# Supplementary material for *Using citizen science to parse climatic and landcover influences on bird occupancy within a tropical biodiversity hotspot*

Vijay Ramesh      Pratik R. Gupte      Morgan W. Tingley      VV Robin      Ruth DeFries

2020-12-23

# Contents

# 1 Introduction

This is supplementary material for a project in preparation that models occupancy for birds in the southern Western Ghats, India. The main project can be found here: https://github.com/pratikunterwegs/eBirdOccupancy.

## 1.1 Attribution

Please contact the following in case of interest in the project.

- Vijay Ramesh (lead author)
    - PhD student, Columbia University
- Pratik Gupte (repo maintainer)
    - PhD student, University of Groningen

# 2 Predicting Species-specific Occupancy

This supplement plots species-specific probabilities of occupancy as a function of significant environmental predictors.

## 2.1 Prepare libraries

```r
# to load data
library(readxl)

# to handle data
library(dplyr)
library(readr)
library(forcats)
library(tidyr)
library(purrr)
library(stringr)

# plotting
library(ggplot2)
library(patchwork)
```

## 2.2 Read data

```
# read data
data <- read_csv("data/results/data_occupancy_predictors.csv")
```

```
# drop na
data <- select(
  data,
  -ci
) %>%
  drop_na() %>%
  nest(data = c(predictor, m_group, seq_x, mean, scale))
```

Figure code is hidden in versions rendered as HTML and PDF. Example output is shown below.

**Figure here**

# 3 Selecting species of interest

This script shows the proportion of checklists that report a particular species across every 25km by 25km grid across the Nilgiris and the Anamalais. Using this analysis, we arrived at a final list of species for occupancy modeling.

We derived this list from inclusion criteria adapted from the State of India's Birds 2020 (Viswanathan et al., 2020). Initially, we considered all 561 species in eBird that occurred within the outlines of our study area. We then considered only those species that had a minimum of 1000 detections each between 2013 and 2019 (reducing to 303 species). Next, the study area was divided into 25 x 25 km cells following (Viswanathan et al., 2020). We then kept only those species that occurred in at least 5% of all checklists across 50% of the 25 x 25 km cells from where they have been reported (reducing to 93 species). We used the above criteria to ensure as much uniform sampling of a species as possible across our study area and to reduce any erroneous associations between environmental drivers and species occupancy. Across our final list of 93 species, we analyzed a total of ~3.2 million detections (presences) between 2013 and 2019.

## 3.1 Prepare libraries

```
# load libraries
library(data.table)
library(readxl)
library(magrittr)
library(stringr)
library(dplyr)
library(tidyr)
library(readr)

library(ggplot2)
library(ggthemes)
library(scico)

# round any function
round_any <- function(x, accuracy = 25000) {
  round(x / accuracy) * accuracy
}
```

## 3.2 Read species of interest

```
1  # add species of interest
2  specieslist <- read.csv("data/species_list.csv")
3  speciesOfInterest <- specieslist$scientific_name
```

## 3.3   Load raw data for locations

```
1   # read in shapefile of the study area to subset by bounding box
2   library(sf)
3   wg <- st_read("data/spatial/hillsShapefile/Nil_Ana_Pal.shp")
4   box <- st_bbox(wg)
5
6   # read in data and subset
7   ebd <- fread("data/01_ebird-filtered-EBD-westernGhats.txt")
8   ebd <- ebd[between(LONGITUDE, box["xmin"], box["xmax"]) &
9     between(LATITUDE, box["ymin"], box["ymax"]), ]
10  ebd <- ebd[year(`OBSERVATION DATE`) >= 2013, ]
11
12  # make new column names
13  newNames <- str_replace_all(colnames(ebd), " ", "_") %>%
14    str_to_lower()
15  setnames(ebd, newNames)
16
17  # keep useful columns
18  columnsOfInterest <- c(
19    "scientific_name", "observation_count", "locality",
20    "locality_id", "locality_type", "latitude",
21    "longitude", "observation_date", "sampling_event_identifier"
22  )
23
24  ebd <- ebd[, ..columnsOfInterest]
```

