

An R Package for Spatio-Temporal Change of Support

Andrew M. Raim

Center for Statistical Research and Methodology

U.S. Census Bureau

andrew.raim@census.gov

2017 Joint Statistical Meetings, Baltimore, Maryland

Joint work with **Scott H. Holan** (U. of Missouri & U.S. Census Bureau),
Jonathan R. Bradley (Florida State U.), and **Christopher K. Wikle** (U. of Missouri)

Disclaimer

This presentation is to inform interested parties of ongoing research and to encourage discussion of work in progress. Any views expressed are those of the authors and not necessarily those of the U.S. Census Bureau.

The American Community Survey (ACS)

- The ACS is an ongoing survey administered by the U.S. Census Bureau to measure key socioeconomic and demographic variables for the U.S. population.
- ACS data is available to the public through the American FactFinder (<http://factfinder.census.gov>) for years 2005 through 2015; 2016 release is planned for mid-Sept.
- Estimates have been released for 1-year, 3-year, or 5-year periods. 3-year estimates were discontinued after 2013.
- Granularity is down to census block-groups. However, estimates for an area are suppressed unless the area meets certain criteria.
- For example, an area must have population > 65,000 for 1-year estimates to be released, but there is no population requirement for 5-year estimates (U.S. Census Bureau, 2016).

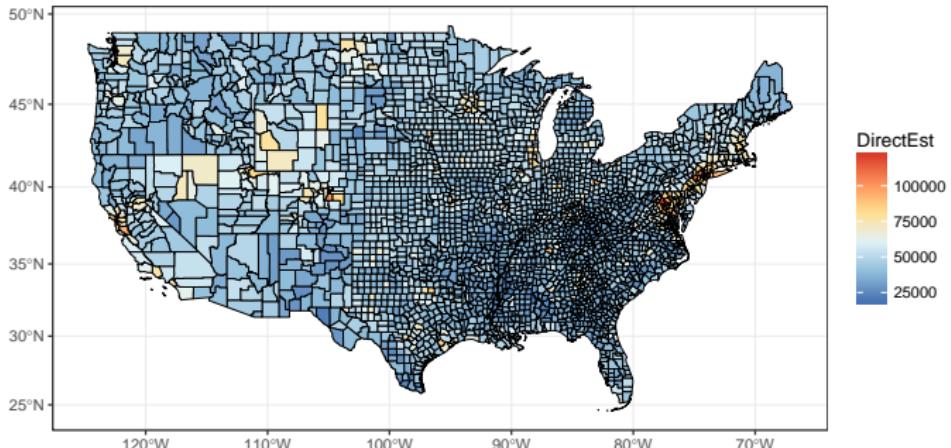
Spatio-Temporal Change of Support in the ACS

- **Spatio-Temporal Change of Support (STCOS) Problem:** using all available ACS releases and their patterns over space and time, provide reasonable estimates for a user-specified support and period.
- This work is based on models developed in Bradley et al. (2015, Stat). We develop the `stcos` R package to make the methodology widely accessible to data users.
- Statistical agencies have direct access to microdata, and can aggregate to any support and period without STCOS methodology.
- The methods and software are not limited to use with ACS data, but were developed with ACS in mind.
- See Bradley et al. (2015, Stat) and the references therein for a review of spatio(-temporal) change of support literature.

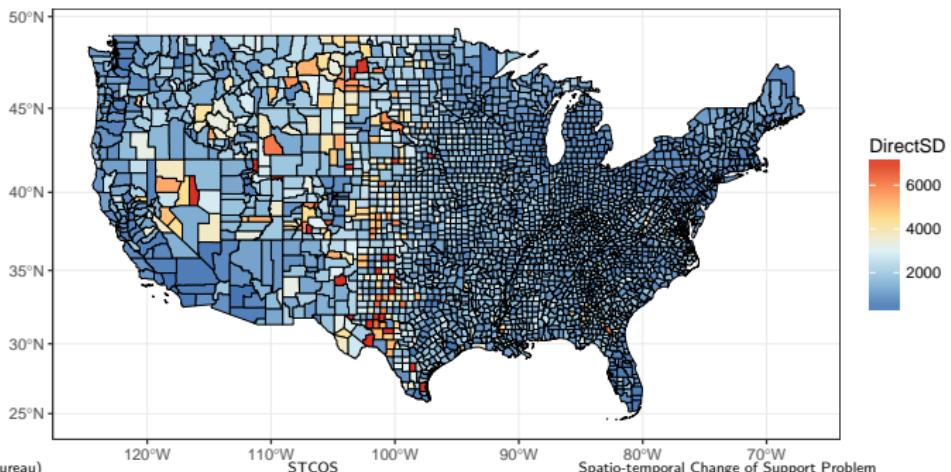
American Indian Reservations

- There are 693 areas in the U.S. designated as American Indian areas/Alaska native areas/Hawaiian home lands (AIANNHs). Of these, there are 397 American Indian reservations within the continental U.S.
- These areas do not necessarily align with census geography (tracts, block groups, counties, etc).
- ACS releases estimates for AIANNHs, which provides a good opportunity to compare with STCOS model estimates.

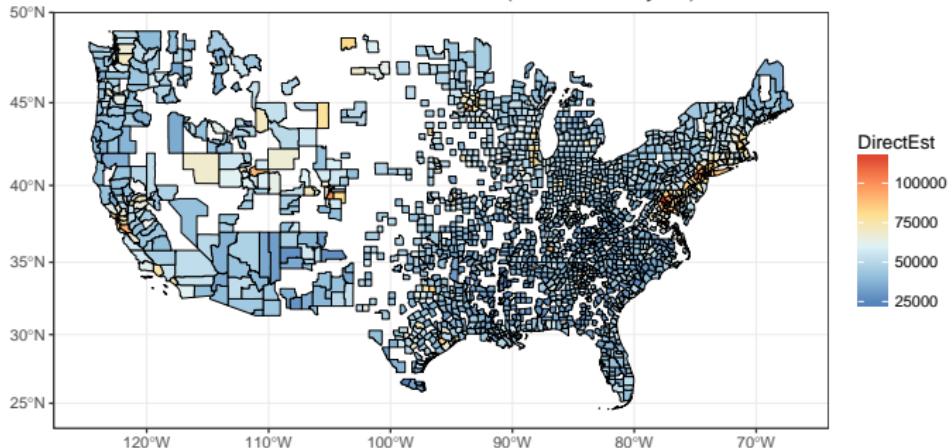
Median Household Income for U.S. Counties (ACS 2013 5-year)



