    1st Qu.: 68.25    1st Qu.:7.666
> Median :124.00            Median : 7.043    Median :106.50    Median :7.716
> Mean     :137.00          Mean     : 7.261    Mean     :124.63    Mean     :7.731
> 3rd Qu.:171.50            3rd Qu.: 8.151    3rd Qu.:161.50    3rd Qu.:7.793
> Max.      :444.00          Max.      :10.588    Max.      :430.00    Max.      :7.971
> NA's :96                  NA's :96        NA's :96        NA's :152
> Salinity PSU Sigma-t Oxygen_ml Oxygen_umol Oxygen_sat
> Min.      :31.30 Mode:logical Min.      :1.184 Mode:logical Mode:logical

```

```

> 1st Qu.:32.37   NA's:174      1st Qu.:1.740   NA's:174      NA's:174
> Median :33.30           Median :2.446
> Mean   :33.05           Mean   :2.998
> 3rd Qu.:33.81           3rd Qu.:3.546
> Max.   :34.00           Max.   :6.939
> NA's    :96             NA's    :96
testthat::expect_equal(unique(sets_raw$"IPHC Reg Area"),
                        as.factor("2B")) # Check just BC
testthat::expect_equal(unique(sets_raw$Year), 2022)
testthat::expect_equal(length(unique(sets_raw$Station)),
                        length(sets_raw$Station))

```

Understand any issues raised above

Uncomment those three `testthat` commands when looking at new data each year. If any of fail then have to comment it out and figure out what it means here.

For 2022 got same results as 2021.

This is for 2020 (check for future years), to look for station(s) that was fished twice. Not really needed for 2021 since that third test passed (in 2022 not quite sure what that's referring to), but `twice_fished` gets used later, so do evaluate here:

```

length(unique(sets_raw$Station))
> [1] 174
length(sets_raw$Station)
> [1] 174
dplyr::count(sets_raw, Station) %>% dplyr::filter(n > 1)
> # A tibble: 0 x 2
> # ... with 2 variables: Station <dbl>, n <int>
twice_fished <- dplyr::count(sets_raw, Station) %>%
  dplyr::filter(n > 1) %>%
  dplyr::select(Station) %>%
  as.numeric()
twice_fished
> [1] NA
# If there's more than a single station then adapt later code
# as.data.frame(dplyr::filter(sets_raw,
#                               Station == twice_fished))

```

Not needed for 2021 or 2022: So Station NA had two vessels fishing the same station (which the code below originally caused a total of four rows for that station, explaining the 200 rows I had in original `setData2020` before fixing the issue). Interestingly the halibut catches were almost double for one vessel than the other (but were 6 days apart):

2020: Note that one of those entries has 'Vessel code' HAN, but HAN only appears once in

the whole data set (as seen in `summary(sets_raw)` above.

For 2021 and 2022, just noting that two vessels were used, and these are different to those in 2020 (for which HAN then got excluded anyway); 2020 used BDP, HAN, VNI and 2022 used BDP and PEN:

```
summary(sets_raw$"Vessel code")
> BDP PEN
> 138 36
summary(sets_raw_2020$"Vessel code")
> BDP HAN VNI
> 139 1 58
```

2020: So given we want to exclude one of the duplicates, makes sense to exclude HAN. (Also, Dana mentioned some gear comparison studies for 2020).

Simplify down to what's needed and rename, based on `iphc2013data.Rnw` (need to include the 'purpose' column, unlike 2013):

```
# sets_simp <- dplyr::filter(sets_raw, `Vessel code` != "HAN") %>%
sets_simp <- dplyr::select(sets_raw,
                           year = Year,
                           station = Station,
                           lat = "MidLat fished",
                           lon = "MidLon fished",
                           avgDepth = "AvgDepth (fm)",
                           skatesHauled = "No. skates hauled",
                           effSkateIPHC = "Effective skates hauled",
                           soakTimeMinutes = "Soak time (min.)", # Joe might want
                           usable = Eff,
                           purpose = "Purpose Code",
                           U32halibut = "U32 Pacific halibut count",
                           O32halibut = "O32 Pacific halibut count") %>%

  arrange(station) %>%
  dplyr::mutate(year = as.integer(year),
                station = as.character(station),
                avgDepth = as.integer(avgDepth),
                usable = as.character(usable))

sets_simp
> # A tibble: 174 x 12
>   year station  lat  lon avgDepth skatesHau~1 effSk~2 soakT~3 usable purpose
>   <int> <chr>   <dbl> <dbl>   <int>      <dbl>   <dbl>   <dbl> <chr>   <fct>
> 1  2022 2001    48.3 -126.     77         8    8.03    614 Y    Standa~
> 2  2022 2003    48.5 -125.     74         8    8.03    426 Y    Standa~
> 3  2022 2004    48.5 -126.     55         8    8.03    511 Y    Standa~
> 4  2022 2005    48.5 -126.     58         8    7.95    433 Y    Standa~
> 5  2022 2006    48.5 -126.    107         8    8.03    738 Y    Standa~
```

```

> 6 2022 2007 48.7 -125. 35 8 8.03 630 Y Standa~
> 7 2022 2008 48.7 -125. 34 8 8.03 512 Y Standa~
> 8 2022 2009 48.7 -126. 55 8 7.95 612 Y Standa~
> 9 2022 2013 48.8 -126. 33 8 8.03 488 Y Standa~
> 10 2022 2015 48.8 -126. 96 8 8.03 476 Y Standa~
> # ... with 164 more rows, 2 more variables: U32halibut <dbl>, O32halibut <dbl>,
> # and abbreviated variable names 1: skatesHauled, 2: effSkateIPHC,
> # 3: soakTimeMinutes

```

Standard grid or not

Need to change purpose to **standard** (Y/N) to match 2018 data (Y for the standard grid). In the raw 2020 data, **Purpose** took three values that we converted to **standard** to save in the package:

```

summary(sets_raw_2020$Purpose)
> Deep expansion Shallow expansion Standard grid
> 3 30 165
summary(setData2020$standard)
> N Y
> 71 126

```

For 2021 and 2022 we have all as Standard Grid, which gets corrected (some stations are non-standard) in the next section.

```

summary(sets_simp$purpose)
> Standard Grid
> 174

sets_simp_std <- dplyr::mutate(sets_simp,
                              standard_tmp = (purpose == "Standard Grid"))
                              # was grid in 2020, Grid in 2021, "Standard Grid"
                              # in 2022

standard <- as.character(sets_simp_std$standard_tmp) # to get the right length
standard[sets_simp_std$standard_tmp] = "Y"
standard[!sets_simp_std$standard_tmp] = "N"
length(standard)
> [1] 174

sets_simp_std <- cbind(sets_simp_std,
                      standard) %>%
  as_tibble() %>%
  dplyr::select(-c("standard_tmp"))
summary(sets_simp_std)
> year station lat lon

```

```

> Min.      :2022   Length:174      Min.      :48.34   Min.      :-133.7
> 1st Qu.:2022   Class :character  1st Qu.:51.35   1st Qu.: -131.1
> Median :2022   Mode  :character  Median :52.52   Median : -130.1
> Mean    :2022      Mean    :52.33   Mean    :-129.8
> 3rd Qu.:2022      3rd Qu.:53.67   3rd Qu.: -128.7
> Max.    :2022      Max.    :55.33   Max.    :-125.1
>   avgDepth      skatesHauled      effSkateIPHC      soakTimeMinutes
> Min.      : 7.00   Min.      :2.000   Min.      :2.110   Min.      :381.0
> 1st Qu.: 45.50   1st Qu.:8.000   1st Qu.:7.950   1st Qu.:466.2
> Median : 66.50   Median :8.000   Median :7.950   Median :548.0
> Mean    : 81.34   Mean    :7.948   Mean    :7.925   Mean    :565.5
> 3rd Qu.:113.50   3rd Qu.:8.000   3rd Qu.:8.030   3rd Qu.:642.8
> Max.    :257.00   Max.    :8.000   Max.    :8.110   Max.    :944.0
>   usable      purpose      U32halibut      032halibut
> Length:174      Standard Grid:174   Min.      : 0.00   Min.      : 0.00
> Class :character      1st Qu.: 2.25   1st Qu.: 10.00
> Mode  :character      Median : 13.00   Median : 22.50
>      Mean    : 27.62   Mean    : 26.07
>      3rd Qu.: 39.00   3rd Qu.: 36.75
>      Max.    :183.00   Max.    :130.00
>   standard
> Length:174
> Class :character
> Mode  :character
>
>
>
unique(sets_simp_std$standard)
> [1] "Y"

```

So they are all classified as standard. For 2020 we stuck with the 2018 definitions of standard, so doing that next.

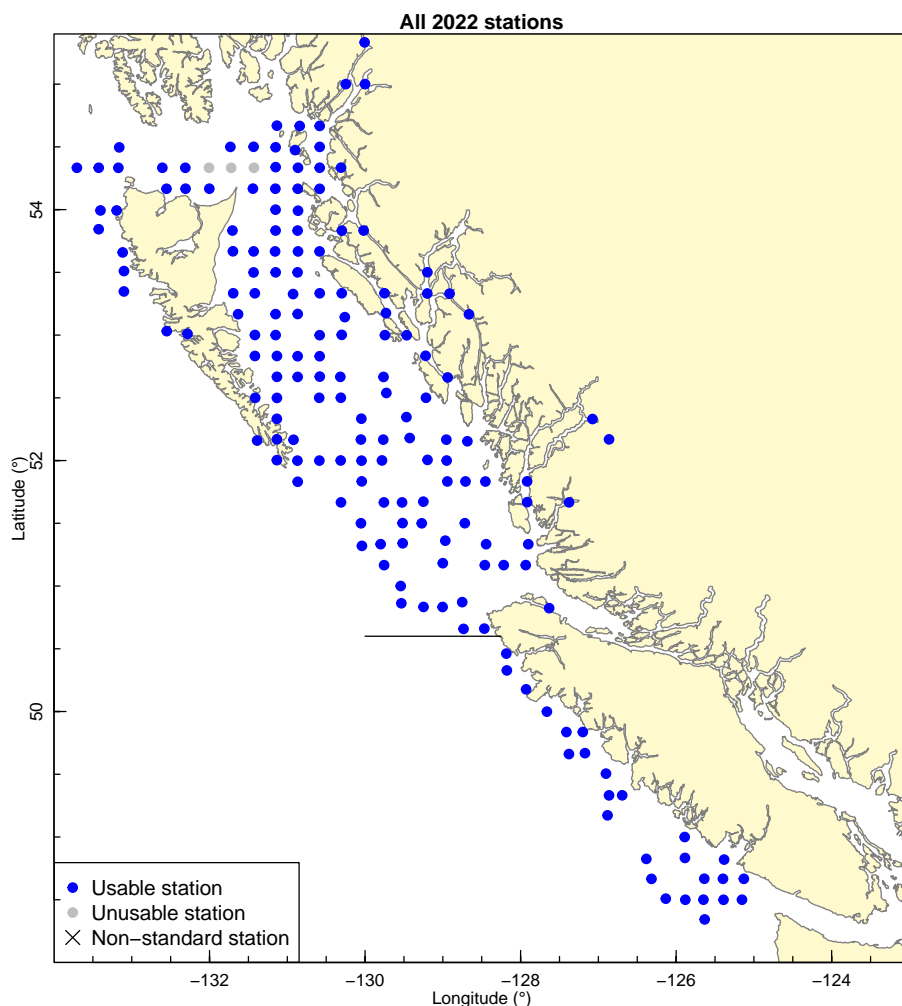
Look at data and show map to understand changing definition of standard station from 2018 to 2020.

The definition of 'standard grid' changed from 2018 (when first needed due to the expanded grid) to 2020 (and 2021). Simply equating them as above is not sufficient. For 2022 we so far have this:

