

ArchaeoGLOBE trend analysis

Nick Gauthier

Last knit on: 30 January, 2019

Sample analysis code for the ArchaeoGlobe database. Here we use Generalized Additive Models (GAMs), a flexible form of nonlinear regression model capable of fitting smooth, time-varying trends to the ordered categorical ArchaeoGLOBE response data.

We model ordered categorical data using a latent variable following a logistic distribution. The model identifies a series of cut points, which correspond to the probabilities of the latent variable falling within each of our categories.

We fit two sets of trends. One trend is fitted to all the data simultaneously, representing the global trend across all archaeological regions. Then we fit region-level trends, which represent the deviation of each region from the global trend. By penalizing the “wiggleness” of the trend lines, we allow regional trends that don’t significantly deviate from the global trend to be penalized to 0, effectively reducing that particular region to the global trend. This is a form of partial pooling, allowing the model to share information between groups and in so doing make the results less sensitive to regions with exceptionally low response rates.

After fitting the model, we can extract the region-specific deviations from the global trend, use a k-means clustering algorithm to group together regions with similar trends, and map the results. We repeat this analysis for both self-reported expertise and perceived data quality.

Setup

Import packages needed for analysis. We’ll use packages from the `tidyverse`, such as `readr`, `dplyr`, and `ggplot2` for data import, processing, and plotting. We’ll also use `mgcv` for fitting nonlinear trends to the data. We’ll use the `sf` package to help us plot shapefiles in a tidy context. Finally, we’ll use `patchwork` to combine multiple ggplots in the same image.

```
library(tidyverse)
library(mgcv)
library(sf)
library(ggplot2)

#install patchwork from github
#devtools::install_github('thomasp85/patchwork')
library(patchwork)
```

Data import

Read in the latest version of the ArchaeoGLOBE database and the regions’ shapefile from the Dataverse repository.

```
library("dataverse")
Sys.setenv("DATAVERSE_SERVER" = "dataverse.harvard.edu")

# get data frame of files on dataverse
ArchaeoGLOBE_Public_Data_DOI <-
  "doi:10.7910/DVN/CNCANQ"
```

```

ArchaeoGLOBE_Public_Data_df <-
  get_dataset(ArchaeoGLOBE_Public_Data_DOI)

# Only download the file we need here
ArchaeoGLOBE_Public_Data_df_files <-
  ArchaeoGLOBE_Public_Data_df$files[grepl("ARCHAEOGLOBE_PUBLIC_DATA|ARCHAEOGLOBE_CONSENSUS_ASSESSMENT",
                                           ArchaeoGLOBE_Public_Data_df$files$filename), ]

# read into local dir
walk(ArchaeoGLOBE_Public_Data_df_files$label,
      ~get_file(.x, ArchaeoGLOBE_Public_Data_DOI) %>%
        writeBin(paste0('data/raw-data/', .x)))

# read into the current environment
archaeoglobe <- read_csv('data/raw-data/ARCHAEOGLOBE_PUBLIC_DATA.tab')
consensus <- read_csv('data/raw-data/ARCHAEOGLOBE_CONSENSUS_ASSESSMENT.tab')

# repeat for shapefile
ArchaeoGLOBE_Regions_DOI <-
  "doi:10.7910/DVN/CQWUBI"

# get data frame of files on DV
ArchaeoGLOBE_Regions_df <-
  get_dataset(ArchaeoGLOBE_Regions_DOI)

# just download the shapefile we want
ArchaeoGLOBE_Regions_df_files <- ArchaeoGLOBE_Regions_df$files[ArchaeoGLOBE_Regions_df$files$filename ==

# read into local dir
walk(ArchaeoGLOBE_Regions_df_files$label,
      ~get_file(.x, ArchaeoGLOBE_Regions_DOI) %>%
        writeBin(paste0('data/raw-data/', .x)))

unzip('data/raw-data/ArchaeGLOBE_Regions.zip',
      overwrite = TRUE,
      exdir = 'data/raw-data/ArchaeGLOBE_Regions')

# read into the current environment, and simplify the polygons for faster plotting
regions_unsimplified <-
  st_read('data/raw-data/ArchaeGLOBE_Regions/ArchaeGLOBE_Regions.shp',
          quiet = TRUE)
regions <- rmapshaper::ms_simplify(regions_unsimplified)

