

Mehrabi et al. 2020. The global divide in data driven farming.
Supplementary Information B: Global farm size map.

Code by Vincent Ricciardi

December 23, 2020

Contents

| | | |
|----------|---|-----------|
| 1 | Aim of document | 2 |
| 2 | Reproducibility | 2 |
| 3 | Data | 2 |
| 3.1 | Country boundaries | 3 |
| 3.2 | Country farm size distributions | 3 |
| 3.3 | Global agricultural area | 4 |
| 3.4 | Global field sizes | 5 |
| 3.4.1 | Read in Data | 5 |
| 3.4.2 | Pre-process Data | 5 |
| 3.4.3 | Interpolate | 5 |
| 3.4.4 | Rasterize | 7 |
| 4 | Create common dataframe | 8 |
| 4.1 | Convert to equal area | 8 |
| 4.2 | Set extents | 8 |
| 4.3 | Create stack | 9 |
| 4.4 | Finalize data set | 9 |
| 5 | Match Field Sizes to Farm Sizes | 11 |
| 5.1 | Algorithm | 11 |
| 5.2 | Resampled Results | 15 |
| 5.3 | Sanity Check | 15 |
| 5.4 | Final Output | 16 |
| | References | 18 |

1 Aim of document

The aim of this document is to create a global map of farm sizes. We do this by merging two currently available sources of information on the distribution of farm sizes across the world: national level data on farm size distributions and subnational data on field size distributions. The resulting map is a "best guess" of where different sizes of farms are distributed globally. We caution that this map should not be used for detailed country level analysis, but is intended for global assessment studies, where results are typically aggregated at the regional or global level. This document explains the underlying data, pre-processing, and algorithm we developed to create this map. We were driven to create this map as a work that can be built upon and improved as higher resolution and more accurate data become available.

2 Reproducibility

Here we call the **R** package `renv` (Ushey 2020). This will create a local library on your computer and install a copy of the packages required by this project as they existed on CRAN by the specified version number, and update the **R** session to use these packages. This helps make our analysis fully reproducible on your machine.

Note, the 'renv.lock' file needs to be in the top level of this project's directory. This code block needs to only be run when initially setting up your project then can be commented out.

```
# Uncomment during first run  
# install.packages('renv')  
# renv::init()
```

We also set the seed of the entire document to ensure the same results when randomly sampling.

```
set.seed(123)
```

Note that the **R** version used here is 3.6.3 (2020-02-29).

For the analysis in this document we will be using the following packages: `R-data.table` (Dowle and Srinivasan 2019), `R-foreach` (Revolution Analytics and Weston 2017), `R-formatR` (Xie 2019a), `R-ggthemes` (**R-ggthemes**), `R-gmodels` (Bolker, Lumley, Johnson, and Johnson 2018), `R-hablar` (Sjoberg 2020), `R-leaflet` (Cheng, Karambelkar, and Xie 2019), `R-lintr` (Hester, Angly, and Hyde 2020), `R-magrittr` (Bache and Wickham 2014), `R-knitr` (Xie 2019b), `R-kknn` (Schliep and Hechenbichler 2016), `R-plyr` (Wickham 2016), `R-raster` (Hijmans 2019), `R-renv` (Ushey 2020), `R-reshape2` (Wickham 2017a), `R-rgdal` (Bivand, Keitt, and Rowlingson 2019), `R-rworldmap` (South 2016), `R-sf` (Pebesma 2020), `R-tidyverse` (Wickham 2017b), `R-viridis` (**R-viridis**), `R-zoo` (Zeileis, Grothendieck, and Ryan 2019).

3 Data

In this analysis we leverage two key datasets. Lowder, Scoet, and Raney 2016 (Lowder, hereafter) contains farm size distributions (in terms of agricultural area and the number of farms) for 100 countries. Lesiv, Laso Bayas, See, Duerauer, Dahlia, Durando, Hazarika, Kumar Sahariah, Vakolyuk, Blyshchyk, Bilous, Perez-Hoyos, Gengler, Prestele, Bilous, Akhtar, Singha, Choudhury, Chetri, Malek, Bungnamei, Saikia, Sahariah, Narzary, Danylo, Sturn, Karner, McCallum, Schepaschenko, Moltchanova, Fraisl, Moorthy, and Fritz 2019 (Lesiv, hereafter) contains a crowd-sourced point data of categorical field size classes (ranging from very small to large fields). We use Lesiv's qualitative field size classes to spatially disaggregate the Lowder farm size classes. Links and access dates to these input data and other ancillary data sets used in our analysis as highlighted below:

1. Global field size data from Lesiv, retrieved from from <http://pure.iiasa.ac.at/id/eprint/15526/> on October 12th, 2018.

2. Country farm size distributions from Lowder, retrieved from <http://iopscience.iop.org/article/10.1088/1748-9326/11/12/124010> on July 12th, 2018.
3. Global cropland and pastureland from Ramankutty, Evan, Monfreda, and Foley 2008 (Ramankutty, hereafter), retrieved from www.earthstat.org on July 12th 2018.
4. Country boundaries, retrieved internally in R through South 2016.

