## From Paper to Digital

Gaurav Sood

June 11, 2015

When we think about paper . . .





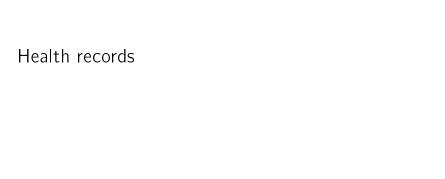




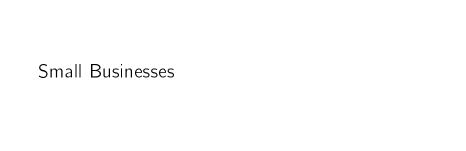
# But paper based storage of information is

common

# Libraries and Archives



Receipts



And it isn't going away	(quickly).

THE Bead Tree Fermi

- Accessible only on location

- Accessible only on location
- Typically needs help of another human, who may in turn want money

- Accessible only on location
- Typically needs help of another human, who may in turn want money
- Hard to copy, distribute

- Accessible only on location
- Typically needs help of another human, who may in turn want money
- Hard to copy, distribute
- Flammable

- Accessible only on location
- Typically needs help of another human, who may in turn want money
- Hard to copy, distribute
- Flammable
- Time consuming to find stuff

- Accessible only on location
- Typically needs help of another human, who may in turn want money
- Hard to copy, distribute
- Flammable
- Time consuming to find stuff
  Google returns average search query in .2 seconds

- Accessible only on location
- Typically needs help of another human, who may in turn want money
- Hard to copy, distribute
- Flammable
- Time consuming to find stuff
  Google returns average search query in .2 seconds
- Hard to analyze, summarize stored information

- Accessible only on location
- Typically needs help of another human, who may in turn want money
- Hard to copy, distribute
- Flammable
- Time consuming to find stuff
  Google returns average search query in .2 seconds
- Hard to analyze, summarize stored information
- Hard to track performance, identify anomalous transactions, identify patterns ...

– Lots of software:

- Lots of software:
  - Adobe Professional

- Lots of software:
  - Adobe Professional
  - Abbyy FineReader

#### – Lots of software:

- Adobe Professional
- Abbyy FineReader
- Tesseract

- Lots of software:
  - Adobe Professional
  - Abbyy FineReader
  - Tesseract
- But ...

#### – Lots of software:

- Adobe Professional
- Abbyy FineReader
- Tesseract

#### But ...

 Still can't handle complex layout, languages other than english etc.

#### – Lots of software:

- Adobe Professional
- Abbyy FineReader
- Tesseract

#### – But

 Still can't handle complex layout, languages other than english etc.

"I found that even native OCR software such as ... the Abbyy Fine Reader proved utterly incapable of extracting words from scanned images of the texts, even when those scanned images were of high quality."

#### – Lots of software:

- Adobe Professional
- Abbyy FineReader
- Tesseract

#### But ...

- Still can't handle complex layout, languages other than english etc.
- No information on how well you do (Quality Metrics).

#### – Lots of software:

- Adobe Professional
- Abbyy FineReader
- Tesseract

#### But ...

- Still can't handle complex layout, languages other than english etc.
- No information on how well you do (Quality Metrics).
- Not scalable

Take images of paper

- Take images of paper
- Within images, find where relevant text is located

- Take images of paper
- Within images, find where relevant text is located
- Find out how the text is laid out

- Take images of paper
- Within images, find where relevant text is located
- Find out how the text is laid out
- Recognize the characters

- Quality of the scan: spine, contrast etc.

# - Quality of the scan: spine, contrast etc.

#### Alabama—Cable Systems

Pay Service 2 Pay Units: 71 (11/30/88). Programming (via satellite): Disney Chan-

Fee: \$7.00 monthly. Pay Service 3

Pay Units: 338 (11/30/88) Programming (via satellite): HR/I Fee: \$7.00 monthly.

Pay Service 4 Pay Units: 173 (11/30/88)

Programming (via satellite): Showtime. Fee: \$7.00 monthly.

Local advertising: Yes, Rates: \$50.00/Minute: \$35,00/30 Seconds Equipment: Scientific-Atlanta headend: C-COR

amplifiers; Times Fiber cable: Scientific-Atlanta satellite antenna; Scientific-Atlanta

satellite receivers. Miles of plant: 80.0 (coaxial). Homes passed:

Manager: Roth Hook: Chief technician: Larry Junkin.

Ownership: Northland Communications Corp.