```

plot_iphc_map(sets_simp_std,
              sp = NULL,
              years = 2022,
              indicate_standard = TRUE)

```



For 2021: So no stations are marked as being outside the standard grid, even though some are clearly new – the ones in the north have never been fished before (see the one-species vignette, though I’ll investigate that here). 2022 the northern inlets ones are there again, as are some other inlets, only 1 in Strait of Georgia.

This next section was to first figure out the twice-fished station 2343 in 2020, and to replicate that original analysis (station ends up being non-standard later), so mostly commented out except first bit which is used later so keeping in case need in future years:

```
hooks_with_bait_revert <- hooks_with_bait

# This should be commented out for 2021 survey analysis in iphc-2021-data.Rmd,
# since the problem is presumably fixed. This is to revert back to the original
# problem, for which 2343 was called standard in 2018 but we changed it. Map on
# page 10 of iphc-2020-data.pdf has this station (second one down off
# north-east tip of Haida Gwaii) as non-standard in 2018 but not 2020.
# hooks_with_bait_revert$set_counts[hooks_with_bait_revert$set_counts$year == 2018 &
#                                   hooks_with_bait_revert$set_counts$station == 2343,
#                                   ]$standard = "Y"
```

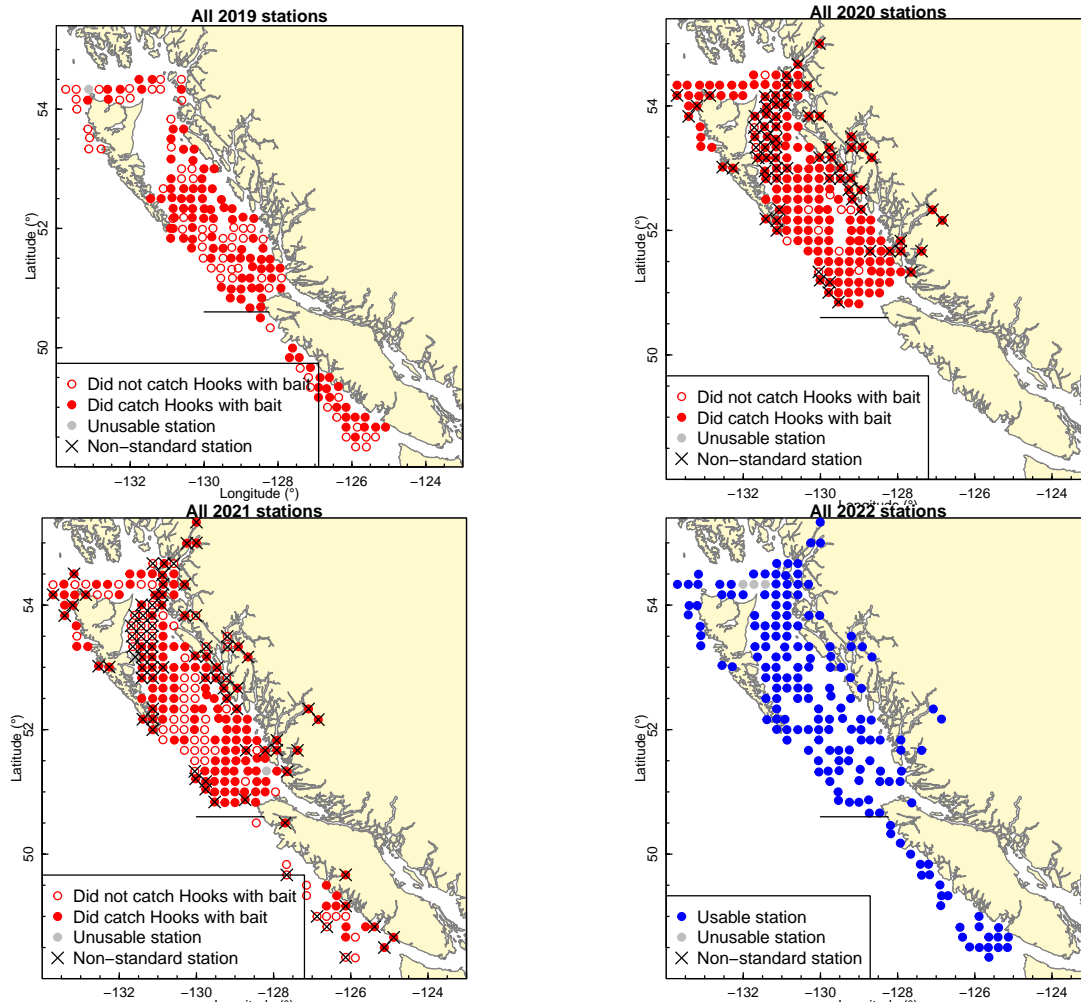
```
#filter(hooks_with_bait$set_counts, year == 2018, station == 2343) %>%
# as.data.frame()          # saved version
#filter(hooks_with_bait_revert$set_counts, year == 2018, station == 2343) %>%
# as.data.frame()          # reverted version
```

Now to figure out standard/non-standard stations. Plotting four years, with crosses showing ‘non-standard’. (2022 is coloured different since no hooks with bait data yet, but the important bit is the crosses).

```
sets_2022 <- dplyr::select(sets_simp_std,
                           -c(U32halibut, O32halibut))
                           # else not the same structure as sets_2018, below

plot_iphc_map_panel(hooks_with_bait$set_counts,
                    sp = "Hooks with bait",
                    years = 2019:2021,    # the latest three years except
                                # current data
                    indicate_standard = TRUE)

plot_iphc_map(sets_2022,
              sp = NULL,
              years = 2022,
              indicate_standard = TRUE)
```



Can see that 2020 has a few less stations just north of Vancouver Island, but not enough to worry about greatly, and 2021 and 2022 have kind of done a few of those. The 2021 and 2022 ones way in in the inlets are not currently flagged as non-standard but will be below (using the 2018 definitions). In fact no stations are flagged for 2021 or 2022 as non-standard. And the other main issue is that 2021 and 2022 is doing a random sample of WCVI stations (some of which will become non-standard). AND for 2021 there are new stations in the north (and maybe elsewhere) that have never been fished before (as I discovered when updating the one-species vignette and redefining the default axes limits for `plot_BC()`; the version before updating that isaved as `iphc-2021-data-all-2021-stations.pdf`). Will examine those shortly.

Need to look and plot values:

```
sets_2018 <- filter(hooks_with_bait_revert$set_counts,
                    year == 2018)
not_std_2018 <- filter(sets_2018,
                      standard == "N")$station

not_std_2022 <- filter(sets_2022,
                      standard == "N")$station
```

```

# Not standard in both:
not_std_2018_and_2022 <- intersect(not_std_2018, not_std_2022)
not_std_2018_and_2022  # Empty for 2022
> character(0)

length(not_std_2018)
> [1] 131
length(not_std_2022)
> [1] 0
length(not_std_2018_and_2022)
> [1] 0

# 2018 has some east of the map, all non-standard:
filter(hooks_with_bait_revert$set_counts, year == 2018, lon > -124)$standard
> [1] N N N N N N N N N N N N N N N
> Levels: Y N
nrow(filter(hooks_with_bait_revert$set_counts, year == 2018, lon > -124))
> [1] 14

std_in_2018_but_not_std_in_2022 <- intersect(filter(sets_2018,
                                                    standard == "Y")$station,
                                                    not_std_2022)

std_in_2018_but_not_std_in_2022
> character(0)

not_std_in_2018_but_std_in_2022 <- intersect(not_std_2018,
                                              filter(sets_2022,
                                                    standard == "Y")$station)

not_std_in_2018_but_std_in_2022
> [1] "2258" "2263" "2265" "2266" "2272" "2275" "2270" "2267" "2290" "2293"
> [11] "2321" "2323" "2331" "2320" "2312" "2314" "2309" "2304" "2302" "2295"
> [21] "2297" "2317" "2315" "2334" "2335" "2333" "2332" "2329" "2328" "2327"
> [31] "2324" "2322" "2318" "2287" "2285" "2288" "2311" "2313" "2289" "2242"
> [41] "2237" "2208" "2213" "2210" "2217" "2278" "2273" "2271" "2274" "2277"
> [51] "2279" "2283" "2307" "2303" "2301" "2294" "2306" "2298" "2291"