3.1 Country boundaries

First we make a raster of the world country data from the `rworldmap` package.

```
world <- rworldmap::getMap(resolution = 'low')
lookup <- as.data.frame(cbind(as.character(
  world@data[, 'ISO3']),
  world@data[, 'ADMIN'],
  as.character(world@data[, 'ADMIN'])))

raster.world <- raster(res = c(0.0833282, 0.0833282))
extent(raster.world) <- extent(world)
world.raster <- rasterize(as(world, 'SpatialPolygons'),
  raster.world,
  field = world@data[, 'ADMIN'],
  fun = 'first')
world.rast <- writeRaster(world.raster,
  'data/tmp/worldrast.tif',
  format = 'GTiff',
  overwrite = TRUE)

world <- raster('data/tmp/worldrast.tif')
```

3.2 Country farm size distributions

Lowder's data contains two variables at the national level, the amount of agricultural area per farm size class, and the number of farms per farm size class. We are interested in the amount of agricultural area per farm size class for our analysis. This data is distributed as the World Census of Agriculture's (WCA) farm size classes, which are: 0-1 ha, 1-2 ha, 2-5 ha, 5-10 ha, 10-20 ha, 20-50 ha, 50-100 ha, 100-200 ha, 200-500 ha, 500-1000 ha, 1000ha. *This data is read in here and columns are relabeled.*

```
# Load Lowder's distribution dataset
lwd <- read.csv('data/lowder/Lowder_2016_dist.csv',
  header = T,
  na.strings = c('', 'NA'))

# Subset only agricultural area
lwd <- lwd[which(lwd$Holdings..agricultural.area != 'Holdings'), ]
lwd <- lwd[, c(1,5:15)]
names(lwd) <- c('country', '0_1', '1_2', '2_5',
  '5_10', '10_20', '20_50', '50_100', '100_200',
  '200_500', '500_1000', '1000_5000')

# Ensure variable type is numeric
for (i in names(lwd)[2:length(names(lwd))]) {
  lwd[[i]] <- as.numeric(lwd[[i]])
}
```

```
}
```

Next we need to calculate the cumulative sum per country across Lowder's farm size classes. To do this we first convert the data from wide to long format, and set classes for countries without data to zero area.

```
lwd <- melt(lwd, id.vars = c('country'))  
lwd[is.na(lwd)] <- 0
```

Then we calculate the proportional agricultural area for each farm size class for each country. We note that Lowder contains some countries without farm size distributions and some countries do not have a distribution for each farm size class, and so we remove those cases here.

```
lwd <- lwd %>%  
  group_by(country) %>%  
  mutate(total = sum(value),  
         perc = value / total) %>%  
  filter(value > 0) %>%  
  select(country, variable, perc)
```

Finally, cast dataframe back into wide form.

```
lwd <- dcast(lwd, country ~ variable)  
write.csv(lwd, 'data/tmp/lwd.csv')
```

3.3 Global agricultural area

We use Ramankutty's 10 km² cropland map (crop, hereafter) and pastureland map (pasture, hereafter) as our reference layer of the distribution of agricultural land at the subnational level. We read in these files here and unionize the crop and pasture maps to create an agricultural land raster (ag, hereafter).

```
crop <- paste0('data/Ramankutty_2008_cropland/',  
              'CroplandPastureArea2000_Geotiff/',  
              'cropland2000_area.tif')  
  
pasture <- paste0('data/Ramankutty_2008_cropland/',  
                 'CroplandPastureArea2000_Geotiff/',  
                 'pasture2000_area.tif')  
  
crop <- raster(crop)  
pasture <- raster(pasture)  
  
# Here we take a union of crop and pasture areas  
# Set NA to 0, ensure same origin, find sum  
crop_tmp <- crop  
pasture_tmp <- pasture  
crop_tmp[is.na(crop_tmp)] <- 0  
pasture_tmp[is.na(pasture_tmp)] <- 0  
origin(crop_tmp) <- origin(pasture_tmp) <- c(0.0, 0.0)  
ag <- mosaic(crop_tmp, pasture_tmp, fun = sum)  
ag[ag[] == 0] <- NA
```

3.4 Global field sizes

Here we interpolate the crowd-sourced field size data point data onto Ramankutty's agricultural area data. While the original sampling frame for this field size was on a unionized cropland map, we extend the distribution here in an attempt to better match and account for the inclusion of non-cropland in Lowder's farm size data (noting that later we will clip the resulting product to cropland only). We use kk-nearest neighbor (kkNN) on a 0.01 degree grid (kkNN is a weighted variant of kNN that we use to try to account as best as possible for Lesiv's sparsely sampled points in Africa).

3.4.1 Read in Data

First, we read in the estimated dominant field sizes point data and subset to not include NA values or values coded as having no fields.

```
fs <- read.csv("data/Fritz_2019/Global Field Sizes/estimated_dominant_field_sizes.csv")
```

3.4.2 Pre-process Data

We recode Lesiv's data according to their coding as detailed below. And then convert this data into a spatial object.