```

Exploratory plots

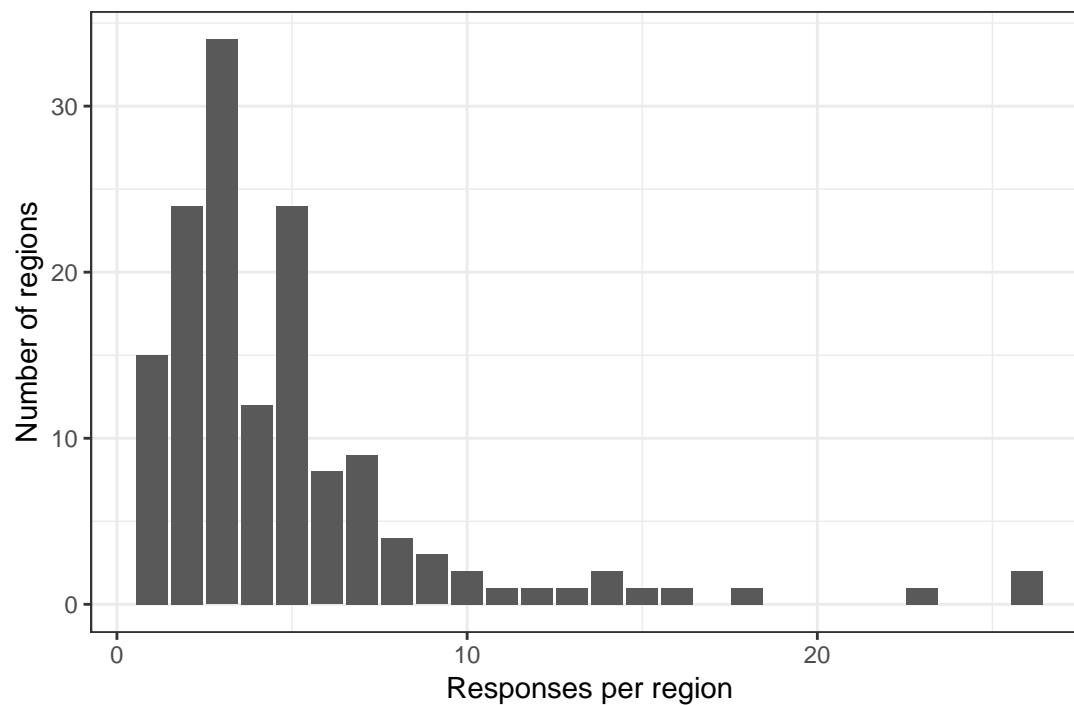
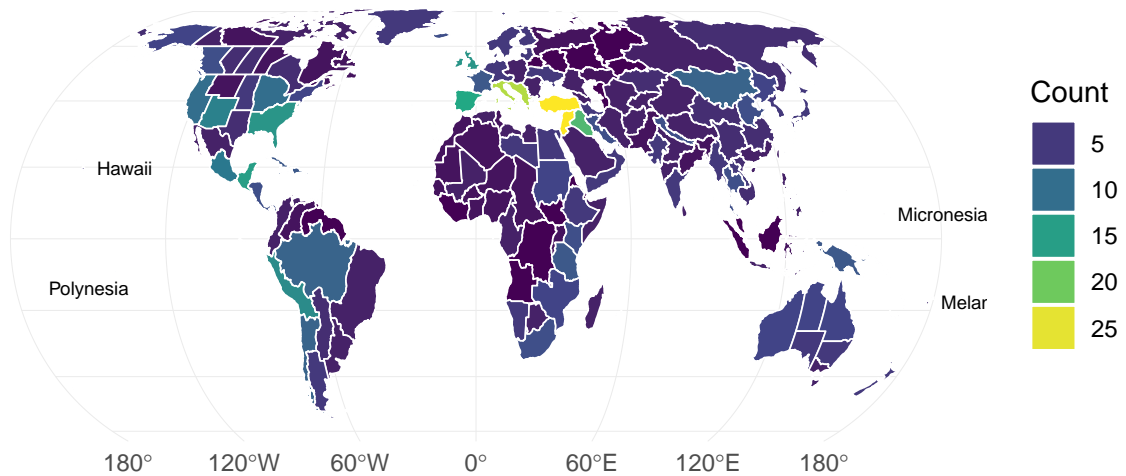
Before running any analyses, let's look at the data. How many responses do we have per region?

```

response_counts <- archaeoglobe %>%
  group_by(REGION_ID) %>%
  count

