

CONFIRM - Clustering Of Noisy Form Images using Robust Metrics

Chris Tensmeyer and Tony Martinez
Machine Learning Laboratory
Brigham Young University

Problem

- Scanned Form Images
- Cluster by Structure
- Noise
- Similar Form Types
- Potentially Large Datasets
 - Not quite solved yet

Reports Control System - (JUST 1001)

UNITED STATES DEPARTMENT OF JUSTICE
Bureau of Prisons
Washington, D.C. 20530
(Rev. 1-2-65)

LIFT NO. _____

LIST OF OUTWARD-BOUND PASSENGERS
(United States Citizens and Nationals)

Sailing from Seattle (Port) 2 March 1969 (Date)

S. S. USAT FUNSTON bound for port of YOKOHAMA (Destination)

Line No.	FAMILY NAME-GIVEN NAME ADDRESS IN UNITED STATES	Age (Years) (P-M)	Sex (M-F)	U.S. PASSPORT No.	DATE AND PLACE OF BIRTH & PLACE OF NATURALIZATION	EXPIRE DATE OF PASSPORT BY WHICH ISSUED	
						(1)	(2)
1	ALEXANDER, RUTH 2110 Harrison St., Evanston, Ill.	24	F	514	Chicago, Ill.	Indef.	
2	Dale L.	5	F		Oakdale, La.		
3	Susan Carol	2	F		Fort Lewis, Washington		
4	Harnett, Susan M.	58	F	513	Chicago, Ill.		
5	ANDERSON, MARY J. MANDEVILLE, La. Box 350	29	F	30020	Mandeville, La.		
6	Henry G.	9	M		Slidell, La.		
7	Jocelyn K.	7	F		Mandeville, La.		
8	Bonita	4	F		Madisonville, La.		
9	Emily L.	3	F		New Orleans, La.		
10	DARTON, DORIS R. 213 N. Penn St., Harboro, Pa.	30	F	1743	Russia, Europe		
11	BENNETT, SANDRA JUANITA 1923 N. Gyre St., Los Angeles, Calif.	6	F	1111	Maywood, Calif.		
12	BERTRAND, HIDIA 320 Surf St., Chicago, Ill. (Friend)	52	F	921	Chicago, Ill.		
13	BYRD, BONNIE M. Saskatchewan, Okla.	23	F		Saskatchewan, Okla.		
14	CAMPBELL, L. FRANCES INSTRUMENTS, INC., Elkhorn	24	F		10/9		
15	Janet Ja.	4	F		Charlottesville, Va.		
16	Patricia J.	2	F		Charlottesville, Va.		
17	CHANDLER, BETTY ANN 1912 N.W. Flagler Miami, Florida	20	F		7/19		
18	Betty Jean	45	F		Hendersonville, N.C.		
19	Patricia Dian	52	F		Kingsport, Tenn.		
20	Majorie Ann	17m	F		Hendersonville, N.C.		
21	CHAPPEL, THIRINA 804 S. Eltony Pomona, Calif.	41	F	1115	Washington, Ind.		
22	Davis, Howard	13	M		Pomona, Calif.		
23	Davis, William	10	M		Pomona, Calif.		
24	Chappel, John	3dm	M		Pomona, Calif.		
25	COHN, TINA 315 E. Peabody Apt. #4 Columbus, Ga.	23	F	1613	Columbus, Ga.	FILE - G.R.V.	

25-vjt

H-44000-1

Form Image Clustering

CENSUS OF ENGLAND AND WALES, 1911.	
 British Schools.)	
SCHEDULE.	
<i>Prepared pursuant to the Census (Great Britain) Act, 1901.</i>	
This space to be filled up by the Enumerator. <hr/> Number of Registration District... <u>582</u> Number of Registration Sub-District <u>582 1</u> Number of Enumeration District... <u>9</u> <hr/> Name of Head of Family or Separate Occupier. <u>W^m Sutton</u> <hr/> Postal Address <u>Mr. Common Chaplain</u>	

NOTICE

This Schedule must be filled up and signed by, or on behalf of, the Head of the Family or other person in occupation, or in

Form Image Clustering

Washington Passenger Lists Dataset

PASSENGER MANIFEST

Owner or operator: NORTHWEST AIRLINES, INC. Flight No.: 501 of 10 Date: 10/2/65

Aircraft No. 43766

Point of Embarkation: Seattle, Washington Point of Disembarkation: Tokyo, Japan

NAME OF PERSON	AGE	SEX	NAME OF PERSON	AGE	SEX	
1. <i>Seaplane</i>	4	M	2. <i>V-897450 To C</i>	1	M	
3. <i>Paganini, John</i>	332		4. <i>V-223344 To Honolulu</i>	2	M	
5. <i>Ketteler, Alondra</i>	370		6. <i>V-223344 To Honolulu</i>	2	M	
7. <i>Romano, Acolino</i>	310		8.			
9.			10.			
11.			12.			
13.			14.			
15.			16.			
17.			18.			
19.			20.			
21.			22.			
23.			24.			
25.			26.			
27. Undocumented Persons:			28. <i>Pat. B. Anderson</i>	3	M	
29.			30.			
31.			32. <i>J. L. Wilson</i>			
33.			34. U. S. CUSTOMS INSPECTOR			
35.			36.			
37.			38.			
39.			40.			
41.			42.			
43.			44.			
45.			46.			
47.			48.			
49.			50.			
51.			52.			
53.			54.			
55.			56.			
57.			58.			
59.			60.			
61.			62.			
63.			64.			
65.			66.			
67.			68.			
69.			70.			
71.			72.			
73.			74.			
75.			76.			
77.			78.			
79.			80.			
81.			82.			
83.			84.			
85.			86.			
87.			88.			
89.			90.			
91.			92.			
93.			94.			
95.			96.			
97.			98.			
99.			100.			
Prepared by: J. L. GISON	1-23-71 U.S.C.					