- 3502 - Very large fields with an area of greater than 100 ha
- 3503 - Large fields with an area between 16 ha and 100 ha
- 3504 - Medium fields with an area between 2.56 ha and 16 ha
- 3505 - Small fields with an area between 0.64 ha and 2.56 ha
- 3506 - Very small fields with an area less than 0.64 ha
- 3507 - no fields
- NA - skipped

```
fs[which(fs$field_size == 3507), 5] <- NA

fs_clean <- fs %>%
  na.omit() %>%
  dplyr::mutate(field_size_fac = as.factor(field_size)) %>%
  dplyr::mutate(fs_verbatim_fac = recode(field_size_fac,
                                        "3502" = "very large",
                                        "3503" = "large",
                                        "3504" = "medium",
                                        "3505" = "small",
                                        "3506" = "very small")) %>%
  dplyr::select(-rowid, -sampilid) %>%
  st_as_sf(coords = c("x", "y"),
           crs = "+proj=longlat +ellps=WGS84")
```

3.4.3 Interpolate

We make a target raster grid to interpolate onto using the agricultural area raster. We use a very simple method of interpolation based on the latitude and longitude.

```

ag_GRO <- ag
ag_GRO[] <- ifelse(ag_GRO[] > 0, 1, NA)
target_pts <- rasterToPoints(ag, spatial = TRUE)
target_pts_agGRO <- rasterToPoints(ag_GRO, spatial = TRUE)

```

We then train the kkNN model using different types of kernels and number of points used in the classification procedure. We use 10 fold cross validation to identify the best parameters to use in the final model (noting that CV does not approximate the test error as would be estimated using an independent hold-out set).

```

fs_train <- data.frame(field_size = fs_clean$fs_verbatim_fac,
                      x = st_coordinates(fs_clean)[,1],
                      y = st_coordinates(fs_clean)[,2])

fs_interp_train <- kknn::train.kknn(field_size ~.,
                                   data = fs_train,
                                   ks = seq(5,51, 5),
                                   distance = 2,
                                   kernel = c("gaussian",
                                             "triangular",
                                             "rectangular",
                                             "epanechnikov",
                                             "optimal"),
                                   kcv = 10)

```

Next we summarize and plot the results of the model fitting. We note the the minimal classifications error is 0.43. Crudely, with 5 classes we would expect a raw misclassification error of 80%, which suggests the model has in the best case helped reduce blind uncertainty by around 53.719%. We also plot the error from the different models below.

Future versions improvements to this model may be made through identifying ancillary features that do a better job of predicting field size classes than simply latitude and longitude alone. Similarly, using an ordinal response, instead of the nominal response defaulted here, is likely to improve predictions also.

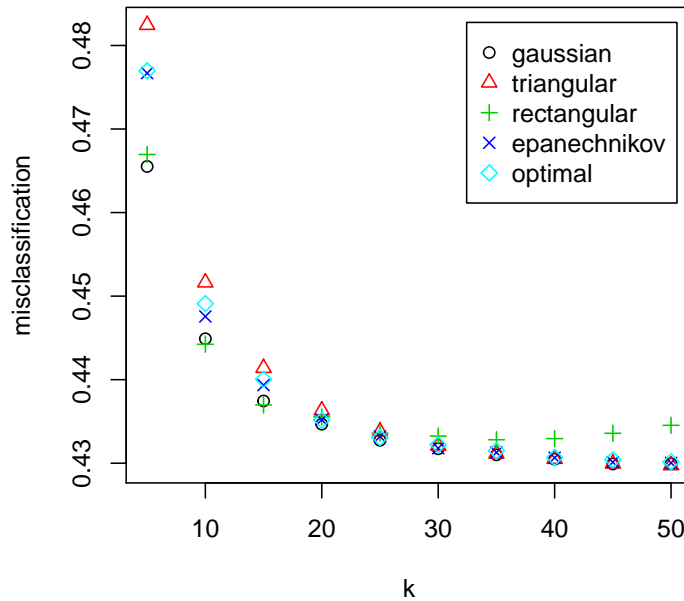


Figure 1: kkNN results of field size interpolation

We use the best fitting model parameters to predict the field size classes across the agricultural area grid, and then extract the fitted values and the associated probabilities.

```
fs_interp_final <- kknm::kknm(field_size ~ .,
                             train = fs_train,
                             test = target_pts_agGRO,
                             kernel = fs_interp_train$best.parameters$kernel,
                             k = fs_interp_train$best.parameters$k)

target_pts_agGRO_res <- target_pts_agGRO %>%
  # extract the interpolated class at each
  # grid cell with the kknm::fitted function
  as.data.frame() %>%
  # only retain the probability of the
  # interpolated size class, discard the others
  mutate(field_size_fit = fitted(fs_interp_final),
         prob = apply(fs_interp_final$prob, 1, function(input) max(input)))
```

3.4.4 Rasterize

Here we convert results back to a raster and save the results of the fitted classes.

```
target_pts_agGRO_res_spdf <- target_pts_agGRO_res
coordinates(target_pts_agGRO_res_spdf) <- ~ x + y

fs_interp_final_ras <- rasterize(target_pts_agGRO_res_spdf,
                                ag,
```

```

as.integer(
  target_pts_agGRO_res_spdf$field_size_fit),
progress = "text")

fs_interp_final_ras_prob <- rasterize(target_pts_agGRO_res_spdf,
  ag,
  target_pts_agGRO_res_spdf$prob,
  progress = "text")

lesiv <- fs_interp_final_ras

```

4 Create common dataframe

To run our algorithm to create a field size map we need to ensure all rasterized datasets (i.e., field size, crop, pasture, ag, and country boundaries) are all in the same resolution, have the same coordinate reference system (crs), and have the same spatial extent, which we do in this section.