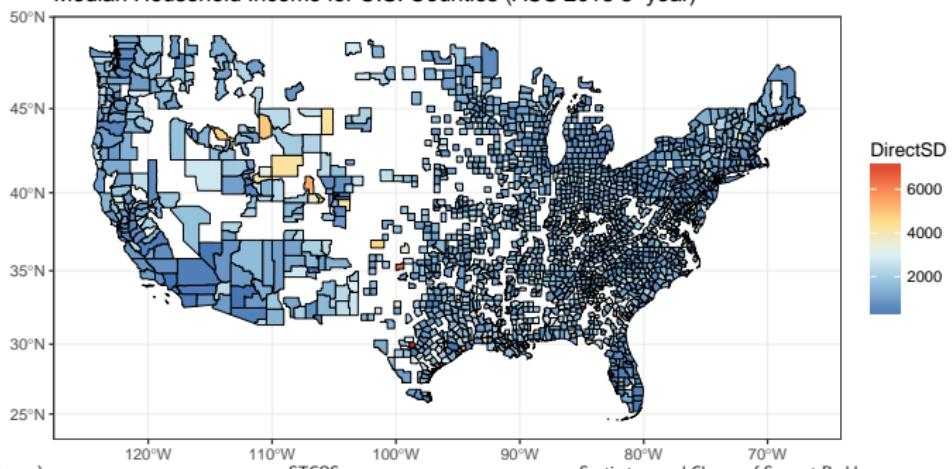
Median Household Income for U.S. Counties (ACS 2013 5-year)



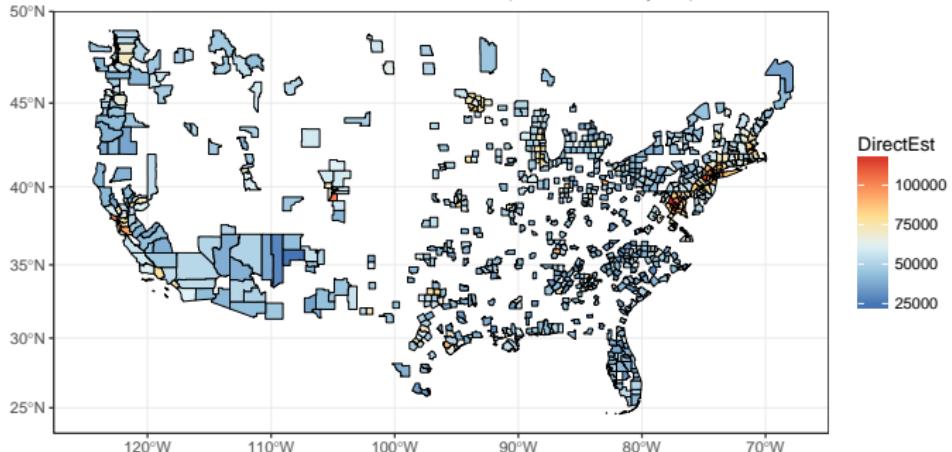
Median Household Income for U.S. Counties (ACS 2013 3-year)



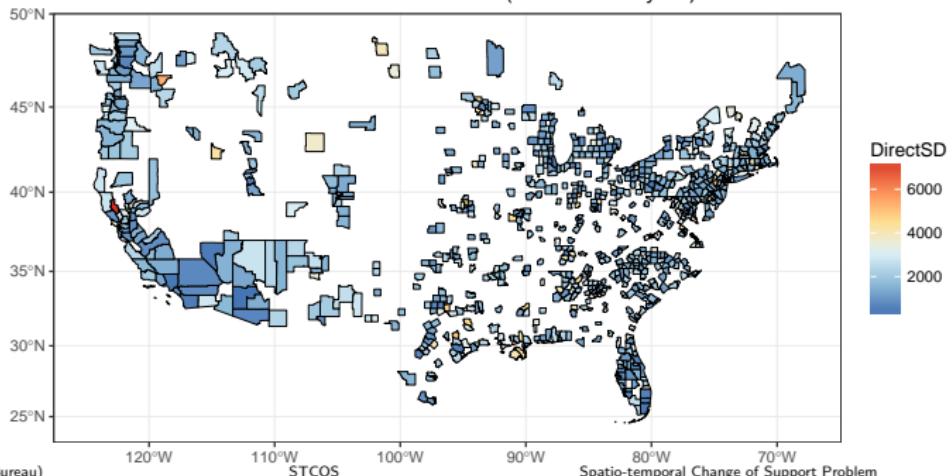
Median Household Income for U.S. Counties (ACS 2013 3-year)



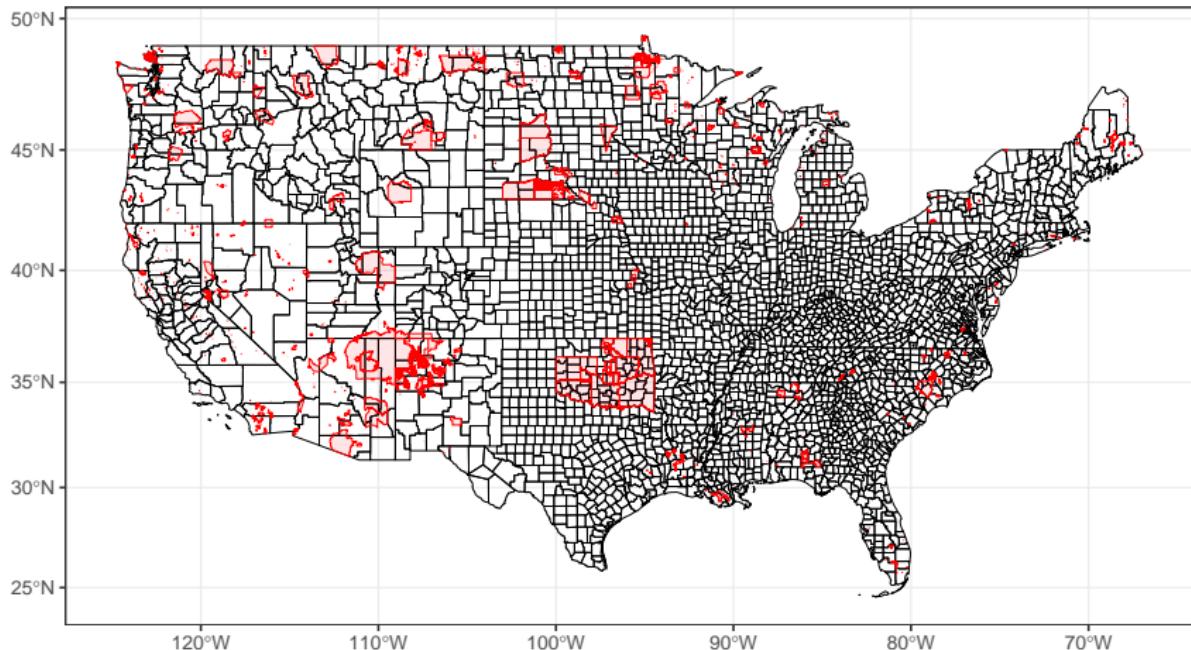
Median Household Income for U.S. Counties (ACS 2013 1-year)



