# The Number of Undocumented Immigrants in the United States

Mohammad M. Fazel-Zarandi, Jonathan S. Feinstein, Edward H. Kaplan

# Motivation

# Research Question and Contribution

- Questions:
    1. How many undocumented immigrants are there in the United States?
    2. Are the current methods of estimating the number of undocumented immigrants adequate?

# Research Question and Contribution

- Questions:
    1. How many undocumented immigrants are there in the United States?
    2. Are the current methods of estimating the number of undocumented immigrants adequate?

- Contribution: Propose a new approach grounded on operational data and mathematical modeling that estimates annual population inflows and outflows from 1990 – 2016

# Research Question and Contribution

- Questions:
    1. How many undocumented immigrants are there in the United States?
    2. Are the current methods of estimating the number of undocumented immigrants adequate?

- Contribution: Propose a new approach grounded on operational data and mathematical modeling that estimates annual population inflows and outflows from 1990 – 2016

- Why: Sets the scale of the issue

# Outline

- Motivation

- Current estimates

- Snapshot of our results

- The model and the simulation

- Results, Receptions, and Policy Implications

# Current Estimates

- Residual Method

  - Passel (2016), Krogstad and Passel (2015), Baker and Rytina (2013), Warren and Warren (2013)

$$\begin{matrix} \text{Estimated Number of} \\ \text{Unauthorized Immigrants} \end{matrix} = \begin{matrix} \text{Estimated Total Foreign} \\ \text{Born Population} \\ (Non-Citizen) \end{matrix} - \begin{matrix} \text{Estimated Lawful} \\ \text{Immigrant Population} \end{matrix}$$

# Current Estimates

- Total Foreign Born Population
  - Based on *surveys* (American Community Survey or Current Population Survey)

# Current Estimates

- Total Foreign Born Population
    - Based on *surveys* (American Community Survey or Current Population Survey)


- Estimate of Lawful Immigrant Population
    - Use Department of Homeland Security data on lawful arrivals

# Current Estimate



Source: Pew Research Center

# Current Estimates

# Hidden Population

- Difficult to locate members of the target population
  (Goel and Salganik (2010), Crawford et al.(2018))

    - Reaching a representative sample of all those born outside of the U.S.
        - Undocumented immigrants are more difficult to locate and survey

    - Accurate responses from survey respondents
        - Undocumented immigrants may misreport their country of origin, citizenship, and number of household residents

# Some Relevant Statistics

ACS response rates

| Year | Response Rate | Refusal | Unable to Locate | No One Home | Temporarily Absent | Language Problem | Insufficient Data | Maximum Contact Attempts Reached | Other |
|------|--------------|---------|------------------|-------------|--------------------|------------------|-------------------|----------------------------------|-------|
| 2017 | 93.7 | 2.7 | 0.0 | 0.9 | 0.1 | 0.1 | 0.4 | 1.1 | 0.9 |
| 2016 | 94.7 | 2.1 | 0.0 | 0.9 | 0.1 | 0.1 | 0.3 | 1.1 | 0.7 |

| Population: Origin and Language [1] | Percent Allocated | | | | |
|-------------------------------------|-------|-------|-------|-------|-------|
| Item | 2017 | 2016 | 2015 | 2014 | 2013 |
| Place of birth total population | 9.3 | 9.1 | 9.2 | 8.8 | 8.6 |
| Sex total population | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |

ACS Question **non**-response rates among respondents

Source:  U.S. Census Bureau

# Some Relevant Statistics

    5% non-response

+ 0.95 × 8% question skippers
_____

    12.6% no clear answer  ≈  40 million people


- Ignoring deliberate misreporting
  - Place birth
  - Number of household residents

# Some Relevant Statistics

5% non-response

+ 0.95 × 8% question skippers
_____

12.6% no clear answer ≈ 40 million people

- Ignoring deliberate misreporting
  - Place birth
  - Number of household residents

- Non-response bias
  - Missing at random vs *missing on purpose*

**Census may be too blunt an instrument to reach a relatively small population with an incentive to remain undetected**

# Census Citizenship Question Debate



The New York Times

**Court Blocks Trump Administration From Asking About Citizenship in Census**

Commerce Secretary Wilbur L. Ross Jr., center, ordered the Census Bureau to add a citizenship question to the 2020 census. Doug Mills/The New York Times
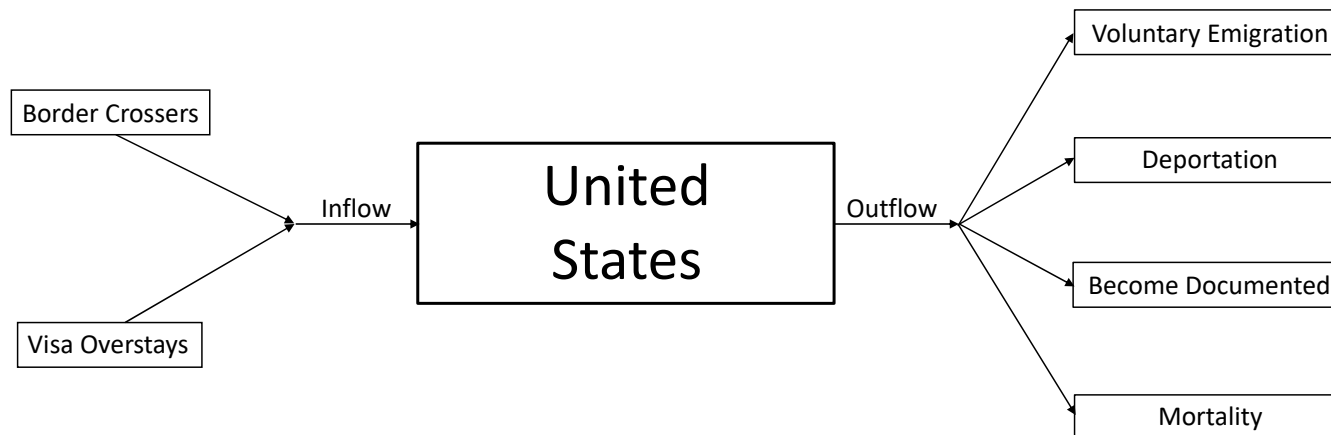
By Michael Wines

Jan. 15, 2019

557

"The result will not only be a decrease in the quality of census data — something Defendants concede — but likely also a net differential undercount (that is, an undercount of certain sectors of the population, including people who live in households containing noncitizens and Hispanics, relative to others)."