LIST OF OUTWARD-BOUND PASSENGERS					
Sailing from		S.E.P.S. (Port)		Jan 5 (Date)	
S. S. USAT REPUBLIC		bound for port of YOKOHAMA		Indef.	
1. HIRRO, VINA E.	31	F	2. <i>Haskogee, Oklahoma</i>	3. <i>Indef.</i>	
3. MOOZ S. JEWOKA			4. <i>Cushing, Oklahoma</i>	5. <i>Cushing, Oklahoma</i>	
6. PAUL N.	12	M	7. <i>Mitchel Field, N. Y.</i>	8. <i>Edwardsville, Illinois</i>	9. <i>Indef.</i>
10. CHRISTOPHER W.	2	M	11. <i>St.Louis, Missouri</i>	12. <i>Granite City, Ilo.</i>	13. <i>Indef.</i>
14. JURITA, EVELYN	39	F	15. <i>Peoria, Illinois</i>	16. <i>Goshen, Indiana</i>	17. <i>Indef.</i>
18. Box 56 Nedora, Illinois			19. <i>Yonkers, New York</i>	20. <i>Gold, Oklahoma</i>	21. <i>Indef.</i>
22. MARGARET	17	F	23. <i>Charles, Oklahoma</i>	24. <i>Enid, Oklahoma</i>	25. <i>Indef.</i>
26. LUCILLE	16	F	27. <i>Panorama, California</i>	28. <i>Detroit, Michigan</i>	29. <i>Indef.</i>
30. JUHES, GERALDINE	29	F	31. <i>Minneapolis, Minn.</i>	32. <i>Detroit, Michigan</i>	33. <i>Indef.</i>
34. 505 Hancock, Peoria, Illinoi			35. <i>Atascadero, California</i>	36. <i>Beaumont, Texas</i>	37. <i>Indef.</i>
38. BUTTERFIELD, L. ELDA	45	F	39. <i>Richmond, California</i>	40. <i>Albany, California</i>	41. <i>Indef.</i>
42. 207 W. High St., Elkhart, Indiana			43. <i>Salem, Oregon</i>	44. <i>Riverside, California</i>	45. <i>Indef.</i>
46. BILLION, C. ERGUS	16	M	47. <i>Louisville, Kentucky</i>	48. <i>Birmingham, Alabama</i>	49. <i>Indef.</i>
50. CLIMENT, HAZEL M.	51	F	51. <i>Louisville, Kentucky</i>	52. <i>Louisville, Kentucky</i>	53. <i>Indef.</i>
54. Huntington Beach, Box 4760, California			55. <i>Indef.</i>	56. <i>Indef.</i>	57. <i>Indef.</i>
58. CHARLES	18	M	60. <i>Indef.</i>	61. <i>Indef.</i>	62. <i>Indef.</i>
62. MARILYN S.	13	F	65. <i>Indef.</i>	66. <i>Indef.</i>	67. <i>Indef.</i>
68. KATHLEEN L.	9	F	70. <i>Indef.</i>	71. <i>Indef.</i>	72. <i>Indef.</i>
72. CAROL, DORIS E.	39	F	75. <i>Indef.</i>	76. <i>Indef.</i>	77. <i>Indef.</i>
76. 3700 34th Ave., Minneapolis, Minnesota			80. <i>Indef.</i>	81. <i>Indef.</i>	82. <i>Indef.</i>
80. KAREN S.	8	F	85. <i>Indef.</i>	86. <i>Indef.</i>	87. <i>Indef.</i>
84. FREDERICK	6	M	90. <i>Indef.</i>	91. <i>Indef.</i>	92. <i>Indef.</i>
88. ROGER E.	4	M	95. <i>Indef.</i>	96. <i>Indef.</i>	97. <i>Indef.</i>
92. GREG, FLORINDA HAS	27	F	100. <i>Indef.</i>	101. <i>Indef.</i>	102. <i>Indef.</i>
96. 5049 Camp Bowie Blvd., Fort Worth, Texas			105. <i>Indef.</i>	106. <i>Indef.</i>	107. <i>Indef.</i>
100. DELIGHT, MARY A.	24	F	110. <i>Indef.</i>	111. <i>Indef.</i>	112. <i>Indef.</i>
104. 6009 Jordan Rd., Richmond, California			115. <i>Indef.</i>	116. <i>Indef.</i>	117. <i>Indef.</i>
108. ERIC H.	17	M	120. <i>Indef.</i>	121. <i>Indef.</i>	122. <i>Indef.</i>
112. DAISY, FLORENCE	45	F	125. <i>Indef.</i>	126. <i>Indef.</i>	127. <i>Indef.</i>
116. 278 Via Calacatta San Lorenzo, California			130. <i>Indef.</i>	131. <i>Indef.</i>	132. <i>Indef.</i>
120. CAROL	20	F	135. <i>Indef.</i>	136. <i>Indef.</i>	137. <i>Indef.</i>
124. DEBORA, LOUISE M.	38	F	140. <i>Indef.</i>	141. <i>Indef.</i>	142. <i>Indef.</i>
128. 1409 Jarvis Lane, Louisville, Kentucky			145. <i>Indef.</i>	146. <i>Indef.</i>	147. <i>Indef.</i>
132. JAMES H. III	7	M	150. <i>Indef.</i>	151. <i>Indef.</i>	152. <i>Indef.</i>

CENSUS OF ENGLAND AND WALES, 1911.