4.1 Convert to equal area

Here we set each dataset to have an equal area (Eckert IV) at a 8.4 km² resolution. To interpolate, we use nearest neighbor for the categorical data and bilinear for the continuous data.

```

rast.list.ngb <- list(lesiv, world)
rast.list.bil <- list(crop, pasture, ag)

resValue <- 8439
rast.list.ngb.eq <- lapply(rast.list.ngb, function(x)
  projectRaster(
    x,
    res = c(resValue, resValue),
    crs = '+proj=eck4 +datum=WGS84 +ellps=WGS84 +towgs84=0,0,0',
    method = 'ngb',
    over = T))

rast.list.bil.eq <- lapply(rast.list.bil, function(x)
  projectRaster(
    x,
    res = c(resValue, resValue),
    crs = '+proj=eck4 +datum=WGS84 +ellps=WGS84 +towgs84=0,0,0',
    method = 'bilinear',
    over = TRUE))

rast.list <- c(rast.list.ngb.eq, rast.list.bil.eq) # Combine data
names(rast.list) <- c('lesiv', 'world', 'crop', 'pasture', 'ag')

```

4.2 Set extents

Next, we need to make sure that all the rasters have the same spatial extents. The country boundaries include Antarctica, while the cropland/pastureland/ag and field size data do not (e.g., crop ymin is -8501789, while the country boundaries ymin is -8503388). To ensure all the datasets match we need to expand the

cropland/pastureland/ag and Lesiv extents to match the country boundaries extent. This means extending the cropland/pastureland/ag and Lesiv data to also include the country boundary grid cells; we do this by adding rows of NA grid cells.

```
ex1 <- extent(rast.list$world)

# Extend
rast.list.ex <- lapply(rast.list,
                      function(x) extend(x, ex1))

# Crop
rast.list.c <- lapply(rast.list.ex,
                    function(x) crop(x, ex1))

# Force equal extent
rast.list.et <- lapply(rast.list.c,
                    function(x) setExtent(x, ex1,
                                           keepres = TRUE))
```

4.3 Create stack

Now that all the layers' resolution, crs, and spatial extent are equal, we can create a raster stack and proceed with the analysis.

```
rast.all <- stack(rast.list.et)
```

4.4 Finalize data set

We now convert the data in the raster stack into a dataframe, which we will use for later analysis.

```
df <- raster::extract(rast.all, 1:ncell(rast.all), df = T)
df <- as.data.frame(df)
df$ID <- 1:nrow(df)
```

Set any crop, pasture value under 0.01 to NA, to account for sparse agricultural areas.

```
df$crop <- ifelse(df$crop >= 0.01, df$crop, NA)
df$pasture <- ifelse(df$pasture >= 0.01, df$pasture, NA)
df$ag <- ifelse(df$ag >= 0.01, df$ag, NA)
```

Next we add back in the country names (presently they are only country ISO3 IDs), then subset the data to only contain countries in the Lowder dataset (this speeds up processing).

```
world <- getMap(resolution = 'low')
lookup <- as.data.frame(
  cbind(
    as.character(world@data[, 'ISO3']),
    world@data[, 'ADMIN'],
    as.character(world@data[, 'ADMIN'])))

matched <- lookup[match(df$world, lookup$V2), ]
df$ISO3 <- matched[[1]]
df$country <- matched[[3]]
```

We then match Lowder's country names to the global country boundary country names.

```
countries <- data.frame(names = sort(unique(df$country)))
lwdNames <- data.frame(names = unique(lwd$country))
```

```

lwdNames$names <- str_trim(lwdNames$names)
lwdNames$indata <- lwdNames$names %in% unique(countries$name)

# List of names in lower that do not match dataset's names
# note Guadeloupe, Martinique, Reunion has no data
# lwdNames[which(lwdNames$indata == FALSE), ]

changeFrom <- c('Korea Rep. of',
               "Lao People's Democratic Republic",
               'Viet Nam',
               'Serbia',
               'Bahamas',
               'Venezuela (Bolivarian Republic of)',
               'St. Kitts & Nevis',
               'Iran (Islamic Republic of)',
               "Cte d'Ivoire",
               'Virgin Islands United States')

changeTo <- c('South Korea',
              'Laos',
              'Vietnam',
              'Republic of Serbia',
              'The Bahamas',
              'Venezuela',
              'Saint Kitts and Nevis',
              'Iran',
              'Ivory Coast',
              'United States Virgin Islands')

require(plyr)
lwd$country <- mapvalues(str_trim(lwd$country),
                       from = changeFrom,
                       to   = changeTo)
detach('package:plyr')

```

Subset the dataset to only the countries containing farm size distributions.

```

lwdNames <- unique(lwd$country)
df_lwd <- df[which(df$country %in% lwdNames), ]
df_not_lwd <- df[which(!df$country %in% lwdNames), ]

```

5 Match Field Sizes to Farm Sizes

5.1 Algorithm

Here are the steps we used to create the farm size map, with uncertainty at the level of the distribution of farm to field size matching. The algorithm makes the assumption that smaller field sizes are owned by smaller farms (although we note that the uncertainty in this assumption grows with farm size, as while smaller farms can not own large fields, large farms can own many small fields).

