Statistical Learning with Big Data

Trevor Hastie Department of Statistics Department of Biomedical Data Science Stanford University



Jean Golding Institute Showcase July 3, 2018

Some Take Home Messages

This talk is about supervised learning: building models from data that predict an outcome using a collection of input features.

- There are some powerful and exciting tools for making predictions from data.
- They are not magic! You should be skeptical. They require good data and proper internal validation.
- Human judgement and ingenuity are essential for their success.
- With big data
 - model fitting takes longer. This might test our patience for model evaluation and comparison.
 - difficult to look at the data; might be contaminated in parts.

Careful subsampling can help with both of these.

Some Definitions

Machine Learning constructs algorithms that can learn from data.

Statistical Learning is a branch of applied statistics that emerged in response to machine learning, emphasizing statistical models and assessment of uncertainty.

Data Science is the extraction of knowledge from data, using ideas from mathematics, statistics, machine learning, computer science, engineering, ...

All of these are very similar — with different emphases.

Some Definitions

Machine Learning constructs algorithms that can learn from data.

Statistical Learning is a branch of applied statistics that emerged in response to machine learning, emphasizing statistical models and assessment of uncertainty.

Data Science is the extraction of knowledge from data, using ideas from mathematics, statistics, machine learning, computer science, engineering, ...

All of these are very similar — with different emphases.

Applied Statistics?

For Statisticians: 15 minutes of fame

- 2009 "I keep saying the sexy job in the next ten years will be statisticians. And I'm not kidding!" Hal Varian, Chief Economist Google
- 2012 "Data Scientist: The sexiest job of the 21st century." Harvard Business Review























The Supervised Learning Paradigm



Training Data Fitting Prediction

The Supervised Learning Paradigm



Training Data Fitting Prediction

Traditional statistics: domain experts work for 10 years to learn good features; they bring the statistician a small clean dataset.

The Supervised Learning Paradigm



Training Data Fitting Prediction

Traditional statistics: domain experts work for 10 years to learn good features; they bring the statistician a small clean dataset.

Today's approach: we start with a large dataset with many features, and use a machine learning algorithm to find the good ones. A huge change.

- IMPORTANT! Don't trust me or anyone who says they have a wonderful machine learning algorithm, unless you see the results of a careful internal validation.
- Eg: divide data into two parts A and B. Run algorithm on part A and then test it on part B. Algorithm must not have seen any of the data in part B.
- If it works in part B, you have (some) confidence in it

- IMPORTANT! Don't trust me or anyone who says they have a wonderful machine learning algorithm, unless you see the results of a careful internal validation.
- Eg: divide data into two parts A and B. Run algorithm on part A and then test it on part B. Algorithm must not have seen any of the data in part B.
- If it works in part B, you have (some) confidence in it

Simple? Yes

- IMPORTANT! Don't trust me or anyone who says they have a wonderful machine learning algorithm, unless you see the results of a careful internal validation.
- Eg: divide data into two parts A and B. Run algorithm on part A and then test it on part B. Algorithm must not have seen any of the data in part B.
- If it works in part B, you have (some) confidence in it

Simple? Yes Done properly in practice? Rarely

- IMPORTANT! Don't trust me or anyone who says they have a wonderful machine learning algorithm, unless you see the results of a careful internal validation.
- Eg: divide data into two parts A and B. Run algorithm on part A and then test it on part B. Algorithm must not have seen any of the data in part B.
- If it works in part B, you have (some) confidence in it

Simple? Yes Done properly in practice? Rarely

In God we trust. All others bring data.

- IMPORTANT! Don't trust me or anyone who says they have a wonderful machine learning algorithm, unless you see the results of a careful internal validation.
- Eg: divide data into two parts A and B. Run algorithm on part A and then test it on part B. Algorithm must not have seen any of the data in part B.
- If it works in part B, you have (some) confidence in it

Simple? Yes Done properly in practice? Rarely

In God we trust. All others bring data.*

* Statistical "proverb" sometimes attributed to W. Edwards Deming.

Big data vary in *shape*. These call for different approaches.



Wide Data

Thousands / Millions of Variables

Hundreds of Samples

Screening and fdr, Lasso, SVM, Stepwise

We have too many variables; prone to overfitting. Need to remove variables, or regularize, or both.