Add a spatial filter and assign grids of 25km x 25km.

```
1   # strict spatial filter and assign grid
2   locs <- ebd[, .(longitude, latitude)]
3
4   # transform to UTM and get 20km boxes
5   coords <- setDF(locs) %>%
6     st_as_sf(coords = c("longitude", "latitude")) %>%
7     `st_crs<-`(4326) %>%
8     bind_cols(as.data.table(st_coordinates(.))) %>%
9     st_transform(32643) %>%
10    mutate(id = 1:nrow(.))
11
12  # convert wg to UTM for filter
13  wg <- st_transform(wg, 32643)
14  coords <- coords %>%
15    filter(id %in% unlist(st_contains(wg, coords))) %>%
16    rename(longitude = X, latitude = Y) %>%
17    bind_cols(as.data.table(st_coordinates(.))) %>%
18    st_drop_geometry() %>%
19    as.data.table()
20
21  # remove unneeded objects
```

```
22  rm(locs)
23  gc()
24
25  coords <- coords[, .N, by = .(longitude, latitude, X, Y)]
26
27  ebd <- merge(ebd, coords, all = FALSE, by = c("longitude", "latitude"))
28
29  ebd <- ebd[(longitude %in% coords$longitude) &
30    (latitude %in% coords$latitude), ]
```

## 3.4 Get proportional obs counts in 25km cells

```
1   # round to 25km cell in UTM coords
2   ebd[, `:=`(X = round_any(X), Y = round_any(Y))]
3
4   # count checklists in cell
5   ebd_summary <- ebd[, nchk := length(unique(sampling_event_identifier)),
6     by = .(X, Y)
7   ]
8
9   # count checklists reporting each species in cell and get proportion
10  ebd_summary <- ebd_summary[, .(nrep = length(unique(
11    sampling_event_identifier
12  ))),
13  by = .(X, Y, nchk, scientific_name)
14  ]
15
16  ebd_summary[, p_rep := nrep / nchk]
17
18  # filter for soi
19  ebd_summary <- ebd_summary[scientific_name %in% speciesOfInterest, ]
20
21  # complete the dataframe for no reports
22  # keep no reports as NA --- allows filtering based on proportion reporting
23  ebd_summary <- setDF(ebd_summary) %>%
24    complete(
25      nesting(X, Y), scientific_name # ,
26      # fill = list(p_rep = 0)
27    ) %>%
28    filter(!is.na(p_rep))
```

## 3.5 Which species are reported sufficiently in checklists?

```
1   # A total of 42 unique grids (of 25km by 25km) across the study area
2   # total number of checklists across unique grids
3
4   tot_n_chklist <- ebd_summary %>%
5     distinct(X, Y, nchk)
6
7   # species-specific number of grids
8   spp_grids <- ebd_summary %>%
9     group_by(scientific_name) %>%
```

```
10      distinct(X, Y) %>%
11      count(scientific_name,
12          name = "n_grids"
13      )
14
15   # Write the above two results
16   write_csv(tot_n_chklist, "data/nchk_per_grid.csv")
17   write_csv(spp_grids, "data/ngrids_per_spp.csv")
18
19   # left-join the datasets
20   ebd_summary <- left_join(ebd_summary, spp_grids, by = "scientific_name")
21
22   # check the proportion of grids across which this cut-off is met for each species
23   # Is it > 90% or 70%?
24   # For example, with a 3% cut-off, ~100 species are occurring in >50%
25   # of the grids they have been reported in
26
27   p_cutoff <- 0.05 # Proportion of checklists a species has been reported in
28   grid_proportions <- ebd_summary %>%
29      group_by(scientific_name) %>%
30      tally(p_rep >= p_cutoff) %>%
31      mutate(prop_grids_cut = n / (spp_grids$n_grids)) %>%
32      arrange(desc(prop_grids_cut))
33
34   grid_prop_cut <- filter(
35      grid_proportions,
36      prop_grids_cut > p_cutoff
37   )
38
39   # Write the results
40   write_csv(grid_prop_cut, "data/chk_5_percent.csv")
41
42   # Identifying the number of species that occur in potentially <5% of all lists
43   total_number_lists <- sum(tot_n_chklist$nchk)
44
45   spp_sum_chk <- ebd_summary %>%
46      distinct(X, Y, scientific_name, nrep) %>%
47      group_by(scientific_name) %>%
48      mutate(sum_chk = sum(nrep)) %>%
49      distinct(scientific_name, sum_chk)
50
51   # Approximately 90 to 100 species occur in >5% of all checklists
52   prop_all_lists <- spp_sum_chk %>%
53      mutate(prop_lists = sum_chk / total_number_lists) %>%
54      arrange(desc(prop_lists))
```

## 3.6   Figure: Checklist distribution

```
1   # add land
2   library(rnaturalearth)
3   land <- ne_countries(
4      scale = 50, type = "countries", continent = "asia",
5      country = "india",
```

```
6      returnclass = c("sf")
7    )
8
9  # crop land
10 land <- st_transform(land, 32643)
```

## 3.7 Prepare the species list

```
1  # write the new list of species that occur in at least 5% of checklists across a minimum of 50% of the grids they h
2
3  new_sp_list <- semi_join(specieslist, grid_prop_cut, by = "scientific_name")
4
5  write_csv(new_sp_list, "data/03_list-of-species-cutoff.csv")
```

# 4    Climate in Relation to Landcover

This script showcases how climate data varies as a function of land cover types across our study area.