1. Group pixels by country
2. Assign pseudo-random number to pixels
3. Reorder the pixels in a nested fashion: first by interpolated field size classes, and then by the pseudo-random numbers
4. Compute pixel wise agricultural area from Ramankutty
5. Compute cumulative agricultural area from Ramankutty
6. Divide cumulative sum of agricultural area by total area in each country to get cumulative proportional sums of agricultural area from Ramankutty
7. Assign pixel fields sizes to farm size classes by matching the cumulative proportional agricultural area from Lowder to the cumulative proportional agricultural area for interpolated Lesiv product where pixels are ordered from smallest to largest fields
8. For countries not in Lowder, hard code farm size based on field sizes.
9. Convert result to a raster
10. Save raster
11. Iterate steps 1 to 10 100x

Steps 1 to 6 are written below.

```
cumsumSkipNA <- function(x) {  
  FUNC <- match.fun('cumsum')  
  x[!is.na(x)] <- FUNC(x[!is.na(x)])  
  return(x)  
}  
  
df_lwd$lesiv <- factor(df_lwd$lesiv)  
levels(df_lwd$lesiv) <- c('5', '4', '3', '2', '1')  
df_lwd$lesiv <- as.integer(as.character(df_lwd$lesiv))  
  
steps_1_6 <- function(df_lwd) {  
  
  dat <- df_lwd %>%  
    mutate(pseudorandom = sample(1:n(), n(), replace = F)) %>%  
    arrange(country, lesiv, pseudorandom) %>%  
    group_by(country) %>%  
    mutate(crop = crop * resValue^2 * 0.0001, # new  
           pasture = pasture * resValue^2 * 0.0001, # new  
           ag = crop + pasture, # new  
           ag_area = ag * resValue^2 * 0.0001, # original was this line only  
           cumsum_area = cumsumSkipNA(ag_area),  
           prop_sums = cumsum_area / max(cumsum_area, na.rm = T))
```

```

return(dat)
}

```

Step 7 is written below.

```

cumsumSkipNA <- function(x) {
  FUNC <- match.fun('cumsum')
  x[!is.na(x)] <- FUNC(x[!is.na(x)])
  return(x)
}

step_7 <- function(dat) {

  tmp <- lwd[, 2:length(lwd)]
  tmp <- t(tmp)
  tmp <- apply(tmp, FUN = function(x) cumsumSkipNA(x), MARGIN = 2)
  tmp <- t(tmp)
  tmp <- as.data.frame(tmp)
  tmp$country <- lwd$country

  dat <- merge(dat, tmp, by = 'country')

  for (x in unique(dat$country)) {

    tmp <- dat[which(dat$country == x), ]

    tmp$farm <- ifelse(tmp$prop_sums <= tmp$`0_1`
                      & !is.na(tmp$`0_1`),
                      '0_1',
                      ifelse(tmp$prop_sums <= tmp$`1_2`
                              & !is.na(tmp$`1_2`),
                              '1_2',
                              ifelse(tmp$prop_sums <= tmp$`2_5`
                                      & !is.na(tmp$`2_5`),
                                      '2_5',
                                      ifelse(tmp$prop_sums <= tmp$`5_10`
                                              & !is.na(tmp$`5_10`),
                                              '5_10',
                                              ifelse(tmp$prop_sums <= tmp$`10_20`
                                                      & !is.na(tmp$`10_20`),
                                                      '10_20',
                                                      ifelse(tmp$prop_sums <= tmp$`20_50`
                                                              & !is.na(tmp$`20_50`),
                                                              '20_50',
                                                              ifelse(tmp$prop_sums <= tmp$`50_100`
                                                                      & !is.na(tmp$`50_100`),
                                                                      '50_100',
                                                                      ifelse(tmp$prop_sums <= tmp$`100_200`
                                                                              & !is.na(tmp$`100_200`),
                                                                              '100_200',
                                                                              ifelse(tmp$prop_sums <= tmp$`200_500`
                                                                                      & !is.na(tmp$`200_500`),
                                                                                      '200_500',
                                                                                      ifelse(tmp$prop_sums <= tmp$`500_1000`
                                                                                          & !is.na(tmp$`500_1000`),
                                                                                          '500_1000',
                                                                                          '500_1000')
                                                                              )
                                                                      )
                                                              )
                                                      )
                                              )
                                      )
                              )
                      )
  }
}

```

```

        & !is.na(tmp$`500_1000`),
        '500_1000',
    ifelse(tmp$prop_sums <= tmp$`1000_5000`
        & !is.na(tmp$`1000_5000`),
        '1000_5000',
        '1000_5000'
    )))))))

    if (x == as.character(unique(dat$country)[1])) {
      out <- tmp
    } else {
      out <- rbind(out, tmp)
    }
  }

  out <- out[,c(2,length(out))]

  dat <- merge(df, out,
              by = 'ID',
              all.x = T,
              all.y = T)

  return(dat)
}

```

Step 8 is written below. For countries not in Lower, we hard code the farm size matched to that field size class. We take the upper end of the given range and match it to Lesiv's codes (e.g., code 3502 becomes 200 ha)

- 1 - 3502 - Very large fields with an area of greater than 100 ha - 100-200 ha
- 2 - 3503 - Large fields with an area between 16 ha and 100 ha - 50-100 ha
- 3 - 3504 - Medium fields with an area between 2.56 ha and 16 ha - 5-10 ha
- 4 - 3505 - Small fields with an area between 0.64 ha and 2.56 ha - 1-2 ha
- 5 - 3506 - Very small fields with an area less than 0.64 ha - 0-1 ha
- NA - 3507 - no fields;

```

step_8 <- function(dat) {

  field_2_farm <- data.frame(
    lesiv = c(1,2,3,4,5),
    # lesiv = c(5,4,3,2,1),
    farm_avg = c('100_200',
                 '50_100',
                 '5_10',
                 '1_2',
                 '0_1'))

  dat <- merge(dat, field_2_farm, by = 'lesiv', all.x = T)

  # Make sure that any inf are NA
  dat2 <- dat %>% rationalize()
}