Judge Jesse M. Furman (United States District Court in Manhattan)
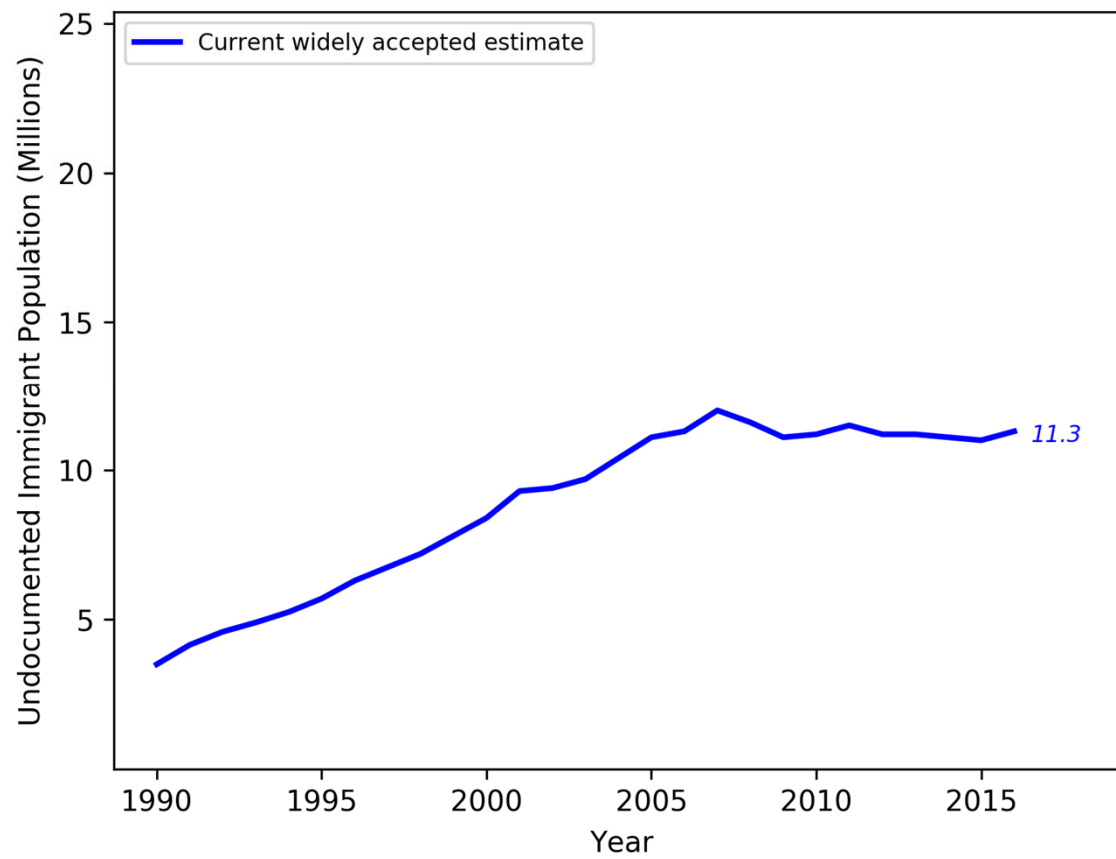
# Our Estimate

- Combine mathematical modeling with data analysis

- Our model tracks and estimates annual inflows and outflows from 1990 - 2016
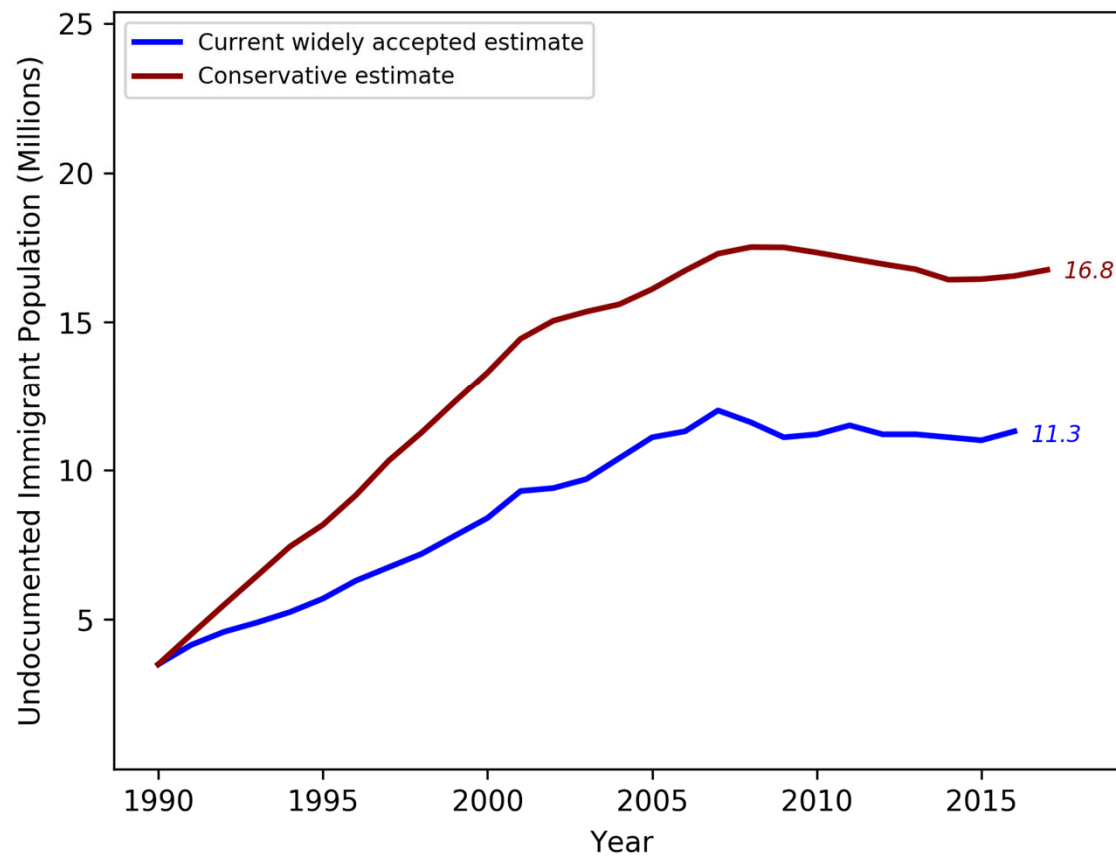
# Our Estimation Result