## 4.1 Prepare libraries

```
1  # load libs
2  library(raster)
3  library(glue)
4  library(purrr)
5  library(dplyr)
6  library(tidyr)
7
8  # plotting options
9  library(ggplot2)
10 library(ggthemes)
11 library(scico)
12
13 # get ci func
14 ci <- function(x) {
15   qnorm(0.975) * sd(x, na.rm = T) / sqrt(length(x))
16 }
```

## 4.2 Prepare environmental data

```
1  # read landscape prepare for plotting
2  landscape <- stack("data/spatial/landscape_resamp01km.tif")
3
4  # get proper names
5  elev_names <- c("elev", "slope", "aspect")
6  chelsa_names <- c("bio_01", "bio_12")
7
8  names(landscape) <- as.character(glue('{c(elev_names, chelsa_names, "landcover")}'))
```

```
1  # make duplicate stack
2  land_data <- landscape[[c("landcover", chelsa_names)]]
3
4  # convert to list
```

Figure 1: Proportion of checklists reporting a species in each grid cell (25km side) between 2013 and 2019. Checklists were filtered to be within the boundaries of the Nilgiris and the Anamalai hills (black outline), but rounding to 25km cells may place cells outside the boundary. Deeper shades of red indicate a higher proportion of checklists reporting a species.

```
5  land_data <- as.list(land_data)
6
7  # map get values over the stack
8  land_data <- purrr::map(land_data, raster::getValues)
9  names(land_data) <- c("landcover", chelsa_names)
10
11 # conver to dataframe and round to 100m
12 land_data <- bind_cols(land_data)
13 land_data <- drop_na(land_data) %>%
14   filter(landcover != 0) %>%
15   pivot_longer(
16     cols = contains("bio"),
17     names_to = "clim_var"
18   ) # %>%
19 # group_by(landcover, clim_var) %>%
20 # summarise_all(.funs = list(~mean(.), ~ci(.)))
```

## 4.3 Climatic variables over landcover

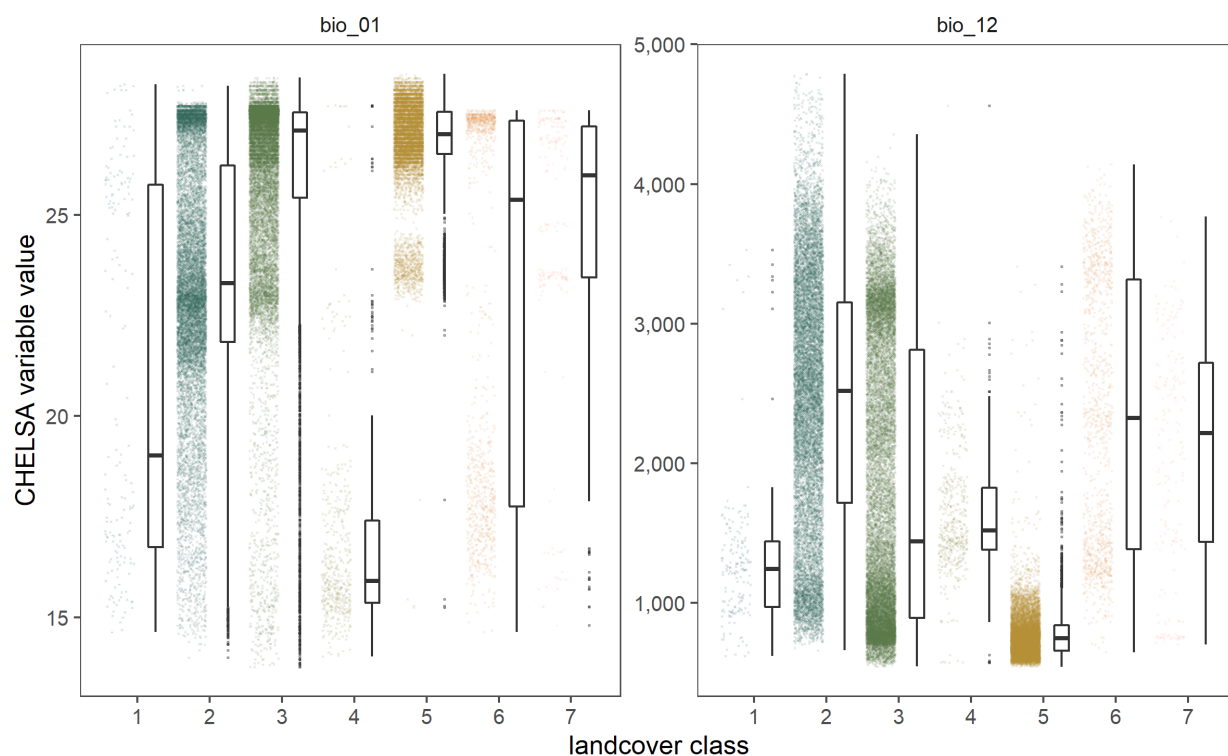Figure code is hidden in versions rendered as HTML and PDF.



