

# MCIS6273 Data Mining / Fall 2017 / Prof. Maull

## LECTURE 1: CLASS POLICIES, TOOLS AND TECHNOLOGIES

Week of 8/30	Lecture Notes
<b>Content</b>	class policies, class tools, introduction, what this course is about, data mining: tools, technologies and techniques
<b>Expected Outcomes</b>	<ul style="list-style-type: none"><li>• overview of course policies</li><li>• overview of data mining concepts, algorithms, methodologies</li><li>• installation of Anaconda and Python 3.6</li><li>• introduction to Jupyter Notebooks</li><li>• creation of Github account</li></ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.1</b></p> <p>» 2014. Leskovec, Jure and Rajaraman, Anand and Ullman, Jeffrey David; <i>Mining of massive datasets</i>. → <b>ch.1</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.1, ch.2</b></p> <p><b>OPTIONAL</b></p> <p>› 2012. Downey, Allen; <i>Think Python</i>. → <b>ch.1-ch.3</b></p> <p>› (website) – 2017; <i>The Periodic Table of Data Science</i>: <a href="https://www.datacamp.com/community/blog/data-science-periodic-table#gs.TF297Gsm">https://www.datacamp.com/community/blog/data-science-periodic-table#gs.TF297Gsm</a>. → <b>Familiarize yourself with the entire table.</b></p>
<b>Homework</b>	<p><b>DUE:</b> Monday, 9/6 - midnight</p> <p>Please see the Blackboard/Github repo for what to turn in.</p>

## LECTURE 2: DATA / REPRESENTATION, PREPARATION AND MANIPULATION

Week of 9/6	Lecture Notes
<b>Content</b>	introduction to core concepts in data; data types and representation of data; data formats including structured and unstructured; concepts in pre-processing data including scaling, sampling, normalizing, binning and imputing
<b>Expected Outcomes</b>	<ul style="list-style-type: none"><li>• understand data types and common formats</li><li>• identify cleaning and adjusting scenarios and apply techniques appropriately</li><li>• utilize and apply the appropriate Python tools (Pandas for data import and cleaning)</li></ul>

Week of 9/6	Lecture Notes
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.1</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.1, ch.2</b></p> <p>» 2012. McKinney, Wes; <i>Python for data analysis: Data wrangling with Pandas, NumPy, and IPython</i>. → <b>ipython/Jupyter notebooks for ch.5, ch.6 and ch.7</b></p> <p>» (website) – 2017; <i>Distance computations (scipy.spatial.distance)</i>: <a href="https://docs.scipy.org/doc/scipy/reference/spatial.distance.html">https://docs.scipy.org/doc/scipy/reference/spatial.distance.html</a>. → <b>euclidean, cosine, correlation, jaccard</b></p> <p><b>OPTIONAL</b></p> <p>› 2012. Downey, Allen; <i>Think Python</i>. → <b>ch.1-ch.3</b></p> <p>› (website) – 2017; <i>Pandas Cookbook</i>: <a href="https://github.com/jvns/pandas-cookbook">https://github.com/jvns/pandas-cookbook</a>. → <b>familiarize yourself with this content of this repo</b></p> <p>› (Michael Kennedy’s Talk Python To Me podcast) – 11-28-2016; <i>Episode #90: Data Wrangling with Python</i>: <a href="http://talkpythontome.fm">http://talkpythontome.fm</a>. → <b>listen to the entire episode</b></p>
<b>Homework</b>	<p><b>DUE:</b> Monday, 9/18 - midnight</p> <p>Please see the Blackboard/Github repo for what to turn in.</p>

---

## LECTURE 3: DATA / DISTANCE, SIMILARITY, STATISTICAL CONCEPTS

Week of 9/13	Lecture Notes
<b>Content</b>	introduction to comparing data using common metrics; introductory concepts in disorder; introductory statistical concepts; intuitions over data dimensionality and common reduction techniques
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• identify common distance metrics and their appropriate contexts</li> <li>• understand similarity (and dissimilarity) in data</li> <li>• develop intuitions of statistical concepts in correlation, distributions and expect value</li> <li>• understand dimensionality reduction via PCA</li> <li>• utilize and apply basic statistical tools in Python (Pandas/Numpy)</li> </ul>

