

**Department of Human Services: Arlington Data System Notes**

**M weights**

System	Last name	First name	gender	ssn	street	zip	birth_yr	birth_mo	birth_dy
Anasazi	0.95	0.9	0.85	0.5	0.4	0.6	0.95	0.95	0.85
Eto	0.95	0.9	0.85	0.9	0.4	0.6	0.95	0.95	0.85
Web Vision	0.95	0.9	0.85	0.9	0.4	0.6	0.96	0.95	0.85
Combined	0.95	0.9	0.85	0.5	0.4	0.6	0.95	0.95	0.85

System	Starting n	# reduced	Final n	Weight cutoff	Notes/observations
<b>Anasazi/CERNER</b> Note to Sayali: make street into address	57065 56558 if you do unique	691	56374	40	Some interesting features: <ul style="list-style-type: none"><li>- case_no = 0 were crazy</li><li>- 'COMBINED A INTO B' formats</li><li>- SSN m-weight set to 0.5 instead of 0.9 to account for missingness (~50%)</li><li>- NO matching IDs (in this case, ID is case number)</li><li>- This system has encounter data; so should integrate it to look for most recent demographics</li><li>- Weights:<ul style="list-style-type: none"><li>- Weight 50-52 has suspected twins of the same gender</li><li>- Weight ~47 has suspected twins of opposite gender</li><li>- Weights in between give funky matches!</li></ul></li></ul>
<b>ETO</b>	25619 25437 if you do unique	460	25159	-37 cut off point. Below 37, SSN matches are either missing or way off: Lower SSN?	Some interesting features <ul style="list-style-type: none"><li>- Eto has merged record ID field that is completely missing</li><li>- Missing SSN is at 28%. Should we lower m weight to 0.7 instead of 0.9?</li><li>- 309 cases of system IDs mismatch whom are the same person</li></ul>

				-For exact matched IDs: 142 system IDs duplicated across the system	- 142 system IDs duplicated across the system
<b>WVS</b>	64927: 58115 if u do, unique:	21595	43332	-50 for non-exact match IDs -Between 40-50: -Cases of twins? Sometimes birthdates are off! -Exact matched IDs went down till 21	- Is it possible that there are twins here as well? Look particularly at 30.33667 - Also, i think we should go down to 31 for non-exact-id-match sets as well
<b>Joint</b>	147,611: Systemwide reduction to: 124,868	4765	120103	First cut: 33 Second cut: 34	- As we improve the matches from each of the three systems, the initial number may change - Again, plenty missing SSNs; so m-weight 0.5 -     asz eto wvs asz 1183 2816 696 eto 0 308 1888 wvs 0 0 1609 - - Off-diagonal entries on the table above tell us how many across-system matches to expect. The current total is 5700 - After unique: 5,000 of them are shared across systems - Around 33.68..., there's definitely a question about twins! - Funnily enough, in this deduplication, the across system matches are significantly better than the within system ones.

## **Coding Steps for merging back**

Web Vision as a case study:

### **Overall Process**

- 1) Group everything by same System ID: Merge the demographic variables (Cut off: 21-80)
- 2) Group everything by row ID: Get the different system id matches and merge the demographic variables (Cut off: 50-80)
- 3) Find the overlaps between #1 and #2. Merge the demo variables there.
- 4) #1 - overlap = Remaining 1
- 5) #2 - overlap = Remaining 2
- 6)  $\text{dedup\_total} = \text{overlap} + \text{Remaining 1} + \text{Remaining 2}$
- 7)  $\text{Undedup\_total} = \text{All data set} - (\text{dedup\_total} + \text{linked system Id pairs from 2})$
- 8) Final data set = dedup\_total + Undedup\_total

Logic that could improve

- 1) Where majority rule fails (Pls refer to email over weekend)
- 2) Cases where : (27, 28,29): (29,31) where 31 is not considered. Code is not deliberately written to take care of it.

### **Remaining Tasks**

For Final Joined Table:

- Complete step 7 and step 8

Create a systems ID pairs table