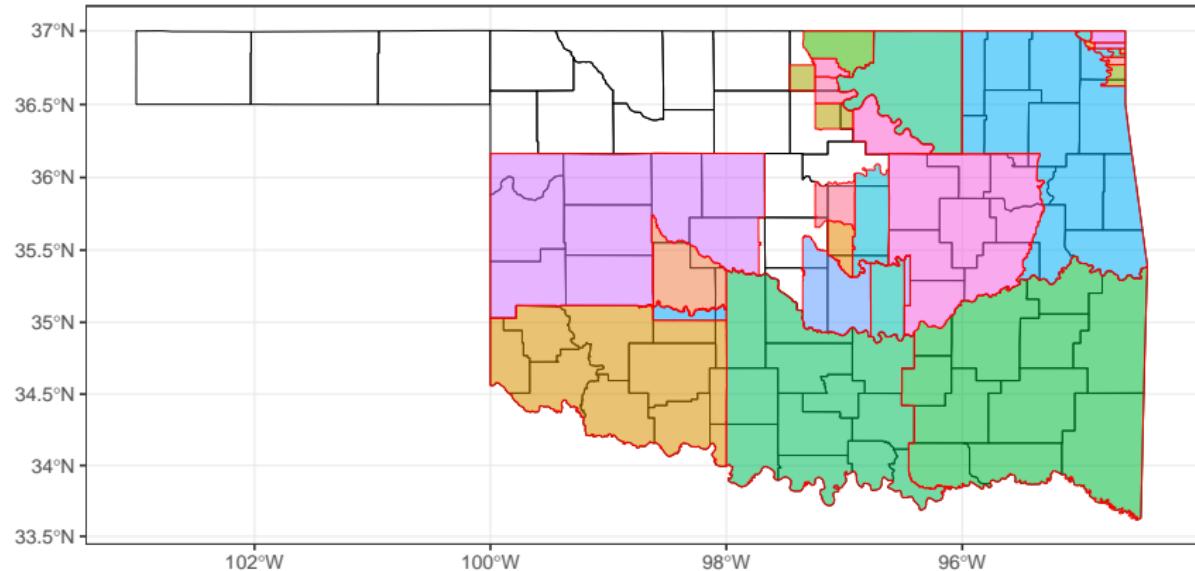
Median Household Income for U.S. Counties (ACS 2013 1-year)



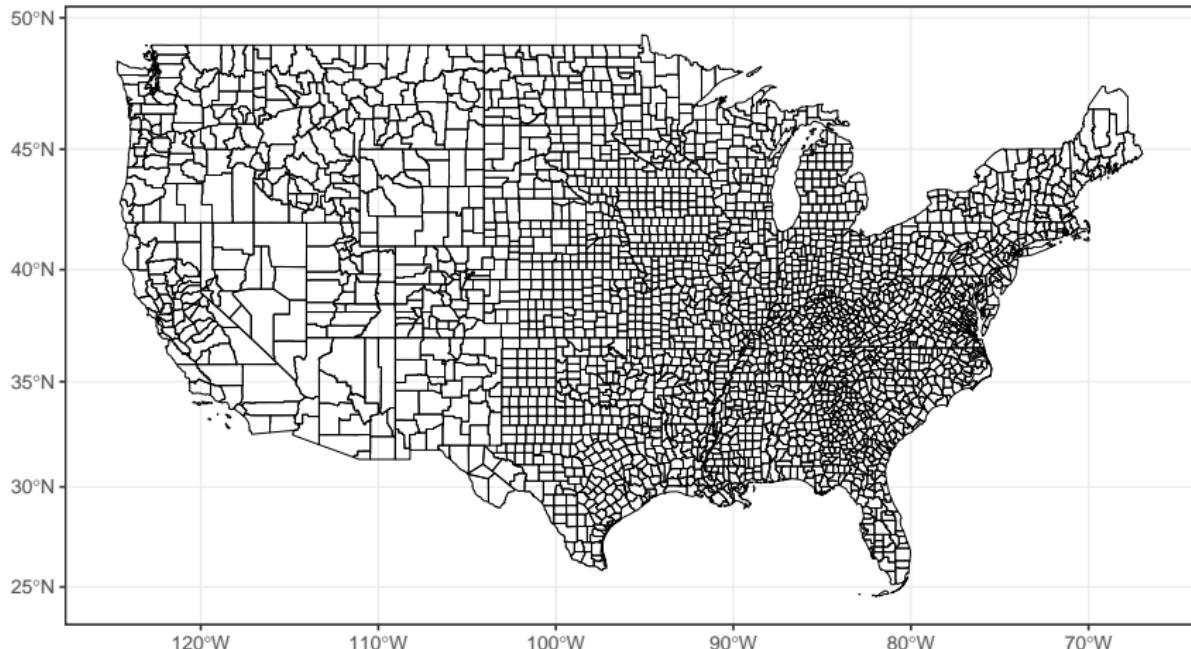
American Indian Reservations in 2015



American Indian Areas within Oklahoma in 2015



U.S. Counties in 2015



The STCOS Model

- $\mathcal{T} = \{T_L, \dots, T_U\}$: times for which direct estimates are available.
- \mathcal{L} : set of lookback periods. For ACS data, $\mathcal{L} = \{1, 3, 5\}$ are possible lookbacks.
- $D_{t\ell}$: source support — collection of areal units with direct estimates — for time $t \in \mathcal{T}$ and period $\ell \in \mathcal{L}$.
- $Z_t^{(\ell)}(A)$ and $\sigma_{t\ell}^2(A)$: direct survey estimate and associated variance for a survey variable of interest, $A \in D_{t\ell}$, $\ell \in \mathcal{L}$, $t \in \mathcal{T}$.
- $D_B = \{B_1, \dots, B_{n_B}\}$ is the fine level support.

STCOS Bayesian Hierarchical Model

- Data Model

$$Z_t^{(\ell)}(A) = Y_t^{(\ell)}(A) + \varepsilon_t^{(\ell)}(A)$$
$$\varepsilon_t^{(\ell)}(A) \stackrel{\text{ind}}{\sim} N(0, \sigma_{t\ell}^2(A))$$

- Process Model

$$Y_t^{(\ell)}(A) = \underbrace{h(A)^\top \mu_B}_{\text{Coarse spatial trend}} + \underbrace{\psi_t^{(\ell)}(A)^\top \eta}_{\text{Fine spatio-temporal trend}} + \xi_t^{(\ell)}(A)$$
$$\xi_t^{(\ell)}(A) \stackrel{\text{iid}}{\sim} N(0, \sigma_\xi^2)$$

- Prior Model

$$\mu_B \sim N(0, \sigma_\mu^2 I), \quad \eta \sim N(0, \sigma_K^2 K),$$
$$\sigma_\mu^2 \sim IG(a_\mu, b_\mu), \quad \sigma_K^2 \sim IG(a_K, b_K), \quad \sigma_\xi^2 \sim IG(a_\xi, b_\xi)$$

Latent Process Model

- Define a continuous-space discrete-time process on $\mathbf{u} \in \bigcup_{i=1}^{n_B} B_i$, $t \in \mathcal{T}$,

$$Y(\mathbf{u}; t) = \delta(\mathbf{u}) + \sum_{j=1}^{\infty} \psi_j(\mathbf{u}; t) \cdot \eta_j,$$

where $\delta(\mathbf{u})$ is a large-scale spatial trend process and $\{\psi_j(\mathbf{u}, t)\}_{j=1}^{\infty}$ is a pre-specified set of spatio-temporal basis functions.