I.  Generate a <span style="color:red">conservative estimate</span> (low-end) of the number of undocumented immigrants

II. Generate probability distribution over the number of undocumented immigrants based on simulating our model over a wide range of assumptions

# Results
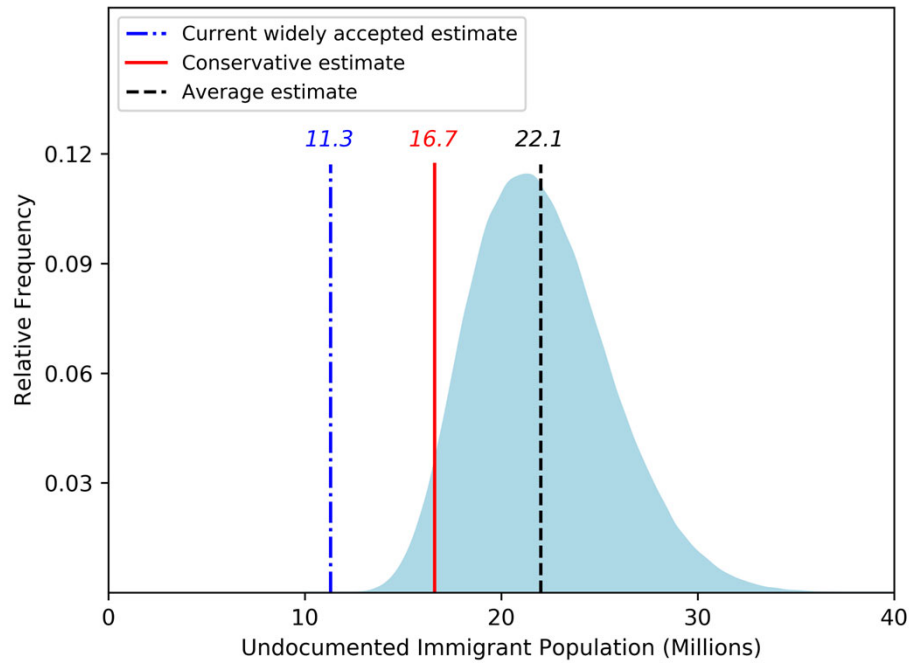
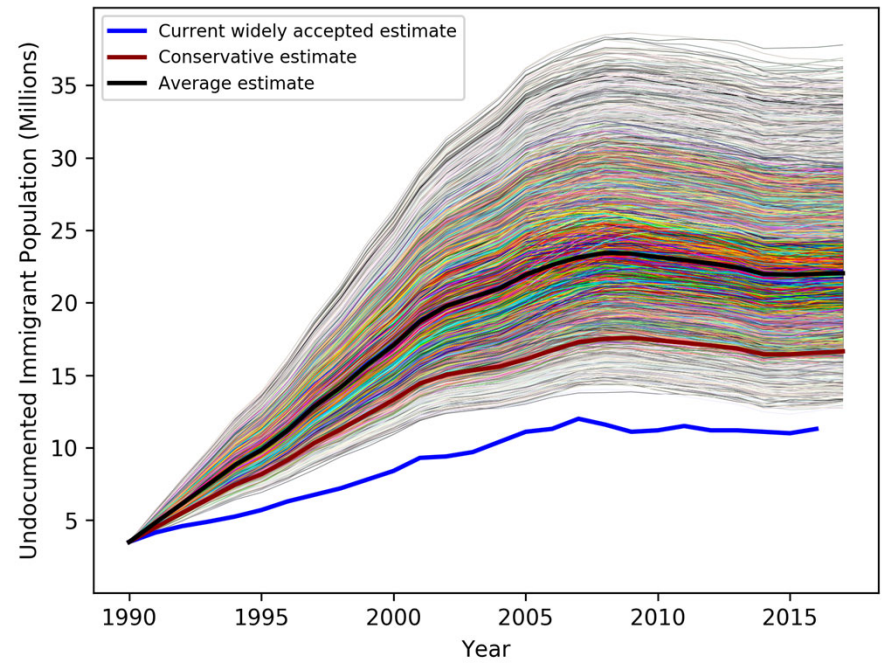# Results – Conservative Estimate

# Results – 1,000,000 Simulation Runs



Probability Distribution (2016)



Trajectories

# The Approach

# The Model

Inflow → United States → Outflow

Number of undocumented immigrants in year $t$ → 

Outflow in year $t$ →

$$N_t = N_{t-1} + I_t - O_t$$

Inflow in year $t$ →

# The Model

Border Crossers

Visa Overstays

Inflow

## United States

Outflow

# The Model

Border Crossers

Visa Overstays

Inflow → United States → Outflow

Voluntary Emigration

Deportation

Become Documented

Mortality

# Estimation Strategy

- Generate a *conservative estimate* (low-end estimate)