# setdiff(x, y) - elements in x but not in y
# setdiff(not_std_2018, not_std_2020) - but 2020 fewer coverage so misleading

```

Plot stations not standard in 2018 but standard in 2022, and vice versa, using each years' lats and lons (to verify that they all still agree – i.e., that station numbers have consistent lats and lons), and show 2019 data to check no 'usual' stations are non-standard in 2018 or 2022. Also (for 2021 and then 2022) adding all stations, since this will clearly show the random sampling off WCVI:

```

plot_BC()
points(lat~lon,
       data = filter(sets_2018,
                     station %in% not_std_in_2018_but_std_in_2022),
       col="red",
       pch = 19)

# Do the same but using 2022 station co-ordinates - should overlap:
points(lat~lon,
       data = filter(sets_2022,
                     station %in% not_std_in_2018_but_std_in_2022),
       col="blue",
       pch = 3)

# And for 2020 showed the single station std in 2018 but not 2020, for 2021 and 2022
# there are none:
points(lat~lon,
       data = filter(sets_2018,
                     station %in% std_in_2018_but_not_std_in_2022),
       col="red",
       pch = 17)
points(lat~lon,
       data = filter(sets_2022,
                     station %in% std_in_2018_but_not_std_in_2022),
       col="blue",
       pch = 1,
       cex = 2)

# Now show all 2019 stations:
points(lat~lon,
       data = filter(hooks_with_bait_revert$set_counts,
                     year == 2019),
       col="darkgreen",
       pch = 0)

# Add all 2022 stations as a small black dot
points(lat~lon,
       data = sets_2022,
       col="black",
       pch = 20,
       cex = 0.8)

legend("bottomleft",

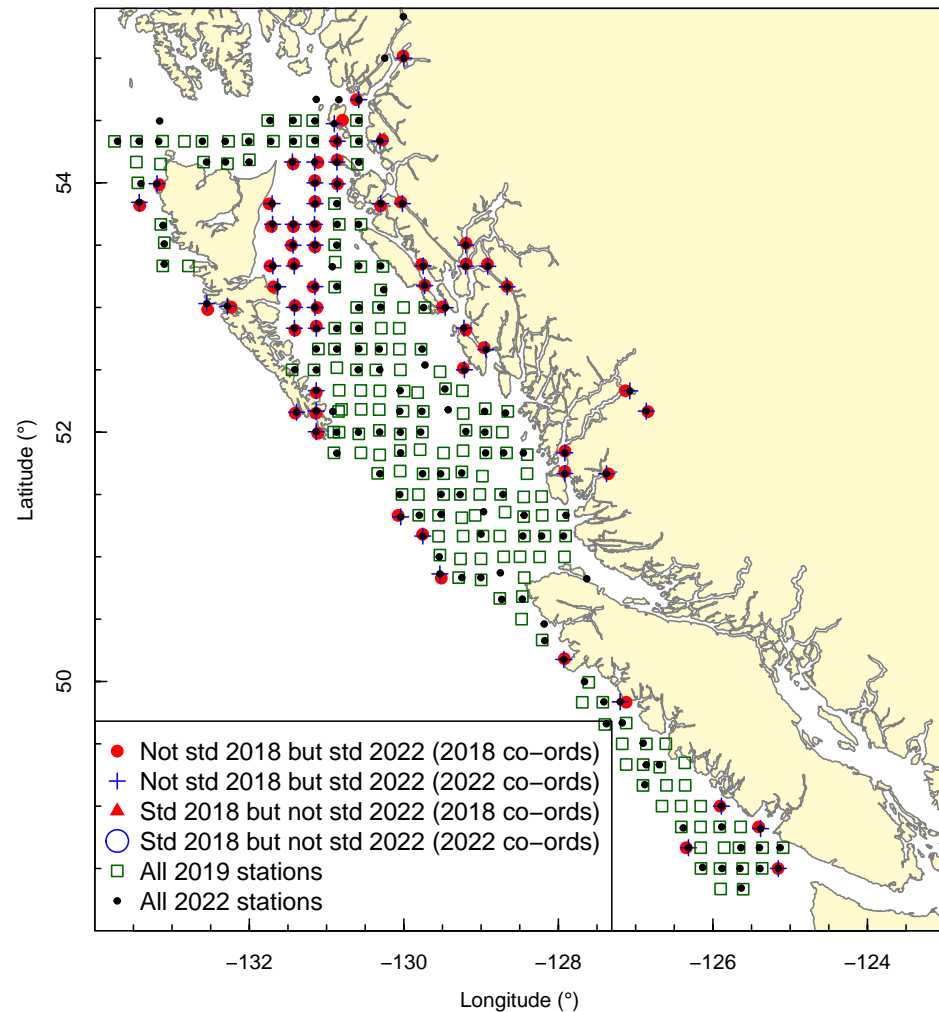
```



```

legend = c("Not std 2018 but std 2022 (2018 co-ords)",
           "Not std 2018 but std 2022 (2022 co-ords)",
           "Std 2018 but not std 2022 (2018 co-ords)",
           "Std 2018 but not std 2022 (2022 co-ords)",
           "All 2019 stations",
           "All 2022 stations"),
pch = c(19, 3, 17, 1, 0, 20),
pt.cex = c(1, 1, 1, 2, 1, 0.8),
col = c("red", "blue", "red", "blue", "darkgreen", "black"))

```



So the co-ordinates look close enough (red circles and blue crosses overlap), none were defined as non-standard in 2022 (or 2021 before) so there are no red triangles or blue circles, and the

green squares for 2019 stations correctly do not overlap with the non-standard 2018 stations. Empty green squares with no black dots (2022 stations) off WCVI clearly shows the reduced coverage there (similar for 2021 and 2022).

2020 only (there were no non-standard stations defined in raw data for 2021 or 2022): Check if the one standard station in 2018 but not in 2020 (not fished at all in 2019) appears in any earlier years:

```
# Fails (so not evaluated here) since empty in 2021 and 2022, and this is corrected
dplyr::filter(hooks_with_bait_revert$set_counts,
              station == std_in_2018_but_not_std_in_2022) %>%
  as.data.frame()
```

For 2020 I worked out it was only fished in 2018 and 2020 so we defined it as non-standard.

So, the conclusion from this section so far is that we should retain the 2018 definitions of standard stations, not the new ones defined in 2021, as we did for 2020. For 2022 this is same as we did for 2021.

Doing that shortly (in `sets_simp_std_corrected`), but first also look for any new 2022 stations. I hadn't expected any in 2021 but saw them when doing the one-species vignette, so had to come back to redo this.

```
# For 2021: yelloweye_rockfish$set_counts is saved in gfiphc, already has 2021 data
# because I had to come back to redo this .pdf after updating the data, hence
# need the <2021 here; station codes do change over time, but I think are
# recently consistent
# For 2022 checking before saving any data into gfiphc. Then rerunning (step 8
# in README) so need the < 2022 here as package contains 2022 data
previous_stations <- dplyr::filter(yelloweye_rockfish$set_counts,
                                  year < 2022)$station %>%
  unique()

stations_in_2022_only <- dplyr::filter(sets_2022,
                                       !(station %in% previous_stations))

stations_in_2022_only
> # A tibble: 2 x 11
>   year station  lat  lon avgDepth skatesHauled effSk~1 soakT~2 usable purpose
>   <int> <chr>   <dbl> <dbl>   <int>         <dbl>   <dbl>   <dbl> <chr>   <fct>
> 1  2022  2248    50.5 -128.     39             8     7.95    480 Y     Standa~
> 2  2022  2256    50.8 -128.    136             8     7.95    428 Y     Standa~
> # ... with 1 more variable: standard <chr>, and abbreviated variable names
> # 1: effSkateIPHC, 2: soakTimeMinutes
```

and plot those stations:

```
plot_BC()
points(lat~lon,
```

```

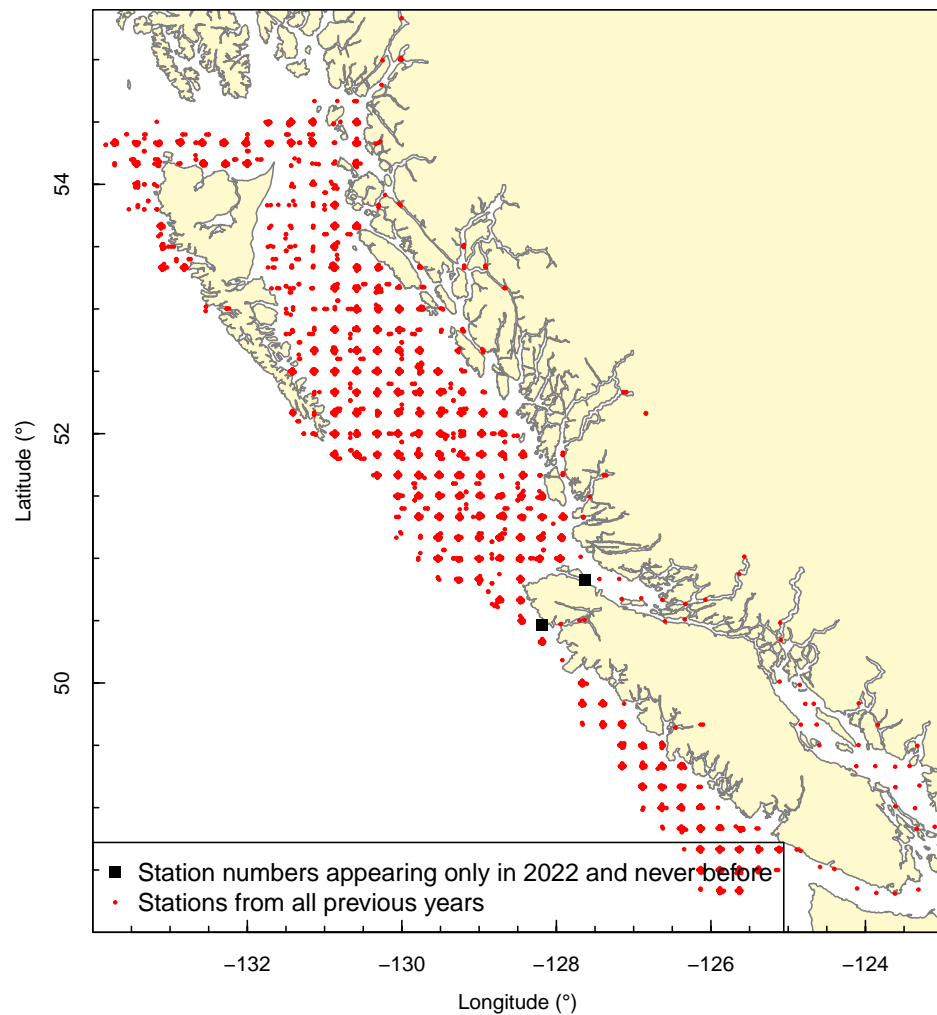
    data = stations_in_2022_only,
    col = "black",
    pch = 15)