```

Total responses per region



Here is the cumulative summary of regions per land-use category based on consensus assessments (Common > 1% to 20% regional land area; Widespread > 20% regional land area).

```
cumsum_landuse_regions <-  
consensus %>%
```

```

select(Region, FHG_10KBP:URBAN_1850CE) %>%
gather(variable, value, - Region) %>%
filter(value %in% c("Widespread", "Common", "Split", "Present")) %>%
separate(variable, into = c("land_use_category", "years_BP"), sep = "_") %>%
mutate(land_use_category = ifelse(str_detect(land_use_category, "AGR"),
                                str_replace_all(land_use_category,
                                                  "AGR",
                                                  "AG"),
                                land_use_category)) %>%
mutate(land_use_category = case_when(
  land_use_category == "FHG" ~ "Foraging",
  land_use_category == "EXAG" ~ "Extensive Agriculture",
  land_use_category == "INAG" ~ "Intensive Agriculture",
  land_use_category == "PAS" ~ "Pastoralism",
  land_use_category == "URBAN" ~ "Urban Centers")
) %>%
mutate(years_BP = ifelse(str_detect(tolower(years_BP), "kbp"),
                        -parse_number(years_BP) * 1000,
                        ifelse(str_detect(tolower(years_BP), "ce"),
                              parse_number(years_BP),
                              -parse_number(years_BP)))) %>%
unite(land_use_category_consensus_assessments,
      c('land_use_category',
        'value'),
      sep = " ") %>%
complete(land_use_category_consensus_assessments,
          nesting(years_BP)) %>%
group_by(land_use_category_consensus_assessments,
          years_BP) %>%
summarise(n = n()) %>%
mutate(perc = n / sum(n) * 100)

# make years label and breaks
years_label <- unique(ifelse(cumsum_landuse_regions$years_BP < 0,
                             str_glue('{-cumsum_landuse_regions$years_BP} BP'),
                             str_glue('{cumsum_landuse_regions$years_BP} CE')
))

years_breaks <- unique(cumsum_landuse_regions$years_BP)

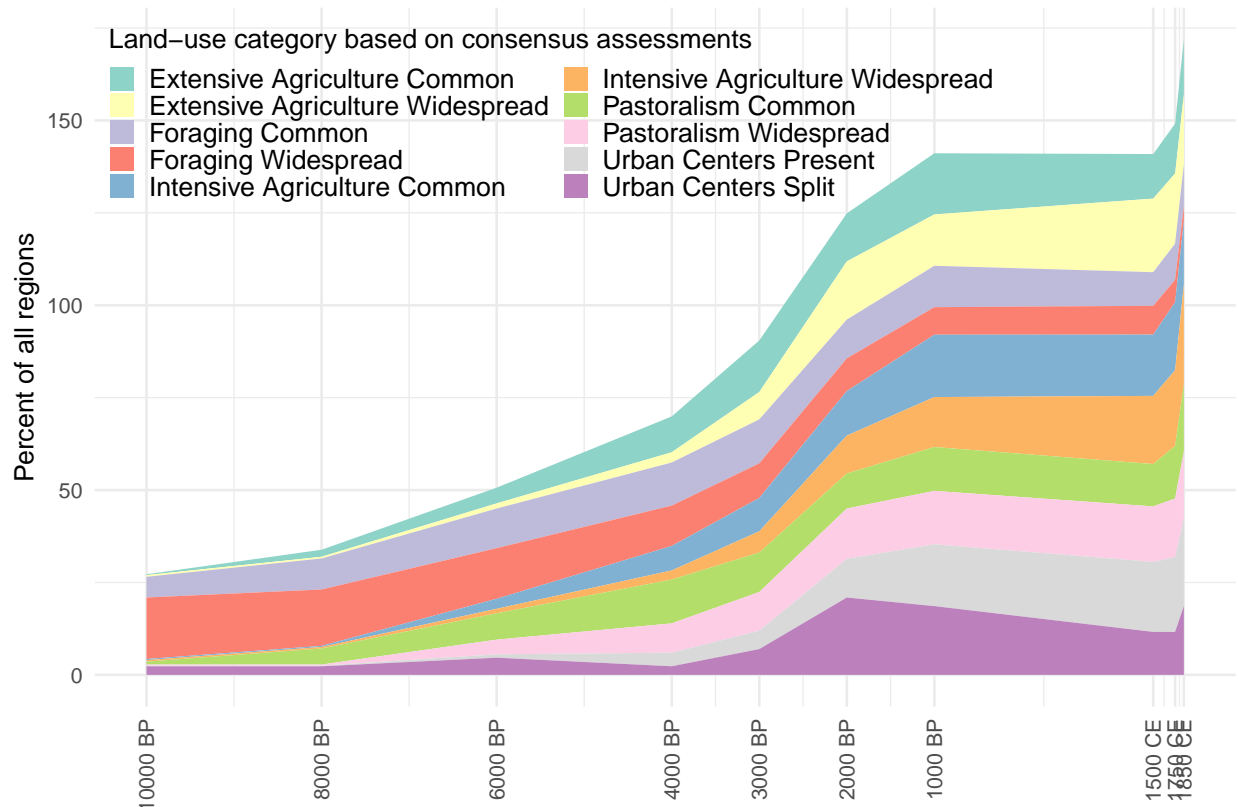
# draw the plot
ggplot(
  cumsum_landuse_regions,
  aes(years_BP,
      perc,
      fill = land_use_category_consensus_assessments)) +
geom_area(position = 'stack') +
scale_fill_brewer(palette = "Set3") +
scale_x_continuous(labels = years_label,
                   breaks = years_breaks) +
theme_minimal(base_size = 10) +
theme(
  axis.text.x = element_text(
    angle = 90,

```

```

    hjust = 1,
    vjust = 0.5),
  # x- and y- offsets from the bottom-left of the plot, ranging from 0 - 1.
  legend.position = c(0.4, 0.85),
  legend.text = element_text(size = 10),
  legend.key.size = unit(0.7, "line")) +
  guides(fill = guide_legend(ncol = 2,
                             title = "Land-use category based on consensus assessments")) +
  xlab("") +
  ylab("Percent of all regions")

```



```

ggsave('figures/cumulative_sum_land_use.png', height = 5, width = 7)

```

Analysis functions

Define some analysis functions that we'll be using repeatedly in the analysis, so that we don't have to keep copying and pasting the same lines of code.

This function subsets the data to highlight a variable of interest, and converts it from a wide to a long "tidy" format to make analysis and plotting easier.

```

preprocess <- function(prefix, categories){
  archaeoglobe %>% # start with the full ArcheoGlobe data
  # drop columns not related to the variable of interest
  select(c(CONTRIBUTR:LAND_AREA, starts_with(prefix))) %>%

```

```

gather(time, value, starts_with(prefix)) %>% # one value per row
mutate(time = parse_number(time) * -1, # convert time period labels to years
       value = ordered(value, levels = categories),
       cat_num = as.numeric(value)) %>%
mutate_if(is.character, as.factor) # convert characters to factors
}

```

This function takes a data frame produced by the above function and fits GAM to the global trend and local deviations for each region, accounting for inter-observer variability. This function takes as arguments a preprocessed data frame containing time slices, regions, contributors, and the ordered categorical response variable transformed to a numeric vector.

```

cores <- max(parallel::detectCores() / 2, 1) # physical cores for parallelization
cl <- parallel::makeCluster(cores)

fit_gam <- function(x, n_cats){
  bam(cat_num ~
    # this spline is for the global trend
    s(time, bs = 'cr', m = 2) +
    # region-specific trends. bs = 'ts' and m = 1
    # help penalize deviation from the global model
    s(time, by = REGION_LAB, bs = 'cs', m = 1) +
    # add back in region-specific intercepts
    REGION_LAB +
    # model contributor as a random effect
    s(CONTRIBUTR, bs = 're'),
    data = x, # data frame to analyze
    family = ocat(R = n_cats), # ordered categorical with n levels
    # final 3 arguments just speed up the model fitting
    method = 'fREML',
    discrete = TRUE,
    cluster = cl)
}