Tens / Hundreds of Variables

Thousands / Millions of Samples GLM, Random Forests, Boosting, Deep Learning

Sometimes simple models (linear) don't suffice. We have enough samples to fit nonlinear models with many interactions, and not too many variables. Good automatic methods for doing this.

Tall Data

Big data vary in *shape*. These call for different approaches.

Tall and Wide Data

Thousands / Millions of Variables

Millions to Billions of Samples

Tricks of the Trade

Exploit sparsity Random projections / hashing Variable screening Subsample rows Divide and recombine Case/ control sampling MapReduce ADMM (divide and conquer)

Big data vary in *shape*. These call for different approaches.

Tall and Wide Data

Thousands / Millions of Variables

Millions to Billions of Samples

Tricks of the Trade

Exploit sparsity Random projections / hashing Variable screening Subsample rows Divide and recombine Case/ control sampling MapReduce ADMM (divide and conquer)

join Google

Also try: english tee store, english tee party, english tee sandwiches

Ads related to: English Tea

Save on English Tea - EnglishTeaStore.com - Try Us www.erglishtestore.com

Largest Selection of Tea Brands, Free Shipping on orders 550 or more. Toware 'Tea Accessories' Tea Sets' British Frod & Bracks Types: Yes Sets, Tespoth, Loose Lest, Tre Bays Find your favorite tealogo by brand at English Tea Store. Tespote. Tes Cittles Linkes - British Food

Harney & Sons Fine Teas | harney.com

hamsy.com/WasterTea/Blenders (931) 209-9963 Use Promo Cotal HCLIDX117 For 2016 Off Al Hamsy & Bons Products: Vid Ship Internationally. - Positive Globel Footprint Styles: Loose, Bulk, Teabaga, Tea Collections, Sachets

What's New + Tea Collections - Hot Clinnamon Spice Tea + All Teas

Rich Tea Biscuits in USA | BritishFoodDepot.com www.BritishFoodDepot.com McWites Rich Tea Cookies Lowest Price Rich Tea Cookies



Also try: english tee store, english tee party, english tee sandwiches

Ads related to: English Tea

Save on English Tea - EnglishTeaStore.com - Try Us www.englishTeaStore.com Largest Selection of Tea Brevich, Free Shoping on onters \$50 or mon. Teawar: Tea Accession: Tea Safe - Whith Prod S Sares Teawar: Tea Accession: Tea Safe - Brith Prod S Sares Teapots - Tea Giffs & Ideas - British Food Teapots - Tea Giffs & Ideas - British Food

Harney & Sons Fine Teas | harney.com hinrey.com/NeterTex/Benet 930 269-993 Use throm Cost HCLIDAY17 For 20% 01 Al Harney & Bona Products We Skip Internstrontly. - Positive Blobal Footprint Syles: Loops, Buk, Tablags, Sie Celeritore, Sochts

What's New + Tea Collections - Hot Clinnamon Spice Tea + All Teas

Rich Tea Biscuits in USA | BritishFoodDepot.com www.BritishFoodDepot.com McMites Rich Tea Cockies Lowest Price Rich Tea Cockies



Click-through rate. Based on the search term, knowledge of this user (IPAddress), and the Webpage about to be served, what is the probability that each of the 30 candidate ads in an ad campaign would be clicked if placed in the sponsored-link locations.

Also try: english tee store, english tee party, english tee sandwiches

Ads related to: English Tea

Save on English Tea - EnglishTeaStore.com - Try Us www.englishTeaStore.com Largest Selection of Tea Branci, Free Stipping on orders 550 or mon. Tawar Tea Accessions - Tas Sate - Stiphi Too S Shuton Tawar - Tea Accessions - Tas Sate - Stiphi Too S Shuton Teapos - Tea Clifte & Ideas - British Food Teapos - Tea Clifte & Ideas - British Food

Harney & Sons Fine Teas | harney.com hinrey.com/NeterTex/Benet 930 269-993 Use throm Cost HCLIDAY17 For 20% 01 Al Harney & Bona Products We Skip Internstrontly. - Positive Blobal Footprint Syles: Loops, Buk, Tablags, Sie Celeritore, Sochts

What's New + Tea Collections - Hot Clinnamon Spice Tea + All Teas

Rich Tea Biscuits in USA | BritishFoodDepot.com www.BritishFoodDepot.com McMites Rich Tea Cockies Lowest Price Rich Tea Cockies



Click-through rate. Based on the search term, knowledge of this user (IPAddress), and the Webpage about to be served, what is the probability that each of the 30 candidate ads in an ad campaign would be clicked if placed in the sponsored-link locations.