- Integrate $Y(\mathbf{u}; t)$ over $\mathbf{u} \in A$ (wrt uniform density) and ℓ lookbacks,

$$\begin{aligned} Y_t^{(\ell)}(A) &= \underbrace{\frac{1}{|A|} \int_A \delta(\mathbf{u}) d\mathbf{u}}_{\text{large-scale spatial trend}} + \underbrace{\frac{1}{\ell|A|} \sum_{k=t-\ell+1}^t \sum_{j=1}^r \int_A \psi_j(\mathbf{u}; k) \cdot \eta_j}_{\text{small-scale spatio-temporal trend}} \\ &\quad + \underbrace{\frac{1}{\ell|A|} \sum_{k=t-\ell+1}^t \sum_{j=r+1}^{\infty} \int_A \psi_j(\mathbf{u}; k) \cdot \eta_j}_{\text{leftovers}} \\ &= \mu(A) + \psi_t^{(\ell)}(A)^{\top} \boldsymbol{\eta} + \xi_t^{(\ell)}(A). \end{aligned}$$

- For the leftovers, assume that $\xi_t^{(\ell)}(A) \stackrel{\text{iid}}{\sim} N(0, \sigma_{\xi}^2)$.

Basis Functions

- We make use of local bisquare basis functions,

$$\psi_j(\mathbf{u}, t) = \left[1 - \frac{\|\mathbf{u} - \mathbf{c}_j\|^2}{w_s^2} - \frac{|t - g_t|^2}{w_t^2} \right]^2 \times \\ I(\|\mathbf{u} - \mathbf{c}_j\| \leq w_s) \cdot I(|t - g_t| \leq w_t).$$

- Spatial knot points \mathbf{c}_j , $j = 1, \dots, r_{\text{space}}$, are selected via a space-filling design on D_B ; see the R `fields` package (Nychka et al., 2015).
- Temporal knot points g_t , $t = 1, \dots, r_{\text{time}}$, are chosen to be equally spaced through \mathcal{T} .
- For area A and lookback period ℓ , we take a Monte Carlo approximation

$$\psi_{jt}^{(\ell)}(A) \approx \frac{1}{\ell Q} \sum_{k=t-\ell+1}^t \sum_{q=1}^Q \psi_j(\mathbf{u}_q, k),$$

using a uniform random sample $\mathbf{u}_1, \dots, \mathbf{u}_Q$ on A .

Change of Support Term

- Suppose for the large-scale spatial trend process that

$$\delta(u) = \sum_{i=1}^{n_B} \mu_i I(u \in A \cap B_i), \quad \text{for a given area } A.$$

- Then, integrating over $u \in A$,

$$\begin{aligned}\mu(A) &= \frac{1}{|A|} \sum_{i=1}^{n_B} \int_{A \cap B_i} \delta(u) du = \frac{1}{|A|} \sum_{i=1}^{n_B} \mu_i \int_{A \cap B_i} du = \sum_{i=1}^{n_B} \mu_i \frac{|A \cap B_i|}{|A|} \\ &= h(A)^\top \boldsymbol{\mu}_B.\end{aligned}$$

- $h(A) = (|A \cap B_1|/|A|, \dots, |A \cap B_{n_B}|/|A|)$ is computed from the source and fine-level supports.
- $\boldsymbol{\mu}_B = (\mu_1, \dots, \mu_{n_B})$ is unknown, to be estimated from the data.

Specification of K

- Suppose the fine-level support behaves according to the process

$$\mathbf{Y}_t^* = \boldsymbol{\mu}_B + \boldsymbol{\nu}_t, \quad \text{for } t \in \mathcal{T}$$

$$\boldsymbol{\nu}_t = \mathbf{M}\boldsymbol{\nu}_{t-1} + \mathbf{b}_t, \quad \mathbf{b}_t \stackrel{\text{iid}}{\sim} N(\mathbf{0}, \sigma_K^2(\mathbf{I} - \mathbf{A})^-).$$

where \mathbf{A} is the adjacency matrix of D_B .

- Let Σ_{y^*} denote the covariance matrix of $(\mathbf{Y}_t^* : t \in \mathcal{T})$.
- Obtain \mathbf{K} by solving

$$\min \|\Sigma_{y^*} - \mathbf{S}\mathbf{C}\mathbf{S}^\top\|_F, \quad \mathbf{C} \text{ is a } r \times r \text{ positive semidefinite matrix}$$

which yields $\mathbf{K} = (\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{S}^\top \Sigma_{y^*} \mathbf{S} (\mathbf{S}^\top \mathbf{S})^{-1}$. The best positive approximant problem is discussed further in Bradley et al. (2015) and Higham (1988).

- We propose several options where $\Sigma_{y^*} = \sigma_K^2 \tilde{\Sigma}_{y^*}$ such that $\tilde{\Sigma}_{y^*}$ is free of unknown parameters and M does not need to be estimated. Here,

$$\mathbf{K} = \sigma_K^2 \tilde{\mathbf{K}}, \quad \tilde{\mathbf{K}} = (\mathbf{S}^\top \mathbf{S})^{-1} \mathbf{S}^\top \tilde{\Sigma}_{y^*} \mathbf{S} (\mathbf{S}^\top \mathbf{S})^{-1}.$$

Specification of K

- **(Independence)** Taking $K = I$ assumes no spatio-temporal covariance in η .
- **(Spatial-only)** Let $\tilde{\Sigma}_{y^*} = \sigma_K^2(I - A)^{-} \otimes I_{|\mathcal{T}|}$ to ignore covariance in time.
- **(Random Walk)** If $M = I$, the process

$$\mathbf{Y}_t^* = \mu_B + M\nu_{t-1} + \mathbf{b}_t, \quad \mathbf{b}_t \stackrel{\text{iid}}{\sim} N(\mathbf{0}, \sigma_K^2(I - A)^{-})$$

is a vector random walk with autocovariance

$$\Gamma(t, h) = \begin{cases} t\sigma_K^2(I - A)^{-} & \text{if } h \geq 0 \\ (t - |h|)\sigma_K^2(I - A)^{-} & \text{if } -t < h < 0. \end{cases}$$

Take

$$\tilde{\Sigma}_{y^*} = \begin{bmatrix} \Gamma(1, 1) & \Gamma(1, 2) & \cdots & \Gamma(1, |\mathcal{T}|) \\ \Gamma(2, 1) & \Gamma(2, 2) & \cdots & \Gamma(2, |\mathcal{T}|) \\ \vdots & \vdots & \ddots & \vdots \\ \Gamma(|\mathcal{T}|, 1) & \Gamma(|\mathcal{T}|, 2) & \cdots & \Gamma(|\mathcal{T}|, |\mathcal{T}|) \end{bmatrix}.$$

Specification of K

- **(Vector Autoregression)** Consider the VAR(1) process with $\mu_B = \mathbf{X}\beta$,

$$\mathbf{Y}_t^* = \mu_B + \mathbf{M}\nu_{t-1} + \mathbf{b}_t, \quad \mathbf{b}_t \stackrel{\text{iid}}{\sim} N(\mathbf{0}, \sigma_K^2 (\mathbf{I} - \mathbf{A})^-).$$