Figure 2: CHELSA climatic variables as a function of landcover class. Grey points in the background represent raw data.

# 5 Distribution of Observer Expertise

This script plots observer expertise over time (2013-2019) as well as across land cover types. ## Prepare libraries

```
1   # load libs
2   library(raster)
3   library(glue)
4   library(purrr)
5   library(dplyr)
6   library(tidyr)
7   library(readr)
8   library(scales)
9
10  # plotting libs
11  library(ggplot2)
12  library(ggthemes)
13  library(scico)
14
15  # get ci func
16  ci <- function(x) {
17    qnorm(0.975) * sd(x, na.rm = T) / sqrt(length(x))
18  }
```

## 5.1   Load observer expertise scores and checklist covariates

```
1   # read in scores and checklist data and link
2   scores <- read_csv("data/03_data-obsExpertise-score.csv")
3   data <- read_csv("data/03_data-covars-perChklist.csv")
4
5   data <- left_join(data, scores, by = c("observer" = "observer"))
6   data <- dplyr::select(data, score, nSp, nSoi, landcover, year) %>%
7     filter(!is.na(score))
```

## 5.2   Species observed in relation to observer expertise

```
1   # summarise data by rounded score and year
2   data_summary01 <- data %>%
3     mutate(score = plyr::round_any(score, 0.2)) %>%
4     dplyr::select(score, year, nSp, nSoi) %>%
5     pivot_longer(
6       cols = c("nSp", "nSoi"),
7       names_to = "variable", values_to = "value"
8     ) %>%
9     group_by(score, year, variable) %>%
10    summarise_at(vars(value), list(~ mean(.), ~ ci(.)))
11
12  # make plot and export
13  fig_nsp_score <-
14    ggplot(data_summary01) +
15    geom_jitter(
16      data = data, aes(x = score, y = nSp),
17      col = "grey", alpha = 0.2, size = 0.1
18    ) +
19    geom_pointrange(aes(
20      x = score, y = mean,
21      ymin = mean - ci, ymax = mean + ci,
```