```

This function extracts the estimated trends for each region, incorporating the global and regional splines as well as the region and contributor specific intercepts. Then it clusters these trends into 6 discrete clusters.

```

extract_trends <-function(mod, n_clusters = 6){
  set.seed(1000) # set seed for reproducibility of clusters
  archaeoglobe %>% # create dummy data for prediction in the following lines
  select(REGION_LAB) %>%
  group_by(REGION_LAB) %>%
  slice(1) %>%
  slice(rep(1:n(), each = 198)) %>%
  ungroup %>%
  mutate(time = rep_len(seq(-10000, -150, 50), n()),
         CONTRIBUTR = 'CYRBU') %>% # select an arbitrary contributor
  mutate(preds = predict(mod, .)) %>% # estimate trend lines
  mutate(preds = plogis(preds)) %>% # transform responses to [0,1] scale
  spread(time, preds) %>%
  # next is the actual kmeans clustering code
  mutate(cluster = kmeans(., -c(1,2), n_clusters, iter.max = 100, nstart = 100)$cluster)
}

```

Analysis

Now we use the functions defined above on the ArchaeoGlobe data. For convenience, first define a data frame that lists the prefixes of the variables we are interested in (e.g. “EXP” for expertise) and the levels of the ordered factors associated with each variable. This will make it easier to quickly focus on a specific variable. The `tribble` command is simply a way to make a data frame by row rather than column, which makes the code easier to read.

```
response_levels <- tribble(
  ~prefix, ~categories,
  'EXP', c('None', 'Low', 'High'),          # Expertise
  'DQ', c('Unknown', 'Low', 'Moderate', 'Good'), # Data Quality
  'HUNT', c('none', 'minimal (<1%)', 'common (1-20%)', 'widespread (>20%)'),
  'EXAG', c('none', 'minimal (<1%)', 'common (1-20%)', 'widespread (>20%)'),
  'INAG', c('none', 'minimal (<1%)', 'common (1-20%)', 'widespread (>20%)'),
  'PAST', c('none', 'minimal (<1%)', 'common (1-20%)', 'widespread (>20%)'),
  'URBN', c('Absent', 'Present')
)
```

Now map each of the above functions to each variable. This allows us to run the analysis for all variables of interest in a single step, and save all the outputs in a tibble format for easy plotting. If you’re running this for the first time, it should take about 40 minutes to run on a Intel NUC with a 5th-gen Intel Core i7-5557U processor and 16gb of RAM running Linux.

```