points(lat~lon,
    data = dplyr::filter(yelloweye_rockfish$set_counts,
                          year < 2022),

    col = "red",
    pch = 20,
    cex = 0.4)

legend("bottomleft",
    legend = c("Station numbers appearing only in 2022 and never before",
               "Stations from all previous years"),
    pch = c(15, 20),
    col = c("black", "red"),
    pt.cex = c(1, 0.4))

```



So there are two 2022 (six in 2021) stations that have never been fished before!

For 2021: That map suggests that we should call the five northern ones non-standard also, to exclude from the standard Series A-F analyses.

2021 (for reference): However, Ann-Marie Huang thinks that these stations may have been fished before but considered as part of Area 2C (Alaskan waters). Some waters around there are claimed by both Canada and the US; there's a clear map and explanation in Canada's Unresolved Maritime Boundaries (clickable), which is linked from this Wikipedia article on Dixon Entrance. So there may be earlier data, which are not in gfiphc because such stations would not have been considered Area 2B, which is the area for which the IPHC sent DFO data in the past (and which I used here for recent years to extract from their website). So there may be data available, and if needed it will have to be obtained. Here we will call those five northern newly-fished stations **non-standard**.

2021: For the sixth station off the northwest of Vancouver Island, zooming in and including the Scott Islands Rockfish Conservation Area (clickable) as a blue rectangle shows (see saved 2021 .pdf, and rectangle in next map).

2022: The two new ones are also off the northwest of Vancouver Island, so zoom in:

```
plot_BC(xlim = c(-130, -127),
        ylim = c(50, 52))

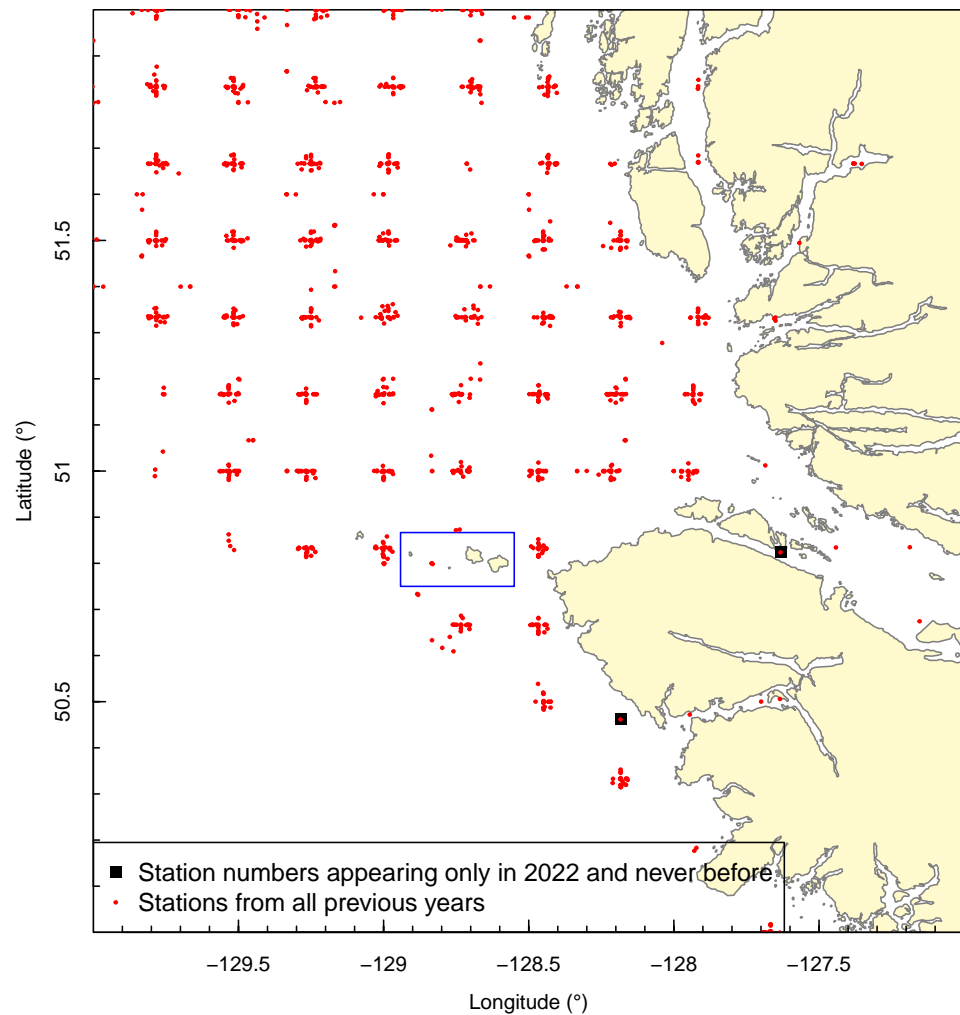
scott_island_RCA_lon <- -c(128 + 56.5/60, 128 + 33/60)
scott_island_RCA_lat <- c(50 + 45/60, 50 + 52/60)

# rect(xleft, ybottom, xright, ytop, density = NULL, angle = 45,
rect(scott_island_RCA_lon[1],
     scott_island_RCA_lat[1],
     scott_island_RCA_lon[2],
     scott_island_RCA_lat[2],
     border = "blue")

points(lat~lon,
       data = stations_in_2022_only,
       col = "black",
       pch = 15)

points(lat~lon,
       data = yelloweye_rockfish$set_counts,
       col = "red",
       pch = 20,
       cex = 0.4)

legend("bottomleft",
      legend = c("Station numbers appearing only in 2022 and never before",
                  "Stations from all previous years"),
      pch = c(15, 20),
      col = c("black", "red"),
      pt.cex = c(1, 0.4))
```



2021: So the new station is just outside the RCA. Presumably in previous years the RCA was avoided as the grid would have put a station in the RCA, close to (or even on) Lanz Island.

2021: It is station 2257 (see above), with a depth of only 40 fathoms, which is not an outlier. For example, for 2013 (depth data for all years is not in gfpnc I don't think):

2022: The two new stations are also not outliers in terms of depth, though are very close to shore:

```
sort(setData2013$avgDepth)
> [1] 18 21 22 24 25 25 26 27 29 32 32 32 35 35 35 36 36 37
> [19] 39 39 40 41 41 42 44 44 44 45 45 46 46 46 47 48 48 48
> [37] 48 48 50 50 51 51 52 52 54 54 54 55 56 56 56 58 58 58
> [55] 58 58 59 61 62 62 62 63 63 64 66 67 67 67 67 67 71 73
```

```

> [73] 74 74 74 75 75 75 76 76 76 77 78 78 78 79 79 81 81 81
> [91] 82 82 87 88 88 88 90 91 92 92 93 93 95 96 96 97 97 98
> [109] 98 98 99 101 101 102 102 102 102 103 103 104 105 105 110 111 112 112
> [127] 113 113 113 114 115 115 116 118 119 120 122 123 123 123 123 124 128 129
> [145] 130 132 135 136 137 139 139 139 140 142 142 144 145 145 150 156 161 183
> [163] 189 190 190 209 215 217 219 256

```

2021: However, since it is a new station and not been used before, we will flag it as **non-standard** (as used for the Series A-F analyses). Also Dana Haggarty says that there is good habitat right close to those islands, but not great further away, and she has used Remotely Operated Vehicles there – it’s all sand/gravel/cobble with massive sand waves from the crazy exposure, but perhaps there are pockets of good habitat. So either way (close to an RCA so may be expected to be good for rockfish at least, or not great rockfish habitat) it shouldn’t really be included for rockfish species, and in general should be excluded since a new station.

2022: Given so close to shore and not close to previous stations, will call these two new 2022 stations non-standard also.

So - retain the 2018 definitions of standard stations (as we did for 2020 and 2021), and call both new 2022 stations non-standard (like we did for the six new 2021 stations):

```

sets_simp_std_corrected <- sets_simp_std
summary(as.factor(sets_simp_std_corrected$standard))
> Y
> 174

sets_simp_std_corrected$standard[sets_simp_std_corrected$station %in%
                                not_std_in_2018_but_std_in_2022] <- "N"
sets_simp_std_corrected$standard[sets_simp_std_corrected$station %in%
                                stations_in_2022_only$station] <- "N"
summary(as.factor(sets_simp_std_corrected$standard))
> N Y
> 61 113
# cbind(sets_simp_std$standard, sets_simp_std_corrected$standard) # to check them

```

Think I hadn’t originally defined them as factors in early code, so keeping them as characters now. Just to verify that none of the 2018 non-standard stations were fished before 2018:

```

dplyr::filter(hooks_with_bait$set_counts,
              station %in% not_std_2018) %>%
  dplyr::select(year) %>%
  unique()
> # A tibble: 4 x 1
>   year
>   <dbl>
> 1 2018