ALLGOOD -- Brookridge Cable Special Purpose Partnership, Suite 404, 7901 Stoneridge Dr., Pleasanton, CA 94588-3600. Phone: 510-463-1919. County: Blount. ICA: ALD165. TV Market Ranking: 40. Franchise award date.

N.A. Franchise expiration date: N.A. Began: February 1, 1989.

Channel capacity: N.A. Channels available but not in use: N.A.

**Basic Service** 

Fee: N.A.

Subscribers: 107 (12/01/90) Programming (received off-airs: WDBB /F) Bessemer; WBIQ (P), WBRC-TV (F), WTTO (I), WVTM-TV (N) Birmingham; WHNT-TV (C) Huntsville-Decatur. Programming (via satellite): WTBS (I) Atlants; WGN-TV (W) Chicago: Tumer Classic Movies.

Programming (via satellite): WGN-TV (W) Chicago: Disney Channel; Family Channel Fee: \$4.65 monthly.

Pay Service 1 Pay Units: 98 (06/01/96) Programming (via satellite): HBO.

Fee: \$11.95 monthly. Pay Service 2 Pay Units: 62 (06/01/96). Programming (via satellite): The Movie

Channel. Fee: \$10.95 monthly. Miles of plant: 33.0 (coaxial)

Manager: Freddy A. Arencibia. Chief technician: Daryl Bunn Franchise fee: 3% of gross Ownership: Falcon Cable TV (MSO).

ANDALUSIA-TV Cable Co. of Andalusia Inc. Box 34, 213 Dunson St. Andalusia Al. 36420-3705. Phone: 334-222-6464. Fax: 334-222-7226. County: Covington. ICA:

AL 0043 TV Market Ranking: Outside TV Markets, Franchise award date: January 1, 1963. Franchise expiration date: September 1, 2017. Began: March 1, 1965

Channel capacity: 40 (not 2-way capable). Channels available but not in use: 4, Rasic Service

Subscribers: 4,150 (06/21/92)

Programming (received off-air): WDHN A), WTVY (C) Dothan: WDIQ (P) Dozier. WEAR-TV (A) Mobile-Pensacola: WAKA (C), WSFA (N) Montgomery-Selma; 1 FM. Programming (via satellite): WTBS (f) Atlanta; WGN-TV (W) Chicago; A & E; BET; C-SPAN: CNBC; CNN; Country Music TV: Discovery Channel; ESPN; Family Channet Headline News; Learning Channel: Lifetime: MTV: Nashville Network: Nickelodeon: Odyssey; The Weather Channel: Turner Network TV: USA Network Current originations: Time-weather: public service announcements. Fac. \$20.00 installation: \$44.00 --

TV Market Ranking: Below 100, Franchise award date: N.A. Franchise expiration date:

N.A. Began: May 1, 1961 Channel capacity: 42 (not 2-way capable) Channels available but not in use: None.

**Basic Service** Subscribers: 35,217 (01/04/96). Programming (received off-air): WJSU-TV (C) Anniston; WGNX (C), WSB-TV (A) Atlanta: WBRC-TV (F), WVTM-TV (N) Bir-

mingham; WNAL-TV (C) Gadsden; WCIQ (P) Mount Cheaha State Park. Programming (via satellite): WTBS (I) At-

Current originations: Channel guide; classified ads. Fee: \$63.89 installation; \$5.19 monthly.

Expanded Basic Service Subscribers: 33,719 (06/27/94). Programming (via satellite): A & E; BET; C-SPAN: CNBC: CNN; Country Music TV: Discovery Channel: El Entertainment TV: ESPN; Family Channel; Headline News; Learning Channel; Lifetime; MTV; Nashwille Network; Nickeladean; Odvssey; QVC: The Weather Channel: Trinity Bosto, Nat-

work; Turner Classic Movies: Turner Net-

work TV: USA Network: VH1. Fee: \$13.39 monthly. Expanded Basic Service 2 Subscribers: N.A.

Programming (via satellite): American Movie Classics; Fox Sports South. Fice: \$1.50 monthly. Pay Service 1 Play Units: 3,398 (06/27/94) Programming (via satellite): Cinemax

Fee: \$17.36 installation; \$10.00 monthly. Pay Service 2 Play Units: 2,353 (06/27/94).