# A simple type of caching...
# Do we want to run the modelling code, or
# load a previously saved result from disk, or
# download a previously saved result from a repository?
# default is not run, then check if there is a saved file, and use that, or download

# devtools::install_github('centerforopenscience/osfr')
library(osfr)

rerun_time_consuming_analysis <- FALSE # FALSE means do not run the modelling code when knitting

if(rerun_time_consuming_analysis) {
  message("running the modelling code, this may take 30-50 min...")
  # go to the next chunk of code
} else {
  # check if there is a local file and if so, load it
  if(file.exists('data/derived-data/trend_dat.rda')) {
    # the file exists on the local disk, so just read it in
    message("Loading previously saved model results from disk...")
    trend_dat <- readRDS('data/derived-data/trend_dat.rda')
  } else {
    # we don't want to run the modelling code, and the result don't exist locally,
    # so download
    message("Downloading previously saved model results, takes 2-3 min...")
    trend_dat <- osf_retrieve_file("kcr2e") %>% osf_download('data/derived-data/trend_dat.rda')
    message("Loading the data downloaded from osf.io...")
    trend_dat <- readRDS('data/derived-data/trend_dat.rda')
    writeLines(paste0('trend_dat.rda downloaded from https://osf.io/kcr2e/ on ', Sys.Date()), con = 'data/')
    message("Done.")
  }
}
```

```

    }
  }

trend_dat <- response_levels %>%
  mutate(data = map2(prefix, categories, ~preprocess(.x,.y)),
         n_cats = map_dbl(categories, length),
         mod = map2(data, n_cats, fit_gam),
         trends = map(mod, extract_trends))

# save to disk
saveRDS(trend_dat,
        file = "data/derived-data/trend_dat.rda")

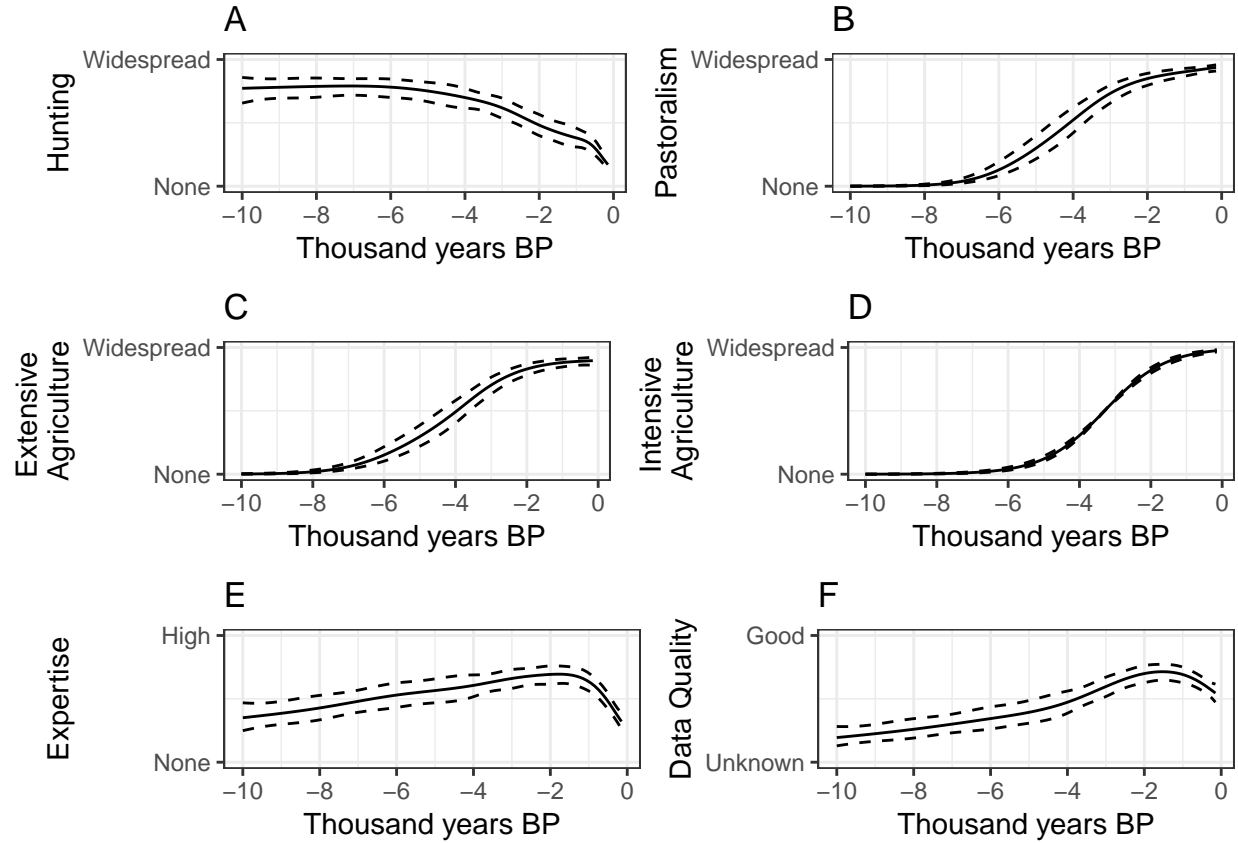
# write a note to indicate the provenance of this file
this_commit <- git2r::revparse_single(git2r::repository('.'), "HEAD")
writeLines(paste0('trend_dat generated on ',
                  Sys.Date(), ' from archaeoglobe.Rmd at git commit ',
                  this_commit$sha, ' made by ',
                  this_commit$author$name, ' on ',
                  this_commit$author$when, " with the message '",
                  this_commit$summary, "'"),
          con = 'data/derived-data/README.md',
          sep = '')
message("Done.")