```

```

# Make sure farm and farm_avg cats are same type
dat$farm <- as.character(dat$farm)
dat$farm_avg <- as.character(dat$farm_avg)

dat$farm_final <- ifelse(dat$ag >= 0.01 & !is.na(dat$farm),
                        dat$farm,
                        ifelse(dat$ag >= 0.01 & !is.na(dat$lesiv),
                              dat$farm_avg, NA))

return(dat)
}

```

Step 9: Generate farm size raster

```

step_9 <- function(dat, x = 'farm_final') {
  # Places farm size values into a raster
  # Params:
  #   df
  # Return:
  #   raster object

  if (x == 'farm_final') {
    dat$rasterValues <- as.numeric(
      str_split_fixed(dat[[x]], '_', 2)[, 2])
  } else {
    dat$rasterValues <- dat[[x]]
  }

  # Clip to cropland area
  dat$rasterValues <- ifelse(is.na(dat$crop), NA, dat$rasterValues)

  dat <- dat[order(dat$ID), ]
  rast_final <- rast.all@layers[[1]]
  rast_final <- setValues(rast_final, dat$rasterValues, dat$ID)

  return(rast_final)
}

```

Step 10: Write raster

```

step_10 <- function(rast_final, i) {
  # Writes raster
  # Params:
  #   rast_final
  #   i: random number
  # Returns
  #   Writes raster to file

  writeRaster(rast_final,
              filename = paste0(
                'raster_out_forCleanRunTest/farmSize_',
                # 'raster_out_forPublication/farmSize_agarea_',
                format(Sys.time(), "%Y%m%d"), '_', i, '.tif'),
              format = 'GTiff',
              overwrite = TRUE)
}

```

```
}
```

5.2 Resampled Results

Step 11: Repeat steps 1-10 100x. Note, 100 runs is illustrative. For more stable results increase to i 1k.

```
# i = 1
# for (i in 1:100) { # use for full run
for (i in 1:2) {    # use for test run
  dat1 <- steps_1_6(df_lwd)
  dat2 <- step_7(dat1)
  dat3 <- step_8(dat2)
  rast <- step_9(dat3)
  step_10(rast, i)
}
```

5.3 Sanity Check

We sanity check that the proportional farm size distributions per country equals Lowders. We use one run as an example. As we can see the data is very close to the one-to-one line (but not exactly as would be expected based on small differences in row matching in our algorithm). Since the sampling scheme randomly shuffles the data, this plot will be slightly different for each sample.

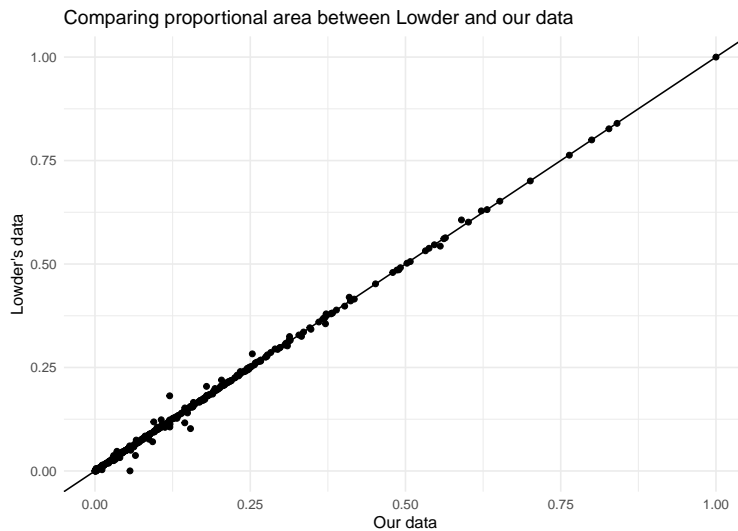


Figure 2: Sanity check that the summed farm size per country equals Lowders

5.4 Final Output

Here is the mode of the 100 farm size maps generated.