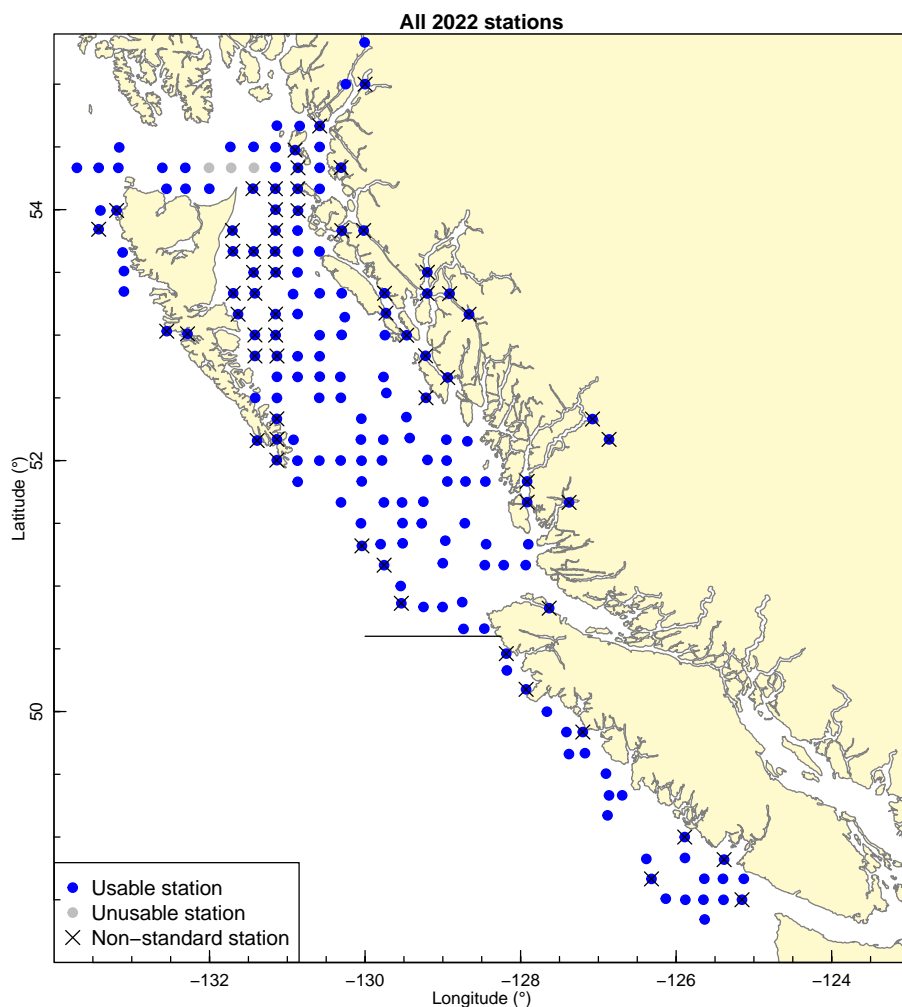
```

```
> 2 2020
> 3 2021
> 4 2022
```

Note 2022 won't show up here until .rda objects are resaved in package, at the end of this .pdf (so it will if this .Rmd has already been run, as it had in 2021).

So here are the final station designations for 2022:

```
plot_iphc_map(sets_simp_std_corrected,
  sp = NULL,
  years = 2022,
  indicate_standard = TRUE)
```



Can see that they did 6 random stations off WCVI that we're calling non-standard (because they were never fished before 2018; was 10 for 2021). Which is a bit of a shame as there are only 19 stations left off WCVI for 2022 (16 for 2022).

2020 (no need to change for 2021 or 2022): So check which functions need changing, since they create a 'standard' column. These do not need changing: `get_iphc_hooks()` and


```
get_iphc_skates_info.
```

2020: Then `get_iphc_sets_info()` does return `standard`, but the `standard` designation is not saved in `GFBio` it is saved in `setDataExpansion` in `gfiphc`. So just need to add a line in `IPHC-stations-expanded.R` and then re-save all `.rda` files. Fixed that, now recreating all `.rda` files, as per the README.

Species counts

Now get the species counts into the desired format (to match `countData2013` shown earlier). First check that the column names and types haven't changed (they did for set data from 2020 to 2021; no change here for 2022):

```
counts_raw_2020 <- readr::read_csv("non-halibut-data-2020.csv") %>%
  dplyr::mutate_if(is.character, factor)
> Rows: 1441 Columns: 13
> -- Column specification -----
> Delimiter: ","
> chr (3): Scientific Name, Species Name, SampleType
> dbl (9): Year, Stlkey, Station, Setno, IPHC Species Code, HooksFished, Hooks...
>
> i Use `spec()` to retrieve the full column specification for this data.
> i Specify the column types or set `show_col_types = FALSE` to quiet this message.

counts_raw <- readr::read_csv("non-halibut-data-2022.csv") %>%
  dplyr::mutate_if(is.character, factor)
> Rows: 1368 Columns: 13
> -- Column specification -----
> Delimiter: ","
> chr (3): Scientific Name, Species Name, SampleType
> dbl (9): Year, Stlkey, Station, Setno, IPHC Species Code, HooksFished, Hooks...
>
> i Use `spec()` to retrieve the full column specification for this data.
> i Specify the column types or set `show_col_types = FALSE` to quiet this message.

testthat::expect_equal(names(counts_raw_2020),
                        names(counts_raw))

testthat::expect_equal(sapply(counts_raw_2020, typeof),
                        sapply(counts_raw, typeof))
```

Great, nothing changed in the structure for 2022.

```
counts_raw
> # A tibble: 1,368 x 13
>   Row numb~1  Year Stlkey Station Setno IPHC ~2 Scien~3 Spec~4 Sampl~5 Hooks~6
```

```

>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <fct> <fct> <fct> <dbl>
> 1      1 2022 2.02e7 2324 1 141 Raja b~ Big Sk~ 20Hook 792
> 2      2 2022 2.02e7 2324 1 143 Raja r~ Longno~ 20Hook 792
> 3      3 2022 2.02e7 2324 1 303 <NA> Hook w~ 20Hook 792
> 4      4 2022 2.02e7 2324 1 304 <NA> Empty ~ 20Hook 792
> 5      5 2022 2.02e7 2324 1 305 <NA> Hook w~ 20Hook 792
> 6      6 2022 2.02e7 2150 2 2 Athere~ Arrowt~ 20Hook 792
> 7      7 2022 2.02e7 2150 2 26 Gadus ~ Pacifi~ 20Hook 792
> 8      8 2022 2.02e7 2150 2 54 Squalu~ Spiny ~ 20Hook 792
> 9      9 2022 2.02e7 2150 2 141 Raja b~ Big Sk~ 20Hook 792
> 10     10 2022 2.02e7 2150 2 143 Raja r~ Longno~ 20Hook 792
> # ... with 1,358 more rows, 3 more variables: HooksRetrieved <dbl>,
> # HooksObserved <dbl>, `Number Observed` <dbl>, and abbreviated variable
> # names 1: `Row number`, 2: `IPHC Species Code`, 3: `Scientific Name`,
> # 4: `Species Name`, 5: SampleType, 6: HooksFished
summary(counts_raw)
> Row number Year Stlkey Station
> Min. : 1.0 Min. :2022 Min. :20220009 Min. :2001
> 1st Qu.: 342.8 1st Qu.:2022 1st Qu.:20220266 1st Qu.:2086
> Median : 684.5 Median :2022 Median :20220358 Median :2146
> Mean : 684.5 Mean :2022 Mean :20220411 Mean :2186
> 3rd Qu.:1026.2 3rd Qu.:2022 3rd Qu.:20220684 3rd Qu.:2275
> Max. :1368.0 Max. :2022 Max. :20220742 Max. :3210
>
> Setno IPHC Species Code Scientific Name
> Min. : 1.00 Min. : 2.0 Squalus suckleyi :148
> 1st Qu.: 47.00 1st Qu.: 54.0 Raja rhina : 94
> Median : 84.00 Median :143.0 Anoplopoma fimbria : 88
> Mean : 76.29 Mean :173.3 Sebastes ruberrimus: 59
> 3rd Qu.:105.00 3rd Qu.:304.0 Ophiodon elongatus : 52
> Max. :138.00 Max. :307.0 (Other) :338
> NA's :589
> Species Name SampleType HooksFished HooksRetrieved
> Empty Hook :174 20Hook:1368 Min. :784.0 Min. :209.0
> Hook with Skin :169 1st Qu.:792.0 1st Qu.:792.0
> Hook with Bait :163 Median :792.0 Median :792.0
> Spiny Dogfish :148 Mean :794.5 Mean :790.1
> Longnose Skate : 94 3rd Qu.:800.0 3rd Qu.:800.0
> Sablefish (Blackcod): 88 Max. :808.0 Max. :808.0
> (Other) :532
> HooksObserved Number Observed
> Min. : 60.0 Min. : 1.00
> 1st Qu.:160.0 1st Qu.: 2.00
> Median :160.0 Median : 5.00

```

```

> Mean      :159.5    Mean      : 18.82
> 3rd Qu.:160.0    3rd Qu.: 23.00
> Max.      :160.0    Max.      :138.00
>
testthat::expect_equal(unique(counts_raw$Year), 2022) # All 2022
testthat::expect_equal(unique(counts_raw$SampleType), as.factor("20Hook")) # All 20Hook