Logistic regression with billions of training observations. Each ad exchange does this, then bids on their top candidates, and if they win, serve the ad — all within 10ms!



Gustaf's Traditional Dutch Soft Liconce Droos 7oz. Tub by Comp Date These \$25,50 + \$2 years

Pase tax explore for Anagon mine In Since

State from any post by Carris Date Rans Carry & GR Stern



Customers Who Viewed This Item Also Viewed



Matjes Herring Tidbits by Skansen (5 ounce) **** (6)



Thick Cut Herring -European Style, 2562 ***** (4) \$8.99



Pickled Henring - 1 Gallon ***** 459.25



Whole Hening - Old Country Style, 26oz *** (2) \$8.99



Guissife Traditional Dutch Soft Liconce Drops 7cz. Tub to Compare Annual Compare Annual The RCS I III Instance The RCS I III Instance The RCS I III Instance

In State. State from any sold by Carris Davis Nette Carry & Gill, State.



Customers Who Viewed This Item Also Viewed



Recommender systems. Amazon online store, online DVD rentals, Kindle books, ...



Guistaf's Traditional Dutch Soft Licentee Drops 7cz. Tub to Cemip See Here: 52:00 - Ni region New Xie organ to Anson Here Common See See

Step from any sold by Carris Drain Terrs Carry & Gill, Step



Customers Who Viewed This Item Also Viewed



Recommender systems. Amazon online store, online DVD rentals, Kindle books, ... Based on my past experiences, and those of others like me, what else would I choose?

• Adverse drug interactions. US FDA (Food and Drug Administration) requires physicians to send in adverse drug reports, along with other patient information, including disease status and outcomes. Massive and messy data.

• Adverse drug interactions. US FDA (Food and Drug Administration) requires physicians to send in adverse drug reports, along with other patient information, including disease status and outcomes. Massive and messy data. Using natural language processing, Stanford BMI researchers found drug interactions associated with good and bad outcomes.

- Adverse drug interactions. US FDA (Food and Drug Administration) requires physicians to send in adverse drug reports, along with other patient information, including disease status and outcomes. Massive and messy data. Using natural language processing, Stanford BMI researchers found drug interactions associated with good and bad outcomes.
- Social networks. Based on who my friends are on Facebook or LinkedIn, make recommendations for who else I should invite. Predict which ads to show me.
Examples of Big Data Learning Problems

- Adverse drug interactions. US FDA (Food and Drug Administration) requires physicians to send in adverse drug reports, along with other patient information, including disease status and outcomes. Massive and messy data. Using natural language processing, Stanford BMI researchers found drug interactions associated with good and bad outcomes.
- Social networks. Based on who my friends are on Facebook or LinkedIn, make recommendations for who else I should invite. Predict which ads to show me. There are more than two billion Facebook members, and two orders of magnitude more connections. Knowledge about friends informs our knowledge about you. Graph modeling is a hot area of research. (e.g. Leskovec lab, Stanford CS.)

The Netflix Recommender

Awesome, glad you enjoyed it! Try these next...



The Netflix Prize — 2006-2009

Ne	tflix Prize	ST/	C	OMPLETED		
Leaderboard Showing Test Score. Click here to show quiz score Display top 20 3 waders.						
Rank	Team Name	Best Test Score	5 Improvement	Best Submit Time		
Grand	Prize - RMSE = 0.8587 - Winning Te	am (BellKor's Pregn	natic Chage			
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:16:28		
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22		
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40		
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31		
5	Vandolay Industries (0.8591	9.81	2009-07-10 00:32:20		
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:58		
7	BelKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09		
8	Date	0.8612	9.59	2009-07-24 17:18:43		
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51		
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59		
11	Opera Solutions	0.8623	9.47	2009-07-24 00:34:07		

41K teams participated! Competition ran for nearly 3 years. Winner "BellKor's Pragmatic Chaos", essentially tied with "The Ensemble".