Before writing on this Schedule please read the Examples and the Instructions given on the other side of the paper, as well as the headings of the Columns. The entries should be written in Ink.

Number of Schedules 103
(To be filled up by the Enumerator after collection.)

The contents of the Schedule will be treated as confidential. Strict care will be taken that no information is disclosed with regard to individual persons. The returns are not to be used for proof of age, as in connection with Old Age Pensions, or for any other purpose than the preparation of Statistical Tables.

NAME AND SURNAME	RELATIONSHIP to Head of Family.	AGE (last Birthday) and SEX.	PARTICULARS as to MARRIAGE.			PROFESSION or OCCUPATION of Persons aged ten years and upwards.						BIRTHPLACE of every person.	NATIONALITY of every Person born in a Foreign Country.	INFIRMITY.	
			For Infants under one year state the number in months as "under one month," "one month," etc.	State whether "Head," or "Wife," or "Son," "Daughter," or other Relative, "Visitor," "Boarder," or "Servant."	Write "Single," "Widower," or "Widow," opposite the names of all persons aged 15 years and upwards.	State, for each Married Woman entered on this Schedule, the number of:— Completed years of the present Marriage has lasted, if less than one year write "under one."	Children born alive to present Marriage. (If no children born alive write "None" in Column 7).	Industry or Service with which worker is connected.	Whether Employer, Worker, or Working on Own Account.	Whether Working at Home.					
1.	2.	3.	4.	5.	6.	7.	8.	9.	10.	11.	12.	13.	14.	15.	16.
1. Frederick Joyce.	Head	62		Married	27	3	2	1	Farmer	100	Employer		Cloford Somerset	140	
2. Elizabeth Hugman Joyce	Wife	57		Married	27								Oldbury on Severn	Glos	
3. Elizabeth Mary Gladys Joyce	Daughter	25		Single					Farmer's Daughter Dairy work	11			Tidenham	Glos	
4. Guendoline Joyce	Daughter	23		Single					Farmer's Daughter Dairy work				Tidenham	Glos	
5.															
6.															
7.															
8.															
9.															
10.															
11.															
12.															
13.															
14.															
15.															

(To be filled up by the Enumerator.)

Total.		
Males.	Females.	Persons.
1	3	4

I certify that—
 (1) All the ages on this Schedule are entered in the proper sex columns.
 (2) I have counted the males and females in Columns 3 and 4 separately, and have compared them with the total number of persons.
 (3) After comparing the figures I have completed all entries on the Schedule which appeared to be defective, and have corrected such as appeared to be erroneous.

Initials of Enumerator *Jay*

(To be filled up by, or on behalf of, the Head of Family or other person in occupation, or in charge, of this dwelling.)

Write below the Number of Rooms in this Dwelling (Garden, Kitchen, Servants' or Apartment). Count the kitchen as a room but do not count scullery, landing, lobby, closet, bathroom; nor warehouse, office, shop.

Fourteen

I declare that this Schedule is correctly filled up to the best of my knowledge and belief.

Signature *Frederick Joyce*

Postal Address *High Hall Farm Tidenham Cheltenham*

103

If any person included in this Schedule is—
 (1) "British subject by parentage,"
 (2) "Totally Deaf," or "Deaf and Dumb,"
 (3) "Blind,"
 (4) "Lunatic,"
 (5) "Imbecile,"
 (6) "Feeble-minded,"
 state the infirmity opposite that person's name, and the age at which he or she became afflicted.

Or

(3) If of foreign nationality, state whether

"French,"

"German,"

"Russian,"

etc.

CENSUS OF ENGLAND AND WALES, 1911.

Before writing on this Schedule please read the Examples and the Instructions given on the other side of the paper, as well as the headings of the Columns. The entries should be written in Ink.

Number of Schedule 57
(To be filled up by the Enumerator
after collection.)

The contents of the Schedule will be treated as confidential. Strict care will be taken that no information is disclosed with regard to individual persons. The returns are not to be used for proof of age, as in connection with Old Age Pensions, or for any other purpose than the preparation of Statistical Tables.

NAME AND SURNAME	RELATIONSHIP to Head of Family.	AGE (last Birthday) and SEX.	PARTICULARS as to MARRIAGE.			PROFESSION or OCCUPATION of Persons aged ten years and upwards.			BIRTHPLACE of every Person	NATIONALITY of every Person Foreign Country.	INFIRMITY.	LANGUAGE SPOKEN.
of every Person, whether Member of Family, Visitor, Boarder or Servant, who			State whether "Head" or "Wife," "Son," "Daughter," or other Relative, "Visitor," "Boarder," or "Servant."	For Infants under one year state the age in months as "under one month," "one month," etc.	Write "Single," "Married," "Widower," or "Widow" opposite the names of all persons aged 15 years and upwards.	State, for each Married Woman entered on this Schedule, the number of:— Completed years of Marriage Less than one year, "under one."	Children born alive to present Marriage. (If no Children born alive write "None" in Column 7.)					
(1) passed the night of Sunday, April 2nd, 1911, in this dwelling and was alive at midnight, or (2) arrived in this dwelling on the morning of Monday, April 3rd, not having been enumerated elsewhere.												
No one else must be included. (For order of entering names see Examples on back of Schedule.)												
1. Thomas Williams	Head	71	Married	38					At Home	Hirnacuton, Mon.	English	
2. Rhoda Ann Williams	Wife	58	Married	38	7	7			Neufport, Mon.	481	English	
3. Eliza Josephine Williams	Daughter	20	Single						At Home	Hirnacuton, Mon.	English	
4.												
5.												
6.												
7.												
8.												
9.												
10.												
11.												
12.												
13.												
14.												
15.												