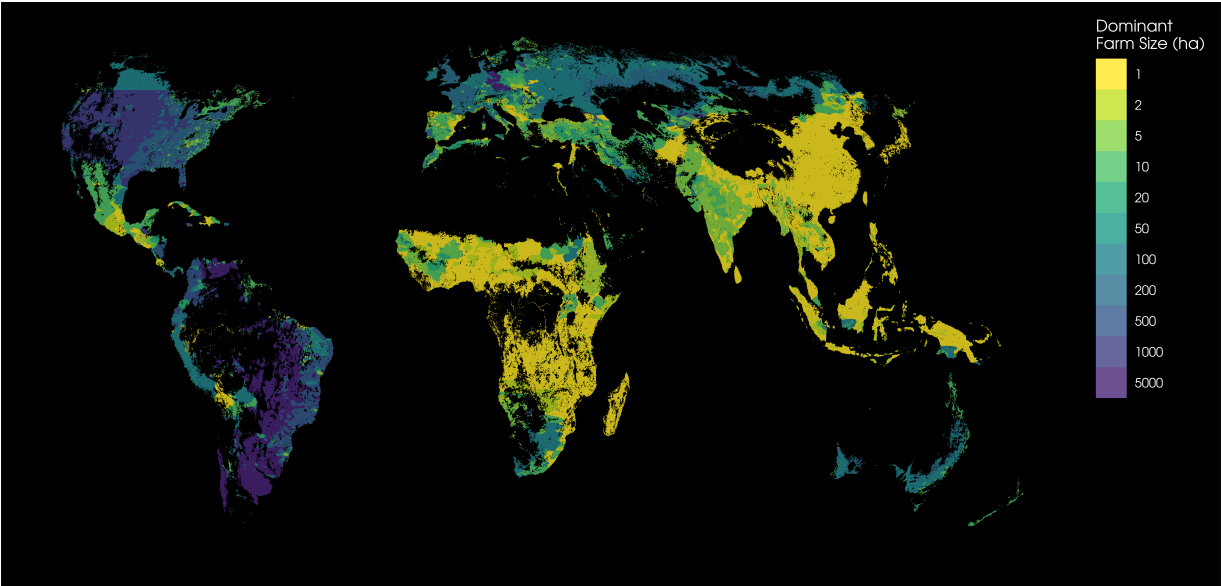


Figure 3: Final map, where dominant farm size (ha) is shown per pixel for the mode of the 100 generated maps.

Session information

```
sessionInfo()

R version 3.6.3 (2020-02-29)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 18.04.5 LTS

Matrix products: default
BLAS: /usr/lib/x86_64-linux-gnu/atlas/libblas.so.3.10.3
LAPACK: /usr/lib/x86_64-linux-gnu/atlas/liblapack.so.3.10.3

locale:
 [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
 [3] LC_TIME=en_US.UTF-8      LC_COLLATE=en_US.UTF-8
 [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
 [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
 [9] LC_ADDRESS=C             LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C

attached base packages:
[1] stats      graphics  grDevices datasets  utils      methods   base

other attached packages:
 [1] zoo_1.8-5      viridis_0.5.1   viridisLite_0.3.0
 [4] forcats_0.5.0 stringr_1.4.0   dplyr_0.8.0.1
 [7] purrr_0.3.4   readr_1.3.1    tidyr_0.8.3
[10] tibble_3.0.1  ggplot2_3.3.2  tidyverse_1.2.1
[13] sf_0.9-3      rworldmap_1.3-6 rgdal_1.4-4
[16] reshape2_1.4.3 raster_2.9-5   sp_1.3-1
[19] kknm_1.3.1    magrittr_1.5   lintr_2.0.1
[22] leaflet_2.0.3 hablar_0.3.0   gmodels_2.18.1
[25] ggthemes_4.2.0 foreach_1.4.4  data.table_1.12.2
[28] renv_0.12.0   formatR_1.6    printr_0.1
[31] knitr_1.22

loaded via a namespace (and not attached):
 [1] colorspace_1.4-1  ellipsis_0.3.1  class_7.3-17
 [4] rprojroot_1.3-2  rstudioapi_0.11 farver_2.0.3
 [7] remotes_2.1.1    fansi_0.4.1     lubridate_1.7.4
[10] xml2_1.2.2       codetools_0.2-16 spam_2.2-2
[13] jsonlite_1.7.0   broom_0.5.2     shiny_1.3.2
[16] compiler_3.6.3   httr_1.4.1      backports_1.1.8
[19] assertthat_0.2.1 Matrix_1.2-17   lazyeval_0.2.2
[22] cli_2.0.2        later_0.8.0     htmltools_0.3.6
[25] tools_3.6.3      igraph_1.2.5    dotCall64_1.0-0
[28] gtable_0.3.0     glue_1.4.1      maps_3.3.0
[31] Rcpp_1.0.5       cellranger_1.1.0 vctrs_0.3.1
[34] gdata_2.18.0     nlme_3.1-139    iterators_1.0.10
[37] crosstalk_1.0.0  xfun_0.6        ps_1.3.3
[40] rvest_0.3.4      mime_0.6        lifecycle_0.2.0
[43] gtools_3.8.2     MASS_7.3-51.4   scales_1.1.1
[46] hms_0.4.2        promises_1.0.1  rex_1.2.0
[49] fields_9.7       gridExtra_2.3   stringi_1.4.3
[52] highr_0.8        mapproj_0.9-5   desc_1.2.0
[55] e1071_1.7-3      cyclocomp_1.1.0 rlang_0.4.6
[58] pkgconfig_2.0.3  evaluate_0.14   lattice_0.20-38
[61] labeling_0.3     htmlwidgets_1.3 processx_3.4.2
[64] tidyselect_0.2.5 plyr_1.8.4      R6_2.4.1
[67] generics_0.0.2  DBI_1.0.0       pillar_1.4.4
[70] haven_2.1.0     foreign_0.8-71  withr_2.2.0
[73] units_0.6-6     modelr_0.1.4    crayon_1.3.4
[76] KernSmooth_2.23-15 grid_3.6.3      readxl_1.3.1
[79] callr_3.4.3     digest_0.6.25  classInt_0.4-3
[82] xtable_1.8-4    httpuv_1.5.1    munsell_0.5.0
```