The Netflix Prize — 2006-2009

Ne	tflix Prize	ST/	C	OMPLETED
Rul	es Leaderboard Update			
Lea	aderboard	Showing Test Score. Click here to show guiz score Display top 20 3 leaders.		
Rank	Team Name	Best Test Score	5 Improvement	Best Submit Time
Grand	Prize - RMSE = 0.8567 - Winning Tr	am: Belikar's Pregn	natic Chage	
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:16:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Orand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Date	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59
	Opera Solutions	0.8523	9.47	2009-07-24 00:34:07
11	Contraction of the second seco			Rept of the option of

41K teams participated! Competition ran for nearly 3 years. Winner "BellKor's Pragmatic Chaos", essentially tied with "The Ensemble". \supset our Lester Mackey \rightarrow



The Netflix Data Set



- Training Data: 480K users, 18K movies, 100M ratings (1–5) (99% ratings missing)
- Goal:
- 1M prize for 10% reduction
- in RMSE over Cinematch
- \cdots $\bullet\,$ BellKor's Pragmatic Chaos
- \cdots declared winners on
 - 9/21/2009
 - Used ensemble of models, an important ingredient being low-rank factorization (SVD) *aka* collaborative filtering

Once the data have been cleaned and organized, we are often left with a massive matrix of observations.

• If data are sparse (lots of zeros or NAs), store using sparse-matrix methods.

Once the data have been cleaned and organized, we are often left with a massive matrix of observations.

• If data are sparse (lots of zeros or NAs), store using sparse-matrix methods. Quantcast example next: fit a sequence of logistic regression models using glmnet in R with 54M rows and 7M predictors. Extremely sparse X matrix, stored in memory (256G) — took 2 hours to fit 100 models of increasing complexity.

Once the data have been cleaned and organized, we are often left with a massive matrix of observations.

- If data are sparse (lots of zeros or NAs), store using sparse-matrix methods. Quantcast example next: fit a sequence of logistic regression models using glmnet in R with 54M rows and 7M predictors. Extremely sparse X matrix, stored in memory (256G) took 2 hours to fit 100 models of increasing complexity.
- If not sparse, use distributed, compressed databases. Many groups are developing fast algorithms and interfaces to these databases.

Once the data have been cleaned and organized, we are often left with a massive matrix of observations.

- If data are sparse (lots of zeros or NAs), store using sparse-matrix methods. Quantcast example next: fit a sequence of logistic regression models using glmnet in R with 54M rows and 7M predictors. Extremely sparse X matrix, stored in memory (256G) took 2 hours to fit 100 models of increasing complexity.
- If not sparse, use distributed, compressed databases. Many groups are developing fast algorithms and interfaces to these databases. For example h2o [CRAN] by h2o.ai interfaces from R to highly compressed versions of data, using Java-based implementations of many of the important modeling tools.

glmnet

Fit regularization paths for a variety of GLMs with lasso and elastic net penalties; e.g. logistic regression

$$\log \frac{\Pr(Y = 1 \mid X = x)}{\Pr(Y = 0 \mid X = x)} = \beta_0 + \sum_{j=1}^p x_j \beta_j$$

- Lasso penalty [Tibshirani, 1996] induces *sparsity* in coefficients: $\sum_{j=1}^{p} |\beta_j| \leq s$. It shrinks them toward zero, and sets many to zero.
- Fit efficiently using coordinate descent. Handles sparse X naturally, and exploits sparsity of solutions, warms starts, variable screening, and includes methods for model selection using cross-validation.

glmnet team: TH, Jerome Friedman, Rob Tibshirani, Noah Simon, Junyang Qian, Balasubramanian Narasimhan.













Example: Large Sparse Logistic Regression

Quantcast is a digital marketing company.* Data are five-minute internet sessions. Binary target is type of family (≤ 2 adults vs adults plus children). 7 million features of session info (web page indicators and descriptors). Divided into training set (54M), validation (5M) and test (5M).