Pay Service 3

Programming (via satellite): Disney Chan-Fee: \$17.36 installation; \$10.00 monthly. Pay Units: 11.431 (06/27/94) Programming (via satelite): HBO

**Basic Service** 

Subscribers: 342 (05/01/96) Programming (received off-air): WDip at Dozier, WALA-TV (F), WEAR-TV WJTC (U), WKRG-TV (C), WPMI Mobile-Pensacola Programming (via satellite): WTBS mile

lanta; WGN-TV (W) Chicago: Tumer () sic Movies. Feer N A

Ownership: Torrence Cable Inc. (MSO).

ARAB Charter Communications Inc., gry Rose Rd., Albertville, AL 35950. Phones 205-878-3802; 800-239-5111. Fax: 206 878-8287. County: Marshall. Also serve Marshall County (portions), Union Grove, ita-

TV Market Ranking: 96, Franchise award rise N.A. Franchise expiration date: N.A. Benze December 12, 1968.

Channel capacity: 38. Channels available by not in use: None.

**Basic Service** Subscribers: 4,019 (02/16/96) Programming (received off-air): WBRC-TV (F), WVTM-TV (N) Birmingham: WT.P.III Gadsden; WAAY-TV (A), WAFF (N), WHO (P), WHNT-TV (C), WZDX (F) Huntsville

Programming (via satellite): WTBS (h.as. lanta: CNN; Cornedy Central; Country Mosic TV; Discovery Channel; ESPN; ESPN2 Family Channel: Fox Sports South: Headline News; Home Shopping Network Learning Channel: Lifetime: MTV: Nashville Network; Nickelodeon; Prevue Charnel; The Weather Channel; Turner Network

TV; USA Network; VH1 Current originations: Public access. Fee: \$20.50 installation; \$10.95 monthly.

Expanded Basic Service Subscribers: 3.907 (02/16/96). Programming: N.A. Fee: \$15.00 months

ennecity: 45 (2-way cap D.way). Channels availab

x: 7.342 (07/14/95 ing (received off-air WAFF IN), WHIQ (P), W (F) Huntsville-Decatur

ning (via satellite): V WIGH-TV (W) Chicago; on Movie Classics; BET; C-S Mr. Country Music TV; Disc el: ESPN; Family Channel: adline News; Learn MTV: Nastville Net Mickelodeon: Nostalgia ( c QVC; The Weather Ch Network; Turner Netw

ment originations: Classific ser \$21.50 installation: \$16

\$1.00 converter; \$32.00 addit

Pay Units: 1,111 (07/14/95). ming (via satellite) ex: \$9.95 monthly

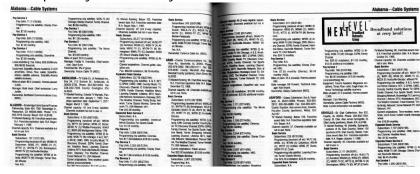
Pay Service 2 w Units: 1.159 (07/14/96) ming (via satellite): ( e \$9.95 monthly. Ty Units: 2.344 (07/14/95) gramming (via satellite):

Fee: \$9.95 monthly Pay Units: 2.267 (07/14/95). Programming (via satellite): I Fee: \$9.95 monthly Pay Service 5 Pay Units: 2.216 (07/14/95)

- Quality of the scan: spine, contrast etc.
- Complexity of the layout

- Quality of the scan: spine, contrast etc.

- Complexity of the layout



- Quality of the scan: spine, contrast etc.
- Complexity of the layout



Network; Turner Network TV; USA ent originations: Classified ads; local

and date

. Began

lable but

/BRC-TV

WIJP®

4), WHID

untsville-

IS (D.AL

ntry Mu-

ESPN2

h; Head-

etwork:

/: Nash-

in Chan-

Network

for \$21,50 installation; \$16.99 monthly; \$1.00 converter; \$32.00 additional installa-

Pay Units: 1,111 (07/14/95)

ming (via satellite): The Movie ex \$9.95 monthly.

Pay Service 2 Pay Units: 1,159 (07/14/95) ming (via satellite): Cinemax Sec \$9.95 monthly.

Pay Units: 2.344 (07/14/95) paramming (via satellite): Disney Chan-

Fee: \$9.95 monthly

ASHFORD-Galaxy Cablevision, Box 8, Headland, AL 36345-0008. Phones: 334-693-2610: 800-365-6988. Fax: 334-693-2291. County: Houston. Also serves Avon, Cowarts, Houston County (portions), Webb, ICA:

State manager: Bill Flowers. Technical man-

Ownership: Galaxy Cablevision (MSO).

ager: Ken Bryant.

not in use: None.