9

```
22      col = as.factor(variable)
23    ),
24    position = position_dodge(width = 0.05)
25    ) +
26    facet_wrap(~year) +
27    scale_y_log10() +
28    #  coord_cartesian(ylim=c(0,50))+
29    scale_colour_scico_d(palette = "cork", begin = 0.2, end = 0.8) +
30    labs(x = "CCI", y = "Number of Species Reported") +
31    theme_few() +
32    theme(legend.position = "none")
33
34  # export figure
35  ggsave(filename = "figs/fig_nsp_score.png", width = 12, height = 7, device = png(), dpi = 300)
36  dev.off()
```
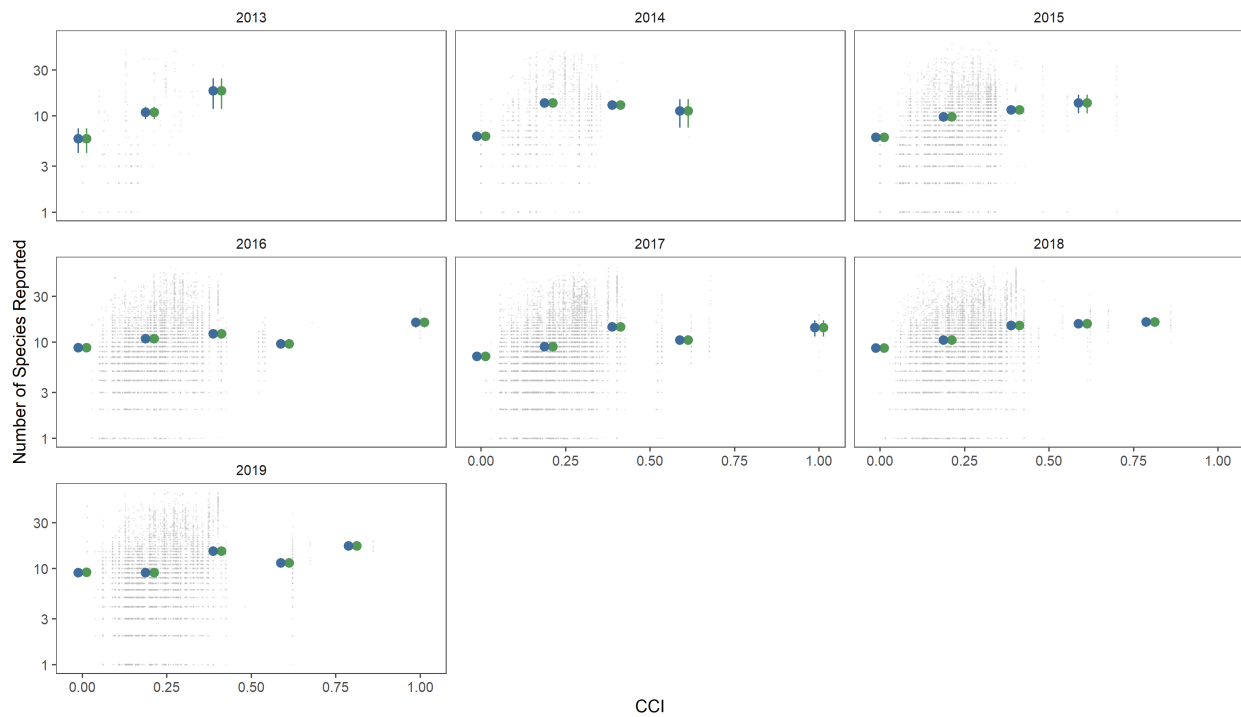


Figure 3: Total number of species (blue) and species of interest to this study (green) reported in checklists from the study area over the years 2013 – 2018, as a function of the expertise score of the reporting observer. Points represent means, with bars showing the 95% confidence intervals; data shown are for expertise scores rounded to multiples of 0.2, and the y-axis is on a log scale. Raw data are shown in the background (grey points).

## 5.3   Observer expertise in relation to landcover

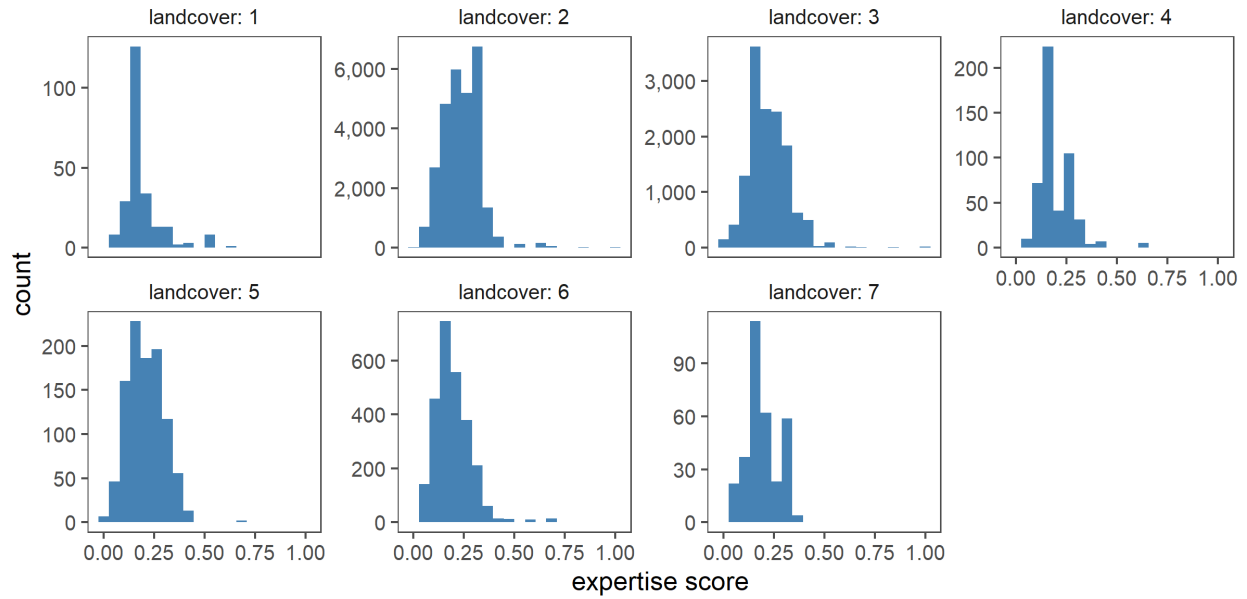Figure code is hidden in versions rendered as HTML or PDF.