- The method of Hughes and Haran (2013) suggests \mathbf{M} as the eigenvectors of $(\mathbf{I} - \mathbf{P}_X)\mathbf{W}(\mathbf{I} - \mathbf{P}_X)$, where $\mathbf{P}_X = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top$ and \mathbf{W} is a pre-specified real-valued matrix.
- We take \mathbf{X} to be a spatial-only bisquare basis expansion of the domain.
- Under VAR(1), the autocovariance becomes

$$\text{vec}(\boldsymbol{\Gamma}(0)) = [\mathbf{I} - \mathbf{M} \otimes \mathbf{M}]^{-1} \text{vec}(\sigma_K^2 (\mathbf{I} - \mathbf{A})^-)$$

$$\boldsymbol{\Gamma}(h) = \begin{cases} \mathbf{M}^h \boldsymbol{\Gamma}(0) & \text{if } h > 0, \\ \boldsymbol{\Gamma}(-h)^\top & \text{if } h < 0. \end{cases}$$

and therefore

$$\tilde{\boldsymbol{\Sigma}}_{y^*} = \begin{bmatrix} \boldsymbol{\Gamma}(0) & \boldsymbol{\Gamma}(-1) & \cdots & \boldsymbol{\Gamma}(-(|\mathcal{T}| - 1)) \\ \boldsymbol{\Gamma}(1) & \boldsymbol{\Gamma}(0) & \cdots & \boldsymbol{\Gamma}(-(|\mathcal{T}| - 2)) \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{\Gamma}(|\mathcal{T}| - 1) & \boldsymbol{\Gamma}(|\mathcal{T}| - 2) & \cdots & \boldsymbol{\Gamma}(0) \end{bmatrix}.$$

Specification of K

- The usual formula

$$\text{vec}(\boldsymbol{\Gamma}(0)) = [\mathbf{I} - \mathbf{M} \otimes \mathbf{M}]^{-1} \text{vec}(\sigma_K^2 (\mathbf{I} - \mathbf{A})^-)$$

for the lag-0 autocovariance is computationally intractable for high dimensional series.

- A more tractable form is

$$\boldsymbol{\Gamma}(0) = \mathbf{V} \cdot \mathcal{M} \left(\text{Diag}(\boldsymbol{\Omega}) \circ \text{vec}(\mathbf{V}^{-1} \sigma_K^2 (\mathbf{I} - \mathbf{A})^- \mathbf{V}^{-\top}) \right) \cdot \mathbf{V}^\top,$$

where \mathbf{V} and $\boldsymbol{\lambda}$ are the eigenvectors/values of \mathbf{M} , \circ is elementwise multiplication, $\mathcal{M} = \text{vec}^{-1}$, and $\boldsymbol{\Omega} = \text{Diag}(\mathbf{1} - \boldsymbol{\lambda} \otimes \boldsymbol{\lambda})^{-1}$.

STCOS Model in Vector Form

We may write

$$\mathbf{Z} = \mathbf{H}\boldsymbol{\mu}_B + \mathbf{S}\boldsymbol{\eta} + \boldsymbol{\xi} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(0, \mathbf{V}),$$

where

$$\mathbf{Z} = \text{vec} \left(Z_t^{(\ell)}(A) : \ell \in \mathcal{L}, t \in \mathcal{T}, A \in \mathcal{D}_{t\ell} \right),$$

$$\mathbf{H} = \text{rbind} \left(h_t^{(\ell)}(A)^T : \ell \in \mathcal{L}, t \in \mathcal{T}, A \in \mathcal{D}_{t\ell} \right),$$

$$\mathbf{S} = \text{rbind} \left(\psi_t^{(\ell)}(A)^T : \ell \in \mathcal{L}, t \in \mathcal{T}, A \in \mathcal{D}_{t\ell} \right),$$

$$\boldsymbol{\xi} = \text{vec} \left(\xi_t^{(\ell)}(A) : \ell \in \mathcal{L}, t \in \mathcal{T}, A \in \mathcal{D}_{t\ell} \right),$$

$$\boldsymbol{\varepsilon} = \text{vec} \left(\varepsilon_t^{(\ell)}(A) : \ell \in \mathcal{L}, t \in \mathcal{T}, A \in \mathcal{D}_{t\ell} \right),$$

$$\mathbf{V} = \text{Diag} \left(\sigma_{t\ell}^2(A) : \ell \in \mathcal{L}, t \in \mathcal{T}, A \in \mathcal{D}_{t\ell} \right),$$

and $h_t^{(\ell)}(A) \equiv h(A)$.

Basis Functions: Dimension Reduction

- The presence of multicollinearity can severely hinder convergence of the Markov-Chain Monte Carlo (MCMC) sampler.
- To protect against multicollinearity, we reduce the $n \times r$ matrix \mathbf{S} using principal components analysis.
- Suppose $\mathbf{U}\mathbf{D}\mathbf{U}^\top$ is the eigendecomposition of $\mathbf{S}^\top\mathbf{S}$, and $\tilde{\mathbf{U}}$ contains the r' columns of \mathbf{U} corresponding the $r' \leq r$ largest magnitude eigenvalues in \mathbf{D} .
- The transformation $T(\mathbf{S}) = \mathbf{S}\tilde{\mathbf{U}}^\top$ is applied to all matrices computed from the basis functions.

Gibbs Sampler

- $[\mu_B \mid \bullet] \sim N(\vartheta_\mu, \Omega_\mu^{-1}),$

$$\vartheta_\mu = \Omega_\mu^{-1} H^\top V^{-1} (Z - S\eta - \xi),$$

$$\Omega_\mu = H^\top V^{-1} H + \sigma_\mu^{-2} I.$$

- $[\eta \mid \bullet] \sim N(\vartheta_\eta, \Omega_\eta^{-1}),$

$$\vartheta_\eta = \Omega_\eta^{-1} S^\top V^{-1} (Z - H\mu_B - \xi),$$

$$\Omega_\eta = S^\top V^{-1} S + \sigma_\eta^{-2} \tilde{K}^{-1}.$$

- $[\xi \mid \bullet] \sim N(\vartheta_\xi, \Omega_\xi^{-1}),$

$$\vartheta_\xi = \Omega_\xi V^{-1} (Z - H\mu_B - S\eta),$$

$$\Omega_\xi^{-1} = V^{-1} + \sigma_\xi^{-2} I.$$

- $[\sigma_\mu^2 \mid \bullet] \sim IG(\alpha_\mu, \beta_\mu), \alpha_\mu = a_\mu + n_B/2 \text{ and } \beta_\mu = b_\mu + \mu_B^\top \mu_B/2.$

- $[\sigma_K^2 \mid \bullet] \sim IG(\alpha_K, \beta_K), \alpha_K = a_K + r/2 \text{ and } \beta_K = b_K + \eta^\top \tilde{K}^{-1} \eta/2.$

- $[\sigma_\xi^2 \mid \bullet] \sim IG(\alpha_\xi, \beta_\xi), \alpha_\xi = a_\xi + N/2 \text{ and } \beta_\xi = b_\xi + \xi^\top \xi/2.$

STCOS R Package

- The `stcos` R package facilitates application of the model.
 1. Preprocess: Prepare Z , V , H , S , and \tilde{K}^{-1} needed to fit the model.
 2. Fit the model via Gibbs sampler.
 3. Postprocess: Make estimates and predictions on user-defined domains using MCMC draws.
- Preprocess once for a given set of source supports. Redo model fit for each survey variable of interest. Redo postprocess for each target support of interest.
- We depend on several other R packages.
 1. Manipulation of shapefiles via the `sf` package (Pebesma, 2017).
 2. Object-oriented programming using the `R6` package (Chang, 2017).
 3. Ability to call C++ code for performance via `Rcpp` (Eddelbuettel, 2013) and `RcppArmadillo` (Eddelbuettel and Sanderson, 2014).
 4. Sparse matrix computations in R via `Matrix` package (Bates and Maechler, 2017).
- Development version of `ggplot2` (Wickham, 2016) can plot `sf` objects.