```

Results

Global Trends

First we plot out the global trends for each land use type. Please refer to the source .rmd file for the plotting code.



Hunting

The global trend in hunting shows constant high prevalence until around 6,000 years ago, after which there is a smooth decline until the present day when it is very rare. Mapping out the clusters reveals a clear east-west divide, which regions in Afro-eurasia seeing hunting earlier than the global mean, and regions in the Americas and Oceania seeing later peaks in hunting.

Extensive Agriculture

The global trends in the prevalence of pastoralism, extensive and intensive agriculture, and urbanism all follow a sigmoidal curve, which means the trend is linear on the scale of the linear predictor (the ordered categorical GAM uses a logit transform as a latent link function). This means that there is a simple increase in the probability of each land use type being prevalent over time.

Regional land–use trends

Hunting

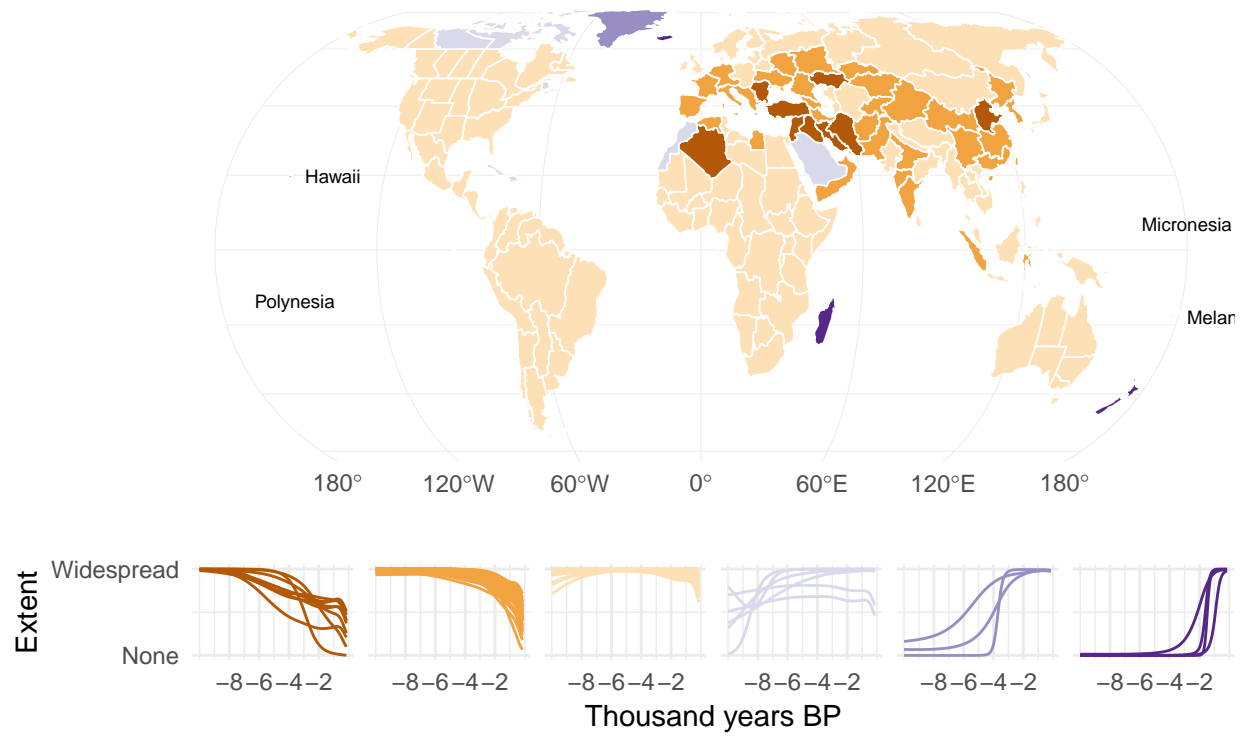
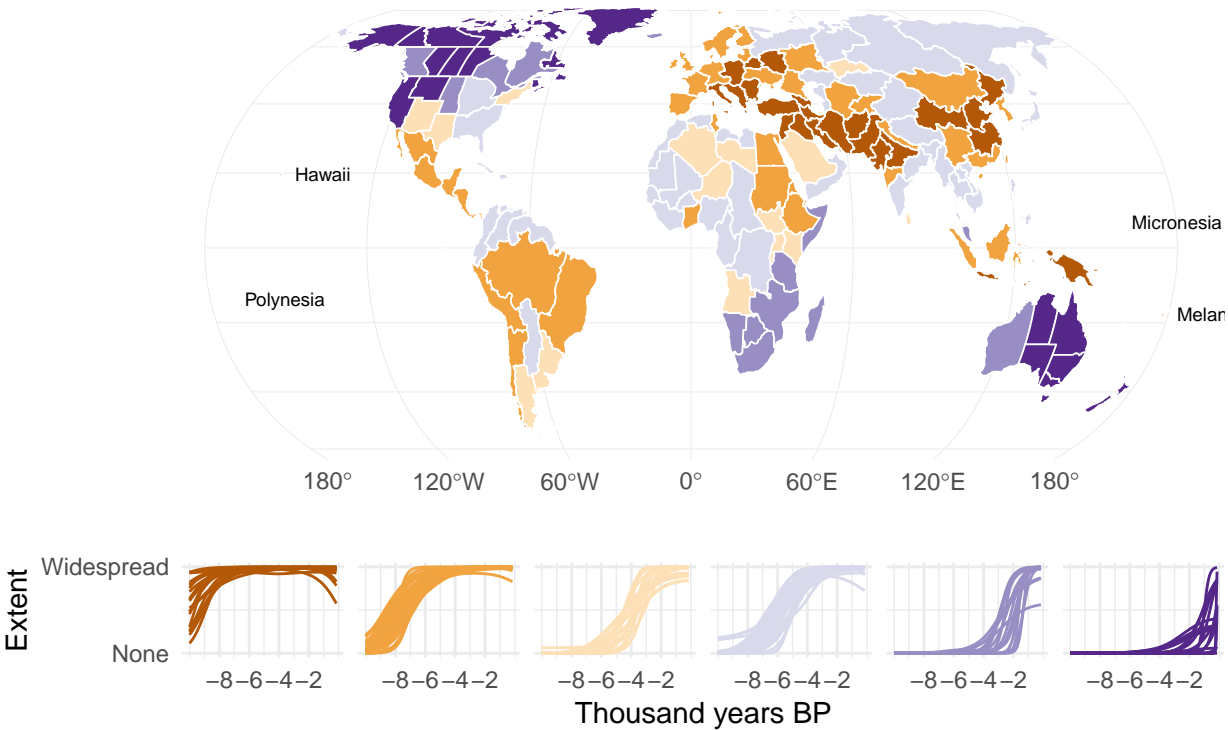


Figure 1: Regional trends in the areal extent of hunting. (A) Global trend (all regions) with 95% confidence interval. (B) Regional deviations from global trend, clustered via k-means. (C) Map of the local deviations from the global trend, same clusters as in B.