Week of 9/13	Lecture Notes
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.7</b></p> <p>» 2014. Leskovec, Jure and Rajaraman, Anand and Ullman, Jeffrey David; <i>Mining of massive datasets</i>. → <b>ch.11</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.2.5, ch.10.4.2</b></p> <p>» 2017. VanderPlas, Jake; <i>Python Data Science Handbook</i>. → <b>ch.5.10 (In-depth Principal Components Analysis notebook)</b></p> <p>» (website) – 2017; <i>sklearn.neighbors.DistanceMetric class</i>:  <a href="http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.DistanceMetric.html">http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.DistanceMetric.html</a>.  → <b>euclidean, cosine, jaccard</b></p> <p><b>OPTIONAL</b></p> <p>› 1997. Charles M. Grinstead, CM and Snell, JL; <i>Introduction to Probability</i>. → <b>nice introductory resource to probability</b></p> <p>› (website) – 2017; <i>Distance computations (scipy.spatial.distance)</i>:  <a href="https://docs.scipy.org/doc/scipy/reference/spatial.distance.html">https://docs.scipy.org/doc/scipy/reference/spatial.distance.html</a>. → <b>cdist, euclidean, cosine, jaccard</b></p> <p>› (O'Reilly Data Show podcast) – 07-06-2017; <i>A framework for building and evaluating data products</i>:  <a href="https://www.oreilly.com/ideas/a-framework-for-building-and-evaluating-data-products">https://www.oreilly.com/ideas/a-framework-for-building-and-evaluating-data-products</a>.  → <b>listen to the entire interview</b></p>
<b>Homework</b>	–

## LECTURE 4: ASSOCIATION RULE MINING, PATTERN MINING

Week of 9/20	Lecture Notes
<b>Content</b>	introduction to concepts for rule and pattern mining; introduction to apriori algorithm for frequent patterns; motivating the market basket analysis context for pattern mining; exploring addition contexts
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand concepts behind frequent patterns</li> <li>• understand association rule mining, apriori algorithm, FP-growth</li> <li>• apply and compute basic patterns by hand</li> <li>• identify the contexts for applying pattern mining</li> </ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.8</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.5</b></p> <p><b>OPTIONAL</b></p> <p>› (PartiallyDerivative.com podcast) – 06-13-2017; <i>The Secret Life Of A Data Scientist</i>:  <a href="http://partiallyderivative.com/podcast/2017/06/13/the-secret-life-of-a-data-scientist">http://partiallyderivative.com/podcast/2017/06/13/the-secret-life-of-a-data-scientist</a>.  → <b>listen to the entire podcast</b></p>

---

Week of 9/20	Lecture Notes
<b>Homework</b>	<b>DUE:</b> Monday, 10/2 - midnight Please see the Blackboard/Github repo for what to turn in.

---



---

## LECTURE 5: UNSUPERVISED TECHNIQUES / INTRODUCTION TO CLUSTERING

---

Week of 9/27	Lecture Notes
<b>Content</b>	introduction to cluster analysis and motivations; introduction to unsupervised clustering algorithms; partitioning (k-means, k-medoids); hierarchical agglomerative methods; model-based (expectation-maximization) neural networks (SOM self-organizing maps); visualizing with voronoi diagrams
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• exposure to unsupervised clustering methods, k-Means</li> <li>• introduction to key clustering algorithms</li> <li>• distinguish between partition and model-based algorithms</li> </ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.13, ch.14, ch.15, ch.17</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.7</b></p> <p>» (website) – 2015; <i>Basic Clustering with k-Means</i>: <a href="https://nbviewer.jupyter.org/github/tmbdev/teaching-mmir/blob/master/30-kmeans.ipynb">https://nbviewer.jupyter.org/github/tmbdev/teaching-mmir/blob/master/30-kmeans.ipynb</a>. → <b>Familiarize yourself with the notebook.</b></p> <p><b>OPTIONAL</b></p> <p>» (LinearDigressions.com podcast) – 04-16-2017; <i>Education Analytics</i>: <a href="http://lineardigressions.com/episodes/2017/4/16/education-analytics">http://lineardigressions.com/episodes/2017/4/16/education-analytics</a>. → <b>listen to the entire podcast</b></p> <p>» (website) – –; <i>Programatically understanding Expectation Maximization</i>: <a href="https://nipunbatra.github.io/blog/2014/em.html">https://nipunbatra.github.io/blog/2014/em.html</a>. → <b>read this practical explanation (with Python code) of the EM algorithm</b></p>
<b>Homework</b>	<b>DUE:</b> Monday, 10/23 - midnight Please see the Blackboard/Github repo for what to turn in.