Source Supports: Shapefiles with Estimates

User provides all supports as shapefiles. Direct estimates and variance estimates should be embedded into source supports.

```
R> library(sf)
R> acs5.2013 <- st_read("county_acs_5yr2013.shp")
R> head(acs5.2013)
Simple feature collection with 6 features and 9 fields
geometry type:  MULTIPOLYGON
dimension:      XY
bbox:           xmin: -9799374 ymin: 3532006 xmax: -9468076 ymax: 4063675
epsg (SRID):   3857
proj4string:   +proj=merc +a=6378137 +b=6378137 +lat_ts=0.0 +lon_0=0.0 +x_0=0.0
                +y_0=0 +k=1.0 +units=m +nadgrids=@null +wktext +no_defs
  GEO_ID STATE COUNTY     NAME    LSAD SHAPE_AREA SHAPE_LEN DirectEst
1 0500000US01001     01    001 Autauga County 2202587903  235761.0    53682
2 0500000US01003     01    003 Baldwin County 5913339907  493065.0    50221
3 0500000US01005     01    005 Barbour County 3262897491  275539.8    32911
4 0500000US01007     01    007 Bibb County   2309817471  223844.6    36447
5 0500000US01009     01    009 Blount County 2454704099  258633.9    44145
6 0500000US01011     01    011 Bullock County 2258507149  233574.8    32033
  DirectVar          geometry
1  831625.9 MULTIPOLYGON((-9675622.416...
2  915703.6 MULTIPOLYGON((-9799373.733...
3  4437638.9 MULTIPOLYGON((-9544726.836...
4  4323124.1 MULTIPOLYGON((-9731765.399...
5  2150305.1 MULTIPOLYGON((-9680577.914...
6  21959826.2 MULTIPOLYGON((-9573382.365...
```

Load Source Supports

```
library(sf)
library(stcos)

# Fine-level support comes from ACS 5-year estimates for 2015
dom.fine <- st_read("shp/county_acs_5yr2015.shp")

# ACS 1-year source supports
acs1.2005 <- st_read("shp/county_acs_1yr2005.shp")
acs1.2006 <- st_read("shp/county_acs_1yr2006.shp")
...
acs1.2015 <- st_read("shp/county_acs_1yr2015.shp")

# ACS 3-year source supports
acs3.2007 <- st_read("shp/county_acs_3yr2007.shp")
acs3.2008 <- st_read("shp/county_acs_3yr2008.shp")
...
acs3.2013 <- st_read("shp/county_acs_3yr2013.shp")

# ACS 5-year source supports
acs5.2009 <- st_read("shp/county_acs_5yr2009.shp")
acs5.2010 <- st_read("shp/county_acs_5yr2010.shp")
...
acs5.2015 <- st_read("shp/county_acs_5yr2015.shp")
```

Prepare Basis

```
1 library(fields)
2
3 # Spatial knots are selected via space-filling design
4 u <- st_sample(dom.fine, size = 5000)
5 M <- matrix(unlist(u), length(u), 2, byrow = TRUE)
6 out <- cover.design(M, 500)
7 knots.sp <- out$design
8
9 # Temporal knots are selected to be evenly spaced
10 knots.t <- c(2005, 2005.5, 2006, 2006.5, 2007, 2007.5, 2008, 2008.5,
11     2009, 2009.5, 2010, 2010.5, 2011, 2011.5, 2012, 2012.5, 2013, 2013.5,
12     2014, 2014.5, 2015)
13
14 # Combined spatio-temporal knots
15 knots <- merge(knots.sp, knots.t)
16 names(knots) <- c("x", "y", "t")
17
18 # Create a Basis object
19 basis <- SpaceTimeBisquareBasis$new(knots[,1], knots[,2], knots[,3], w.s = 1, w.t = 1)
```

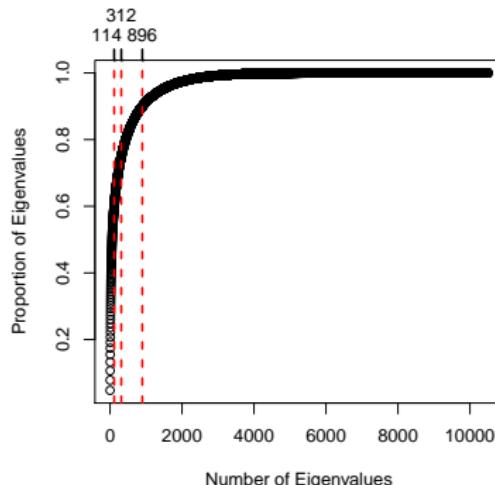
Construct an STCOSPrep Object

```
sp <- STCOSPrep$new(fine_domain = dom.fine, fine_domain_geo_name = "GEO_ID",
  basis = basis, basis_mc_reps = 500)
sp$add_obs(acs1.2015, period = 2015, estimate_name = "DirectEst",
  variance_name = "DirectVar", geo_name = "GEO_ID")
sp$add_obs(acs3.2013, period = 2011:2013, estimate_name = "DirectEst",
  variance_name = "DirectVar", geo_name = "GEO_ID")
sp$add_obs(acs5.2013, period = 2009:2013, estimate_name = "DirectEst",
  variance_name = "DirectVar", geo_name = "GEO_ID")
...
Z <- sp$get_Z()
V <- sp$get_V()
H <- sp$get_H()
S <- sp$get_S()
```

```
R> sp$add_obs(acs1.2015, period = 2015, estimate_name = "DirectEst",
  variance_name = "DirectVar", geo_name = "GEO_ID")
2017-07-13 15:03:18 - Begin adding observed space-time domain
2017-07-13 15:03:18 - Computing overlap matrix using field 'GEO_ID'
2017-07-13 15:03:22 - Computing basis functions
2017-07-13 15:03:30 - Computing basis for area 100 of 812
2017-07-13 15:03:37 - Computing basis for area 200 of 812
...
2017-07-13 15:04:12 - Computing basis for area 700 of 812
2017-07-13 15:04:19 - Computing basis for area 800 of 812
2017-07-13 15:04:20 - Extracting survey estimates from field 'DirectEst'
  and variance estimates from field 'DirectVar'
2017-07-13 15:04:20 - Finished adding observed space-time domain
```

Dimension Reduction for S

```
1 eig <- eigen(t(S) %*% S)
2 rho <- eig$values
3
4 idx.S <- which(cumsum(rho) / sum(rho) < 0.6)
5 Tx.S <- t(eig$vectors[idx.S,])
6 f <- function(S) { S %*% Tx.S }
7 sp$set_basis_reduction(f)
8
9 S.reduced <- sp$get_reduced_S()
```