- All but 1.1M features could be screened because ≤ 3 nonzero values.
- Fit 100 models in 2 hours in R using glmnet.
- Richest model had 42K nonzero coefficients, and explained 10% deviance (like R-squared).

* TH on SAB

54M train, 5M val, 5M test



% Deviance Explained on Training Data



Hour of Day

H2O Billion Row Machine Learning Benchmark GLM Logistic Regression



Compute Hardware: AWS EC2 c3.2xlarge - 8 cores and 15 GB per node, 1 GbE interconnect

Airline Dataset 1987-2013, 42 GB CSV, 1 billion rows, 12 input columns, 1 outcome column 9 numerical features, 3 categorical features with cardinalities 30, 376 and 380

* TH on SAB

• Online (stochastic) learning algorithms are popular — need not keep data in memory.

- Online (stochastic) learning algorithms are popular need not keep data in memory.
- Subsample if possible!

- Online (stochastic) learning algorithms are popular need not keep data in memory.
- Subsample if possible! When modeling click-through rate, there is typically 1 positive example per 10,000 negatives. You do not need all the negatives, because beyond some point the variance comes from the paucity of positives. 1 in 15 is sufficient.



Will Fithian and TH (2014, Annals of Statistics) Local Case-Control Sampling: Efficient Subsampling in Imbalanced Data Sets

- Online (stochastic) learning algorithms are popular need not keep data in memory.
- Subsample if possible! When modeling click-through rate, there is typically 1 positive example per 10,000 negatives. You do not need all the negatives, because beyond some point the variance comes from the paucity of positives. 1 in 15 is sufficient.



Will Fithian and TH (2014, Annals of Statistics) Local Case-Control Sampling: Efficient Subsampling in Imbalanced Data Sets

• Think out of the box!

- Online (stochastic) learning algorithms are popular need not keep data in memory.
- Subsample if possible! When modeling click-through rate, there is typically 1 positive example per 10,000 negatives. You do not need all the negatives, because beyond some point the variance comes from the paucity of positives. 1 in 15 is sufficient.



Will Fithian and TH (2014, Annals of Statistics) Local Case-Control Sampling: Efficient Subsampling in Imbalanced Data Sets

• Think out of the box! How much accuracy do you need? Timeliness can play a role, as well as the ability to explore different approaches. Explorations can be done on subsets of the data.

Thinking out the Box: Spraygun



Work with Brad Efron



Beer ratings 1.4M ratings 0.75M vars (sparse document features)

Lasso regression path: 70 mins. Split data into 25 parts, distribute, and average: 30 secs. In addition, free prediction standard errors and CV error.

Predicting the Pathogenicity of Missense Variants

Goal: prioritize list of candidate genes for prostate cancer

Joint work with Epidemiology colleagues Weiva Sieh, Nilah Monnier Ioannidis, Joe Rothstein, Alice Whittemore, \cdots



REVEL — rare exome variant ensemble learner

Ioannidis, N., ..., Sieh, W. (Oct 2016) Amer. J. Human Genetics

Approach

- A number of existing scores for disease status do not always agree (e.g SIFT, MutPred).
- Idea is to use a Random Forest algorithm to integrate these scores into a single consensus score for predicting disease.
- We will use existing functional prediction scores, conservation scores, etc as features 12 features in all.
- Data acquired through Human Gene Mutation Database, SwissVar and ClinVar.

	Neutral	Disease
Train	123,706	6,182
Test	$2,\!406$	1,953

Correlation of Features





Trees use the features to create subgroups in the data to refine the estimate of disease.



Trees use the features to create subgroups in the data to refine the estimate of disease. Shallow trees are too coarse/inaccurate.

Random Forests



Leo Breiman (1928–2005)

- Deep trees (fine subgroups) are more accurate, but very noisy.
- Idea: fit many (1000s) different and very-deep trees, and average their predictions to reduce the noise.
- How to get different trees?
 - Grow trees to bootstrap subsampled versions of the data.
 - Randomly ignore variables as candidates for splits.

Random Forests are very effective and give accurate predictions. They are automatic, and give good CV estimates of prediction error (for free!). R package RandomForest.

Results for REVEL



Performance evaluated on independent test set, and REVEL compared with 7 other ensemble competitors.