Figure 4: Distribution of expertise scores in the seven landcover classes present in the study site.

# 6 Spatial Autocorrelation of Climatic Predictors

## 6.1 Load libraries

```
1  # load libs
2  library(raster)
3  library(gstat)
4  library(stars)
5  library(purrr)
6  library(tibble)
7  library(dplyr)
8  library(tidyr)
9  library(glue)
10 library(scales)
11 library(gdalUtils)
12 library(sf)
13
14 # plot libs
15 library(ggplot2)
16 library(ggthemes)
17 library(scico)
18 library(gridExtra)
19 library(cowplot)
20 library(ggspatial)
21
22 #' make custom functiont to convert matrix to df
23 raster_to_df <- function(inp) {
24
25   # assert is a raster obj
26   assertthat::assert_that("RasterLayer" %in% class(inp),
27     msg = "input is not a raster"
```

```
28    )
29
30    coords <- coordinates(inp)
31    vals <- getValues(inp)
32
33    data <- tibble(x = coords[, 1], y = coords[, 2], value = vals)
34
35    return(data)
36  }
```

## 6.2 Prepare data

```
1   # list landscape covariate stacks
2   landscape_files <- "data/spatial/landscape_resamp01_km.tif"
3   landscape_data <- stack(landscape_files)
4
5   # get proper names
6   elev_names <- c("elev", "slope", "aspect")
7   chelsa_names <- c("bio_01", "bio_12")
8   names(landscape_data) <- c(elev_names, chelsa_names, "landcover")
9
10
11  # get chelsa rasters
12  chelsa <- landscape_data[[chelsa_names]]
13  chelsa <- purrr::map(as.list(chelsa), raster_to_df)
```

## 6.3 Calculate variograms of environmental layers

```
1   # prep variograms
2   vgrams <- purrr::map(chelsa, function(z) {
3     z <- drop_na(z)
4     vgram <- gstat::variogram(value ~ 1, loc = ~ x + y, data = z)
5     return(vgram)
6   })
7
8   # save temp
9   save(vgrams, file = "data/chelsa/chelsaVariograms.rdata")
10
11  # get variogram data
12  vgrams <- purrr::map(vgrams, function(df) {
13    df %>% select(dist, gamma)
14  })
15  vgrams <- tibble(
16    variable = chelsa_names,
17    data = vgrams
18  )
```

```
1   wg <- st_read("data/spatial/hillsShapefile/Nil_Ana_Pal.shp") %>%
2     st_transform(32643)
3   bbox <- st_bbox(wg)
4
5   # add lamd
6   library(rnaturalearth)
```

```
7   land <- ne_countries(
8     scale = 50, type = "countries", continent = "asia",
9     country = "india",
10    returnclass = c("sf")
11  )
12
13  # crop land
14  land <- st_transform(land, 32643)
```

