

1 Estimating the Eligible to Naturalize Population: Developing a 2 Robust Method

3 Katherine M. Condon, U.S. Citizenship & Immigration Services

4 Ronald E. Wilson, U.S. Citizenship & Immigration Services

5
6 *This paper is being presented to inform interested parties of ongoing research on immigrant population*
7 *estimation. The intent of the presentation is to encourage discussion on method improvements and the*
8 *use of unique data sources. Any views expressed on technical or operational issues are those of the*
9 *authors and not that of the U.S. Citizenship and Immigration Services.*

10 11 INTRODUCTION

12 Estimating the eligible to naturalize (EtN) population accurately and properly requires a detailed and
13 thorough process. Currently, multiple divisions within the U.S. Citizenship and Immigration Services
14 (USCIS), as well as the Office of Immigration Statistics (OIS) of the Department of Homeland Security
15 (DHS) produce EtN estimates from administrative records of the lawful permanent resident (LPR) living
16 in the United States, with the purpose of assessing resource allocation in operational planning. While
17 each division's method produces similar estimates, none completely capitalize on USCIS's data sources
18 in refining the method to obtain the most accurate estimate possible. We present a method in this
19 paper that maximizes USCIS own data to refine estimates of EtNs. The range of detailed immigrant
20 beneficiary characteristics across multiple beneficiary types allow for using sound social science
21 methods and research to bring a high level of accuracy to the EtN estimates. To do so, the estimation
22 must be done in several stages.¹

23 Refining the estimates facilitate the fine tuning of resource allocation within USCIS, but also in the
24 estimation of LPRs who most likely will naturalize², as well as any other estimations based off of EtN
25 population estimates.

26 BACKGROUND

27 Multiple agencies, inside and outside of USCIS, use Federal government administrative data sources to
28 estimate the level of immigration to the United States. Such sources are the U.S. (i) Census Bureau's
29 decennial censuses, (ii) Department of Housing & Urban Development's American Housing Survey, and
30 (iii) Department of Labor, which include characteristics on place of birth, citizenship, and year of entry
31 into the United States. These data provide essential information on the total foreign-born population,
32 naturalized citizens, and non-U.S. citizens residing in the United States. While these data sources are
33 used in support of their overall missions for operational planning and forecasting, the broad categories
34 of non-U.S. citizen data do not help overcome the limits of error introduced into estimates from using
35 non-immigration specific data.

36 Estimation by USCIS is guided by the Immigration and Nationality Act (INA), which produces detailed
37 data sets of immigrant characteristics for foreign nationals the allow us to specifically identify the

¹ We envision the full paper will be divided into the following parts that describes each stage. These parts include: (i) concept behind the stage; (ii) any underlying empirical social science research that underpins the stage; (iii) the data used; and (iv) statistical mechanics in creating the estimates at that stage.

² Not every LPR will naturalize, to which there will be variation in the patterns of timing to naturalize.

38 number of LPRs in the U.S. that serves as solid foundation for improving the estimation of LPRs eligible
39 to naturalize. The following is an overview of the data sources and method for producing a highly
40 refined estimate of the EtN from the national level.

41 DATA AND METHOD

42 The main objective of this work is to refine the process of estimating with respect to both the data
43 sources, as well as the methods. While there are a number of differences between other estimation
44 methods of EtN population, it is important to state that there is a basic overlap between the current
45 method to estimate and the one proposed in this paper.

46 Further, we propose that separate population estimates of eligible to naturalize LPRs be developed by
47 year of admission and class of admission. With regard to year of admission: (a) LPRs who were admitted
48 to the United States before 1973 and (b) LPRs who were admitted after 1973. This was done, because
49 USCIS computerization of their administrative data did not start until 1973. With regard to the class of
50 admission, there will be separate estimates for (a) LPRs admitted for humanitarian reasons, e.g.,
51 refugees and asylees; and (b) all other. These estimate sets are added together to obtain the overall EtN
52 population estimates as of July 1, 2015. We will also take into consideration out-migration and
53 mortality.

54 The following are the general stages used to produce our estimates: 1) Standardizing the data; 2)
55 Assembling of base data set; 3) Establishing a Refined Clock Start Date; 4) Identifying Eligible LPRs; 5)
56 Separating the Data Set by Class of Admission; and 6) Accounting for Out-migration and Mortality.