TV Market Ranking: Below 100. Franchise award date: N.A. Franchise expiration date: N.A. Becar: N.A. Channel capacity: 37. Channels available but

Basic Service Subscribers: 235 (07/01/96) Programming (received off-air): WLTZ (N). WRBL (C), WTVM (A) Columbus; WDHN (A) WTVY (C) Dothan: WDIO (P) Dozier:

Fee: \$8.00 installation: \$10.00 monthly. Local advertising: No.

Foulthment Tocom headend: Magnavox amplifiers: Times Fiber cable; Eagle traps; Antenna Technology satellite receivers. Miles of plant: 27.0 (cooxial). Homes passed:

1.300 Manager: Cary Manning. Ownership: James Cable Partners (MSO). Note: Current information not available.

ASHVILLE-St. Clair Cablevision, Box 932, Fayette, AL 35555. Phone: 205-932-7264. County: St. Clair. Also serves Springville, St. Clair County (portions), Steele. ICA: ALD168. TV Market Ranking: 40 (Ashville, Springville,

portions of St. Clair County); Below 100 (portions of St. Clair County, Steele). Franchise award date: N.A. Franchise expiration date: N.A. Began: January 1, 1988. Channel capacity: 42. Channels available but not in use: 4.

TV Market Ranking: 96. Franchise award date: N.A. Franchise expiration date: N.A. Began:

Channel capacity: 40. Channels available but

Programming (received off-air): WAAY-TV (A). WAFF (N). WHIQ (P), WHNT-TV (C). WZDX (F) Huntsville-Decatur; allband FM. Programming (via satellite): A & E: American Movie Classics; BET; C-SPAN; CNBC; Comedy Central; Country Music TV; ESPN; Family Channel; Fox Sports South; Home Shapping Network; Knowledge TV; Learning Channel: Lifetime: MTV: Nickelodeon; QVC; The Weather Channel: Travel Channel: Trinity Bosto, Network: Turner Network TV; USA Network: VH1

Current originations: Time-weather; newsticker; stock ticker; bulletin board; message

Fee: \$19.99 monthly **Expanded Basic Service** Subscribers: 7.707 (06/01/96). Programming (via satellite): CNN; Discovery Channel; Headline News.

Fee: \$1.65 monthly **Expanded Basic Service 2** Subscribers: 7,543 (06/01/96).

- Quality of the scan: spine, contrast etc.
- Complexity of the layout
- Font

- Quality of the scan: spine, contrast etc.
- Complexity of the layout
- Font
- Language

- Quality of the scan: spine, contrast etc.
- Complexity of the layout
- Font
- Language
- Hardware and Software (duh!)

Make images

- Make images
- Detect Text

- Make images
- Detect Text
  Alabama—Cable Systems

- Make images
- Detect Text
  Alabama Cable Systems
- Segment "Characters"

- Make images
- Detect Text
  Alabama—Cable Systems
- Segment "Characters"
  Alabama Cable Systems

- Make images
- Detect Text
  Alabama—Cable Systems
- Segment "Characters"
  Alabama Cable Systems
- Classify "Characters"

- Make images
- Detect Text
  Alabama—Cable Systems
- Segment "Characters"
  Alabama—Cable Systems
- Classify "Characters"A → A → a

Detect Text

- Detect Text
  - Supervised Learning

- Detect Text
  - Supervised Learning
  - Blobs with text, Blobs without

#### Detect Text

- Supervised Learning
- Blobs with text, Blobs without
- But size of a blob is an issue

- Detect Text
  - Supervised Learning
  - Blobs with text, Blobs without
  - But size of a blob is an issue
- Character Segmentation

- Detect Text
  - Supervised Learning
  - Blobs with text, Blobs without
  - But size of a blob is an issue
- Character Segmentation
  - Supervised Learning

#### Detect Text

- Supervised Learning
- Blobs with text, Blobs without
- But size of a blob is an issue

## Character Segmentation

- Supervised Learning
- Letters (and Ligatures) versus Splits

#### Detect Text

- Supervised Learning
- Blobs with text. Blobs without
- But size of a blob is an issue

## Character Segmentation

- Supervised Learning
- Letters (and Ligatures) versus Splits

# Classify Characters (and Ligatures)

- Supervised Learning
- A versus B versus C...

Classified (training) data

Classified (training) data

– Estimate a model

- Classified (training) data
- Estimate a model  $logit[p(spam)] = \alpha + f'\beta$  where f is frequencies.