## 6.4   Visualise variograms of environmental data

```
1   # make ggplot of variograms
2   yaxis <- c("semivariance", "")
3   xaxis <- c("", "distance (km)")
4   fig_vgrams <- purrr::pmap(list(vgrams$data, yaxis, xaxis), function(df, ya, xa) {
5     ggplot(df) +
6       geom_line(aes(x = dist / 1000, y = gamma), size = 0.2, col = "grey") +
7       geom_point(aes(x = dist / 1000, y = gamma), col = "black") +
8       scale_x_continuous(labels = comma, breaks = c(seq(0, 100, 25))) +
9       scale_y_log10(labels = comma) +
10      labs(x = xa, y = ya) +
11      theme_few() +
12      theme(
13        axis.text.y = element_text(angle = 90, hjust = 0.5, size = 8),
14        strip.text = element_blank()
15      )
16  })
17  # fig_vgrams <- purrr::map(fig_vgrams, ggplot2::ggplotGrob)
18
19  # make ggplot of chelsa data
20  chelsa <- as.list(landscape_data[[chelsa_names]]) %>%
21    purrr::map(stars::st_as_stars)
22
23  # colour palettes
24  pal <- c("bilbao", "davos")
25  title <- c(
26    "a Annual Mean Temperature",
27    "b Annual Precipitation"
28  )
29  direction <- c(1, 1)
30  lims <- list(
31    range(values(landscape_data$bio_01), na.rm = T),
32    range(values(landscape_data$bio_12), na.rm = T)
33  )
34  fig_list_chelsa <-
35    purrr::pmap(
36      list(chelsa, pal, title, direction, lims),
37      function(df, pal, t, d, l) {
38        ggplot() +
39          stars::geom_stars(data = df) +
40          geom_sf(data = land, fill = NA, colour = "black") +
41          geom_sf(data = wg, fill = NA, colour = "black", size = 0.3) +
42          scale_fill_scico(
```

```
43          palette = pal, direction = d,
44          label = comma, na.value = NA, limits = l
45        ) +
46        coord_sf(
47          xlim = bbox[c("xmin", "xmax")],
48          ylim = bbox[c("ymin", "ymax")]
49        ) +
50        ggspatial::annotation_scale(location = "tr", width_hint = 0.4, text_cex = 1) +
51        theme_few() +
52        theme(
53          legend.position = "top",
54          title = element_text(face = "bold", size = 8),
55          legend.key.height = unit(0.2, "cm"),
56          legend.key.width = unit(1, "cm"),
57          legend.text = element_text(size = 8),
58          axis.title = element_blank(),
59          axis.text.y = element_text(angle = 90, hjust = 0.5),
60          panel.background = element_rect(fill = "lightblue"),
61          legend.title = element_blank()
62        ) +
63        labs(x = NULL, y = NULL, title = t)
64    }
65  )
66 # fig_list_chelsa <- purrr::map(fig_list_chelsa, ggplotGrob)
```

```
1  # fig_list_chelsa <- append(fig_list_chelsa, fig_vgrams)
2  # lmatrix <- matrix(c(c(1, 2, 3, 4, 5), c(1, 2, 3, 4, 5), c(6, 7, 8, 9, 10)),
3  #   nrow = 3, byrow = T
4  # )
5  # plot_grid <- grid.arrange(grobs = fig_list_chelsa, layout_matrix = lmatrix)
6  #
7  # ggsave(
8  #   plot = plot_grid, filename = "figs/fig_chelsa_variograms.png",
9  #   dpi = 300, width = 12, height = 6
10 # )
11 # dev.off()
12
13 library(patchwork)
14 fig_variogram <- wrap_plots(append(fig_list_chelsa, fig_vgrams))
15 ggsave(fig_variogram,
16   filename = "figs/fig_chelsa_variograms.png",
17   dpi = 300,
18   width = 6, height = 6
19 )
```