Specification for K

- Random Walk

```
K.inv <- sp$get_Kinv(2005:2015)
```

- Spatial-only

```
K.inv <- sp$get_Kinv(2005:2015, autoreg = FALSE)
```

- Independence

```
K.inv <- diag(x = 1, nrow = ncol(S.reduced))
```

- VAR(1)

```
1 # Take a spatial-only basis expansion of the fine-level support, and project
2 # orthogonally to this design matrix.
3 sp.basis <- SpatialBisquareBasis$new(knots.sp[,1], knots.sp[,2], w = 1)
4 X <- compute_sp_basis_mc(basis = sp.basis, domain = dom.fine,
5   R = 500, report.period = 100)
6 K.inv <- sp$get_Kinv(2005:2015, X)
```

Gibbs Sampler

```
1 # Std'ize before MCMC
2 D <- Diagonal(n = length(Z), x = 1/sd(Z))
3 Z.scaled <- (Z - mean(Z)) / sd(Z)
4 V.scaled <- V / var(Z)
5
6 # Use MLE as initial value for MCMC
7 mle.out <- mle.stcos(Z.scaled, S.reduced, V.scaled, H, init = list(sig2xi = 1))
8 init <- list(
9   sig2xi = mle.out$sig2xi.hat,
10  mu_B = mle.out$mu.hat,
11  eta = mle.out$eta.hat
12 )
13
14 # Gibbs Sampler
15 gibbs.out <- gibbs.stcos.raw(Z.scaled, S.reduced, V.scaled, K.inv, H, R = 10000,
16   report.period = 100, burn = 1000, thin = 10, init = init)
```

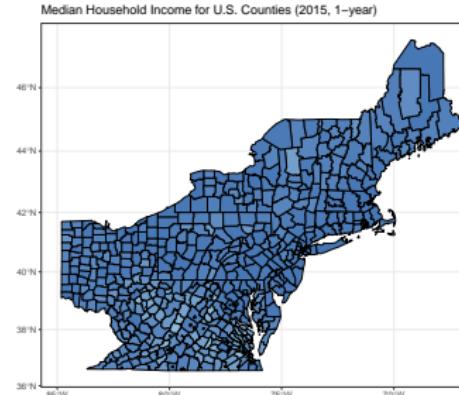
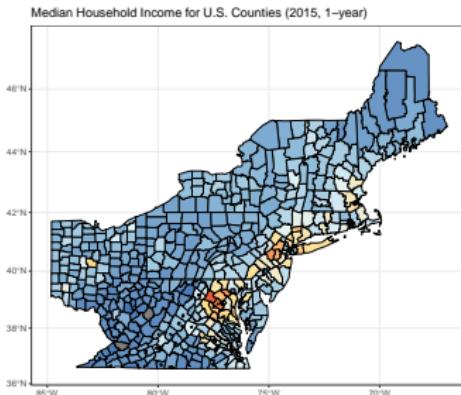
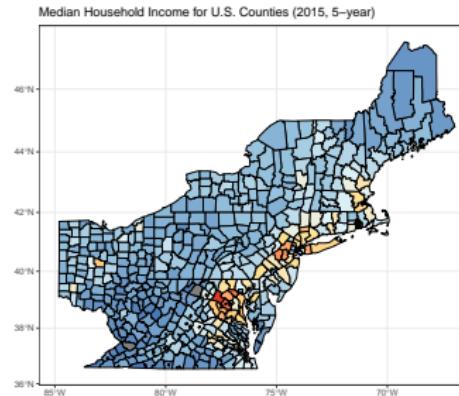
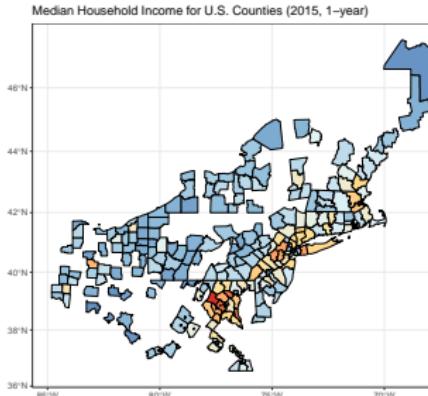
```
2017-07-06 13:21:58 - Begin Gibbs sampler
2017-07-06 13:24:09 - Begin iteration 100
2017-07-06 13:26:17 - Begin iteration 200
...
2017-07-06 16:46:22 - Begin iteration 10000
2017-07-06 16:46:23 - Finished Gibbs sampler
```

Estimation & Prediction on Target Support

```
1 # Load a target support and transform to fine-level support's projection
2 aiannh.2015 <- st_read("shp/cb_2015_us_aiannh_500k.shp")
3 dom <- st_transform(aiannh.2015, crs = st_crs(dom.fine))
4
5 # Compute H and S matrices
6 target.out <- sp$domain2model(dom, period = 2015, geo_name = "AFFGEOID")
7
8 # Posterior distribution for E(Y)
9 E.hat.scaled <- fitted(gibbs.out, target.out$H, target.out$S.reduced)
10 E.hat <- sd(Z) * E.hat.scaled + mean(Z) # Uncenter and unscale
11 dom$E.mean <- colMeans(E.hat) # Point estimates
12 dom$E.sd <- apply(E.hat, 2, sd) # SDs
13 dom$E.lo <- apply(E.hat, 2, quantile, prob = 0.025) # Credible interval lo
14 dom$E.hi <- apply(E.hat, 2, quantile, prob = 0.975) # Credible interval hi
15
16 # Posterior predictive distribution of Y
17 Y.pred.scaled <- predict(gibbs.out, target.out$H, target.out$S.reduced)
18 Y.pred <- sd(Z) * Y.pred.scaled + mean(Z) # Uncenter and unscale
19 dom$PP.mean <- colMeans(Y.pred) # Point estimates
20 dom$PP.sd <- apply(Y.pred, 2, sd) # SDs
21 dom$PP.lo <- apply(Y.pred, 2, quantile, prob = 0.025) # Prediction interval lo
22 dom$PP.hi <- apply(Y.pred, 2, quantile, prob = 0.975) # Prediction interval hi
23
24 > dom
```

	AFFGEOID	NAME	ALAND	AWATER	E.mean	PP.mean	...
84	2500000US9865	Sappony	111355719	10556265	34672.74	34596.95	...
92	2500000US9260	Pamunkey	4422053	1958365	66197.45	66215.14	...
323	2500000US9230	Mattaponi	279160	31696	64608.35	64619.12	...
470	2500000US9580	Chickahominy	134292954	64589	49225.30	49279.00	...
673	2500000US9675	Eastern Chickahominy	5781366	0	70456.80	70523.77	...