(To be filled up by the Enumerator.)

Total.		
Males.	Females.	Persons.
1	2	3

- I certify that:—
 (1) All the entries in this Schedule are entered in the proper sequence.
 (2) I have counted the males and females in Columns 3 and 4 separately, and have compared their sum with the total number of persons.
 (3) After making the necessary inquiries I have completed all entries on the Schedule which were found to be defective, and have corrected such as appeared to be erroneous.

Initials of Enumerator J.P.W.

(To be filled up by, or on behalf of, the Head of Family or other person in occupation, or in charge, of this dwelling.)

Write below the Number of Rooms in this Dwelling (House, Tenement or Apartment). Count the kitchen as a room but do not count scullery, landing, lobby, closet, bathroom; nor warehouse, office, shop.
Four 4

I declare that this Schedule is correctly filled up to the best of my knowledge and belief.

Signature Thomas Williams

Postal Address Rose Cottage, Earlwood, Llanfair M. Cheshire, Mon

CENSUS OF ENGLAND AND WALES, 1911.

Number of Schedule 163
(To be filled up by the Enumerator
after collection.)

Before writing on this Schedule please read the Examples and the Instructions given on the other side of the paper, as well as the headings of the Columns. The entries should be written in Ink.

The contents of the Schedule will be treated as confidential. Strict care will be taken that no information is disclosed with regard to individual persons. The returns are not to be used for proof of age, as in connection with Old Age Pensions, or for any other purpose than the preparation of Statistical Tables.

NAME AND SURNAME	RELATIONSHIP to Head of Family.	AGE (last Birthday) and SEX.	PARTICULARS as to MARRIAGE.						PROFESSION or OCCUPATION of Persons aged ten years and upwards.						BIRTHPLACE of every person.	NATIONALITY of every Person born in a Foreign Country.	INFIRMITY.
			For Infants under one year state the age in months as "under one month," "one month," etc.	Write "Single," "Married," "Widower," or "Widow." Opposite the names of all persons aged 15 years and upwards.	State, for each Married Woman entered on this Schedule, the number of:—	Completed years the present Marriage has lasted.	Children born alive to present Marriage. (If no children born alive write "None" in Column 7).	Total Children Born Alive.	Children still Living.	Children who have Died.	Personal Occupation.	Industry or Service with which worker is connected.	Whether Employer, Worker, or Working on Own Account.	Whether Working at Home.			
of every Person, whether Member of Family, Visitor, Boarder, or Servant, who (1) passed the night of Sunday, April 2nd, 1911, in this dwelling and was alive at midnight, or (2) arrived in this dwelling on the morning of Monday, April 3rd, not having been enumerated elsewhere. No one else must be included. (For order of entering names see Examples on back of Schedule.)																	
L.	2.	3.	4.	5.	6.	7.	8.	9.	10.	11.	12.	13.	14.	15.	16.		

CENSUS OF ENGLAND AND WALES, 1911.

Number of Schedule 57
(To be filled up by the Enumerator
after collection.)

Before writing on this Schedule please read the Examples and the Instructions given on the other side of the paper, as well as the headings of the Columns. The entries should be written in Ink.

The contents of the Schedule will be treated as confidential. Strict care will be taken that no information is disclosed with regard to individual persons. The returns are not to be used for proof of age, as in connection with Old Age Pensions, or for any other purpose than the preparation of Statistical Tables.

NAME AND SURNAME	RELATIONSHIP to Head of Family.	AGE (last Birthday) and SEX.	PARTICULARS as to MARRIAGE.						PROFESSION or OCCUPATION of Persons aged ten years and upwards.						BIRTHPLACE of every person.	NATIONALITY of every Person born in a Foreign Country.	INFIRMITY.	LANGUAGE SPOKEN.
			For Infants under one year state the age in months as "under one month," "one month," etc.	Write "Single," "Married," "Widower," or "Widow." Opposite the names of all persons aged 15 years and upwards.	State, for each Married Woman entered on this Schedule, the number of:—	Completed years the present Marriage has lasted.	Children born alive to present Marriage. (If no children born alive write "None" in Column 7).	Total Children Born Alive.	Children still Living.	Children who have Died.	Personal Occupation.	Industry or Service with which worker is connected.	Whether Employer, Worker, or Working on Own Account.	Whether Working at Home.				
of every Person, whether Member of Family, Visitor, Boarder, or Servant, who (1) passed the night of Sunday, April 2nd, 1911, in this dwelling and was alive at midnight, or (2) arrived in this dwelling on the morning of Monday, April 3rd, not having been enumerated elsewhere. No one else must be included. (For order of entering names see Examples on back of Schedule.)																		
L.	2.	3.	4.	5.	6.	7.	8.	9.	10.	11.	12.	13.	14.	15.	16.			