# This mismatches for 2020, not for 2021 or 2022:
testthat::expect_equal(length(unique(counts_raw$Station)),
                        length(sets_raw$Station))

unique(counts_raw$"Species Name")
> [1] Big Skate                      Longnose Skate
> [3] Hook with Skin                  Empty Hook
> [5] Hook with Bait                  Arrowtooth Flounder
> [7] Pacific Cod                     Spiny Dogfish
> [9] Brittle Star                    Sea Anemone
> [11] Bent/Broken/Missing             Sea Pen
> [13] unident. Starfish               Basketstar
> [15] Quillback Rockfish              Copper Rockfish
> [17] Oregon Rock Crab                Soupfin Shark
> [19] Spotted Ratfish                 Sablefish (Blackcod)
> [21] Redbanded Rockfish              Lingcod
> [23] Yelloweye Rockfish              Shortspine Thornyhead
> [25] Aleutian Skate                  Sleeper Shark
> [27] Blackspotted Rockfish           Canary Rockfish
> [29] unident. Sculpin                Silvergray Rockfish
> [31] Fish-eating Star                Red Irish Lord
> [33] Rougheye Rockfish               Shortraker Rockfish
> [35] Bocaccio                       Yellowmouth Rockfish
> [37] Sandpaper Skate                 Rosethorn Rockfish
> [39] Petrale Sole                    Greenstriped Rockfish
> [41] Walleye Pollock                 unident. thornyhead (Idiot)
> [43] unident. Invertebrate           China Rockfish
> [45] Blue Shark                      Cabezon
> [47] unident. Coral                  Stylaster campylecus (coral)
> [49] Unknown/Unspecified            Wolf-Eel
> [51] Sunflower Sea Star              unident. Sponge
> [53] Grenadier (Rattails)
> 53 Levels: Aleutian Skate Arrowtooth Flounder ... Yellowmouth Rockfish

```

Here's what was seen in 2020 but not 2022, and vice versa:

```

# Seen in 2020 not 2022
setdiff(unique(counts_raw_2020$"Species Name"),

```

```

      unique(counts_raw$"Species Name"))
> [1] "Sand Dab"          "Glass Sponge"          "Sea Urchin"
> [4] "Octopus"           "Gastropod"             "Tiger Rockfish"
> [7] "Red Tree Coral"     "Giant Pacific Octopus"
# Seen in 2022 not 2020
setdiff(unique(counts_raw$"Species Name"),
        unique(counts_raw_2020$"Species Name"))
> [1] "Red Irish Lord"      "Sandpaper Skate"
> [3] "Rosethorn Rockfish"  "Greenstriped Rockfish"
> [5] "Walleye Pollock"     "unident. Invertebrate"
> [7] "China Rockfish"      "Cabezon"
> [9] "Stylaster campylecus (coral)" "Unknown/Unspecified"
> [11] "Grenadier (Rattails)"

```

2021: Presumably Sun Sea Star and Sunflower Sea Star are the same. Will mention this later on.

2022: Updated the csv file to include Rosethorn and Red Irish Lord, the rest are all dealt with.

Note that halibut are not included in these counts:

```

dplyr::filter(counts_raw, "Species Name" == "Pacific Halibut")
> # A tibble: 0 x 13
> # ... with 13 variables: Row number <dbl>, Year <dbl>, Stlkey <dbl>,
> #   Station <dbl>, Setno <dbl>, IPHC Species Code <dbl>, Scientific Name <fct>,
> #   Species Name <fct>, SampleType <fct>, HooksFished <dbl>,
> #   HooksRetrieved <dbl>, HooksObserved <dbl>, Number Observed <dbl>
# Should be: dplyr::filter(counts_raw, `Species Name` == as.character("Pacific
#                               Halibut")) %>% as.data.frame()
# Still 0 in 2021 and 2022

```

which I presume explains why total number of counts for a station does not add up to HooksObserved. See later for halibut calculations.

2020 only: Need to remove the HAN records for the twice-fished station, which turns out to be set number 4 for station 2104:

```

dplyr::filter(counts_raw, Station == twice_fished) %>%
  dplyr::select(c("Station", "Setno", "Species Name",
                  "Number Observed")) %>%
  as.data.frame()
> [1] Station      Setno          Species Name    Number Observed
> <0 rows> (or 0-length row.names)

dplyr::filter(sets_raw, Station == twice_fished)
> # A tibble: 0 x 44

```

```

> # ... with 44 variables: Row number <dbl>, Year <dbl>, Stlkey <dbl>,
> #   Vessel code <fct>, Station <dbl>, Setno <dbl>, IPHC Reg Area <fct>,
> #   IPHC Stat Area <dbl>, IPHC Charter Region <fct>, Purpose Code <fct>,
> #   Date <fct>, Eff <fct>, Ineffcde <fct>, BeginLat <dbl>, BeginLon <dbl>,
> #   BeginDepth (fm) <dbl>, EndLat <dbl>, EndLon <dbl>, EndDepth (fm) <dbl>,
> #   MidLat fished <dbl>, MidLon fished <dbl>, AvgDepth (fm) <dbl>,
> #   Lat - Grid target <dbl>, Lon - Grid target <dbl>, ...

```

So for 2020 had to use that here to remove the species counts for that vessel (note that vessel code is not in counts_raw), just commenting that part out for 2021 and 2022:

```

dplyr::filter(counts_raw,
              Station == twice_fished & Setno == 4)
> # A tibble: 0 x 13
> # ... with 13 variables: Row number <dbl>, Year <dbl>, Stlkey <dbl>,
> #   Station <dbl>, Setno <dbl>, IPHC Species Code <dbl>, Scientific Name <fct>,
> #   Species Name <fct>, SampleType <fct>, HooksFished <dbl>,
> #   HooksRetrieved <dbl>, HooksObserved <dbl>, Number Observed <dbl>

# So just keep these:
# dplyr::filter(counts_raw,
#               !(Station == twice_fished & Setno == 4))
#countData2020_no_halibut <- dplyr::filter(counts_raw,
#                                           !(Station == twice_fished & Setno == 4)) %>%
# Seems that can't just keep using that even if twice_fished = NA
countData2022_no_halibut <- counts_raw %>%
  dplyr::select(year = Year,
                station = Station,
                spNameIPHC = "Species Name",
                specCount = "Number Observed") %>%
  arrange(station) %>%
  dplyr::mutate(year = as.integer(year),
                station = as.character(station),
                spNameIPHC = as.character(spNameIPHC),
                specCount = as.integer(specCount))

testthat::expect_equal(names(countData2013), names(countData2022_no_halibut))
countData2022_no_halibut
> # A tibble: 1,368 x 4
>   year station spNameIPHC      specCount
>   <int> <chr>   <chr>          <int>
> 1  2022 2001    Spiny Dogfish         50
> 2  2022 2001    Yelloweye Rockfish      1
> 3  2022 2001    Empty Hook          108

```

```

> 4 2022 2003 Arrowtooth Flounder 1
> 5 2022 2003 Sablefish (Blackcod) 33
> 6 2022 2003 Spiny Dogfish 2
> 7 2022 2003 Longnose Skate 1
> 8 2022 2003 Hook with Skin 1
> 9 2022 2003 Empty Hook 120
> 10 2022 2003 Hook with Bait 2
> # ... with 1,358 more rows
summary(countData2022_no_halibut)
>      year      station      spNameIPHC      specCount
> Min.   :2022   Length:1368      Length:1368      Min.    : 1.00
> 1st Qu.:2022   Class :character  Class :character  1st Qu.:  2.00
> Median :2022   Mode  :character  Mode  :character  Median :  5.00
> Mean    :2022                                     Mean    :18.82
> 3rd Qu.:2022                                     3rd Qu.:23.00
> Max.    :2022                                     Max.    :138.00

```

Hooks observed and retrieved

Now, obtain the numbers of hooks observed and retrieved from `counts_raw`, to then merge into the set details:

```

# hook_details <- dplyr::filter(counts_raw,
#                               !(Station == twice_fished & Setno == 4)) %>%
hook_details <- counts_raw %>%
  dplyr::group_by(Station) %>%
  dplyr::summarise(year = unique(Year),
                   hooksRetr = unique(HooksRetrieved),
                   hooksObs = unique(HooksObserved)) %>%
  dplyr::rename(station = Station) %>%
  dplyr::ungroup() %>%
  arrange(station) %>%
  dplyr::mutate(year = as.integer(year),
                station = as.character(station))

hook_details
> # A tibble: 174 x 4
>   station year hooksRetr hooksObs
>   <chr>   <int>   <dbl>   <dbl>
> 1 2001    2022     800     160
> 2 2003    2022     800     160
> 3 2004    2022     800     160
> 4 2005    2022     792     160
> 5 2006    2022     800     160
> 6 2007    2022     800     160

```

```

> 7 2008      2022      800      160
> 8 2009      2022      792      160
> 9 2013      2022      800      160
> 10 2015     2022      800      160
> # ... with 164 more rows

```

```
testthat::expect_equal(sets_simp_std_corrected$station, hook_details$station)
```

So now need to get the hook details into the set details, and keep columns as for setData2013 but also with standard, and may as well keep hooksRetr and hooksObs:

```

setData2022 <- dplyr::left_join(sets_simp_std_corrected,
                                hook_details,
                                by = c("year", "station")) %>%
  dplyr::mutate(E_it20 = effSkateIPHC * hooksObs / hooksRetr) %>%
  dplyr::select(year,
                station,
                lat,
                lon,
                avgDepth,
                effSkateIPHC,
                E_it20,
                usable,
                standard,
                hooksRetr,
                hooksObs) %>%
  dplyr::mutate(year = as.integer(year),
                station = as.character(station),
                avgDepth = as.integer(avgDepth),
                usable = as.character(usable),
                standard = as.factor(standard))

setData2022
> # A tibble: 174 x 11
>   year station  lat  lon avgDepth effSkateIPHC E_it20 usable stand~1 hooks~2
>   <int> <chr>   <dbl> <dbl>    <int>      <dbl>   <dbl> <chr>   <fct>    <dbl>
> 1  2022 2001    48.3 -126.     77      8.03    1.61 Y      Y      800
> 2  2022 2003    48.5 -125.     74      8.03    1.61 Y      Y      800
> 3  2022 2004    48.5 -126.     55      8.03    1.61 Y      Y      800
> 4  2022 2005    48.5 -126.     58      7.95    1.61 Y      Y      792
> 5  2022 2006    48.5 -126.    107      8.03    1.61 Y      Y      800
> 6  2022 2007    48.7 -125.     35      8.03    1.61 Y      Y      800
> 7  2022 2008    48.7 -125.     34      8.03    1.61 Y      Y      800
> 8  2022 2009    48.7 -126.     55      7.95    1.61 Y      Y      792
> 9  2022 2013    48.8 -126.     33      8.03    1.61 Y      Y      800
> 10 2022 2015    48.8 -126.     96      8.03    1.61 Y      Y      800

```

