

# Data Analysis of Data Science Job Ads

with suggestions

Liuhua Peng Lecturer, School of Mathematics and Statistics





#### **Thanks** !

#### **Fiona Simpson**

Careers & Industry Consultant Academic Engagement

Faculty of Science

The University of Melbourne

#### **Peter Karutz**

Senior Advisor Experiential Learning, Industry

Global Leadership and Employability, Academic Services

University Services

The University of Melbourne

#### **Kathy Ryan**

Academic Services and Registrar University Services The University of Melbourne

Copyright:





## Number of Job Ads with data science OR data scientist in the Title



Title with : data science OR data scientist

Last 365 Days: May 03, 2017 – May 02, 2018

Occupations	• Job Titles • Employers
Skills	<ul> <li>Baseline Skills</li> <li>Software and Programming Skills</li> <li>Specialized Skills</li> </ul>
Education and Experience	<ul><li>Degree</li><li>Experience</li></ul>
Salaries	
Salaries	



### **Trend in Advertised Occupations**

Occupation (Job Title)	# Job Ads in 2016	# Job Ads in 2017	# Job Ads in Last 365 Days
Data Science	616	1030	1165
Data Analyst	13	38	32
Software Programmers	8	4	4
Database Administrator	5	9	17
University Lecturer	5	14	14
Software Engineer	3	7	8
Business Intelligence Analyst	2	9	8
Business Manager	1	2	4







Levels of Education and Levels of Experience Requested





#### **Top Skills: Baseline and Specialized Skills**



7



Percent of Job Ads with specific software and programming skills requested



■ 2016 ■ 2017 ■ Last 365 Days



2016



Salary Range 35 31 28 30 25 20 16 15 9 10 5 0 0 535,000 to 549,000 550,00°0574,999 575,00000599,999 5100,00 205149,999 Morethans150,000

> Number of job ads within each salary category

Mean salary: \$110,000 Median salary: \$100,000



Number of job ads within each salary category

Mean salary: \$109,000 Median salary: \$100,000 Last 365 Days



Number of job ads within each salary category

Mean salary: \$112,000 Median salary: \$108,000



### Salaries !!!!!!!

2016

Salary Range



- \$50,000 to \$74,999
- \$75,000 to \$99,999
- \$100,000 to \$149,999
- More than \$150,000

Mean salary: \$110,000 Median salary: \$100,000





- \$35,000 to \$49,999
- \$50,000 to \$74,999
- \$75,000 to \$99,999
- \$100,000 to \$149,999
- More than \$150,000

Mean salary: \$109,000 Median salary: \$100,000 Last 365 Days

Salary Range



- \$35,000 to \$49,999
- \$50,000 to \$74,999
- \$75,000 to \$99,999
- \$100,000 to \$149,999
- More than \$150,000

Mean salary: \$112,000 Median salary: \$108,000



# **Suggestions** and Tips



#### **The Data Science Process**



Schutt and O'Neil (2013), Doing Data Science. O'Reilly.



#### **Useful Resources**

- 1. Kaggle <u>https://www.kaggle.com/</u>
- 2. Coursera https://www.coursera.org/
- 3. Amazon Web Services <u>https://aws.amazon.com/</u>

The Elements of Statistical Learning https://web.stanford.edu/~hastie/ElemStatLearn/



An Introduction to Statistical Learning with Applications in R <u>http://www-bcf.usc.edu/~gareth/ISL/</u>





### **Top Python Data Analysis Libraries**

- 1. Numpy <u>http://www.numpy.org/</u>
- 2. Scipy <u>https://www.scipy.org/</u>
- 3. Pandas <u>https://pandas.pydata.org/</u>
- 4. Matplotlib <u>https://matplotlib.org/</u>
- 5. Seaborn <u>https://seaborn.pydata.org/</u>
- 6. Bokeh <u>https://bokeh.pydata.org/</u>
- 7. Scikit-learn http://scikit-learn.org/
- 8. Theano <u>http://deeplearning.net/software/theano/</u>
- 9. Tensorflow <a href="https://www.tensorflow.org/">https://www.tensorflow.org/</a>
- 10.BeautifulSoup4, Urllib2, Selenium, Scrapy

Beautiful is better than ugly. Explicit is better than implicit. Simple better than complex. Complex is better than complicated. Flat is better than nested. Sparse is better than dense. Readability counts. Special cases aren't special enough to break the rules. icality beats purity. Errors should never less explicitly silenced. In the face of jasoup jo ajour

Although practicality beats purity. Errors should never pass silently. Unless explicitly silenced. In the face of ambiguity, refuse the temptation to guess. There should be one – and preferably only one – obvious way to do it. Although that way may not be obvious at first *unless you're Dutch*. Now is better than never. Although never is often better than right now. If the implementation is hard to explain, it's a bad as. If the implementation

nay be a good idea

Namespaces are

more of those!

Beautiful is better than ugly. Explicit is better than complex. Complex is better is better than complex. Complex is better than complex de. Fals is better than dense. than complexed. Fals is better than dense. Recail and the set of the counts. Special cases aren't special cases aren't Although practicality beats purity. Errors should never pass silently. Unless **explicitly** silenced. In the face of many gury, refuse the temptation to guess. There should never ambigury, refuse the temptation to guess. There should be one and you do it. Now is pass silently. Unloss explicitly silenced. In the face of may may not be obvious at first unless you're Dutch. Now is better than never. Although never is often better than right mow. If the implementation is hard to explain, it's a bad mow. If the implementation is hard to explain, it's a bad



- Anaconda is a free and open source distribution of the Python and R for data science.
- 2. It aims to simplify package management and deployment.
- 3. Use Jupyter Notebook as IDE
- Anaconda Distribution is used by over 6 million users, and it includes more than 250 popular data science packages suitable for Windows, Linux, and MacOS.
- 5. <u>https://www.anaconda.com/download/</u>



DATA SCIENCE LIBRARIES



...and many more!



## **High Performance Computing in R**

- 1. *Rcpp, RcppArmadillo, RcppEigen*: integrate R with C++
- 2. Parallel Computing in R:
  - ✓ *parallel, foreach*: execute the for loop in parallel
  - ✓ SNOW, doSNOW
  - $\checkmark$  H20: facilitates machine learning (e.g. random forests, GBM) in a parallel environment
  - ✓ *Rhadoop*, *RHIPE*: interface between R and Hadoop
  - ✓ *Rmpi, pbdMPI*: interface MPI in R
  - ✓ *rgpu, gcbd, OpenCl*: GPU programming
  - ✓ *data.table, ff, bigmemory*: large memory and out-of-memory data

https://cran.r-project.org/web/views/HighPerformanceComputing.html



## **Useful Packages in R/Rstudio**

- 1. tidyverse <a href="https://www.tidyverse.org/">https://www.tidyverse.org/</a>
- 2. rvest
- 3. dplyr, tidyr, stringr, purrr
- 4. lubridate
- 5. ggplot2, maptools, rgdal
- 6. leaflet
- 7. caret, e1071, kernlab, nnet, rpart, xgboost, tensorflow

https://cran.r-project.org/web/views/MachineLearning.html

- 8. shiny <u>http://shiny.rstudio.com/</u>
- 9. knitr





# Thank you

**Questions?** 

Liuhua Peng Lecturer, School of Mathematics and Statistics