  - Underestimate Inflows

  - Overestimate Outflows

# Inflows

# The Model

Border Crossers

Visa Overstays

Inflow

## United States

Outflow

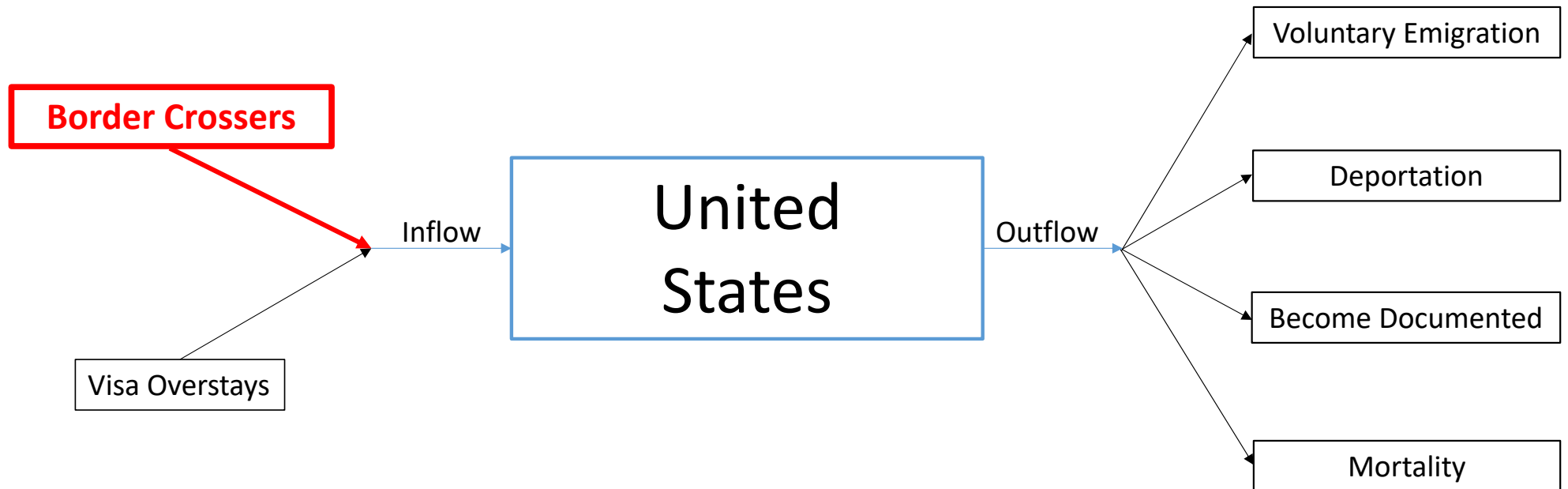Voluntary Emigration

Deportation

Become Documented

Mortality

# Illegal Border Crossers

- Really hard component of the model

- To estimate this we use apprehension data on the Mexico - U.S. border

- Our model of the border crossers is based on standard repeated trials (Espenshade (1995)), i.e. a *Bernoulli process model*

- The model estimates the number of border crossings consistent with observed apprehensions that are happening in the data

# Number of Apprehensions



Data Source: U.S. Border Patrol

# Frequency of Apprehension (2005)



Data Source: U.S. Department of Homeland Security

# Frequency of Apprehension (2005)



Data Source: U.S. Department of Homeland Security

# The Border Crossing Model



$1 - p_t$

$p_t$

$C_t$

Probability of
apprehension
at year t

Number of individuals
attempting to cross the
border at year $t$

# The Border Crossing Model

$1 - p_t$

Number of
apprehensions

$p_t$

1

$C_t$

$C_t p_t$

# The Border Crossing Model



$1 - p_t$

$p_t$

$1 - d_t$

1

$C_t$

$C_t p_t$

$d_t$

Deterrence due to apprehension (estimated from Encuesta sobre Migración en la Frontera Norte y Sur de México (EMIF))

# The Border Crossing Model



$$1 - p_t$$

$$C_t p_t (1 - d_t)$$

$$p_t$$

$$1 - d_t$$

$$1$$

$$C_t$$

$$C_t p_t$$

$$d_t$$

# The Border Crossing Model

# The Border Crossing Model

# The Border Crossing Model



Diagram: $C_t \xrightarrow{p_t} 1$ (with $1 - p_t$ branch up, labeled arrow), then $1 \xrightarrow{1 - d_t}$ with $d_t$ branch down. $C_t p_t$. Then $C_t p_t (1 - d_t) \xrightarrow{p_t} 2$ (with $1 - p_t$ branch up), $C_t p_t^2 (1 - d_t)$, then $2 \xrightarrow{1 - d_t}$ with $d_t$ branch down $\dots$

Number of Apprehensions in Year $t$ $\longrightarrow$ $A_t = \underbrace{C_t p_t}_{First\ Attempt} + \underbrace{C_t p_t^2 (1 - d_t)}_{Second\ Attempt} + \cdots$

# The Border Crossing Model



$$C_t \xrightarrow{p_t} 1 \xrightarrow{1-d_t} \quad C_t p_t (1-d_t) \xrightarrow{p_t} 2 \xrightarrow{1-d_t}$$

$1 - p_t$

$C_t p_t$

$C_t p_t^2 (1-d_t)$

$d_t$

$\ldots$

$$A_t = \underbrace{C_t p_t}_{First\ Attempt} + \underbrace{C_t p_t^2 (1-d_t)}_{Second\ Attempt} + \cdots$$

Expected Number of Apprehensions in Year $t$ $\longrightarrow$

$$A_t = \underbrace{C_t}_{\substack{Number \\ of \\ attempters}} \times \underbrace{\frac{1}{1 - p_t(1-d_t)}}_{\substack{Expected \\ number\ of \\ attempts}} \times \underbrace{p_t}_{\substack{Probability\ of \\ apprehension \\ in\ each\ attempt}}$$

# Mean Apprehensions Over
# All Attempted Border Crossers

- Let $p$ = Pr{Apprehension per crossing attempt}
- Let $d$ = Pr{Deterrence | Apprehension}
- Let $a$ = mean number of apprehensions over all attempted border crossers

Number of Apprehensions

$1 - p$

0

Successful Crossing
(Fail to Apprehend)

$a =$

$d$

1

$p$

Deterred
Quit

Apprehension

$1 - d$

$1 + a$

Not Deterred
Try again

$a = (1 - p) \times 0 + pd \times 1 + p(1 - d) \times (1 + a)$

$$a = \frac{p}{1 - p(1 - d)}$$

If $C$ persons are attempting to cross, then

$$E(Apprehensions) = Ca = C\,\frac{p}{1 - p(1 - d)}$$

# What Does the Model Imply So Far?

- Number of apprehensions at year $t$: $\qquad\qquad A_t = C_t \dfrac{p_t}{1 - p_t(1 - d_t)}$

- Number of attempters at year $t$: $\qquad\qquad C_t = A_t \dfrac{1 - p_t(1 - d_t)}{p_t}$

- Number of undocumented border crossers at year $t$: $B_t = C_t - Q_t$ $\longleftarrow$ Individuals who give up