```

> # ... with 164 more rows, 1 more variable: hooksObs <dbl>, and abbreviated
> #   variable names 1: standard, 2: hooksRetr
testthat::expect_equal(names(setData2013), names(setData2022)[1:ncol(setData2013)])
summary(setData2022)
>      year      station      lat      lon
> Min.   :2022   Length:174   Min.   :48.34   Min.   : -133.7
> 1st Qu.:2022   Class :character 1st Qu.:51.35   1st Qu.: -131.1
> Median :2022   Mode  :character Median :52.52   Median : -130.1
> Mean   :2022                Mean  :52.33   Mean   : -129.8
> 3rd Qu.:2022                3rd Qu.:53.67   3rd Qu.: -128.7
> Max.   :2022                Max.   :55.33   Max.   : -125.1
>      avgDepth      effSkateIPHC      E_it20      usable      standard
> Min.   : 7.00   Min.   :2.110   Min.   :0.6057   Length:174   N: 61
> 1st Qu.: 45.50   1st Qu.:7.950   1st Qu.:1.6060   Class :character  Y:113
> Median : 66.50   Median :7.950   Median :1.6061   Mode  :character
> Mean   : 81.34   Mean   :7.925   Mean   :1.6001
> 3rd Qu.:113.50   3rd Qu.:8.030   3rd Qu.:1.6061
> Max.   :257.00   Max.   :8.110   Max.   :1.6067
>      hooksRetr      hooksObs
> Min.   :209.0   Min.   : 60.0
> 1st Qu.:792.0   1st Qu.:160.0
> Median :792.0   Median :160.0
> Mean   :789.5   Mean   :159.4
> 3rd Qu.:800.0   3rd Qu.:160.0
> Max.   :808.0   Max.   :160.0

```

Pacific Halibut counts

As noted above, the data extraction for the counts is for all non-halibut species. We still want the halibut counts for just the first 20 hooks – the `data_for_all_species` vignette (for data up to 2019) shows that the 20-hook and full hook counts (Series A and B) are very similar when rescaled, and the rescaling is miniscule with $G_A/G_B = 1.005$. So this justifies sticking with 20-hook counts for halibut, even though the full data are available for all sets, given it is a halibut survey. (Using all hooks for all years could be done, but would be a lot of new code).

There are two options for getting halibut counts for the first 20 hooks (given we don't have hook-by-hook data, though it could probably be obtained just maybe not from the IPHC website).

Option 1.

Take the halibut counts for all the hooks (which we have in `sets_raw` and subsequent objects) and create $N_{it20_halibut_est} = E_{it20} / E_{it} * N_{it}$, or equivalently just $N_{it20_halibut_est} = \text{hooksObs} / \text{hooksRetr} * N_{it}$. Note that observed refers to ob-

served for non-halibut species (presumably `hooksRetr` works for halibut). Not strictly the first 20 hooks, but is a rescaling. But will not guarantee integer values.

```
setData2022_and_halibut <-
  dplyr::left_join(setData2022,
                    dplyr::select(sets_simp_std_corrected,
                                   c(station,
                                       U32halibut,
                                       O32halibut)),
                    by = "station") %>%
  dplyr::mutate(N_it_halibut = U32halibut + O32halibut,
                 N_it20_halibut_opt_1 = hooksObs / hooksRetr * N_it_halibut)
setData2022_and_halibut %>% dplyr::select(station,
                                           N_it_halibut,
                                           N_it20_halibut_opt_1)

> # A tibble: 174 x 3
>   station N_it_halibut N_it20_halibut_opt_1
>   <chr>      <dbl>      <dbl>
> 1 2001         2         0.4
> 2 2003         3         0.6
> 3 2004        19         3.8
> 4 2005         6         1.21
> 5 2006         1         0.2
> 6 2007        37         7.4
> 7 2008       113        22.6
> 8 2009        23         4.65
> 9 2013        54        10.8
> 10 2015         9         1.8
> # ... with 164 more rows
```

Option 2.

Add all the 20-hook counts for a set (which include `Hook with Skin` etc.) and compare with `hooksObs`. The latter is higher (or equal), and the difference is halibut (as the only **non** non-halibut species). Compare with the results from option 1. If close then use option 2, since it will be just be halibut counts and gives an integer number, and is based on the first 20 hooks.

Add counts for each set:

```
counts_20 <- countData2022_no_halibut %>%
  dplyr::group_by(station) %>%
  dplyr::summarise(non_halibut = sum(specCount)) %>%
  dplyr::ungroup()
counts_20
> # A tibble: 174 x 2
```

```

> station non_halibut
> <chr> <int>
> 1 2001 159
> 2 2003 160
> 3 2004 156
> 4 2005 159
> 5 2006 159
> 6 2007 151
> 7 2008 141
> 8 2009 155
> 9 2013 147
> 10 2015 159
> # ... with 164 more rows

```

Now join the two options together to calculate `N_it20_halibut_opt_2` and then compare the two estimates of `N_it20_halibut`:

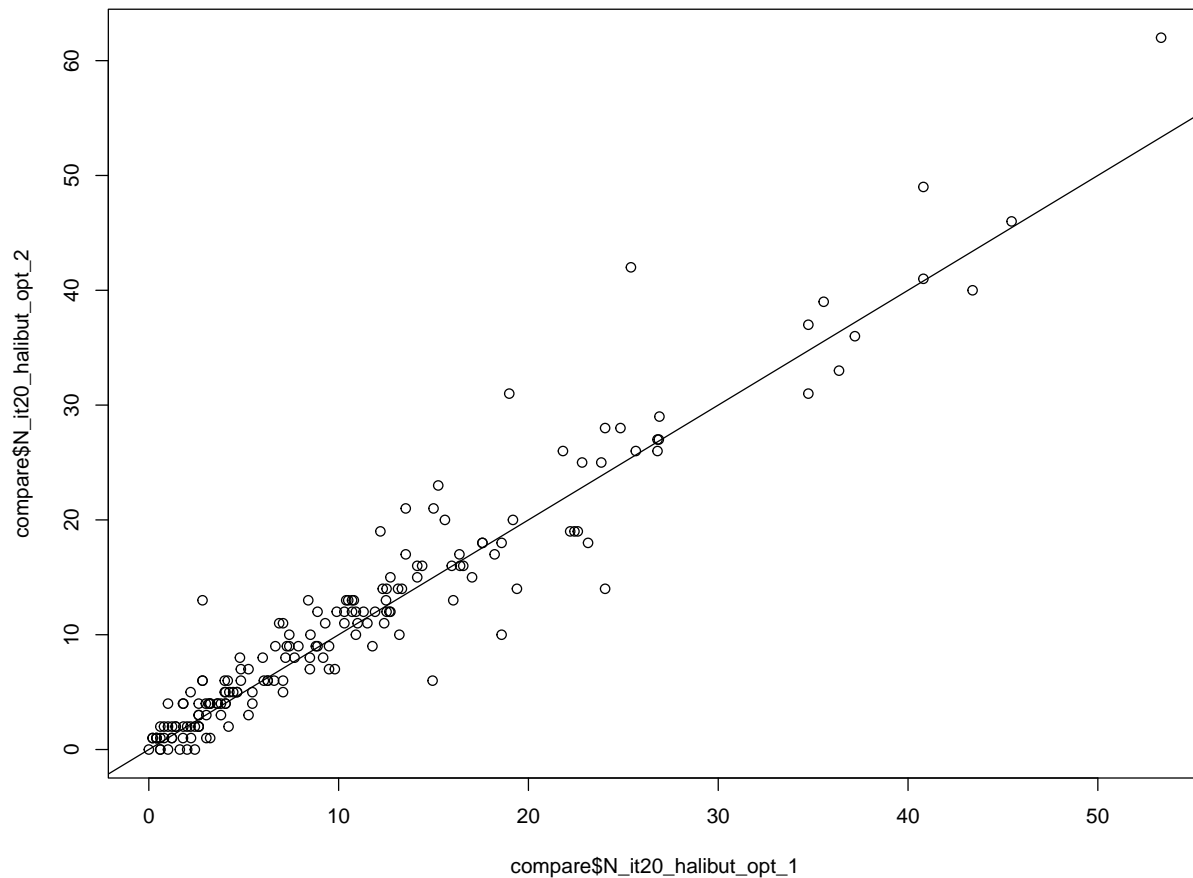
```

compare <-
  dplyr::left_join(setData2022_and_halibut,
                    counts_20,
                    by = "station") %>%
  dplyr::mutate(N_it20_halibut_opt_2 = hooksObs - non_halibut,
                N_it20_opt_1_over_opt_2 = N_it20_halibut_opt_1 / N_it20_halibut_opt_2) %>%
  dplyr::select(year,
                station,
                usable,
                N_it20_halibut_opt_1,
                N_it20_halibut_opt_2,
                N_it20_opt_1_over_opt_2)
compare$spNameIPHC <- "Pacific Halibut"
compare
> # A tibble: 174 x 7
>   year station usable N_it20_halibut_opt_1 N_it20_halibut_opt_2 N_it20_opt_1_over_opt_2 spNameIPHC
>   <int> <chr>   <chr>           <dbl>           <dbl>           <dbl>           <chr>
> 1  2022 2001     Y             0.4             1             0.4           Pacifi~
> 2  2022 2003     Y             0.6             0             Inf            Pacifi~
> 3  2022 2004     Y             3.8             4             0.95          Pacifi~
> 4  2022 2005     Y             1.21            1             1.21          Pacifi~
> 5  2022 2006     Y             0.2             1             0.2           Pacifi~
> 6  2022 2007     Y             7.4             9             0.822         Pacifi~
> 7  2022 2008     Y            22.6            19             1.19          Pacifi~
> 8  2022 2009     Y             4.65            5             0.929         Pacifi~
> 9  2022 2013     Y            10.8            13             0.831         Pacifi~
> 10 2022 2015     Y             1.8             1             1.8           Pacifi~
> # ... with 164 more rows, and abbreviated variable names

```