Regional land–use trends
Extensive Agriculture

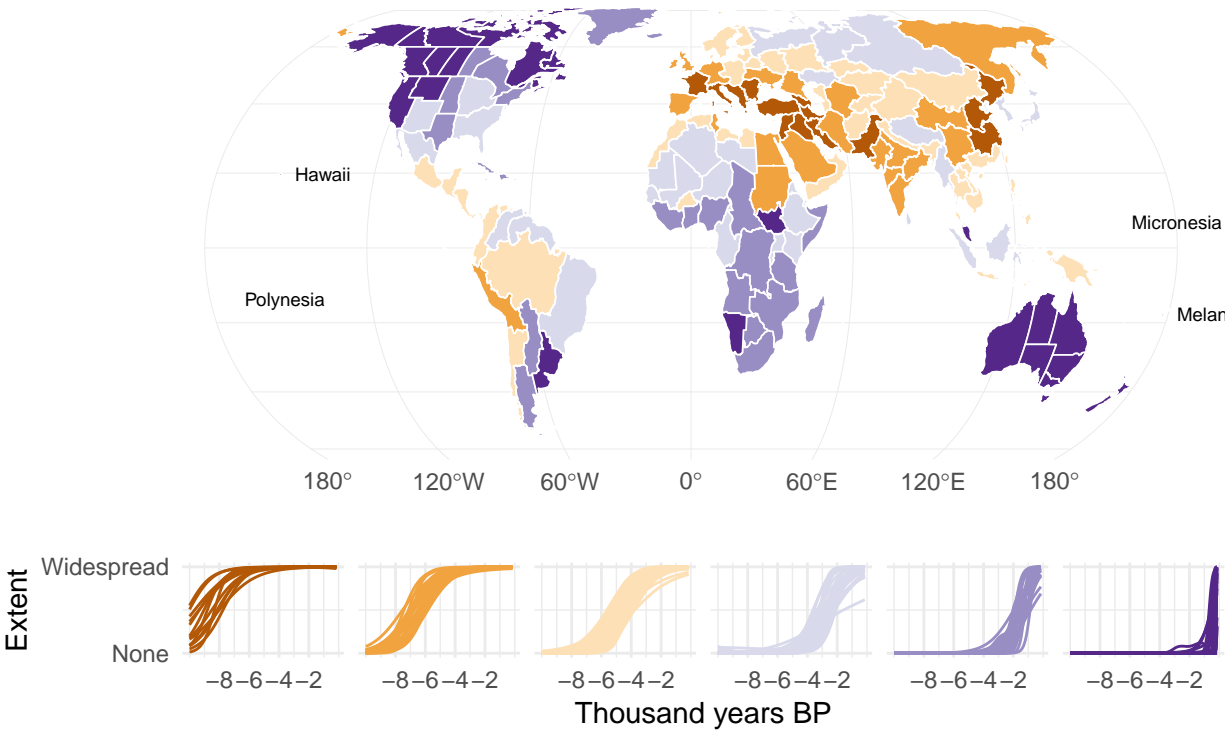


Intensive Agriculture

See above.

Regional land–use trends

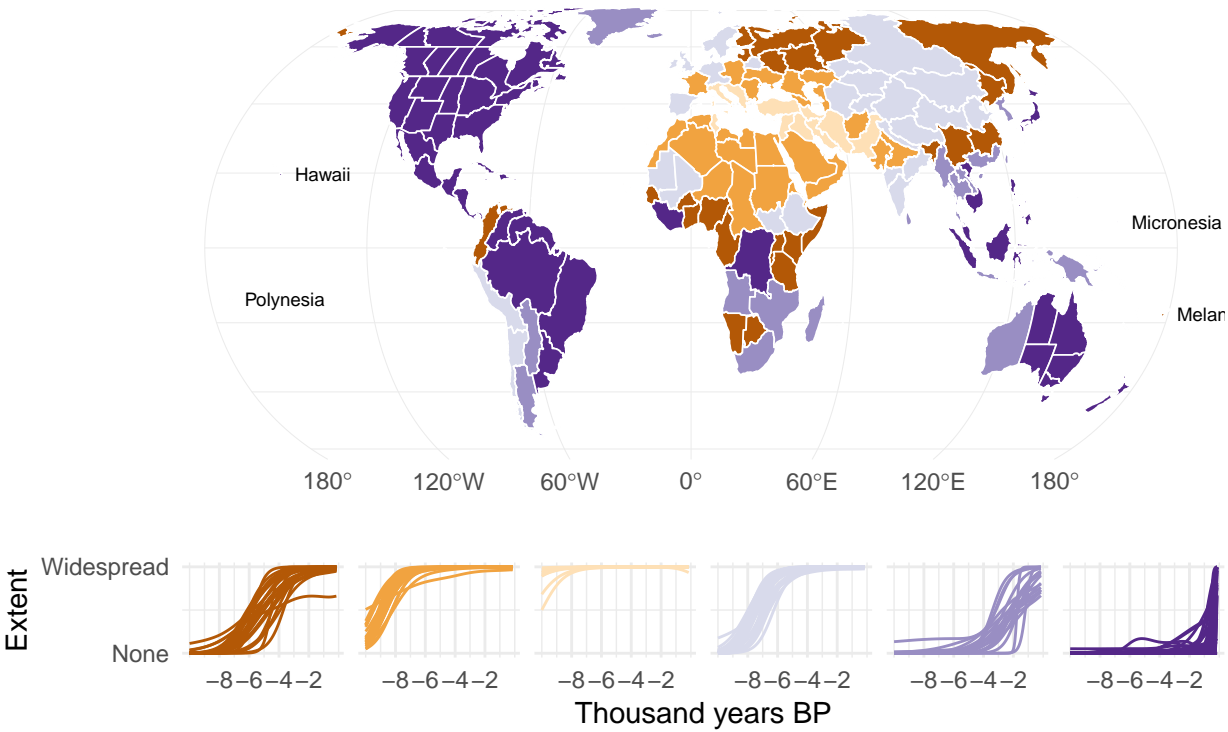
Intensive Agriculture



Pastoralism

See above.