- Number of individuals deterred at year $t$: $\qquad Q_t = A_t d_t$

- Number of undocumented border crossers at year $t$: $\boldsymbol{B_t = A_t \dfrac{1 - p_t}{p_t}}$

# We Need The Apprehension Probability $p$

- Define $\bar{A}$ = *recidivist* apprehensions (that is, # 2nd or higher apprehension)
- Recall $A$ = total apprehensions
- All apprehensions = First time apprehensions + recidivist apprehensions
- $A = Cp + \bar{A}$
- Recall $C = A\dfrac{1-p(1-d)}{p}$
- ==> $\bar{A} = A - Cp = A - A(1 - p(1 - d)) = Ap\,(1 - d)$
- Solve for $p$:  $p = \dfrac{\bar{A}}{A(1-d)}$
- Note: $A(1 - d)$ is total undeterred apprehensions, which equals the number of crossing opportunities for *recidivist* apprehensions

# What Does the Model Imply So Far?

- Number of apprehensions at year $t$:

$$A_t = C_t \frac{p_t}{1 - p_t(1 - d_t)}$$

- Number of attempters at year $t$:

$$C_t = A_t \frac{1 - p_t(1 - d_t)}{p_t}$$

- Number of individuals deterred at year $t$:

$$Q_t = A_t d_t$$

- Number of undocumented border crossers at year $t$:
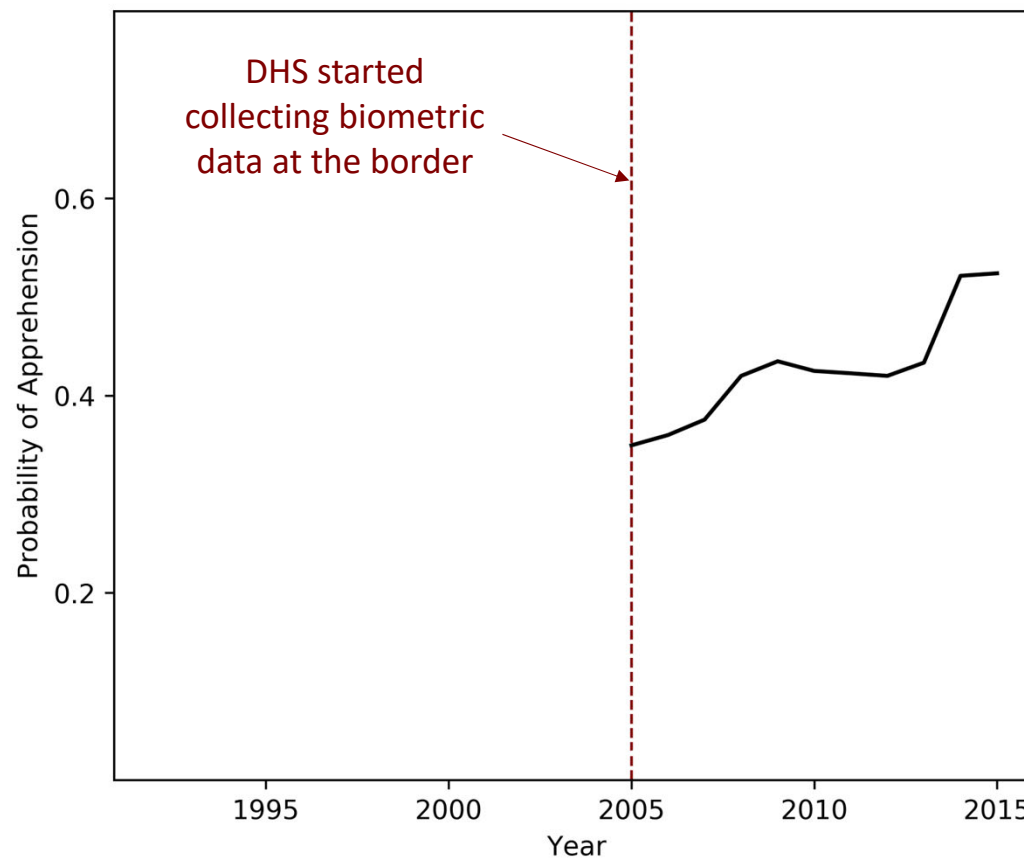
$$B_t = A_t \frac{1 - p_t}{p_t}$$

- The probability of apprehension at year $t$:

Number of recidivist apprehensions

$$p_t = \frac{\bar{A}_t}{A_t(1 - d_t)}$$

# Probability of Apprehension

# Probability of Apprehension

# Probability of Apprehension

- What does the literature says?

  - Massey et al. (2016):
    - Start from the low 20% range in the 1990s ranging upward to approximately 30% in the earlier 2000s

  - Wein and Motskin (2009):
    - Estimate the 2003 rate at around 20%

# Probability of Apprehension

# The Model

Border Crossers

**Visa Overstays**

Inflow → United States → Outflow

Voluntary Emigration

Deportation

Become Documented

Mortality

# Visa Overstays

- Non-immigrants who are admitted to the U.S. lawfully, but do not leave after the period during which they have been allowed to remain in the U.S. legally ends

- Comprehensively measured by Department of Homeland Security starting in 2016

# Number of Non-Immigrant Visas Issued



Data Source: U.S. Department of State

# Visa Overstays

$V_j$ : Number of visas issued in year $j$

$S_j$ : Number of visa overstays in year $j$

$r_j$: Visa overstay rate in year $j$



$$S_j = r_j \times V_j$$

# Visa Overstays

$V_j$ : Number of visas issued in year $j$

$S_j$ : Number of visa overstays in year $j$

$r_j$: Visa overstay rate in year $j$



$(1 - r_j)$

$V_j \longrightarrow \boxed{S_j}$

2016 visa overstay rate

$$S_j = r_j \times V_j \qquad r_j = \frac{S_{2016}}{V_{2016}}, \qquad \forall j$$