---



---

## LECTURE 6: UNSUPERVISED TECHNIQUES / MORE CLUSTERING

---

Week of 10/4	Lecture Notes
<b>Content</b>	continued clustering, hierarchical algorithms (agglomerative), introduction to density-based algorithms (DBSCAN)
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand hierarchical and density-based algorithms</li> <li>• develop intuitions for choosing algorithms in various contexts</li> <li>• utilize algorithms on real-world data</li> </ul>

---

Week of 10/4	Lecture Notes
<b>Readings &amp; Supplemental Homework</b>	No assigned readings. Please complete readings from previous week if not current. –

## LECTURE 7: SUPERVISED TECHNIQUES / CLASSIFICATION AND PREDICTION

Week of 10/11	Lecture Notes
<b>Content</b>	classification and prediction; understanding decision trees, concepts and theory; probabilistic approaches to classification - naïve bayes; introduction to bayesian belief networks
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand and explain decision trees</li> <li>• develop probabilistic models of classification using naïve Bayes</li> <li>• identify BBNs and their application context</li> <li>• utilize naïve Bayes in real-world applications</li> </ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.18, ch.19</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.6.3, ch.6.4</b></p> <p><b>OPTIONAL</b></p> <p>› (DataSkeptic.com podcast) – 08-04-2017; <i>MINI: Bayesian Belief Networks</i>: <a href="https://dataskeptic.com/blog/episodes/2017/bayesian-belief-networks">https://dataskeptic.com/blog/episodes/2017/bayesian-belief-networks</a>. → <b>explore this light discussion of BBNs</b></p> <p>› 2012. Barber, D.; <i>Bayesian Reasoning and Machine Learning</i>. → <b>explore ch.3 for in a deeper theoretical treatment of BBNs</b></p>
<b>Homework</b>	<b>DUE:</b> Monday, 11/3 - midnight Please see the Blackboard/Github repo for what to turn in.

## LECTURE 8: SUPERVISED TECHNIQUES / CLASSIFICATION AND PREDICTION

Week of 10/18	Lecture Notes
<b>Content</b>	linear regression models for prediction; logistic regression models for prediction; introduction to generalized linear models
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand and develop linear regression models</li> <li>• understand and interpret logistic regression models</li> <li>• exposure to generalized linear models</li> </ul>

Week of 10/18	Lecture Notes
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.20</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.6.11</b></p> <p><b>OPTIONAL</b></p> <p>› (DataSkeptic.com podcast) – 01-27-2017; <i>MINI: Logistic Regression on Audio Data</i>: <a href="https://dataskeptic.com/blog/episodes/2017/logistic-regression-on-audio-data">https://dataskeptic.com/blog/episodes/2017/logistic-regression-on-audio-data</a>. → <b>listen to the entire podcast</b></p> <p>› (website) – –; <i>Building a logistic regression classifier from the ground up</i>: <a href="http://inmachineswetrust.com/posts/building-logistic-regression/">http://inmachineswetrust.com/posts/building-logistic-regression/</a>. → <b>this is a nice explanation (and code) in Python</b></p>
<b>Homework</b>	–

## LECTURE 9: SUPERVISED TECHNIQUES / CLASSIFICATION AND MODEL EVALUATION

Week of 10/25	Lecture Notes
<b>Content</b>	support vector machines; neural networks and the basic NN model and its relation to learning algorithms; evaluating models and applying techniques to model validation
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand support vector machines and their strengths</li> <li>• understand neural networks, their basic theory and application</li> <li>• identify and develop intuition around model evaluation and validation</li> </ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i>. → <b>ch.21</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.6.6, ch.6.7</b></p> <p><b>OPTIONAL</b></p> <p>› (DataSkeptic.com podcast) – 05-27-2017; <i>Data Science at eHarmony</i>: <a href="https://dataskeptic.com/blog/episodes/2016/data-science-at-eharmony">https://dataskeptic.com/blog/episodes/2016/data-science-at-eharmony</a>. → <b>listen to the entire podcast</b></p>
<b>Homework</b>	<p><b>DUE:</b> Monday, 11/30 - midnight</p> <p>Please see the Blackboard/Github repo for what to turn in.</p>

## LECTURE 10: ENSEMBLE METHODS

Week of 11/1	Lecture Notes
<b>Content</b>	ensemble methods; introduction to boosting, bagging, random forests and related methods

Week of 11/1	Lecture Notes
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand and identify the need for ensembles</li> <li>• identify and develop intuition around ensemble model evaluation and validation</li> </ul>
<b>Readings &amp; Supplemental</b>	<b>REQUIRED</b> » 2014. Zaki, Mohammed J and Meira Jr, Wagner; <i>Data mining and analysis: fundamental concepts and algorithms</i> . → <b>ch.22</b> » 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i> . → <b>ch.6.12, ch.6.13, ch.6.14, ch.6.15</b>
<b>Homework</b>	–