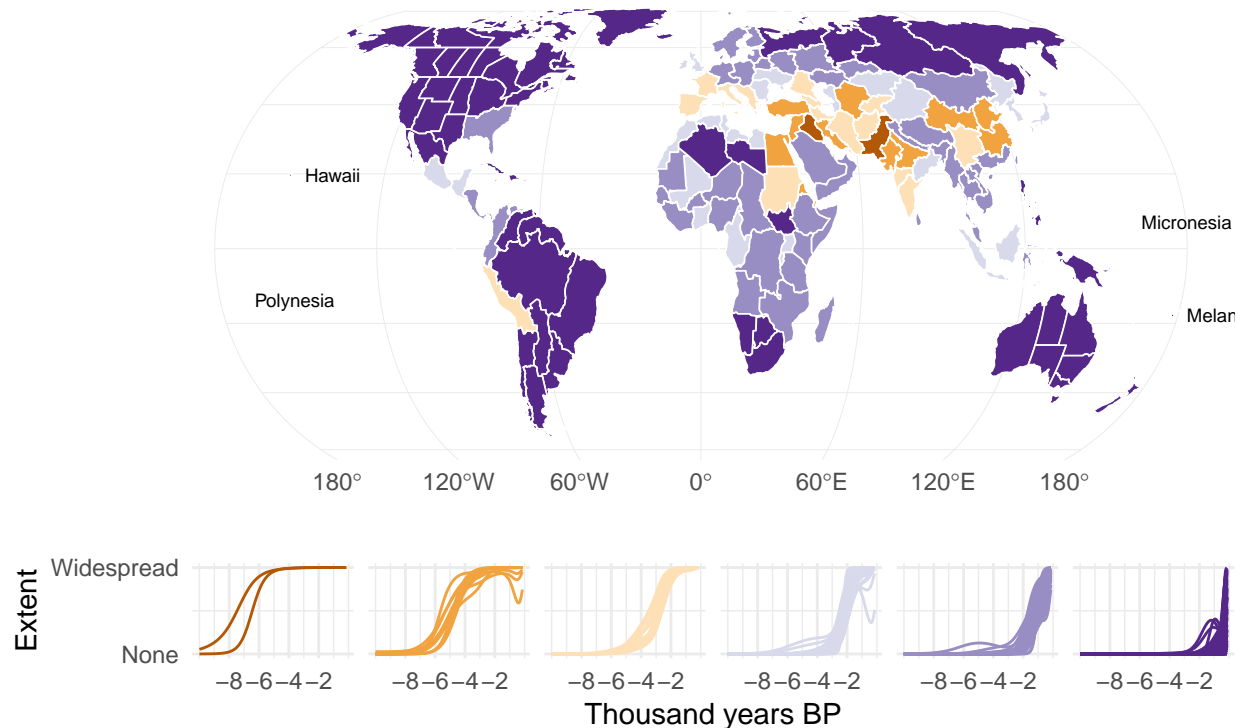
Regional land–use trends
Pastoralism



Urbanism

Regional land-use trends

Urbanism



See above.

Expertise and Data Quality

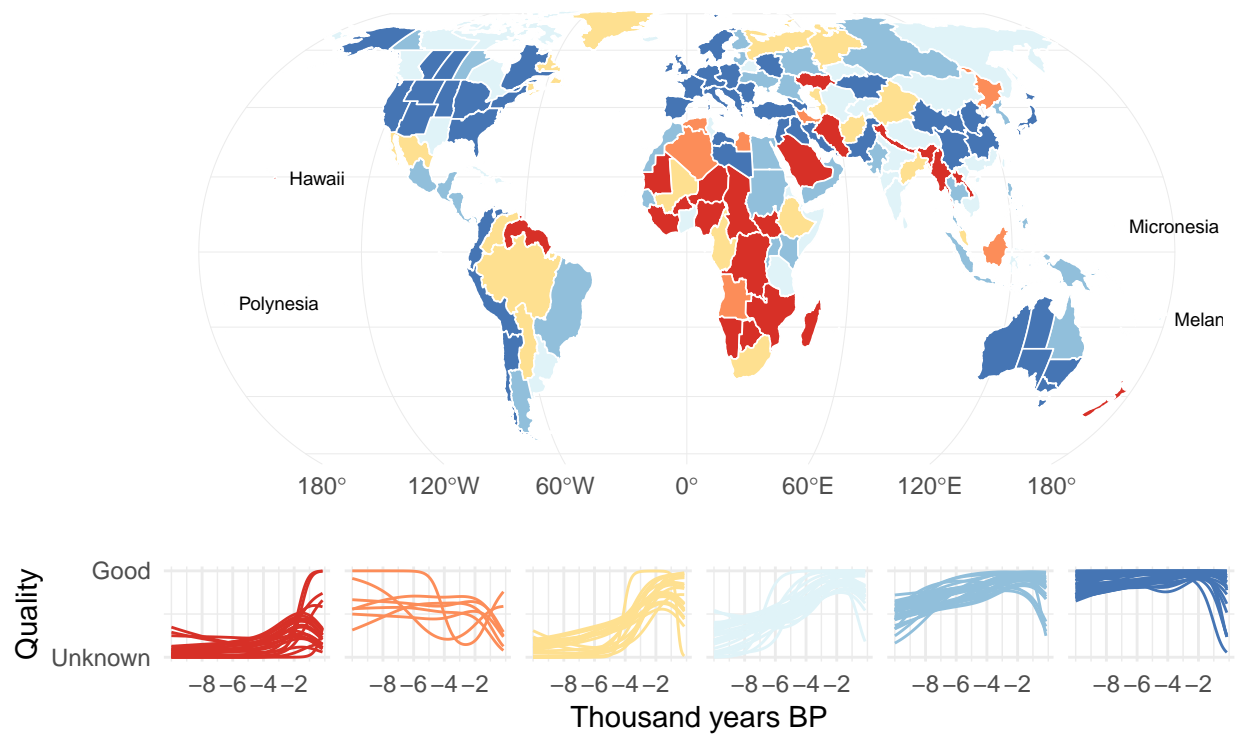
How does self-professed level of expertise vary in each region over time? The global trend is a roughly linear increase in self-reported expertise from 10ka BP up to 2ka BP, then a falloff continuing to the present day. The present day expertise values are approximately the same as at 10ka BP. This makes sense, as it points to both the increased frequency of preserved archaeological materials with time as well as the reduction in archaeological attention in periods with extensive historical records.

Now we cluster together the local deviations from the global trend using a k-means algorithm. The selection of 6 clusters is somewhat arbitrary, and is made simply based on visual comparisons of different cluster solutions with the goal making the results visually interpretable. The trajectories in these clusters are deviations from the global trend, so a horizontal line would indicate no deviation from the global trend.

The global trend in data quality is more or less the same as the expertise data, with the peak in data quality occurring more recently than for expertise and with a less dramatic falloff leading to the present day. Unlike expertise, which reaches the same values at 10ky BP and present, data quality in the present day remains high in spite of the falloff in the last 2 millennia. Also note the confidence interval for the global trend is generally wider than for the expertise responses.

Archaeological trends

Data Quality



Archaeological trends

Expertise