# Visa Overstays

$V_j$ : Number of visas issued in year $j$

$S_j$ : Number of visa overstays in year $j$
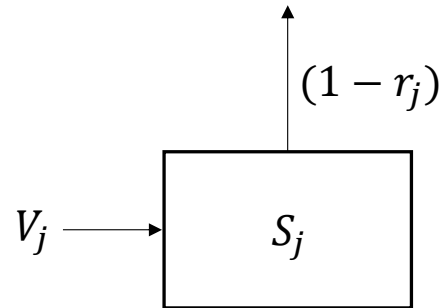
$r_j$: Visa overstay rate in year $j$

$(1 - r_j)$

$V_j \longrightarrow \boxed{S_j}$

2016 visa overstay rate

$$S_j = r_j \times V_j \qquad r_j = \frac{S_{2016}}{V_{2016}}, \qquad \forall j$$

**Calibration:**

Fraction of undocumented immigrants arriving in year $j$ still in the US in 2015

$$\sum_{j=1}^{t=26} S_j \times Pr\{\tau_j \geq t - j\} \quad < \quad$$

Number of overstayers in the current widely accepted estimate for 2015

Model estimate of the population of overstayers in 2015

# Outflows

# The Model

# Outflows

- Voluntary Emigration

- Deportations

- Become Documented
  - Including DACA in the outflows

- Mortality

# Outflows

- Voluntary Emigration

- Deportations

- Become Documented
  - Including DACA in the outflows

- Mortality

# Voluntary Emigration

- Emigration decreases with time spent in the country
  - Van Hook and Zhang (2011), Bhaskar et al. (2013), Warren and Warren (2013)

- For each undocumented immigrant we must keep track of duration in the country

  - Duration-dependent emigration rate

# Voluntary Emigration

Undocumented Immigrant Population → Visitors → Less than 1 year

Undocumented Immigrant Population → Residents → Less than 10 years

Residents → More than 10 years

# Voluntary Emigration – Simplified Example

Visitors → Less than 1 year → $\mu_s$

Undocumented Immigrant Population

Residents → Less than 10 years → $\mu_m$

Residents → More than 10 years → $\mu_l$

$I_t(1-\mu_s)(1-\mu_m)^9(1-\mu_l)$

$I_t(1-\mu_s)(1-\mu_m)^2$ … $I_t(1-\mu_s)(1-\mu_m)^9$

$I_t(1-\mu_s)(1-\mu_m)$

$I_t(1-\mu_s)$

$I_t$

Year 1   Year 2   Year 3   …   Year 10   Year 11   …

# Voluntary Emigration – Simplified Example

# Voluntary Emigration - Data

- Residents: use *largest* values in published academic and government sources

  o Ahmed and Robinson (1994), Mulder (2003), Van Hook and Zhang (2011), Baker and Rytina (2013), Warren and Warren (2013), Bhaskar et al. (2013)

- Visitors: use data on first-year exit rate for visa overstayers
  o The assumed emigration rates are an upper bound for border crossers (Massey et al (2002), Massey et al. (2016))

- The emigration rates used significantly overestimate outflows

# The Final Formula

$$N_t = N_{t-1} + I_t - O_t$$

$$N_t = N_0 \times Pr\{\tau_0 > t\} + \sum_{j=1}^{t} \left( \left( I_j(1 - \mu_s) - D_j \right) \times Pr\{\tau_j > t\} \right)$$

$$I_t = S_t + B_t = r \times V_t + A_t \frac{1 - p_t}{p_t} \quad , \quad p_t = \frac{\bar{A}_t / A_t}{(1 - d_t)}$$

$$Pr\{\tau_j > t\} = \begin{cases} (1 - \mu_m - \delta)^{10}(1 - \mu_l - \delta)^{t-10}, & j = 0 \\ (1 - \mu_m - \delta)^9(1 - \mu_l - \delta)^{t-j-9}, & 0 < j \leq t - 10 \\ (1 - \mu_m - \delta)^{t-j}, & j > t - 10 \end{cases}$$

# Results – Conservative Estimate

# Incorporate Uncertainty Into The Model

- Build uncertainty into the model to take into account variability

  - Produce probability distribution over the number of undocumented immigrants

- Main source of uncertainty

  - Parameter uncertainty

# Incorporate Uncertainty Into The Model

- Parameter Uncertainty:

$$\left\{ \begin{array}{c} Overstay \\ Rate \end{array}, \begin{array}{c} Apprehension \\ Probabilities \end{array}, \begin{array}{c} Emigration \\ Rates \end{array}, \begin{array}{c} Mortality \\ Rate \end{array} \right\}$$

$\underbrace{\phantom{xxxxxxxxxxxxxxxxx}}$ Inflow Parameters   $\underbrace{\phantom{xxxxxxxxxxxxxxx}}$ Outflow Parameters

- Taking into account:

    - Circular flows which link apprehension probability with emigration rates (Massey (2004, 2013), Massey and Pren (2012))

    - Cohort dependent emigration rates

# Incorporate Uncertainty Into The Model

- Parameter Uncertainty:

$$\left\{ \begin{matrix} Overstay \\ Rate \end{matrix}, \begin{matrix} Apprehension \\ Probabilities \end{matrix}, \begin{matrix} Emigration \\ Rates \end{matrix}, \begin{matrix} Mortality \\ Rate \end{matrix} \right\}$$

Inflow Parameters          Outflow Parameters

$$\Downarrow$$

$$N_t = N_0 \times Pr\{\tau_0 > t\} + \sum_{j=1}^{t} \left( \left( I_j (1 - \mu_s) - D_j \right) \times Pr\{\tau_j > t\} \right)$$

$$I_t = S_t + B_t = r \times V_t + A_t \frac{1 - p_t}{p_t} \quad , \qquad p_t = \frac{\bar{A}_t / A_t}{(1 - d_t)}$$

$$Pr\{\tau_j > t\} = \begin{cases} (1 - \mu_m - \delta)^{10}(1 - \mu_l - \delta)^{j-10}, & j = 0 \\ (1 - \mu_m - \delta)^{9}(1 - \mu_l - \delta)^{t-j-9}, & 0 < j \le t - 10 \\ (1 - \mu_m - \delta)^{t-j}, & j > t - 10 \end{cases}$$