References

- [1] K. Ushey, *Renv: Project environments*, R package version 0.12.0, 2020. [Online]. Available: <https://CRAN.R-project.org/package=renv>.
- [2] M. Dowle and A. Srinivasan, *Data.table: Extension of 'data.frame'*, R package version 1.12.2, 2019. [Online]. Available: <https://CRAN.R-project.org/package=data.table>.
- [3] Revolution Analytics and S. Weston, *foreach: Provides foreach looping construct for r*, R package version 1.4.4, 2017. [Online]. Available: <https://CRAN.R-project.org/package=foreach>.
- [4] Y. Xie, *Formatr: Format r code automatically*, R package version 1.6, 2019. [Online]. Available: <https://CRAN.R-project.org/package=formatR>.
- [5] G. R. W. and Ben Bolker, T. Lumley, R. C. Johnson, and R. C. Johnson, *Gmodels: Various r programming tools for model fitting*, R package version 2.18.1, 2018. [Online]. Available: <https://CRAN.R-project.org/package=gmodels>.
- [6] D. Sjoberg, *Hablar: Non-astonishing results in r*, R package version 0.3.0, 2020. [Online]. Available: <https://CRAN.R-project.org/package=hablar>.
- [7] J. Cheng, B. Karambelkar, and Y. Xie, *Leaflet: Create interactive web maps with the javascript 'leaflet' library*, R package version 2.0.3, 2019. [Online]. Available: <https://CRAN.R-project.org/package=leaflet>.
- [8] J. Hester, F. Angly, and R. Hyde, *LintR: A 'linter' for r code*, R package version 2.0.1, 2020. [Online]. Available: <https://CRAN.R-project.org/package=lintr>.
- [9] S. M. Bache and H. Wickham, *Magrittr: A forward-pipe operator for r*, R package version 1.5, 2014. [Online]. Available: <https://CRAN.R-project.org/package=magrittr>.
- [10] Y. Xie, *Knitr: A general-purpose package for dynamic report generation in r*, R package version 1.22, 2019. [Online]. Available: <https://CRAN.R-project.org/package=knitr>.
- [11] K. Schliep and K. Hechenbichler, *Kknn: Weighted k-nearest neighbors*, R package version 1.3.1, 2016. [Online]. Available: <https://CRAN.R-project.org/package=kknn>.
- [12] H. Wickham, *PlyR: Tools for splitting, applying and combining data*, R package version 1.8.4, 2016. [Online]. Available: <https://CRAN.R-project.org/package=plyr>.
- [13] R. J. Hijmans, *Raster: Geographic data analysis and modeling*, R package version 2.9-5, 2019. [Online]. Available: <https://CRAN.R-project.org/package=raster>.
- [14] H. Wickham, *Reshape2: Flexibly reshape data: A reboot of the reshape package*, R package version 1.4.3, 2017. [Online]. Available: <https://CRAN.R-project.org/package=reshape2>.
- [15] R. Bivand, T. Keitt, and B. Rowlingson, *Rgdal: Bindings for the 'geospatial' data abstraction library*, R package version 1.4-4, 2019. [Online]. Available: <https://CRAN.R-project.org/package=rgdal>.
- [16] A. South, *Rworldmap: Mapping global data*, R package version 1.3-6, 2016. [Online]. Available: <https://CRAN.R-project.org/package=rworldmap>.
- [17] E. Pebesma, *Sf: Simple features for r*, R package version 0.9-3, 2020. [Online]. Available: <https://CRAN.R-project.org/package=sf>.
- [18] H. Wickham, *Tidyverse: Easily install and load the 'tidyverse'*, R package version 1.2.1, 2017. [Online]. Available: <https://CRAN.R-project.org/package=tidyverse>.
- [19] A. Zeileis, G. Grothendieck, and J. A. Ryan, *Zoo: S3 infrastructure for regular and irregular time series (z's ordered observations)*, R package version 1.8-5, 2019. [Online]. Available: <https://CRAN.R-project.org/package=zoo>.
- [20] S. K. Lowder, J. Scoet, and T. Raney, "The number, size, and distribution of farms, smallholder farms, and family farms worldwide," *World Development*, vol. 87, pp. 16–29, 2016.
- [21] M. Lesiv, J. C. Laso Bayas, L. See, M. Duerauer, D. Dahlia, N. Durando, R. Hazarika, P. Kumar Sahariah, M. Vakolyuk, V. Blyshchyk, A. Bilous, A. Perez-Hoyos, S. Gengler, R. Prestele, S. Bilous, I. u. H. Akhtar, K. Singha, S. B. Choudhury, T. Chetri, i. Malek, K. Bungnamei, A. Saikia, D. Sahariah, W. Narzary, O. Danylo, T. Sturn, M. Karner, I. McCallum, D. Schepaschenko, E. Moltchanova, D. Fraisl, I. Moorthy, and S. Fritz, "Estimating the global distribution of field size using crowdsourcing," *Global Change Biology*, vol. 25, no. 1, pp. 174–186, 2019. DOI: 10.1111/gcb.14492. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/gcb.14492>. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1111/gcb.14492>.

- [22] N. Ramankutty, A. T. Evan, C. Monfreda, and J. A. Foley, “Farming the planet : 1 . Geographic distribution of global agricultural lands in the year 2000,” vol. 22, no. August 2007, pp. 1–19, 2008. DOI: 10.1029/2007GB002952.