---

## LECTURE 11: DATA VISUALIZATION: INTRODUCTORY CONCEPTS

Week of 11/8	Lecture Notes
<b>Content Expected Outcomes</b>	introduction to data visualization; building data narratives <ul style="list-style-type: none"> <li>• understand core social mining algorithms</li> <li>• understand concepts in network analysis</li> </ul>
<b>Readings &amp; Supplemental</b>	<b>REQUIRED</b> » 2015. Knafllic, Cole Nussbaumer; <i>Storytelling with data: A data visualization guide for business professionals</i> . → <b>ch.8</b> » (website) – 2017; <i>D3.js: Data-Driven Documents</i> : <a href="http://d3js.org">http://d3js.org</a> . → <b>familiarize yourself with some of the visualizations and capabilities of D3.js</b>  <b>OPTIONAL</b> › 2014. B\^orner, Katy and Polley, David E; <i>Visual insights: A practical guide to making sense of data</i> . → <b>ch.5</b> › (website) – 2017; <i>Analyzing Scrabble Games</i> : <a href="http://rpubs.com/jalapic/scrabblr">http://rpubs.com/jalapic/scrabblr</a> . → <b>This is a very interesting exploration in analysis and visualization.</b> › (website) – 2017; <i>World Population Growth</i> : <a href="https://ourworldindata.org/world-population-growth/">https://ourworldindata.org/world-population-growth/</a> . → <b>explore some of the data and visualizations</b> › (website) – 2017; <i>RAWGraphs: The missing link between spreadsheets and data visualization</i> : <a href="http://rawgraphs.io/">http://rawgraphs.io/</a> . → <b>explore this site and its galleries</b> › (website) – 2016; <i>Rio 2016 Medals Race: An analysis of the 2016 Olympic Medals</i> : <a href="http://timesofoman.com/extra/rio_2016_medal_tally/index.html">http://timesofoman.com/extra/rio_2016_medal_tally/index.html</a> . → <b>explore this visualization</b>
<b>Homework</b>	–

---

## LECTURE 12: INTRODUCTION TO SOCIAL MINING

Week of 11/15	Lecture Notes
<b>Content</b>	introduction to social mining; introduction to recommendation systems, collaborative and content-based filtering

Week of 11/15	Lecture Notes
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand core social mining algorithms</li> <li>• understand concepts in network analysis</li> <li>• understand core recommender system concepts</li> </ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2014. Leskovec, Jure and Rajaraman, Anand and Ullman, Jeffrey David; <i>Mining of massive datasets</i>. → <b>ch.10</b></p> <p>» 2015. Grus, Joel; <i>Data science from scratch: First principles with Python</i>. → <b>ch.22</b></p> <p><b>OPTIONAL</b></p> <p>› 2014. Leskovec, Jure and Rajaraman, Anand and Ullman, Jeffrey David; <i>Mining of massive datasets</i>. → <b>ch.9</b></p> <p>› 2014. B\^ornier, Katy and Polley, David E; <i>Visual insights: A practical guide to making sense of data</i>. → <b>ch.5</b></p>
<b>Homework</b>	–

## LECTURE 13: INTRODUCTION TO TEXT MINING

Week of 11/29	Lecture Notes
<b>Content</b>	introduction to text mining; concepts in document preparation pipeline (tokenizing, stemming, etc.); TFIDF, cosine similarity; corpus selection
<b>Expected Outcomes</b>	<ul style="list-style-type: none"> <li>• understand introductory concepts in text mining and information retrieval</li> <li>• understand document preparation tools</li> <li>• apply basic concepts to real-world data</li> </ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; <i>Data mining: concepts and techniques</i>. → <b>ch.10.4</b></p> <p>» 2008. Manning, Christopher D and Raghavan, Prabhakar and Sch\^utze, Hinrich; <i>Introduction to information retrieval</i>. → <b>ch.6</b></p> <p><b>OPTIONAL</b></p> <p>› 2008. Manning, Christopher D and Raghavan, Prabhakar and Sch\^utze, Hinrich; <i>Introduction to information retrieval</i>. → <b>ch.13</b></p> <p>› (O'Reilly Data Show podcast) – 07-06-2017; <i>Language understanding remains one of AI's grand challenges</i>: <a href="https://www.oreilly.com/ideas/language-understanding-remains-one-of-ais-grand-challenges">https://www.oreilly.com/ideas/language-understanding-remains-one-of-ais-grand-challenges</a>. → <b>listen to the entire interview</b></p> <p>› (LinearDigressions.com podcast) – 04-30-2017; <i>Word2Vec</i>: <a href="http://lineardigressions.com/episodes/2017/4/30/word2vec">http://lineardigressions.com/episodes/2017/4/30/word2vec</a>. → <b>listen to the entire podcast</b></p>
<b>Homework</b>	–