57 *Standardizing the Data Sets*

58 Since the data come from different data systems in both OIS and USCIS, the naming conventions of the
59 variables used differ. We established a set of common variable names to combine the same information
60 from different field names across the systems into one common field. Most of the required
61 characteristics across the data sets were straight-forward when combining into a single standardized
62 field, such as (i) Class of Admission, (ii) Sex, (iii) Date of Birth, (iv) Date of Naturalization, (v) Country of
63 Birth, (vii) State, and (viii) ZIP Code. However, the standardization of the date on which eligibility begins
64 is based on the path used in becoming an LPR.³

65 There are two main paths followed to become an LPR, which are an adjustment of status from already
66 being in the U.S., and being a new arrival from a foreign country. An (ix) Approval Date was created to
67 capture either date an immigrant becomes an LPR is based on which path is used. If the immigrant is an
68 adjustment, then the date eligibility is based off is the approval date of their adjustment. If a new arrival,
69 then the date is based off the entry date into the U.S. The OIS data contained Admission variable that
70 was put into the Approval Date variable. But for both paths, there is a subcategory humanitarian
71 immigrants (asylees and refugees) whose eligibility dates are not the adjustment approval or arrival
72 dates. The approval dates for the immigrants need to be adjusted for determining when they are eligible
73 to naturalize. The details about the adjustment of this subcategory of immigrants for establishing
74 eligibility date for is discussed in more detail below in regards to how we established the “Clock Start
75 Date” on which being eligible to naturalize begins.

76 We needed to create additional variables to not only ensure the back tracing of a record to the original
77 data source, but also needed information toward estimating the EtN. Finally one data issue we

³ Additional variables were standardized, such as Marital Status, Name, Received Date, Form Number, and Receipt Number for the purpose of managing the processing of the data. But the variables listed in the narrative are the primary variables used in the estimation.

78 encountered is with regard to a small number of records that either have the sex of the LPR as unknown
79 or is missing.⁴

80 *Data Set Assembly*

81 With both the USCIS and OIS data sets standardized, each data set needed to be truncated and
82 combined. The OIS data set contains all LPRs from 1973 to 2015. The USCIS data ranged from 1993 to
83 the present.⁵ Because USCIS data systems only came on line after its creation, it took time for the
84 agency's data systems to mature and contain all the records necessary conducting any analysis. In our
85 analysis, we identified the year 2010 to be the first year to contain complete data.

86 LPRs who arrived before 1973 and who have not naturalized are estimated using the aggregated data
87 from the most recent 5-year estimate from the Census Bureau's American Community Survey (ACS) to
88 supplement the estimates derived from individual records after the processing the OIS and USCIS data.
89 These data, though, were not combined with the OIS and USCIS data because they are aggregated and
90 could not be processed by individual records and only added in aggregate at the end.

91 As such, we combined the USCIS and OIS data sets at July 1, 2010 to make a comprehensive data set
92 that contained all the information needed to filter the data, identify existing naturalizations, and identify
93 eligibility start dates based on class of admission. The begin date of this data set was truncated at July 1,
94 1973 and truncated at June 30, 2015. This July 1 date aligns with the Census Bureau's population
95 estimate dates. Using this date ensures consistency when combining USCIS/OIS EtN estimates with the
96 EtN estimates from ACS data prior to 1973, as well as aligning with the yearly mortality rate dates.

97 Finally, because the USCIS and OIS data contained all LPRs across time, we matched and removed the
98 records for LPRs who had naturalized. We did this in two stages. First, since the OIS data already
99 identified those LPRs who had naturalized, we removed those records that had a naturalization date.
100 Second, we then matched the naturalization records in USCIS databases to the entire data set for those
101 who had naturalized with (1) the N-400 (Application for Naturalization) and (2) those who applied for
102 citizenship certificate with N-600 (Application for Certificate of Citizenship).⁶ We applied the N-400 and
103 N-600 records to not only identify those who naturalized in the USCIS data, but also to pick up any
104 additional naturalization in the OIS that might have not been recorded correctly.

105 *Establishing an Eligibility Clock Start Date and Identifying Eligible LPRs*

106 Depending on the pathway into LPR status, the start date towards eligibility to naturalize varies by the
107 initial classification for LPR status. All LPRs have a waiting period before and LPR can become eligible for
108 naturalization. In general, whether an adjustment of status or new arrival, the waiting period is five
109 years from first date the immigrant officially becomes an LPR. The exceptions are those LPRs who
110 married a U.S. Citizen (USC) or are protected under the Violence Against Women Act (VAWA), to which
111 the waiting period is three years. Asylees, refugees, parolees, or who have a cancellation or removal are

⁴ While these records make up less than a fraction of LPRs use in our estimation, we, nonetheless, impute the sex for those records to preserve as many records as possible. We impute sex for these records by proportionally assign the 1.1% sex ratios, which migration research shows that it is about 1.01 to 1.02 males to females. We use simple random sampling to draw a 51.5% of cases with either missing or unknown sex and assign those cases to be male, with remaining 48.5% to be female.