CONFIRM

- Form Specific Features
 - OCR Text
 - Rule Lines
- Matching algorithms
- Pick T random forms as templates
- Encode forms by matching with templates
- Learn form-form similarity metric

Extract Text Features

CENSUS OF ENGLAND AND WALES,

NAME AND SURNAME

Children born alive to

Signature

OCR

Noisy
Transcription of Text

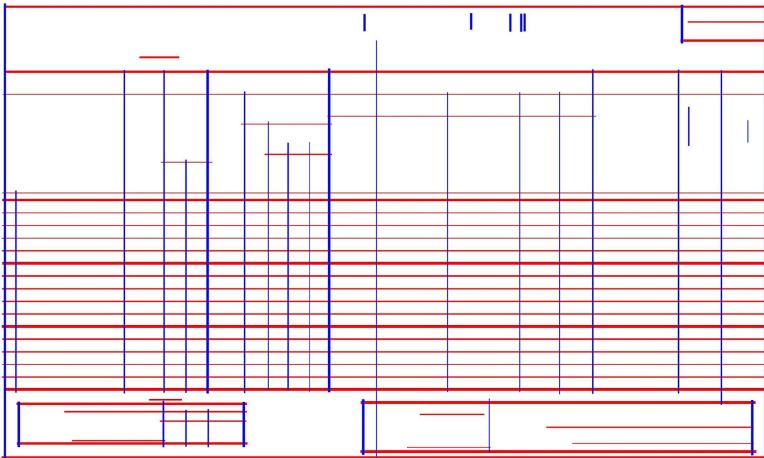
“CENSUS OF ENGLAND
AND WALES”
“NAME AND SURNAME”
“Children bom alive to”
“Signature”

Text Matching

- Correspondence problem
- Do greedy matching
 - Similar location and content
 - Prefix/Suffix matching
- Edit Distance
 - Number of character edits to transform one string into another.

Extract Rule Lines

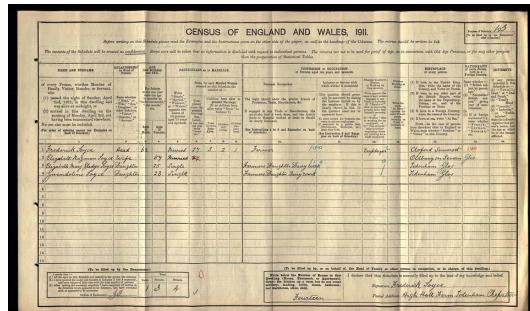
Line Extraction



Rule Line Matching

- Horz/Vert as sequences of parallel lines
- Edit Distance between sequences
 - Also shows how to merge two sequences
- Operations include
 - Match Deletion
 - Contains Overlap
 - Transpose Connect

Matching Vector



Query Form

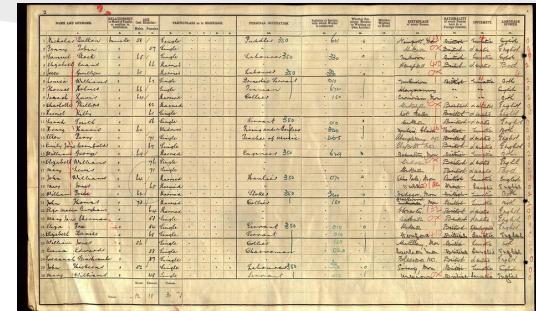


Match Vectors

$$(0, 1, 1, \dots)$$

(1, 1, 0.5, 0.75, ...)

(0, 0, 0.25, 0, ...)



Template Form

$$(0, 1, 1, \dots, 1, 1, 0.5, 0.75, \dots, 0, 0, 0.25, 0, \dots)$$

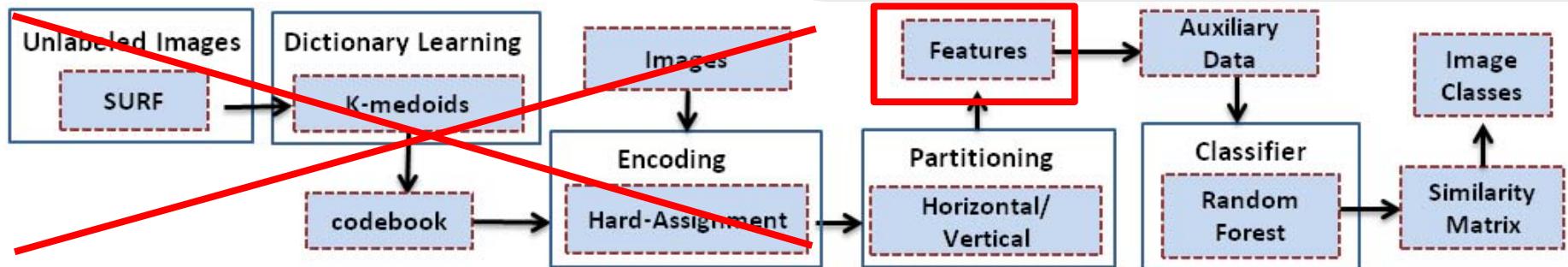
New Feature Vector

Text

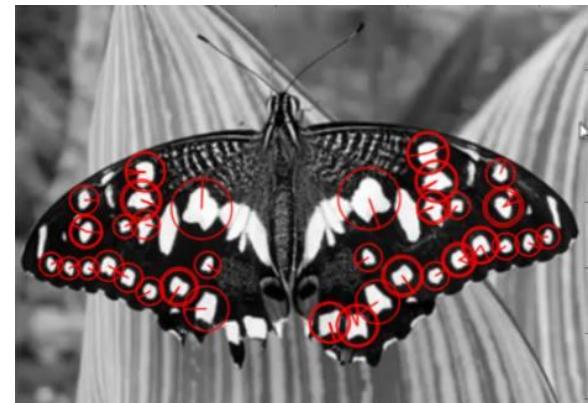
Horizontal Lines

Vertical Lines

Kumar and Doermann



- Create codebook of SURFs
- Encode images as histogram of codewords
- Learn similarity matrix using Random Forest