## LECTURE 14: OPEN DATA, ETHICS IN DATA MINING, THE FUTURE OF DATA SCIENCE

Week of 12/6	Lecture Notes
<b>Content</b>	open data portals, APIs, tools and technologies; ethics in data mining; anonymization, privacy and data considerations; data science and the future
<b>Expected Outcomes</b>	<ul style="list-style-type: none"><li>• exposure to open data portals and open data technologies</li><li>• exposure to open APIs and tools for open data access</li><li>• understand data mining ethics and why ethics (and privacy) are critically important</li><li>• the future to data science, analytics and intelligent systems built on big data</li></ul>
<b>Readings &amp; Supplemental</b>	<p><b>REQUIRED</b></p> <p>» (DataStori.es podcast) – 05-18-2016; 74 - <i>Data Ethics and Privacy with Eleanor Saitta</i>: <a href="http://datastori.es/74-data-ethics-and-privacy-with-eleanor-saitta/">http://datastori.es/74-data-ethics-and-privacy-with-eleanor-saitta/</a>. → <b>listen to the entire podcast</b></p> <p>» (website) – 2017; <i>ProgrammableWeb.com: The Journal of the API Economy</i>: <a href="https://www.programmableweb.com/">https://www.programmableweb.com/</a>. → <b>familiarize yourself with this site and some APIs</b></p> <p><b>OPTIONAL</b></p> <p>» (LinearDigressions.com podcast) – 08-13-2017; <i>Curing Cancer with Machine Learning is Super Hard</i>: <a href="http://lineardigressions.com/episodes/2017/8/13/curing-cancer-with-machine-learning-is-super-hard">http://lineardigressions.com/episodes/2017/8/13/curing-cancer-with-machine-learning-is-super-hard</a>. → <b>listen to the entire podcast</b></p>
<b>Homework</b>	–

## RESOURCES

1. 2014. Zaki, Mohammed J and Meira Jr, Wagner; *Data mining and analysis: fundamental concepts and algorithms*.
2. 2014. Leskovec, Jure and Rajaraman, Anand and Ullman, Jeffrey David; *Mining of massive datasets*.
3. 1997. Charles M. Grinstead, CM and Snell, JL; *Introduction to Probability*.
4. 2011. Yau, Nathan; *Visualize this: the FlowingData guide to design, visualization, and statistics*.
5. 2014. B{"o}rner, Katy and Polley, David E; *Visual insights: A practical guide to making sense of data*.
6. 2012. Downey, Allen; *Think Python*.
7. 2012. Conway, Drew and White, John; *Machine learning for hackers*.
8. 2015. Grus, Joel; *Data science from scratch: First principles with Python*.
9. (website) – 2017; *The Periodic Table of Data Science*: <https://www.datacamp.com/community/blog/data-science-periodic-table#gs.TF297Gsm>.
10. 2011. Han, Jiawei and Pei, Jian and Kamber, Micheline; *Data mining: concepts and techniques*.
11. 2012. McKinney, Wes; *Python for data analysis: Data wrangling with Pandas, NumPy, and IPython*.