⁵ USCIS and OIS were established in 2002 under the Homeland Security Act of 2002 out of what had previously been known as Immigration and Naturalization Service (INS). (<https://www.dhs.gov/office-immigration-statistics>) The Department of Homeland Security's Office of Immigration Statistics (OIS) leads the collection and dissemination to Congress and the public of statistical information and analysis useful in evaluating the social, economic, environmental, and demographic impact of immigration laws, migration flows, and immigration enforcement. While the Department of Homeland Security's U.S. Citizenship and Immigration Services (USCIS).

⁶ The N-600 was used to capture any derivative LPRs who became eligible and officially obtained documentation of becoming a naturalized citizen. Not all derivatives will apply for a certificate, making for an incomplete filter on all naturalizations. Identifying this group remains difficult to estimate because there is no source to identify these LPRs.

112 all subject to having their approval dates “rolled back” to an earlier date. Asylees have their clock start
113 date adjusted to one year prior to the approval of the application for asylum. Refugees, parolees, and
114 those with removal cancellations all have their dates rolled back to either the date of entry into the U.S.
115 or 10 years prior to the time of their approval, whichever date occurs earlier.. The approval dates for
116 each LPR are transferred into a Clock Start Date field, adjusting for special immigrant classes, then
117 projecting forward the waiting period to establish the date at which the LPR is eligible to naturalize.⁷

118 *Accounting for Out-migration and Mortality*

119 Neither USCIS nor OIS track deaths or outmigration as part of an LPR’s administrative record. This leaves
120 the need to apply probabilities of dying or leaving the U.S. to the LPR data set to reduce the number to a
121 more realistic number of LPRs in the U.S. who might naturalize. For estimating out-migration⁸ to reduce
122 the data set to only those who remain in the U.S. in order to then calculate mortality. We use USCIS data
123 from I-407 (Record of Abandonment of Lawful Permanent Resident Status⁹) records in aggregate to
124 produce probabilities of out-migration by age to reduce the number LPRs who can apply for
125 naturalization.¹⁰

126 For each year, starting in 1973, the out-migration probabilities are applied to each age group, and then
127 those who “stayed” are survived forward one year using NCHS life tables to age them forward. This will
128 be done as an iterative annual process.

129 **CALCULATING THE FINAL NATIONAL-LEVEL ESTIMATE**

130 After processing the data, we summarize the estimates in the third to last year (2012) of the final year
131 (2015) of the outmigration and mortality data sets across the age groups, leaving out those under age
132 18. The year 2012 year was used because no LPR within three years of the last recorded year are eligible
133 to naturalize. We then combined the three summarized estimates for 1973 to 2012 with the American
134 Community Survey estimates from 1973 and prior to create the final EtN national-level estimate.

135 Local-level estimates can also be derived from the process, particularly for state and metropolitan areas.
136 The process can be followed up to the stage applying the out migration and mortality probabilities.
137 Estimates at that stage can be split out by the desired local-level geography, which the remainder of the
138 process can be followed on the estimates for those geographies to produce more localized estimates for
139 more detailed planning.

140 While outside agencies cannot make use of USCIS administrative data, our hope is that our method will
141 prompt creative thinking from outside agencies toward finding data sources that allow them to follow
142 this method and improve their estimation efforts.

⁷ We keep the under age 18 population during the production because they are essential for using our cascading approach to estimating mortality and out migration.

⁸ Prior to this step for accounting of out-migration and mortality, we separated LPRs by class of admission, i.e., asylees, refugees, parolees, and those with removal cancellations, to only apply the mortality probabilities. Given the nature of their immigrating to the U.S., it is unlikely that many will want to return to their country of birth/nationality. As such, very few of these immigrants are likely to out-migrate and we do not include them in our out-migration calculations. We, then, summarize the individual records from the base data set into a new data set of age groups and apply the out-migration probabilities derived from the I-407 population only to LPRs with class of admission as “Other”.

⁹ The purpose of the I-407 form is to record in USCIS data systems that an LPR has decided to voluntarily abandon their status as an LPR from having left the country and not returned after a given period of time.

¹⁰ The distribution of ages in these records appears consistent with migration research in accordance with risk of out-migration based on (i) age to do so, (ii) financial means to do so by age, and (iii) intent/motivation to do so by age. This distribution can be used to create a series of emigration probabilities by age.