AUC by Allele Frequency



Feature Importance



Two New(ish) Methods

GLINTERNET

With past PhD student Michael Lim (JCGS 2014). Main effect + two-factor interaction models selected using the *group lasso*.



GAMSEL

With past Ph.D student Alexandra Chouldechova, using *overlap group lasso*. Automatic, *sticky* selection between zero, linear or nonlinear terms in GAMs:



$$\eta(x) = \sum_{j=1}^{p} f_j(x_j)$$

GLINTERNET

Example: GWAS with $p=27K~{\rm Snps}$, each a 3-level factor, and a binary response, N=3500.

- Let X_j be $N \times 3$ indicator matrix for each Snp, and $X_{j:k} = X_j \star X_k$ be the $N \times 9$ interaction matrix.
- We fit model

$$\log \frac{\Pr(Y=1|X)}{\Pr(Y=0|X)} = \alpha + \sum_{j=1}^{p} X_j \beta_j + \sum_{j < k} X_{j:k} \theta_{j:k}$$

- note: $X_{j:k}$ encodes main effects and interactions.
- Maximize group-lasso penalized likelihood:

$$\ell(\mathbf{y}, \mathbf{p}) - \lambda \left[\sum_{j=1}^p \|\beta_j\|_2 + \sum_{j < k} \|\theta_{j:k}\|_2 \right]$$

• Solutions map to traditional hierarchical main-effects/interactions model (with effects summing to zero).

GLINTERNET (continued)

- Strong rules for feature filtering essential here parallel and distributed computing useful too. GWAS search space of 729M interactions!
- Formulated for all types of interactions, not just categorical variables.
- GLINTERNET very fast two-orders of magnitude faster than competition, with similar performance.

Example: Mining Electronic Health Records for Synergistic Drug Combinations

Using Oncoshare database (EHR from Stanford Hospital and Palo Alto Medical Foundation) looked for synergistic effects between 296 drugs in treatment of 9,945 breast cancer patients.

Used **GLINTERNET** to discover three potential synergies. Joint work with Yen Low, Michael Lim, TH, Nigam Shah and others.



Low, Y., $\cdots,$ Shah, N. (May 2017) J Am Med Inform Assoc



GAMSEL: Generalized Additive Model Selection

 $\underset{\alpha_{0}, \{f_{j}\}_{1}^{p}}{\text{minimize}} \frac{1}{N} \sum_{i=1}^{N} L[y_{i}, \sum_{i=1}^{p} f_{j}(x_{ij})] + \lambda \sum_{i=1}^{p} P_{j}(f_{j})$

GAMSEL: Generalized Additive Model Selection

$$\underset{\alpha_{0}, \{f_{j}\}_{1}^{p}}{\text{minimize}} \frac{1}{N} \sum_{i=1}^{N} L[y_{i}, \sum_{j=1}^{p} f_{j}(x_{ij})] + \lambda \sum_{j=1}^{p} P_{j}(f_{j})$$

Here $P_j(f_j)$ is an overlap group lasso penalty that enables state selection, as well as degree of roughness if nonlinear.


38 / 39

All the tools I described are implemented in R, which is wonderful free software that gets increasingly more powerful as it interfaces with other systems. R can be found on CRAN: http://cran.us.r-project.org

R user conference held in Brussels July 2017; > 1100 attendees!

All the tools I described are implemented in R, which is wonderful free software that gets increasingly more powerful as it interfaces with other systems. R can be found on CRAN: http://cran.us.r-project.org

R user conference held in Brussels July 2017; > 1100 attendees!

 \cdots and now for some cheap marketing ...

All the tools I described are implemented in R, which is wonderful free software that gets increasingly more powerful as it interfaces with other systems. R can be found on CRAN: http://cran.us.r-project.org

R user conference held in Brussels July 2017; > 1100 attendees!

 \cdots and now for some cheap marketing \ldots





All the tools I described are implemented in R, which is wonderful free software that gets increasingly more powerful as it interfaces with other systems. R can be found on CRAN: http://cran.us.r-project.org

R user conference held in Brussels July 2017; > 1100 attendees!

 \cdots and now for some cheap marketing \ldots