## 7   Climatic raster resampling

### 7.1   Prepare landcover

```
1  # read in landcover raster location
2  landcover <- "data/landUseClassification/classifiedImage-UTM.tif"
3  # get extent
4  e <- bbox(raster(landcover))
```
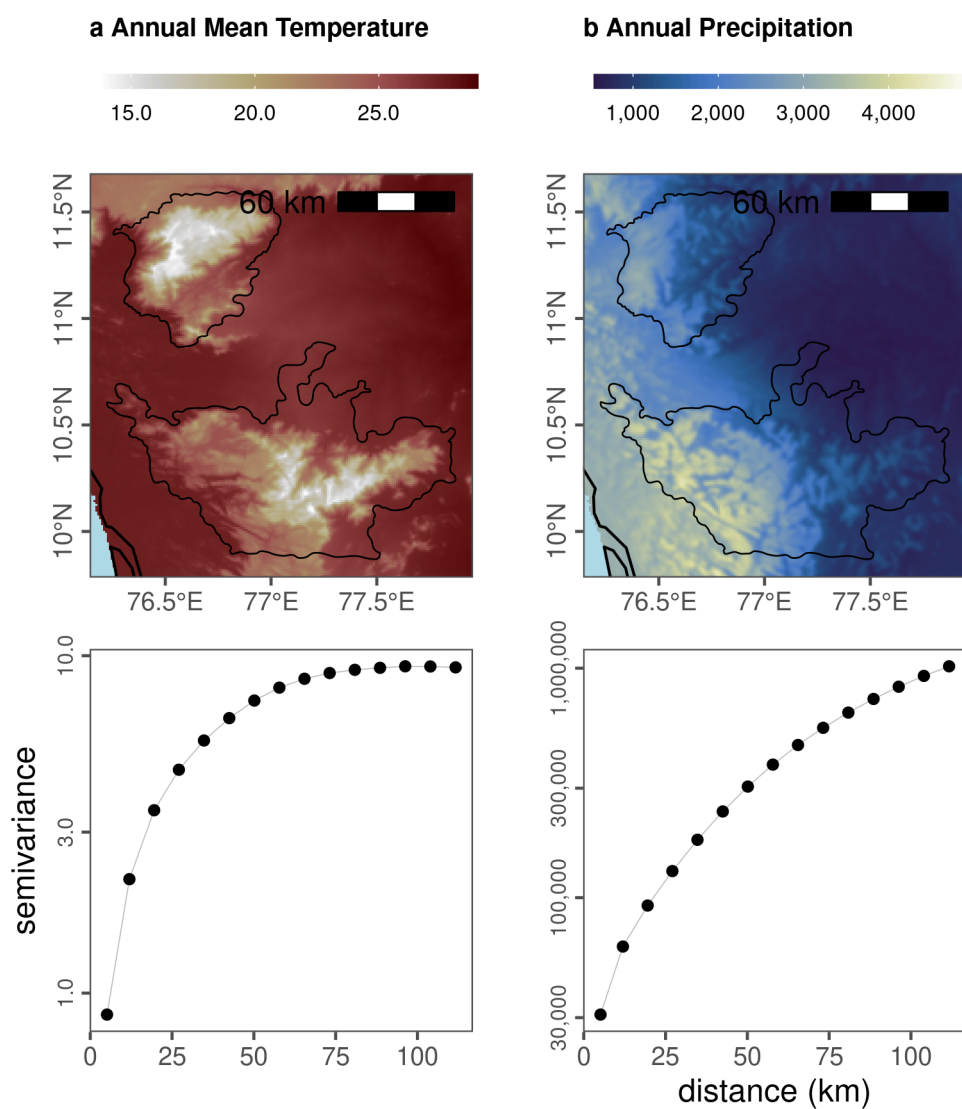
14

Figure 5: CHELSA rasters with study area outline, and associated semivariograms. Semivariograms are on a log-transformed y-axis.

```r
5
6  # init resolution
7  res_init <- res(raster(landcover))
8  # res to transform to 1000m
9  res_final <- map(c(100, 250, 500, 1e3, 2.5e3), function(x) {
10    x * res_init
11 })
12
13 # use gdalutils gdalwarp for resampling transform
14 # to 1km from 10m
15 for (i in 1:length(res_final)) {
16   this_res <- res_final[[i]]
17   this_res_char <- stringr::str_pad(this_res[1], 5, pad = "0")
18   gdalUtils::gdalwarp(
19     srcfile = landcover,
20     dstfile = as.character(glue("data/landUseClassification/lc_{this_res_char}m.tif")),
21     tr = c(this_res), r = "mode", te = c(e)
22   )
23 }
```

```r
1  # read in resampled landcover raster files as a list
2  lc_files <- list.files("data/landUseClassification/", pattern = "lc", full.names = TRUE)
3  lc_data <- map(lc_files, raster)
```

## 7.2  Prepare spatial extent

```r
1  # load hills
2  library(sf)
3  hills <- st_read("data/spatial/hillsShapefile/Nil_Ana_Pal.shp")
4  hills <- st_transform(hills, 32643)
5  buffer <- st_buffer(hills, 3e4) %>%
6    st_transform(4326)
7  bbox <- st_bbox(hills)
```

## 7.3  Prepare CHELSA rasters

```r
1  # list chelsa files
2  chelsaFiles <- list.files("data/chelsa/", full.names = TRUE, pattern = "*.tif")
3
4  # gather chelsa rasters
5  chelsaData <- purrr::map(chelsaFiles, function(chr) {
6    a <- raster(chr)
7    crs(a) <- crs(buffer)
8    a <- crop(a, as(buffer, "Spatial"))
9    return(a)
10 })
11
12 # stack chelsa data
13 chelsaData <- raster::stack(chelsaData)
14 names(chelsaData) <- c("chelsa_bio10_01", "chelsa_bio10_12")
```