- Classified (training) data
- Estimate a model
  Predict class (e.g. Blobs with or without text) using features
  (pixel by pixel rgb)
  Use cross-validation to tune the parameters

- Classified (training) data
- Estimate a model
- Predict classes of unseen data (groups of pixels)

# Paper to Digital Pipeline

- Take images of paper
- Within images, find where relevant text is located
- Find out how the text is laid out
- Recognize the characters

# Paper to Digital Pipeline

- Take images of paper
- Within images, find where relevant text is located
- Find out how the text is laid out
- Recognize the characters
- Every step is error prone

Optimize all steps w.r.t final error rate.

Optimize all steps w.r.t final error rate. How to deal with errors that remain

How confident are you that...

- How confident are you that...
  - An area has relevant text

- How confident are you that...
  - An area has relevant text
  - Split is correct

- How confident are you that...
  - An area has relevant text
  - Split is correct
  - Right character (or ligature) is recognized

#### How to Fix Errors

- How confident are you that...
  - An area has relevant text
  - Split is correct
  - Right character (or ligature) is recognized
- Flag low confidence areas, splits, characters...

#### How to Fix Errors

- How confident are you that...
  - An area has relevant text
  - Split is correct
  - Right character (or ligature) is recognized
- Flag low confidence areas, splits, characters...
- Get humans to identify the correct classes

#### How to Fix Errors

- How confident are you that...
  - An area has relevant text
  - Split is correct
  - Right character (or ligature) is recognized
- Flag low confidence areas, splits, characters...
- Get humans to identify the correct classes
- Use that knowledge to fix other errors

Search and Replace

- Search and Replace
- OCR makes certain kinds of errors (| is mistaken for an |)

- Search and Replace
- OCR makes certain kinds of errors (| is mistaken for an |)
- Compare against a corpora (dictionary) and replace

- Search and Replace
- OCR makes certain kinds of errors (| is mistaken for an I)
- Compare against a corpora (dictionary) and replace
- But replace with what?

- Search and Replace
- OCR makes certain kinds of errors (| is mistaken for an |)
- Compare against a corpora (dictionary) and replace
- But replace with what?
- standd -> strand, stand, stood, or sand?

– How similar are two strings?

- How similar are two strings?
- Typically refers to minimum edit distance

- How similar are two strings?
- Typically refers to minimum edit distance
- Minimum number of editing operations (Insertion, Deletion, Substitution) to convert one string to another.

- How similar are two strings?
- Typically refers to minimum edit distance
- Minimum number of editing operations (Insertion, Deletion, Substitution) to convert one string to another.
- Levenshtein Distance, substitution cost = 2

- How similar are two strings?
- Typically refers to minimum edit distance
- Minimum number of editing operations (Insertion, Deletion, Substitution) to convert one string to another.
- Levenshtein Distance, substitution cost = 2
- You can implement this at word level so Microsoft Corp. is 1 away from Microsoft.

 But edit distance isn't context aware. Use surrounding words.

- But edit distance isn't context aware. Use surrounding words.
- How likely is a certain word within a phrase?

- But edit distance isn't context aware. Use surrounding words.
- How likely is a certain word within a phrase?
- $-\sim$  Contemporary spelling correction algorithms

- But edit distance isn't context aware. Use surrounding words.
- How likely is a certain word within a phrase?
- $-\sim$  Contemporary spelling correction algorithms
- A bigram model of language: given previous word, probability of next word

- But edit distance isn't context aware. Use surrounding words.
- How likely is a certain word within a phrase?
- $-\sim$  Contemporary spelling correction algorithms
- A bigram model of language: given previous word, probability of next word
- But good training data is paramount.

Training data is 'similar data' (topic model)

and data from human computation

Training data is 'similar data' (topic model)

and data from human computation

Estimate a model based on similar data

- Training data is 'similar data' (topic model) and data from human computation
- Estimate a model based on similar data
- Use stochastic gradient descent to continue to tweak parameters based on human computation

- Training data is 'similar data' (topic model)
  and data from human computation
- Estimate a model based on similar data
- Use stochastic gradient descent to continue to tweak parameters based on human computation
- Human computation parallelized, data for costlier (most duplicated low confidence strings, errors in recognition correlated) errors prioritized

- Training data is 'similar data' (topic model)
  and data from human computation
- Estimate a model based on similar data
- Use stochastic gradient descent to continue to tweak parameters based on human computation
- Human computation parallelized, data for costlier (most duplicated low confidence strings, errors in recognition correlated) errors prioritized
- Calculate error rate against trained random