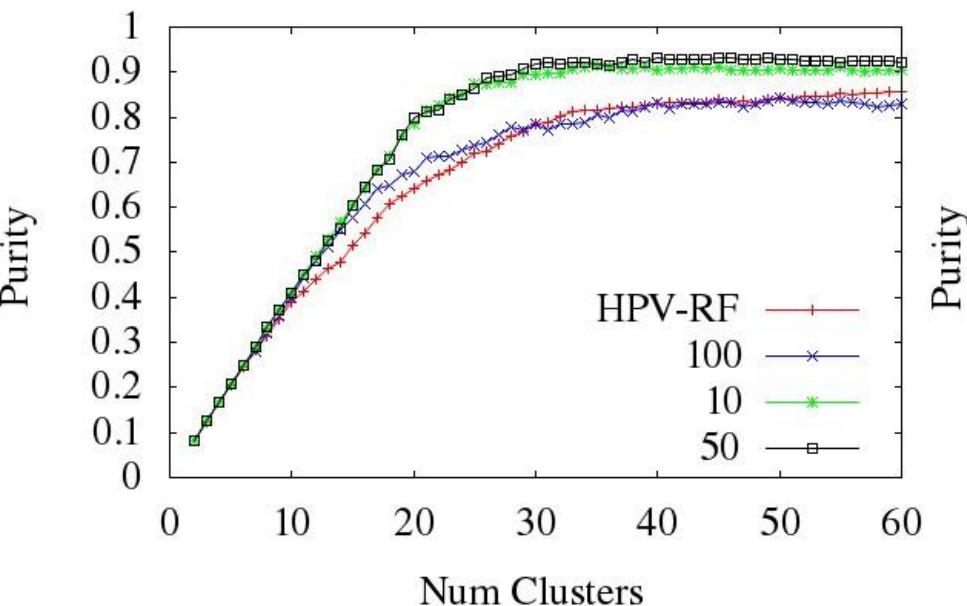


Experiment

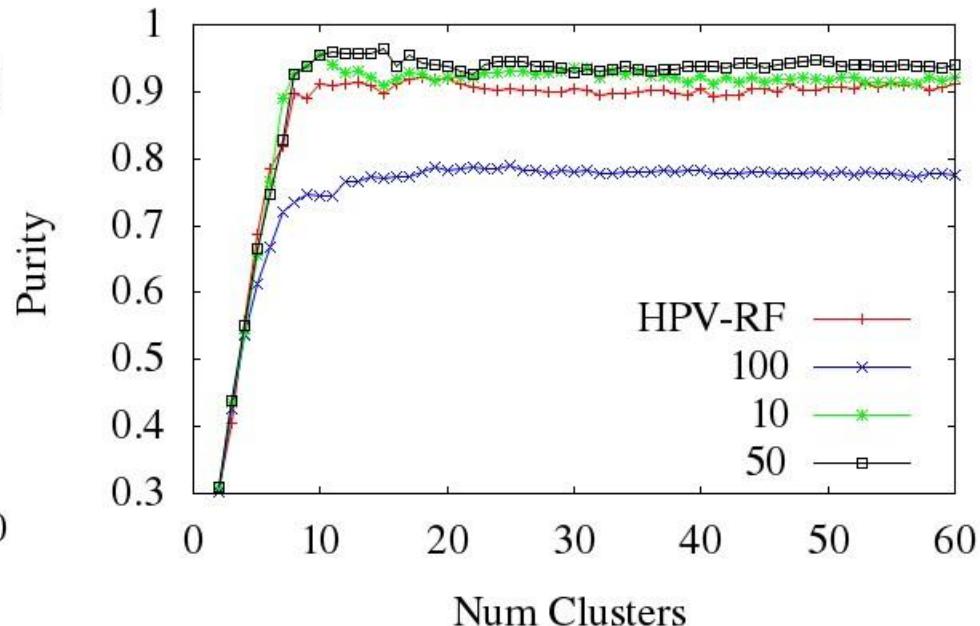
- CONFIRM vs Kumar & Doermann
 - Vary number of random templates
 - Vary types among random templates
- Historical Datasets
 - Census
 - Death Certificates
 - Passenger Lists
- Metric: Purity