# Incorporate Uncertainty Into The Model

- Parameter Uncertainty:

$$\left\{ \begin{array}{c} \textcolor{green}{Overstay} \\ \textcolor{green}{Rate} \end{array} , \begin{array}{c} \textcolor{green}{Apprehension} \\ \textcolor{green}{Probabilities} \end{array} , \begin{array}{c} \textcolor{blue}{Emigration} \\ \textcolor{blue}{Rates} \end{array} , \begin{array}{c} \textcolor{blue}{Mortality} \\ \textcolor{blue}{Rate} \end{array} \right\}$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\textcolor{green}{\text{Inflow Parameters}}} \qquad \underbrace{\qquad\qquad\qquad\qquad}_{\textcolor{blue}{\text{Outflow Parameters}}}$$
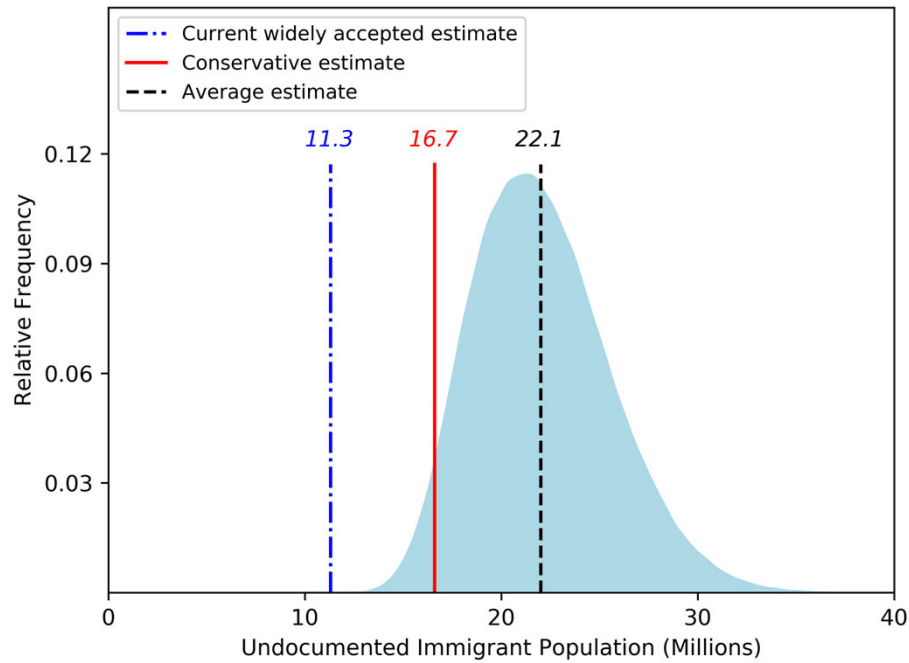
$$\Downarrow$$

$$N_t = N_0 \times Pr\{\tau_0 > t\} + \sum_{j=1}^{t} \left( \left( I_j(1 - \mu_s) - D_j \right) \times Pr\{\tau_j > t\} \right)$$

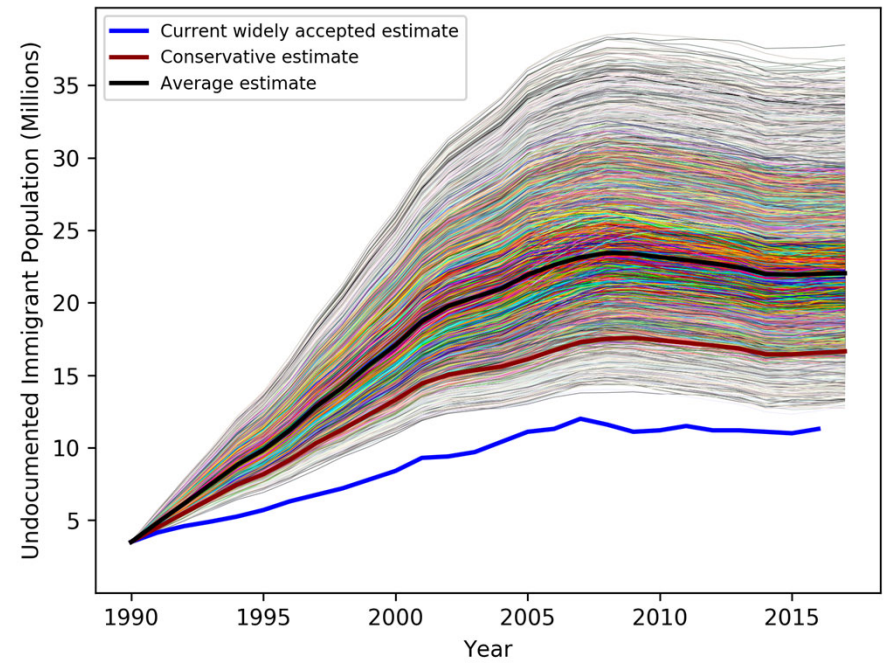$$I_t = S_t + B_t = r \times V_t + A_t \frac{1 - p_t}{p_t} \quad , \quad p_t = \frac{\bar{A}_t / A_t}{(1 - d_t)}$$

$$Pr\{\tau_j > t\} = \begin{cases} (1 - \mu_m - \delta)^{10}(1 - \mu_l - \delta)^{j-10}, & j = 0 \\ (1 - \mu_m - \delta)^{9}(1 - \mu_l - \delta)^{t-j-9}, & 0 < j \le t - 10 \\ (1 - \mu_m - \delta)^{t-j}, & j > t - 10 \end{cases}$$

- **Poisson structure with mean dependent upon the underlying parameter**

# Results – 1,000,000 Simulation Runs



Probability Distribution (2016)



Trajectories

# Reception

# Reception

Some of the issues in dispute:

1. The undercount implied by the new model is too high.

2. The range of model-estimated populations is too large to be useful for policy purposes, while the residual method gives a much smaller range of uncertainty.

3. The voluntary emigration rates of undocumented immigrants employed in the model are too low.

# Reception

Some of the issues in dispute:

1. **The census undercount implied by the new model is too high.**

2. The range of model-estimated populations is too large to be useful for policy purposes, while the residual method gives a much smaller range of uncertainty.