12. 2008. Manning, Christopher D and Raghavan, Prabhakar and Sch{"u"}tze, Hinrich; *Introduction to information retrieval*.
13. 2015. Knafllic, Cole Nussbaumer; *Storytelling with data: A data visualization guide for business professionals*.
14. 2016. Rose, Doug; *Data Science: Create Teams That Ask the Right Questions and Deliver Real Value*.
15. (website) – 2013; *Mining the Social Web: Data Mining Facebook, Twitter, LinkedIn, Google+, GitHub, and More*: <https://github.com/ptwobrussell/Mining-the-Social-Web-2nd-Edition/>.
16. 2017. Wexler, Steve and Shaffer, Jeffrey and Cotgreave, Andy; *The Big Book of Dashboards: Visualizing Your Data Using Real-World Business Scenarios*.
17. 2017. VanderPlas, Jake; *Python Data Science Handbook*.
18. (website) – 2015; *Basic Clustering with k-Means*: <https://nbviewer.jupyter.org/github/tmbdev/teaching-mmir/blob/master/30-kmeans.ipynb>.
19. (website) – 2017; *Distance computations (scipy.spatial.distance)*: <https://docs.scipy.org/doc/scipy/reference/spatial.distance.html>.
20. (website) – 11-15-2016; *Jupyter Notebook Tutorial: The Definitive Guide*: <https://www.datacamp.com/community/tutorials/tutorial-jupyter-notebook#gs.zExWvMw>.
21. (website) – 2017; *Pandas Cookbook*: <https://github.com/jvns/pandas-cookbook>.
22. (website) – 2017; *sklearn.neighbors.DistanceMetric class*: <http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.DistanceMetric.html>.
23. ({Michael Kennedy’s Talk Python To Me} podcast) – 11-28-2016; *Episode #90: Data Wrangling with Python*: <http://talkpythontome.fm>.
24. ({O’Reilly Data Show} podcast) – 07-06-2017; *A framework for building and evaluating data products*: <https://www.oreilly.com/ideas/a-framework-for-building-and-evaluating-data-products>.
25. ({O’Reilly Data Show} podcast) – 07-06-2017; *Language understanding remains one of AI’s grand challenges*: <https://www.oreilly.com/ideas/language-understanding-remains-one-of-ais-grand-challenges>.
26. (PartiallyDerivative.com podcast) – 06-13-2017; *The Secret Life Of A Data Scientist*: <http://partiallyderivative.com/podcast/2017/06/13/the-secret-life-of-a-data-scientist>.
27. (LinearDigressions.com podcast) – 04-16-2017; *Education Analytics*: <http://lineardigressions.com/episodes/2017/4/16/education-analytics>.
28. (LinearDigressions.com podcast) – 06-04-2017; *PageRank*: <http://lineardigressions.com/episodes/2017/6/4/pagerank>.
29. (LinearDigressions.com podcast) – 08-13-2017; *Curing Cancer with Machine Learning is Super Hard*: <http://lineardigressions.com/episodes/2017/8/13/curing-cancer-with-machine-learning-is-super-hard>.
30. (LinearDigressions.com podcast) – 04-30-2017; *Word2Vec*: <http://lineardigressions.com/episodes/2017/4/30/word2vec>.
31. (DataStori.es podcast) – 05-18-2016; *74 - Data Ethics and Privacy with Eleanor Saitta*: <http://datastori.es/74-data-ethics-and-privacy-with-eleanor-saitta/>.
32. (website) – 2017; *ProgrammableWeb.com: The Journal of the API Economy*: <https://www.programmableweb>.

com/.

33. (website) – 2017; *Analyzing Scrabble Games*: <http://rpubs.com/jalapic/scrabblr>.
34. (website) – 2017; *GSS Data Explorer*: <https://gssdataexplorer.norc.umd.edu/>.
35. (website) – 2017; *World Population Growth*: <https://ourworldindata.org/world-population-growth/>.
36. (website) – 2017; *RAWGraphs: The missing link between spreadsheets and data visualization*: <http://rawgraphs.io/>.
37. (website) – 2016; *Rio 2016 Medals Race: An analysis of the 2016 Olympic Medals*: [http://timesofoman.com/extra/rio\\_2016\\_medal\\_tally/index.html](http://timesofoman.com/extra/rio_2016_medal_tally/index.html).
38. (website) – 2017; *D3.js: Data-Driven Documents*: <http://d3js.org>.
39. (DataSkeptic.com podcast) – 08-04-2017; *MINI: Bayesian Belief Networks*: <https://dataskeptic.com/blog/episodes/2017/bayesian-belief-networks>.
40. (DataSkeptic.com podcast) – 01-27-2017; *MINI: Logistic Regression on Audio Data*: <https://dataskeptic.com/blog/episodes/2017/logistic-regression-on-audio-data>.
41. (DataSkeptic.com podcast) – 05-27-2017; *Data Science at eHarmony*: <https://dataskeptic.com/blog/episodes/2016/data-science-at-eharmony>.
42. (website) – –; *Programatically understanding Expectation Maximization*: <https://nipunbatra.github.io/blog/2014/em.html>.
43. (website) – –; *Building a logistic regression classifier from the ground up*: <http://inmachineswetrust.com/posts/building-logistic-regression/>.
44. 2012. Barber, D.; *{Bayesian Reasoning and Machine Learning}*.