Results - Census

Balanced

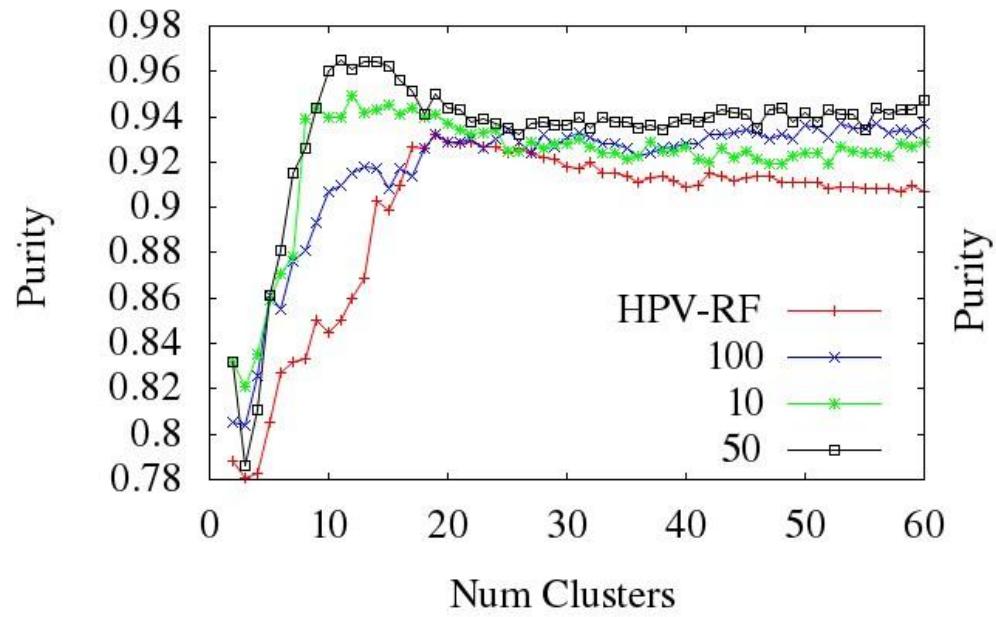


Skewed

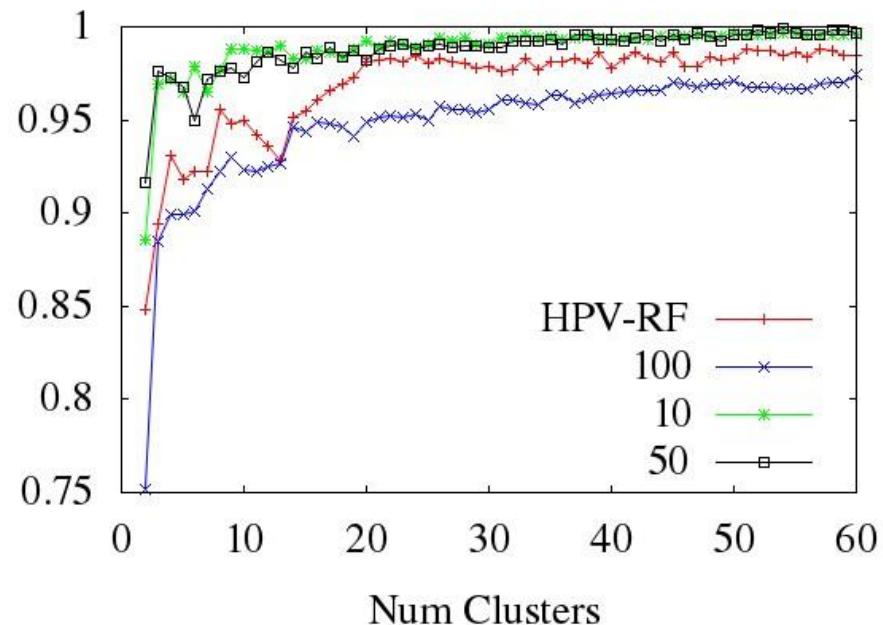


Results - Cont

Death Certificates

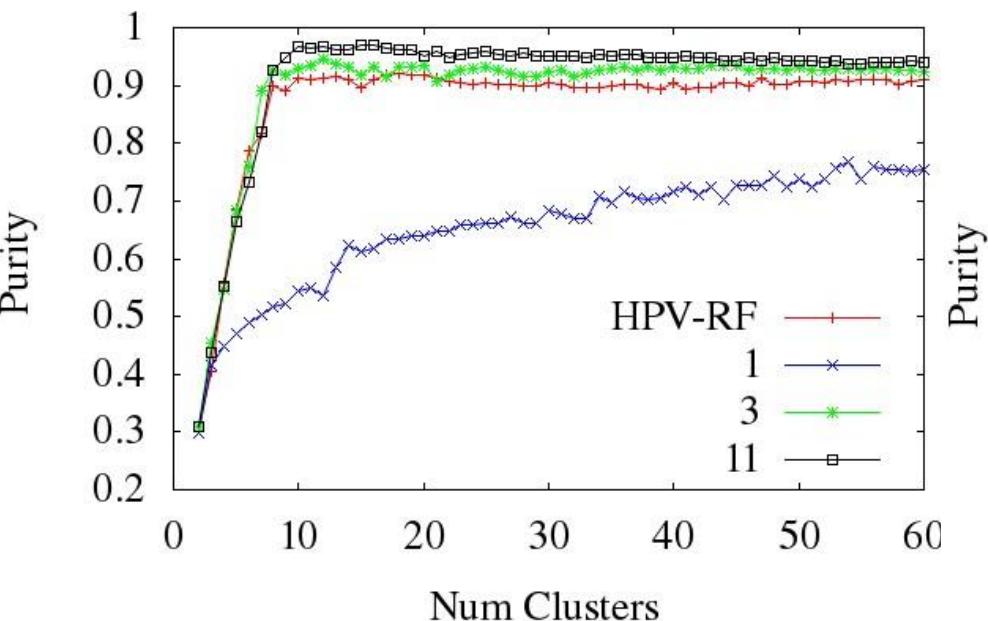


Passenger Lists

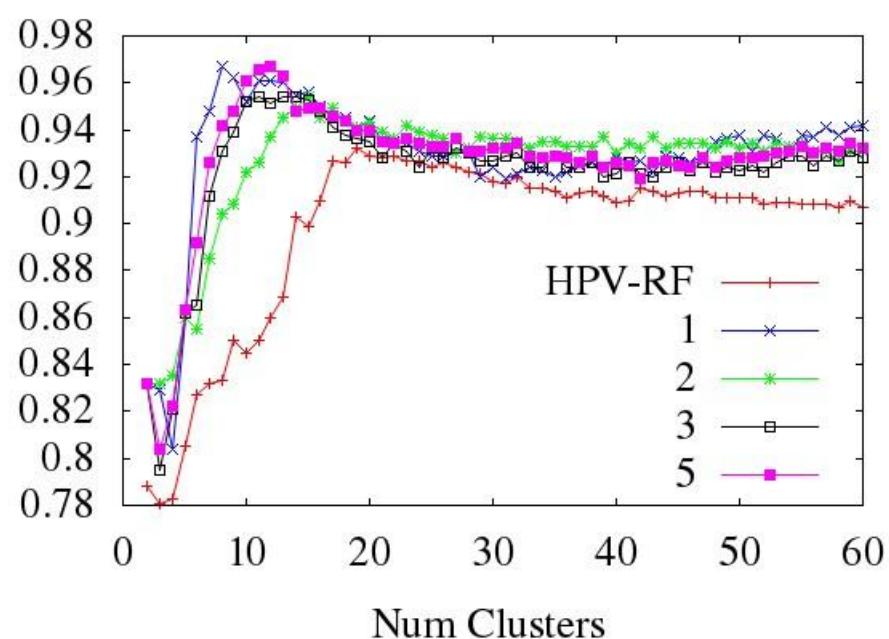


Results - Num Types