3. The voluntary emigration rates of undocumented immigrants employed in the model are too low.

# Census/American Community Survey Undercount

- For around 40 million there is no clear answer to place of birth (ignored deliberate misreporting)

- Fill in the blanks using "hot deck" allocation

  - Missing at random vs *missing on purpose*

  - Due to *missing on purpose*, the number of "donors" to the hot deck will be disproportionately US born

  - Imputed value for the place of birth variable will disproportionately point to "born in the USA"

  - Undercounting number foreign born ⇒ undercount in number of undocumented immigrants

# Reception

Some of the issues in dispute:

1. The undercount implied by the new model is too high.

2. **The range of model-estimated populations is too large to be useful for policy purposes, while the residual method gives a much smaller range of uncertainty.**

3. The voluntary emigration rates of undocumented immigrants employed in the model are too low.

# Reception

Some of the issues in dispute:

1. The undercount implied by the new model is too high.

2. **The range of model-estimated populations is too large to be useful for policy purposes, while the residual method gives a much smaller range of uncertainty.**

- **Precise estimate of wrong quantity: undocumented immigrants who are located and answering truthfully**

- **Small variability stems from the sampling variation that accompanies large samples**

# Reception

Some of the issues in dispute:

1. The undercount implied by the new model is too high.

2. The range of model-estimated populations is too large to be useful for policy purposes, while the residual method gives a much smaller range of uncertainty.

3. **The voluntary emigration rates of undocumented immigrants employed in the model are too low.**

# Mexican Migration Project Emigration Rate

- Data source critics use to deduce the new emigration rates is the Mexican Migration Project

# Mexican Migration Project Emigration Rate

- Data source critics use to deduce the new emigration rates is the Mexican Migration Project

- The main problem with using this data to estimating emigration rates is ***physical*** *sampling bias*:
  - Nearly everyone in the survey are in Mexico at the time of sampling (Massey et al (2016), Lessem (2018))
  - Misses those still in the United States at the time of the survey

# Mexican Migration Project Emigration Rate

- Data source critics use to deduce the new emigration rates is the Mexican Migration Project

- The main problem with using this data to estimating emigration rates is physical sampling bias:
  - Nearly everyone in the survey are in Mexico at the time of sampling (Massey et al (2016), Lessem (2018))
  - Misses those still in the United States at the time of the survey

- After accounting for the sampling bias, we find that the resulting emigration rates are consistent with our presumed rates

# Snapshot Models of Undocumented Immigration

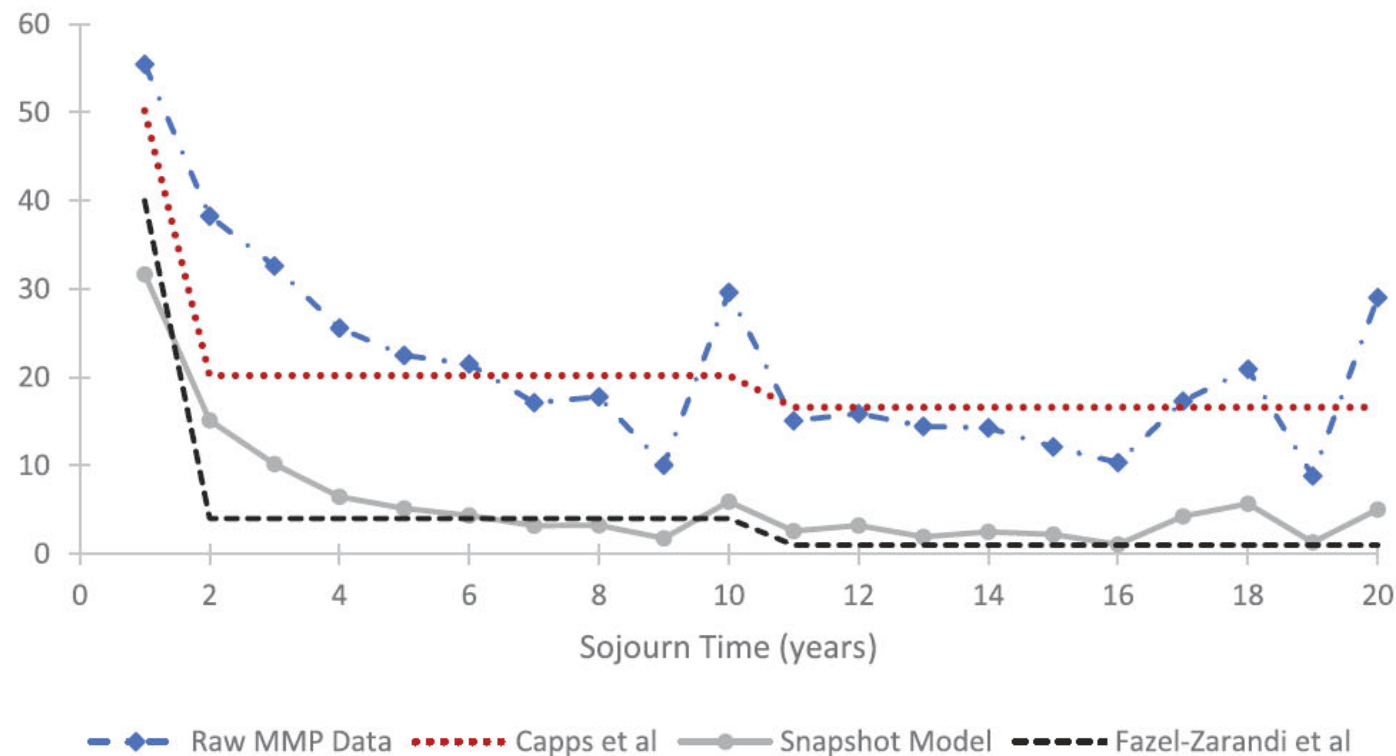## Scott Rodilitz [iD][1,*] and Edward H. Kaplan [iD][2]

Fig 8. Annual emigration probabilities from the United States.

# Policy Implications

- **Social services**:
  - agencies that have been working off of the previous estimate should recognize that the resources they have allocated for this population may be too low

- **Border control**:
  - In the last few years, the majority of new undocumented entrants are visa overstays
  - The number of illegal border crossers has substantially decreased

- **Crime:**
  - Further calls into question the claim of elevated risks of criminality surrounding undocumented immigrants.
  - The crime rate is much smaller