```
> # 1: N_it20_halibut_opt_2, 2: N_it20_opt_1_over_opt_2, 3: spNameIPHC

plot(compare$N_it20_halibut_opt_1, compare$N_it20_halibut_opt_2)
abline(a = 0, b = 1)
```



```
cor(compare$N_it20_halibut_opt_1,
     compare$N_it20_halibut_opt_2)
> [1] 0.9591653
```

So this is the right approach and correlation coefficient is high, though numbers not quite as close as may have thought. But these data are used for aggregating across all stations in a year (and any further analyses on halibut for management purposes should be done using the full halibut data anyway – we wouldn't really need that). And the means aren't too bad:

```
mean(compare$N_it20_halibut_opt_1)
> [1] 10.84046
mean(compare$N_it20_halibut_opt_2)
> [1] 11.47126
```

(9.6 and 10.5 in 2021; 10.8 and 11.5 in 2022).

So either of these would work. So use option 2 since gives an integer count:

```
compare$N_it20_halibut_opt_2
> [1] 1 0 4 1 1 9 19 5 13 1 4 6 19 4 14 28 17 20 19 16 2 5 16 11 2
> [26] 6 13 27 6 17 12 6 2 15 37 33 5 6 12 11 14 26 62 49 12 2 12 16 9 20
> [51] 28 4 7 8 12 11 2 12 8 23 9 41 2 13 7 6 8 0 12 4 25 21 7 7 1
> [76] 25 6 6 18 18 36 6 18 4 11 10 27 16 40 5 42 26 39 9 46 12 9 1 8 3
> [101] 3 14 11 29 14 31 5 1 0 0 6 11 15 12 2 14 13 3 12 1 8 7 17 2 4
> [126] 10 2 4 13 6 13 1 13 3 1 5 8 10 0 2 5 4 4 0 1 5 2 4 2 13
> [151] 5 9 4 1 2 31 3 18 21 15 16 26 14 11 1 19 9 10 4 10 9 4 0 2
countData2022_halibut <- dplyr::select(compare,
                                     year,
                                     station,
                                     spNameIPHC,
                                     specCount = N_it20_halibut_opt_2) %>%
  dplyr::mutate(specCount = as.integer(specCount))
countData2022 <- rbind(countData2022_no_halibut,
                       countData2022_halibut) %>%
  dplyr::arrange(station)
# First time running, called the above countData2020_NEW to check remaining data didn't
# expect_equal(countData2020, filter(countData2020_NEW, spNameIPHC !=
#                                     "Pacific Halibut"))
```

Note that for 2021 and 2022 this does give zeros for Pacific Halibut (the only species that will have a zero, because we have a value for each station because zero counts are in the original sets_raw):

```
summary(dplyr::filter(countData2022,
                      spNameIPHC == "Pacific Halibut"))
>      year      station      spNameIPHC      specCount
> Min.   :2022  Length:174      Length:174      Min.    : 0.00
> 1st Qu.:2022  Class :character  Class :character  1st Qu.: 4.00
> Median :2022  Mode  :character  Mode  :character  Median : 9.00
> Mean    :2022                                Mean    :11.47
> 3rd Qu.:2022                                3rd Qu.:15.00
> Max.    :2022                                Max.    :62.00
unique(dplyr::filter(countData2022, specCount == 0)$spNameIPHC)
> [1] "Pacific Halibut"
```

Check species names

The file `inst/extdata/iphc-spp-names.csv` contains species common names (as used for gfsynopsis, and a few extra like unidentified skate) and the IPHC common name. The function `check_iphc_spp_name()` has a list of non-groundfish species that are automatically

ignored. These first results are from running these functions *before* updating anything, so the results are hardwired here (chunks are not evaluated); so set `eval=TRUE` then back to `eval=FALSE`; then we update the species list and re-run the functions.

These are IPHC names that are not given in `iphc-spp-names.csv` (automatically ignoring obvious ones that are listed in the function), for years up to 2020 (since not updated code yet):

```
check_iphc_spp_name()
## [1] "Unidentified Shark"           "Unident. Rockfish"
## [3] "unident. thornyhead (Idiot)"  "Grenadier (Rattails)"
## [5] "Miscellaneous Shark"         "Eelpout"
## [7] "unident. Roundfish"          "unident. Sculpin"
## [9] "Unident. Flatfish"           "Greenland Turbot"
## [11] "unident. Hagfish"            "Starry Skate"
## [13] "Black Skate"                 "Brittle Star"
## [15] "Glass Sponge"               "Basketstar"
## [17] "Blackspotted Rockfish"
```

These are the ones just for the new 2022 data:

```
check_iphc_spp_name(countData2022)
## [1] "Basketstar"                  "unident. thornyhead (Idiot)"
## [3] "Sandpaper Skate"             "Sea Whip"
## [5] "Stylaster campylecus (coral)" "Brittle Star"
## [7] "unident. Sculpin"            "Glass Sponge"
## [9] "Sun Sea Star"                "Salmon Shark"
## [11] "Jellyfish"                   "Great Sculpin"
## [13] "Cabezon"                     "Unident. Salmon"
## [15] "Unident. Rockfish"           "unident. organic matter"
## [17] "Dungeness Crab"              "Blackspotted Rockfish"
```

There were only six for 2020 though (a lot more for 2021):

```
check_iphc_spp_name(countData2020)
## [1] "unident. thornyhead (Idiot)" "Brittle Star"
## [3] "Glass Sponge"               "Basketstar"
## [5] "Blackspotted Rockfish"      "unident. Sculpin"
```

For 2020 I said that only the Thornyhead and Blackspotted Rockfish are likely of interest (Issues #17 and #18). And the sharks from the earlier list. So look at just the new ones in 2022 that aren't in 2020 or any previous year (switch this to `eval=TRUE`, paste results in, then set back to `eval=FALSE`, or maybe try leaving as it should automatically give `character(0)` once fixed (as it does for 2021):

```
setdiff(check_iphc_spp_name(countData2022),
        check_iphc_spp_name())
> character(0)
```

```
# Before updating anything this gives:
# [1] "Rosethorn Rockfish" "Red Irish Lord"
```

2021: Of these, Sandpaper Skate, Salmon Shark, Great Sculpin, and Cabezon are in gfsynopsis but have not been designated an `iphc_common_name` in `iphc-spp-names.csv` (have to do that manually). Though Sandpaper Skate, Salmon Shark, and Great Sculpin do show up as having IPHC data in 2019 gfsynopsis report, but looks like only data from GFBio, looking carefully at the `data_for_all_species` vignette for 2020: http://htmlpreview.github.io/?https://github.com/pbs-assess/gfiphc/blob/master/vignettes/data_for_all_species.html They did not have 2020 IPHC data, but do for 2021 (GS had 1995 and 1996 as zeros; don't think others did). Cabezon has no previous data.

So, in 2021 added those species to `iphc-spp-names.csv`, which may discover some old data for those years when I redo the vignettes, as it seems strange that they never seem to show up in the 20-hook-only data, just in GFBio.

2021: Also add these to the `ignore_obvious` list in `check_iphc_spp_name()`:

“Sea Whip”, “Stylaster campylecus (coral)”, “Sun Sea Star”, “Jellyfish”, “Unident. Salmon”, “unident. organic matter”, “Dungeness Crab”

That list already had Sunflower Sea Star in it, presumably the same as Sun Sea Star.

2022: The `setdiff()` just gave "Rosethorn Rockfish" "Red Irish Lord" which do appear in latest synopsis report with older IPHC data, which I presume is why I just need to update `iphc-spp-names.csv` for which we have NA and ***** as their IPHC names. Doing that, and rerunning the `setdiff()` now correctly gives an empty result (above and then here also):

Then redoing those above commands with updated code gives this, where some species are returned because they are not non-groundfish ones (or Brittle Star or Glass Sponge which we also kept in the past) that we want to automatically ignore:

```
check_iphc_spp_name(countData2022)
> [1] "unident. thornyhead (Idiot)" "Basketstar"
> [3] "Brittle Star"                "Grenadier (Rattails)"
> [5] "unident. Sculpin"            "Blackspotted Rockfish"
# That still retains some we may want to think about further at some point, but
# these are all in the overall list for all years:
setdiff(check_iphc_spp_name(countData2022),
        check_iphc_spp_name())
> character(0)
check_iphc_spp_name()
> [1] "Unidentified Shark"          "Unident. Rockfish"
> [3] "unident. thornyhead (Idiot)" "Grenadier (Rattails)"
> [5] "Miscellaneous Shark"        "Eelpout"
> [7] "unident. Roundfish"         "unident. Sculpin"
> [9] "Unident. Flatfish"          "Greenland Turbot"
> [11] "unident. Hagfish"           "Starry Skate"
```

```
> [13] "Black Skate"           "Brittle Star"
> [15] "Glass Sponge"         "Basketstar"
> [17] "Blackspotted Rockfish"
```

Save data sets

```
usethis::use_data(countData2022,
                    overwrite = TRUE)
> v Saving 'countData2022' to 'data/countData2022.rda'
> * Document your data (see 'https://r-pkgs.org/data.html')

usethis::use_data(setData2022,
                    overwrite = TRUE)
> v Saving 'setData2022' to 'data/setData2022.rda'
> * Document your data (see 'https://r-pkgs.org/data.html')
```

Add descriptions for new years in R/data.R.