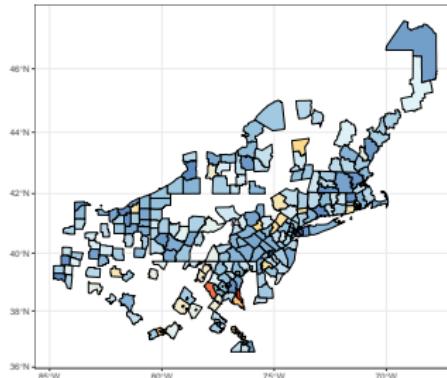
Direct vs. Model Estimates (2015)



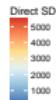
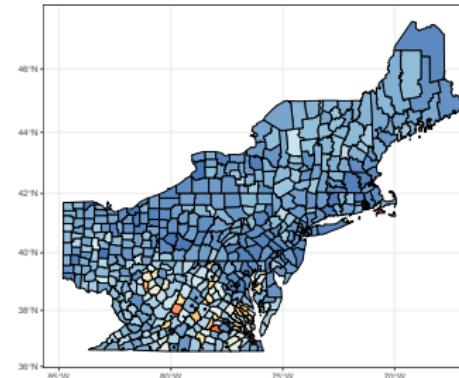
Model results shown are based on MCMC draws of $E(\mathbf{Y} | \boldsymbol{\theta}) = \mathbf{H}\boldsymbol{\mu}_B + \mathbf{S}\boldsymbol{\eta}$.

Direct vs. Model SDs (2015)

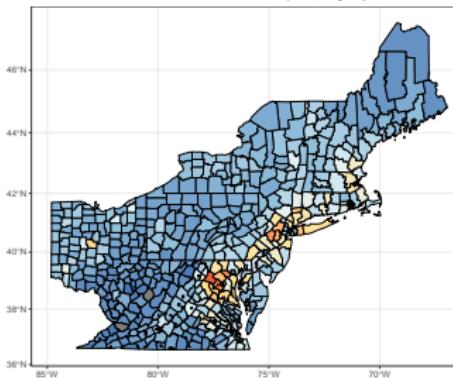
Median Household Income for U.S. Counties (2015, 1-year)



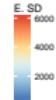
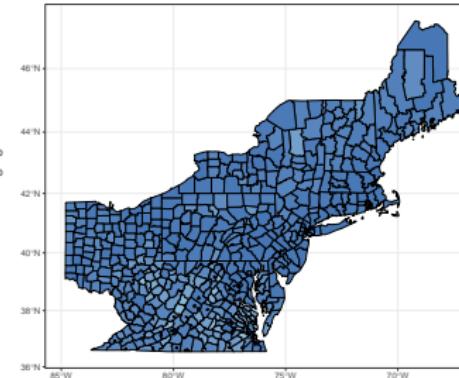
Median Household Income for U.S. Counties (2015, 5-year)



Median Household Income for U.S. Counties (2015, 1-year)

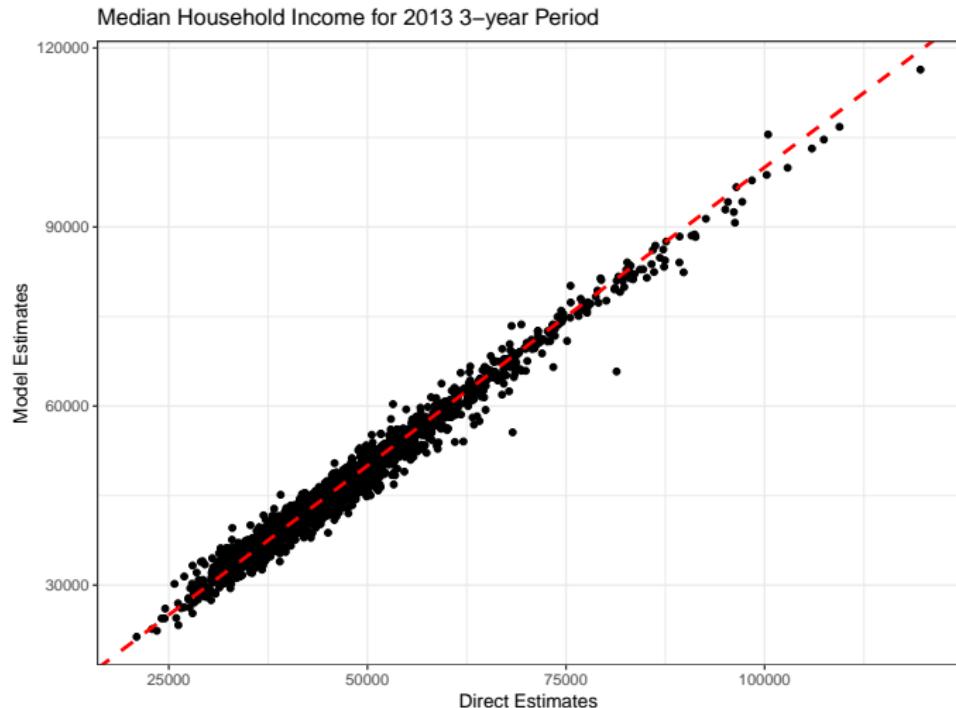


Median Household Income for U.S. Counties (2015, 1-year)



Model results shown are based on MCMC draws of $E(\mathbf{Y} | \boldsymbol{\theta}) = \mathbf{H}\boldsymbol{\mu}_B + \mathbf{S}\boldsymbol{\eta}$.

Direct vs. Model Estimates (2013, 3-year, Counties)



2013 3-year source support was excluded from model fit (Leave One Out).

Conclusions

- We reviewed the STCOS methodology from Bradley et al. (2015) and demonstrated the `stcos` R package.
- Improvements are underway to lower programming burden for users, and to increase performance (speed, memory usage) where possible.
- We expect a CRAN submission and a companion article to be available within a few months.
- We hope that this software will empower users to explore official statistics on custom geographies and time periods.

Contact: andrew.raim@census.gov



References I

- Douglas Bates and Martin Maechler. *Matrix: Sparse and Dense Matrix Classes and Methods*, 2017. Available online at
<https://cran.r-project.org/package=Matrix>. R package version 1.2-10.
- Jonathan R. Bradley, Christopher K. Wikle, and Scott H. Holan. Spatio-temporal change of support with application to American Community Survey multi-year period estimates. *Stat*, 4(1):255–270, 2015.
- Winston Chang. *R6: Classes with Reference Semantics*, 2017. Available online at
<https://cran.r-project.org/package=R6>. R package version 2.2.2.
- Dirk Eddelbuettel. *Seamless R and C++ Integration with Rcpp*. Springer, 2013.
- Dirk Eddelbuettel and Conrad Sanderson. Rcpparmadillo: Accelerating r with high-performance C++ linear algebra. *Computational Statistics and Data Analysis*, 71:1054–1063, March 2014.
- Nicholas J. Higham. Computing a nearest symmetric positive semidefinite matrix. *Linear Algebra and its Applications*, 103:103–118, 1988.
- John Hughes and Murali Haran. Dimension reduction and alleviation of confounding for spatial generalized linear mixed models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(1):139–159, 2013.

References II

- Douglas Nychka, Reinhard Furrer, John Paige, and Stephan Sain. *fields: Tools for spatial data*. University Corporation for Atmospheric Research, Boulder, CO, USA, 2015. Available online at www.image.ucar.edu/fields. R package version 9.0.
- Edzer Pebesma. *sf: Simple Features for R*, 2017. Available online at <https://cran.r-project.org/package=sf>. R package version 0.5-1.
- U.S. Census Bureau. American Community Survey data suppression, September 2016. Available online at <https://www.census.gov/programs-surveys/acs/technical-documentation/data-suppression.html>.
- Hadley Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York, 2016.