Census - 11 total types



Death Certs - 5 total types

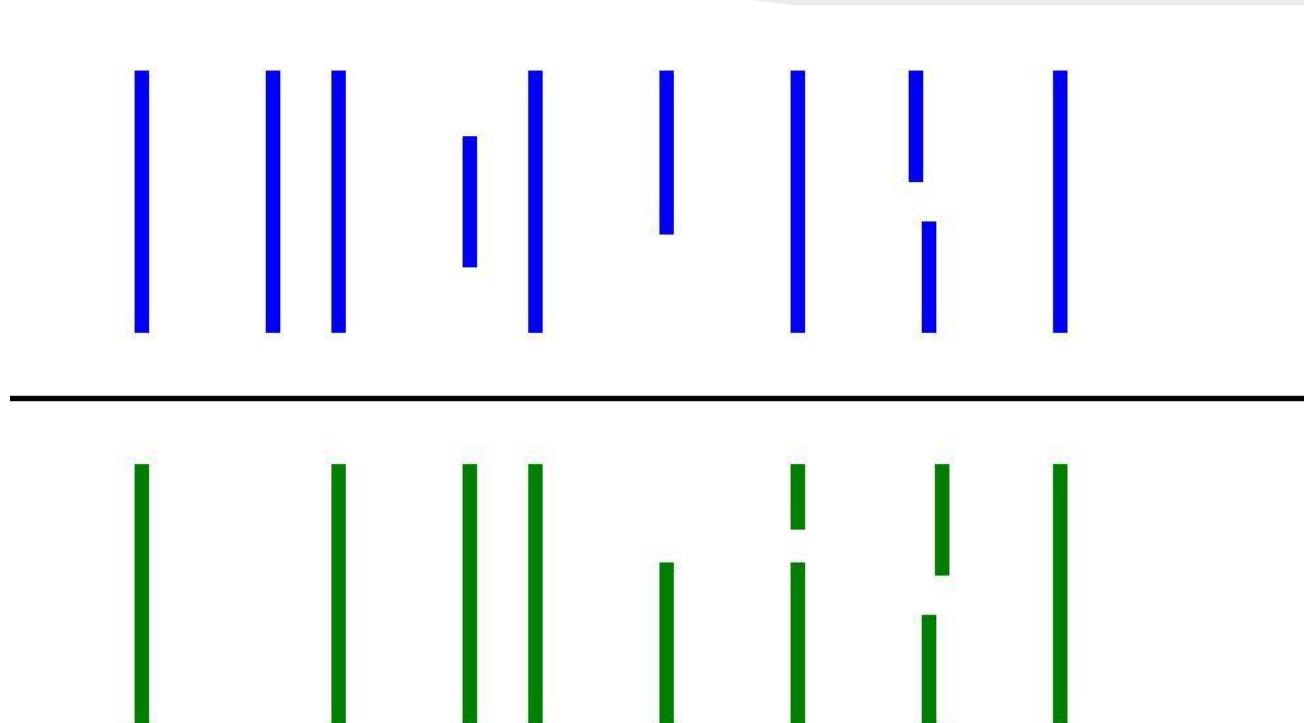


Future Work

- Detect and split impure clusters
- Bootstrap to classification
- Scale to millions of images
- Other features

Questions

Line Edit Distance - Aligned



Line Edit Distance - Offset

